

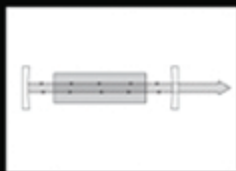
Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume II

*Design, Fabrication, and Testing; Sources
and Detectors; Radiometry and Photometry*



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is a professor at the University at Albany, and director of the Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume II
Design, Fabrication, and Testing;
Sources and Detectors;
Radiometry and Photometry

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Virendra N. Mahajan Associate Editor
*The Aerospace Corporation
El Segundo, California*

Eric Van Stryland Associate Editor
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

Copyright © 2010 by The McGraw-Hill Companies, Inc. All rights reserved. Except as permitted under the United States Copyright Act of 1976, no part of this publication may be reproduced or distributed in any form or by any means, or stored in a database or retrieval system, without the prior written permission of the publisher.

ISBN 978-0-07162927-0, MHID 0-07-162927-0

The material in this eBook also appears in the print version of this title. ISBN: P/N 978-0-07-163600-1 of set 978-0-07-149890-6. MHID: P/N 0-07-163600-5 of set 0-07-149890-7.

All trademarks are trademarks of their respective owners. Rather than put a trademark symbol after every occurrence of a trademarked name, we use names in an editorial fashion only, and to the benefit of the trademark owner, with no intention of infringement of the trademark. Where such designations appear in this book, they have been printed with initial caps.

McGraw-Hill eBooks are available at special quantity discounts to use as premiums and sales promotions, or for use in corporate training programs. To contact a representative please e-mail us at bulk-sales@mcgraw-hill.com.

TERMS OF USE

This is a copyrighted work and The McGraw-Hill Companies, Inc. (“McGraw-Hill”) and its licensors reserve all rights in and to the work. Use of this work is subject to these terms. Except as permitted under the Copyright Act of 1976 and the right to store and retrieve one copy of the work, you may not decompile, disassemble, reverse engineer, reproduce, modify, create derivative works based upon, transmit, distribute, disseminate, sell, publish or sublicense the work or any part of it without McGraw-Hill’s prior consent. You may use the work for your own noncommercial and personal use; any other use of the work is strictly prohibited. Your right to use the work may be terminated if you fail to comply with these terms.

THE WORK IS PROVIDED “AS IS.” MCGRAW-HILL AND ITS LICENSORS MAKE NO GUARANTEES OR WARRANTIES AS TO THE ACCURACY, ADEQUACY OR COMPLETENESS OF OR RESULTS TO BE OBTAINED FROM USING THE WORK, INCLUDING ANY INFORMATION THAT CAN BE ACCESSED THROUGH THE WORK VIA HYPERLINK OR OTHERWISE, AND EXPRESSLY DISCLAIM ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. McGraw-Hill and its licensors do not warrant or guarantee that the functions contained in the work will meet your requirements or that its operation will be uninterrupted or error free. Neither McGraw-Hill nor its licensors shall be liable to you or anyone else for any inaccuracy, error or omission, regardless of cause, in the work or for any damages resulting therefrom. McGraw-Hill has no responsibility for the content of any information accessed through the work. Under no circumstances shall McGraw-Hill and/or its licensors be liable for any indirect, incidental, special, punitive, consequential or similar damages that result from the use of or inability to use the work, even if any of them has been advised of the possibility of such damages. This limitation of liability shall apply to any claim or cause whatsoever whether such claim or cause arises in contract, tort or otherwise.

Information contained in this work has been obtained by The McGraw-Hill Companies, Inc. (“McGraw-Hill”) from sources believed to be reliable. However, neither McGraw-Hill nor its authors guarantee the accuracy or completeness of any information published herein, and neither McGraw-Hill nor its authors shall be responsible for any errors, omissions, or damages arising out of use of this information. This work is published with the understanding that McGraw-Hill and its authors are supplying information but are not attempting to render engineering or other professional services. If such services are required, the assistance of an appropriate professional should be sought.

COVER ILLUSTRATIONS

Left: Telescope such as used by Galileo to discover lunar craters and Jupiter's moons. The basic design is still used in opera and sports glasses. See Chap. 1.

Middle: Simplified schematic of a laser showing the gain medium which amplifies the light, and the resonator which defines the light's direction and spatial distribution. The third critical part, the source to excite the gain medium, is not shown. See Chap. 16.

Right: Zernike circle polynomial representing balanced astigmatism with a standard deviation of one wave illustrated as an isometric plot on the top, interferogram on the left, and point-spread function on the right. See Chap. 11.

This page intentionally left blank.

CONTENTS

Contributors	xvii
Brief Contents of All Volumes	xix
Editors' Preface	xxv
Preface to Volume II	xxvii
Glossary and Fundamental Constants	xxix

Part 1. Design

Chapter 1. Techniques of First-Order Layout	<i>Warren J. Smith</i>	1.3
<hr/>		
1.1	Glossary / 1.3	
1.2	First-Order Layout / 1.4	
1.3	Ray-Tracing / 1.4	
1.4	Two-Component Systems / 1.5	
1.5	Afocal Systems / 1.7	
1.6	Magnifiers and Microscopes / 1.8	
1.7	Afocal Attachments / 1.8	
1.8	Field Lenses / 1.8	
1.9	Condensers / 1.10	
1.10	Zoom or Varifocal Systems / 1.11	
1.11	Additional Rays / 1.12	
1.12	Minimizing Component Power / 1.13	
1.13	Is It a Reasonable Layout? / 1.13	
1.14	Achromatism / 1.14	
1.15	Athermalization / 1.15	
Chapter 2. Aberration Curves in Lens Design	<i>Donald C. O'Shea and Michael E. Harrigan</i>	2.1
<hr/>		
2.1	Glossary / 2.1	
2.2	Introduction / 2.1	
2.3	Transverse Ray Plots / 2.2	
2.4	Field Plots / 2.4	
2.5	Additional Considerations / 2.5	
2.6	Summary / 2.6	
2.7	References / 2.6	
Chapter 3. Optical Design Software	<i>Douglas C. Sinclair</i>	3.1
<hr/>		
3.1	Glossary / 3.1	
3.2	Introduction / 3.2	
3.3	Lens Entry / 3.2	
3.4	Evaluation / 3.8	
3.5	Optimization / 3.16	
3.6	Other Topics / 3.21	
3.7	Buying Optical Design Software / 3.22	
3.8	Summary / 3.24	
3.9	References / 3.24	

Chapter 4. Optical Specifications *Robert R. Shannon* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.1
- 4.3 Preparation of Optical Specifications / 4.5
- 4.4 Image Specifications / 4.6
- 4.5 Element Description / 4.8
- 4.6 Environmental Specifications / 4.10
- 4.7 Presentation of Specifications / 4.10
- 4.8 Problems with Specification Writing / 4.11
- 4.9 References / 4.12

Chapter 5. Tolerancing Techniques *Robert R. Shannon* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 Wavefront Tolerances / 5.3
- 5.4 Other Tolerances / 5.7
- 5.5 Starting Points / 5.8
- 5.6 Material Properties / 5.9
- 5.7 Tolerancing Procedures / 5.9
- 5.8 Problems in Tolerancing / 5.11
- 5.9 References / 5.11

Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.* 6.1

- 6.1 Glossary / 6.1
- 6.2 Introduction and Summary / 6.1
- 6.3 Mounting Individual Rotationally Symmetric Optics / 6.2
- 6.4 Multicomponent Lens Assemblies / 6.5
- 6.5 Mounting Windows and Domes / 6.11
- 6.6 Mounting Small Mirrors and Prisms / 6.11
- 6.7 Mounting Moderate-Sized Mirrors / 6.17
- 6.8 Contact Stresses in Optics / 6.21
- 6.9 Temperature Effects on Mounted Optics / 6.21
- 6.10 References / 6.25

Chapter 7. Control of Stray Light *Robert P. Breault* 7.1

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.1
- 7.3 Concepts / 7.2
- 7.4 Optical Software for Stray Light Analysis / 7.24
- 7.5 Methods / 7.27
- 7.6 Conclusion / 7.30
- 7.7 Sources of Information on Stray Light and Scattered Light / 7.31
- 7.8 References / 7.32

Chapter 8. Thermal Compensation Techniques
Philip J. Rogers and Michael Roberts 8.1

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.2
- 8.3 Homogeneous Thermal Effects / 8.2
- 8.4 Tolerable Homogeneous Temperature Change (No Compensation) / 8.5
- 8.5 Effect of Thermal Gradients / 8.6
- 8.6 Intrinsic Athermalization / 8.7
- 8.7 Mechanical Athermalization / 8.8
- 8.8 Optical Athermalization / 8.12
- 8.9 References / 8.15

Part 2. Fabrication

Chapter 9. Optical Fabrication *Michael P. Mandina* 9.3

- 9.1 Introduction / 9.3
- 9.2 Material Forms of Supply / 9.3
- 9.3 Basic Steps in Spherical Optics Fabrication / 9.4
- 9.4 Plano Optics Fabrication / 9.7
- 9.5 Asphere Optics Fabrication / 9.7
- 9.6 Crystalline Optics / 9.8
- 9.7 Purchasing Optics / 9.9
- 9.8 Conclusion / 9.9
- 9.9 References / 9.9

Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans* 10.1

- 10.1 Glossary / 10.1
- 10.2 Introduction / 10.1
- 10.3 The Diamond-Turning Process / 10.2
- 10.4 The Advantages of Diamond Turning / 10.2
- 10.5 Diamond-Turnable Materials / 10.4
- 10.6 Comparison of Diamond Turning and Traditional Optical Fabrication / 10.6
- 10.7 Machine Tools for Diamond Turning / 10.6
- 10.8 Basic Steps in Diamond Turning / 10.8
- 10.9 Surface Finish of Diamond-Turned Optics / 10.9
- 10.10 Metrology of Diamond-Turned Optics / 10.12
- 10.11 Conclusions / 10.13
- 10.12 References / 10.14

Part 3. Testing

Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan* 11.3

- Abstract / 11.3
- 11.1 Glossary / 11.3
- 11.2 Introduction / 11.4
- 11.3 Orthonormal Polynomials / 11.5
- 11.4 Zernike Circle Polynomials / 11.6
- 11.5 Zernike Annular Polynomials / 11.13
- 11.6 Hexagonal Polynomials / 11.21
- 11.7 Elliptical Polynomials / 11.21
- 11.8 Rectangular Polynomials / 11.27
- 11.9 Square Polynomials / 11.30
- 11.10 Slit Polynomials / 11.30
- 11.11 Aberration Balancing and Tolerancing, and Diffraction Focus / 11.30
- 11.12 Isometric, Interferometric, and PSF Plots for Orthonormal Aberrations / 11.36
- 11.13 Use of Circle Polynomials for Noncircular Pupils / 11.37
- 11.14 Discussion and Conclusions / 11.39
- 11.15 References / 11.40

Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández* 12.1

- 12.1 Glossary / 12.1
- 12.2 Introduction and Definitions / 12.2

- 12.3 Length and Straightness Measurements / 12.2
- 12.4 Angle Measurements / 12.10
- 12.5 Curvature and Focal Length Measurements / 12.17
- 12.6 References / 12.25

Chapter 13. Optical Testing *Daniel Malacara-Hernández* **13.1**

- 13.1 Glossary / 13.1
- 13.2 Introduction / 13.1
- 13.3 Classical Noninterferometric Tests / 13.1
- 13.4 Interferometric Tests / 13.7
- 13.5 Increasing the Sensitivity of Interferometers / 13.13
- 13.6 Interferogram Evaluation / 13.14
- 13.7 Phase-Shifting Interferometry / 13.18
- 13.8 Measuring Aspherical Wavefronts / 13.23
- 13.9 References / 13.28

Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant* **14.1**

- 14.1 Glossary / 14.1
- 14.2 Introduction / 14.1
- 14.3 Plotting CGHs / 14.3
- 14.4 Interferometers Using Computer-Generated Holograms / 14.4
- 14.5 Accuracy Limitations / 14.6
- 14.6 Experimental Results / 14.7
- 14.7 Discussion / 14.9
- 14.8 References / 14.9

Part 4. Sources

Chapter 15. Artificial Sources *Anthony LaRocca* **15.3**

- 15.1 Glossary / 15.3
- 15.2 Introduction / 15.3
- 15.3 Radiation Law / 15.4
- 15.4 Laboratory Sources / 15.7
- 15.5 Commercial Sources / 15.13
- 15.6 References / 15.53

Chapter 16. Lasers *William T. Silfvast* **16.1**

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Laser Properties Associated with the Laser Gain Medium / 16.4
- 16.4 Laser Properties Associated with Optical Cavities or Resonators / 16.19
- 16.5 Special Laser Cavities / 16.25
- 16.6 Specific Types of Lasers / 16.29
- 16.7 References / 16.37

Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman* **17.1**

- 17.1 Glossary / 17.1
- 17.2 Introduction / 17.2
- 17.3 Light-Generation Processes / 17.2
- 17.4 Light Extraction / 17.6
- 17.5 Device Structures / 17.8

- 17.6 Material Systems / 17.15
- 17.7 Substrate Technology / 17.20
- 17.8 Epitaxial Technology / 17.21
- 17.9 Wafer Processing / 17.23
- 17.10 Led Quality and Reliability / 17.25
- 17.11 Led-Based Products / 17.29
- 17.12 References / 17.35

Chapter 18. High-Brightness Visible LEDs **18.1**
Winston V. Schoenfeld

- 18.1 The Materials Systems / 18.1
- 18.2 Substrates and Epitaxial Growth / 18.2
- 18.3 Processing / 18.3
- 18.4 Solid-State Lighting / 18.4
- 18.5 Packaging / 18.5

Chapter 19. Semiconductor Lasers *Pamela L. Derry,
Luis Figueroa, and Chi-Shain Hong* **19.1**

- 19.1 Glossary / 19.1
- 19.2 Introduction / 19.3
- 19.3 Applications for Semiconductor Lasers / 19.3
- 19.4 Basic Operation / 19.4
- 19.5 Fabrication and Configurations / 19.6
- 19.6 Quantum Well Lasers / 19.9
- 19.7 High-Power Semiconductor Lasers / 19.18
- 19.8 High-Speed Modulation / 19.30
- 19.9 Spectral Properties / 19.36
- 19.10 Surface-Emitting Lasers / 19.39
- 19.11 Conclusion / 19.41
- 19.12 References / 19.43

Chapter 20. Ultrashort Optical Sources and Applications **20.1**
Jean-Claude Diels and Ladan Arissian

- 20.1 Introduction / 20.1
- 20.2 Description of Optical Pulses and Pulse Trains / 20.2
- 20.3 Pulse Evolution toward Steady State / 20.9
- 20.4 Coupling Circulating Pulses Inside a Cavity / 20.12
- 20.5 Designs of Cavities with Two Circulating Pulses / 20.15
- 20.6 Analogy of a Two-Level System / 20.22
- 20.7 Conclusion / 20.28
- 20.8 References / 20.28

Chapter 21. Attosecond Optics *Zenghu Chang* **21.1**

- 21.1 Glossary / 21.1
- 21.2 Introduction / 21.2
- 21.3 The Driving Laser / 21.4
- 21.4 Attosecond Pulse Generation / 21.6
- 21.5 Attosecond Pulse Characterization / 21.8
- 21.6 Acknowledgments / 21.10
- 21.7 References / 21.10

Chapter 22. Laser Stabilization *John L. Hall,
Matthew S. Taubman, and Jun Ye* **22.1**

- 22.1 Introduction and Overview / 22.1
- 22.2 Servo Principles and Issues / 22.5

- 22.3 Practical Issues / 22.12
- 22.4 Summary and Outlook / 22.23
- 22.5 Conclusions and Recommendations / 22.24
- 22.6 Acknowledgments / 22.24
- 22.7 References / 22.24

Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy* **23.1**

- 23.1 Glossary / 23.1
- 23.2 Introduction / 23.5
- 23.3 Some History of the Photon Concept / 23.6
- 23.4 Quantum Theory of the Laser / 23.14
- 23.5 The Laser Phase-Transition Analogy / 23.35
- 23.6 Exotic Masers and Lasers / 23.40
- 23.7 Acknowledgments / 23.45
- 23.8 References / 23.46

Part 5. Detectors

Chapter 24. Photodetectors *Paul R. Norton* **24.3**

- 24.1 Scope / 24.3
- 24.2 Thermal Detectors / 24.4
- 24.3 Quantum Detectors / 24.6
- 24.4 Definitions / 24.10
- 24.5 Detector Performance and Sensitivity / 24.13
- 24.6 Other Performance Parameters / 24.18
- 24.7 Detector Performance / 24.21
- 24.8 References / 24.101
- 24.9 Suggested Readings / 24.102

Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen* **25.1**

- 25.1 Glossary / 25.1
- 25.2 Introduction / 25.2
- 25.3 Principle of Operation / 25.3
- 25.4 Applications / 25.11
- 25.5 Reliability / 25.13
- 25.6 Future Photodetectors / 25.15
- 25.7 Acknowledgment / 25.17
- 25.8 References / 25.18
- 25.9 Additional Reading / 25.19

Chapter 26. High-Speed Photodetectors *J. E. Bowers and Y. G. Wey* **26.1**

- 26.1 Glossary / 26.1
- 26.2 Introduction / 26.3
- 26.3 Photodetector Structures / 26.3
- 26.4 Speed Limitations / 26.5
- 26.5 *p-i-n* Photodetectors / 26.10
- 26.6 Schottky Photodiode / 26.16
- 26.7 Avalanche Photodetectors / 26.17
- 26.8 Photoconductors / 26.20

- 26.9 Summary / 26.24
 26.10 References / 26.24

Chapter 27. Signal Detection and Analysis *John R. Willison* **27.1**

- 27.1 Glossary / 27.1
 27.2 Introduction / 27.1
 27.3 Prototype Experiment / 27.2
 27.4 Noise Sources / 27.3
 27.5 Applications Using Photomultipliers / 27.6
 27.6 Amplifiers / 27.10
 27.7 Signal Analysis / 27.12
 27.8 References / 27.15

Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse* **28.1**

- 28.1 Glossary / 28.1
 28.2 Thermal Detector Elements / 28.1
 28.3 Arrays / 28.7
 28.4 References / 28.13

Part 6. Imaging Detectors

Chapter 29. Photographic Films *Joseph H. Altman* **29.3**

- 29.1 Glossary / 29.3
 29.2 Structure of Silver Halide Photographic Layers / 29.4
 29.3 Grains / 29.5
 29.4 Processing / 29.5
 29.5 Exposure / 29.5
 29.6 Optical Density / 29.6
 29.7 The D-Log H Curve / 29.8
 29.8 Spectral Sensitivity / 29.11
 29.9 Reciprocity Failure / 29.11
 29.10 Development Effects / 29.12
 29.11 Color Photography / 29.12
 29.12 Microdensitometers / 29.15
 29.13 Performance of Photographic Systems / 29.16
 29.14 Image Structure / 29.17
 29.15 Acutance / 29.17
 29.16 Graininess / 29.19
 29.17 Sharpness and Graininess Considered Together / 29.22
 29.18 Signal-to-Noise Ratio and Detective Quantum Efficiency / 29.22
 29.19 Resolving Power / 29.24
 29.20 Information Capacity / 29.24
 29.21 List of Photographic Manufacturers / 29.25
 29.22 References / 29.25

Chapter 30. Photographic Materials *John D. Baloga* **30.1**

- 30.1 Introduction / 30.1
 30.2 The Optics of Photographic Films and Papers / 30.2
 30.3 The Photophysics of Silver Halide Light Detectors / 30.7
 30.4 The Stability of Photographic Image Dyes toward Light Fade / 30.10
 30.5 Photographic Spectral Sensitizers / 30.13

- 30.6 General Characteristics of Photographic Films / 30.18
- 30.7 References / 30.28

Chapter 31. Image Tube Intensified Electronic Imaging **31.1**
C. Bruce Johnson and Larry D. Owen

- 31.1 Glossary / 31.1
- 31.2 Introduction / 31.2
- 31.3 The Optical Interface / 31.3
- 31.4 Image Intensifiers / 31.7
- 31.5 Image Intensified Self-Scanned Arrays / 31.19
- 31.6 Applications / 31.27
- 31.7 References / 31.30

Chapter 32. Visible Array Detectors **32.1**
Timothy J. Tredwell

- 32.1 Glossary / 32.1
- 32.2 Introduction / 32.2
- 32.3 Image Sensing Elements / 32.2
- 32.4 Readout Elements / 32.12
- 32.5 Sensor Architectures / 32.21
- 32.6 References / 32.35

Chapter 33. Infrared Detector Arrays **33.1**
Lester J. Kozlowski and Walter F. Kosonocky

- 33.1 Glossary / 33.1
- 33.2 Introduction / 33.3
- 33.3 Monolithic FPAs / 33.10
- 33.4 Hybrid FPAs / 33.14
- 33.5 Performance: Figures of Merit / 33.23
- 33.6 Current Status and Future Trends / 33.28
- 33.7 References / 33.31

Part 7. Radiometry and Photometry

Chapter 34. Radiometry and Photometry **34.3**
Edward F. Zalewski

- 34.1 Glossary / 34.3
- 34.2 Introduction / 34.5
- 34.3 Radiometric Definitions and Basic Concepts / 34.7
- 34.4 Radiant Transfer Approximations / 34.13
- 34.5 Absolute Measurements / 34.20
- 34.6 Photometry / 34.37
- 34.7 References / 34.44

Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection **35.1**
James M. Palmer

- 35.1 Glossary / 35.1
- 35.2 Introduction and Terminology / 35.2
- 35.3 Transmittance / 35.3
- 35.4 Absorptance / 35.4
- 35.5 Reflectance / 35.4
- 35.6 Emittance / 35.7
- 35.7 Kirchhoff's Law / 35.7
- 35.8 Relationship between Transmittance, Reflectance, and Absorptance / 35.7
- 35.9 Measurement of Transmittance / 35.8

-
- 35.10 Measurement of Absorptance / 35.10
 - 35.11 Measurement of Reflectance / 35.10
 - 35.12 Measurement of Emittance / 35.14
 - 35.13 References / 35.16
 - 35.14 Further Reading / 35.23

Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer* 36.1

- 36.1 Glossary / 36.1
- 36.2 Introduction and Background / 36.2
- 36.3 Symbols, Units, and Nomenclature in Radiometry / 36.4
- 36.4 Symbols, Units, and Nomenclature in Photometry / 36.5
- 36.5 Conversion of Radiometric Quantities to Photometric Quantities / 36.11
- 36.6 Conversion of Photometric Quantities to Radiometric Quantities / 36.12
- 36.7 Radiometric/Photometric Normalization / 36.14
- 36.8 Other Weighting Functions and Conversions / 36.17
- 36.9 References / 36.17
- 36.10 Further Reading / 36.18

Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno* 37.1

- 37.1 Introduction / 37.1
- 37.2 Basis of Physical Photometry / 37.1
- 37.3 Photometric Base Unit—the Candela / 37.3
- 37.4 Quantities and Units in Photometry and Radiometry / 37.3
- 37.5 Principles in Photometry and Radiometry / 37.8
- 37.6 Practice in Photometry and Radiometry / 37.11
- 37.7 References / 37.12

Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis* 38.1

- 38.1 Introduction / 38.1
- 38.2 Definitions, Calculations, and Figures of Merit / 38.1
- 38.3 General Features of Spectroradiometry Systems / 38.7
- 38.4 Typical Spectroradiometry System Designs / 38.13
- 38.5 References / 38.19

Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly* 39.1

- 39.1 Introduction / 39.1
- 39.2 Basic Calculations / 39.2
- 39.3 Software Modeling of Nonimaging Systems / 39.6
- 39.4 Basic Building Blocks / 39.8
- 39.5 Concentration / 39.12
- 39.6 Uniformity and Illumination / 39.22
- 39.7 Acknowledgments / 39.41
- 39.8 References / 39.41

Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel* 40.1

- 40.1 Glossary / 40.1
- 40.2 Introduction / 40.1
- 40.3 Vision Biology and Perception / 40.3
- 40.4 The Science of Lighting Design / 40.6

40.5	Luminaires	/	40.24
40.6	Lighting Measurements	/	40.51
40.7	Lighting Application Areas	/	40.54
40.8	Acknowledgments	/	40.71
40.9	References	/	40.72

Index	I.1
--------------	------------

CONTRIBUTORS

- Joseph H. Altman** *Institute of Optics, University of Rochester, Rochester, New York* (CHAP. 29)
- Ladan Arissian** *Texas A&M University, College Station, Texas, and National Research Council of Canada, Ottawa, Ontario, Canada* (CHAP. 20)
- John D. Baloga** *Imaging Materials and Media, Eastman Kodak Company, Rochester, New York* (CHAP. 30)
- János A. Bergou** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Physics and Astronomy, Hunter College of the City University of New York, New York, New York* (CHAP. 23)
- John E. Bowers** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (CHAP. 26)
- Robert P. Breault** *Breault Research Organization, Tucson, Arizona* (CHAP. 7)
- William Cassarly** *Optical Research Associates, Pasadena, California* (CHAP. 39)
- Zenghu Chang** *Department of Physics, Kansas State University, Cardwell Hall, Manhattan, Kansas* (CHAP. 21)
- M. George Craford** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- Katherine Creath** *Optineering, Tucson, Arizona, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 14)
- Pamela L. Derry** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- Jean-Claude Diels** *Departments of Physics and Electrical Engineering, University of New Mexico, Albuquerque, New Mexico* (CHAP. 20)
- Berthold-Georg Englert** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Abteilung Quantenphysik der Universität Ulm, Ulm, Germany* (CHAP. 23)
- Chris J. Evans** *Zygo Corporation, Middlefield, Connecticut* (CHAP. 10)
- Luis Figueroa** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- Anurag Gupta** *Optical Research Associates, Tucson, Arizona* (CHAP. 40)
- Roland H. Haitz** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- John L. Hall** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Michael E. Harrigan** *Harrigan Optical Design, Victor, New York* (CHAP. 2)
- Chi-Shain Hong** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- C. Bruce Johnson** *Johnson Scientific Group, Inc., Phoenix, Arizona* (CHAP. 31)
- Abhay M. Joshi** *Discovery Semiconductors, Inc., Cranbury, New Jersey* (CHAP. 25)
- R. John Koschel** *Photon Engineering LLC, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 40)
- Walter F. Kosonocky*** *New Jersey Institute of Technology, University Heights, Newark, New Jersey* (CHAP. 33)
- Lester J. Kozlowski** *Altasens, Inc., Westlake Village, California* (CHAP. 33)
- Paul W. Kruse** *Consultant, Edina, Minnesota* (CHAP. 28)
- Anthony LaRocca†** *General Dynamics, Advanced Information Systems, Ypsilanti, Michigan* (CHAP. 15)

*Deceased.

†Retired.

- Melvin Lax*** *Department of Physics, City College of the City University of New York, New York, New York* (CHAP. 23)
- Virendra N. Mahajan** *The Aerospace Corporation, El Segundo, California* (CHAP. 11)
- Zacarias Malacara** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (CHAP. 12)
- Daniel Malacara-Hernández** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (CHAPS. 12, 13)
- Michael P. Mandina** *Brandon Light, Optimax Systems, Inc., Ontario, New York* (CHAP. 9)
- Paul R. Norton** *U.S. Army Night Vision and Electronics Directorate, Fort Belvoir, Virginia* (CHAP. 24)
- Donald C. O'Shea** *Georgia Institute of Technology, School of Physics, Atlanta, Georgia* (CHAP. 2)
- Yoshi Ohno** *Optical Technology Division, National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 37)
- Gregory H. Olsen** *Sensors Unlimited, Inc., Princeton, New Jersey* (CHAP. 25)
- Larry D. Owen** *NuOptics International, Phoenix, Arizona* (CHAP. 31)
- James M. Palmer*** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAPS. 35, 36)
- Richard L. Rhorer** *National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 10)
- Michael Roberts** *Pilkington Optronics, Wales, United Kingdom* (CHAP. 8)
- Philip J. Rogers** *Pilkington Optronics, Wales, United Kingdom* (CHAP. 8)
- Winston V. Schoenfeld** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 18)
- Marian O. Scully** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany* (CHAP. 23)
- Robert R. Shannon†** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAPS. 4, 5)
- Carolyn J. Sher DeCusatis** *Pace University, White Plains, New York* (CHAP. 38)
- William T. Silfvast** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 16)
- Douglas C. Sinclair** *Sinclair Optics, Inc., Fairport, New York* (CHAP. 3)
- Warren J. Smith*** *Kaiser Electro-Optics, Inc., Carlsbad, California* (CHAP. 1)
- Matthew S. Taubman** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Timothy J. Tredwell** *Sensor Systems Division, Imager Systems Development Laboratory, Eastman Kodak Company, Rochester, New York* (CHAP. 32)
- Herbert Walther*** *Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Sektion Physik der Universität München, Garching bei München, Germany* (CHAP. 23)
- Robert H. Weissman** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- Yih G. Wey** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (CHAP. 26)
- John R. Willison** *Stanford Research Systems, Inc., Sunnyvale, California* (CHAP. 27)
- William L. Wolfe** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 28)
- James C. Wyant** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 14)
- Jun Ye** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Paul R. Yoder, Jr.** *Consultant in Optical Engineering, Norwalk, Connecticut* (CHAP. 6)
- Edward F. Zalewski** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 34)
- M. Suhail Zubairy** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Electronics, Quaid-i-Azam University, Islamabad, Pakistan* (CHAP. 23)

*Deceased.

†Retired.

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENT AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes, and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*

- Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
- Chapter 31. Optical Spectrometers *Brian Henderson*
- Chapter 32. Interferometers *Parameswaran Hariharan*
- Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
- Chapter 34. Xerographic Systems *Howard Stark*
- Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
- Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
- Chapter 3. Optical Design Software *Douglas C. Sinclair*
- Chapter 4. Optical Specifications *Robert R. Shannon*
- Chapter 5. Tolerancing Techniques *Robert R. Shannon*
- Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
- Chapter 7. Control of Stray Light *Robert P. Breault*
- Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
- Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
- Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
- Chapter 13. Optical Testing *Daniel Malacara-Hernández*
- Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
- Chapter 16. Lasers *William T. Silfvast*
- Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
- Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
- Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
- Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arissian*
- Chapter 21. Attosecond Optics *Zenghu Chang*
- Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
- Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
- Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
- Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
- Chapter 27. Signal Detection and Analysis *John R. Willison*
- Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
- Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design *B. Ralph Chou*
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*
 Chapter 4. Properties of Metals *Roger A. Paquin*

- Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
- Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
- Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
- Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
- Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
- Chapter 11. Coherent Optical Transients *Paul R. Berman and D. G. Steel*
- Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
- Chapter 13. Optical Limiting *David J. Hagan*
- Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
- Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and M. Bashkansky*
- Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
- Chapter 17. Continuous-Wave Optical Parametric Oscillators *M. Ebrahim-Zadeh*
- Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
- Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
- Chapter 21. Strong Field Physics *Todd Ditmire*
- Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
- Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
- Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
- Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
- Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
- Chapter 7. Electro-Optic Modulators *Georgianne M. Purvinis and Theresa A. Maldonado*
- Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
- Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
- Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and G. J. Pearce*
- Chapter 12. Infrared Fibers *James A. Harrington*
- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
- Chapter 14. Optical Fiber Amplifiers *John A. Buck*
- Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*

- Chapter 16. Fiber-Based Couplers *Daniel Nolan*
 Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
 Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
 Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
 Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
 Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
 Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
 Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
 Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
 Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn A. MacDonald*
 Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
 Chapter 28. Requirements for X-Ray diffraction *Scott T. Misture*
 Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
 Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
 Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
 Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
 Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
 Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
 Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and M. Sanchez del Rio*
 Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
 Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells*
 Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
 Chapter 40. Zone Plates *Alan Michette*
 Chapter 41. Multilayers *Eberhard Spiller*
 Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
 Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Reflective Optics *James Harvey*
 Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
 Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
 Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
 Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
 Chapter 49. Pore Optics *Marco Beijersbergen*
 Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
 Chapter 51. The Schwarzschild Objective *Franco Cerrina*
 Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
 Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter M. Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
 Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
 Chapter 56. Laser Generated Plasmas *Alan Michette*

- Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter M. Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, nonlinear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

PREFACE TO VOLUME II

Volume II of the *Handbook of Optics* is a continuation of Volume I. It starts with optical system design and covers first-order layout, aberration curves, design software, specifications and tolerances, component mounting, stray light control, and thermal compensation techniques. Optical fabrication and testing are discussed next. A new chapter on the use of orthonormal polynomials in optical design and testing has been added. Such a polynomial representing balanced astigmatism is illustrated on the cover. The section on sources includes different types of lasers, laser stabilization, laser theory, and a discussion of ultrashort laser sources. Light-emitting diodes including the new “high-brightness” LEDs are presented. Artificial sources of light for both the laboratory and field are described along with a discussion of light standards calibration. The section on detectors includes high-speed and thermal detectors along with an analysis of signal detection. Imaging using film, detector arrays, and image tubes is discussed. This volume ends with a section on radiometry and photometry. Two new chapters have been added in this area. One is on spectroradiometry and the other is on lighting and applications.

Every effort was made to contact all the authors of chapters in the second edition that would appear in this edition so that they could update their chapters. However, the authors of several chapters could not be located or were not available. Their chapters are reproduced without update. Every effort has been made to ensure that such chapters have been correctly reproduced.

There are many other chapters in this edition of the *Handbook* that could have been included in Volumes I and II. However, page limitations prevented that. For example, in Volume V there is a section on Atmospheric Optics. It consists of three chapters, one on transmission through the atmosphere, another on imaging through atmospheric turbulence, and a third on adaptive optics to overcome some of the deleterious effects of turbulence.

The chapters are generally aimed at the graduate students, though practicing scientists and engineers will find them equally suitable as references on the topics discussed. Each chapter has sufficient references for additional and/or further study.

Virendra N. Mahajan
The Aerospace Corporation
Eric Van Stryland
CREOL, The College of Optics and Photonics
Associate Editors

This page intentionally left blank.

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the *Handbook* represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire *Handbook*, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the *Handbook* also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary unit arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets of terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance, exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript ν indicates a quantity normalized to the response of the eye. Thereby, L_ν is luminance, E_ν illuminance, and M_ν and I_ν luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

<i>Symbol</i>	<i>Name</i>	<i>Common name</i>	<i>Exponent of ten</i>
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f	Fermi occupation function, conduction band
f_v	Fermi occupation function, valence band

FN	focal ratio (<i>f</i> /number) [—]
<i>g</i>	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
<i>h</i>	height [m]
<i>I</i>	irradiance (see also <i>E</i>) [Wm^{-2}]
<i>I</i>	radiant intensity [Wsr^{-1}]
<i>I</i>	nuclear spin quantum number [—]
<i>I</i>	current [A]
<i>i</i>	$\sqrt{-1}$
Im()	imaginary part of
<i>J</i>	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{s}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
<i>k</i>	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
<i>k</i>	extinction coefficient [—]
<i>L</i>	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
<i>L</i>	inductance [h, $\text{m}^2 \text{kg C}^2$]
<i>L</i>	laser cavity length
<i>L, M, N</i>	direction cosines [—]
<i>M</i>	angular magnification [—]
<i>M</i>	radiant exitance [radiant emittance] [Wm^{-2}]
<i>m</i>	linear magnification [—]
<i>m</i>	effective mass [kg]
MTF	modulation transfer function [—]
<i>N</i>	photon flux [s^{-1}]
<i>N</i>	carrier (number) density [m^{-3}]
<i>n</i>	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
<i>P</i>	macroscopic polarization [C m^{-2}]
Re()	real part of [—]
<i>R</i>	resistance [Ω]
r	position vector [m]
<i>S</i>	Seebeck coefficient [VK^{-1}]
<i>s</i>	spin quantum number [—]
<i>s</i>	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
<i>t</i>	time [s]
<i>t</i>	thickness [m]
<i>u</i>	slope of ray with the optical axis [rad]
<i>V</i>	Abbe reciprocal dispersion [—]
<i>V</i>	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
<i>x, y, z</i>	rectangular coordinates [m]
<i>Z</i>	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\mathfrak{R}	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
College of Optical Sciences
University of Arizona
Tucson, Arizona

PART

1

DESIGN

This page intentionally left blank.

1

TECHNIQUES OF FIRST-ORDER LAYOUT

Warren J. Smith*

*Kaiser Electro-Optics, Inc.
Carlsbad, California*

1.1 GLOSSARY

A, B	scaling constants
d	distance between components
f	focal length
h	image height
I	invariant
j, k	indices
l	axial intercept distance
M	angular magnification
m	linear, lateral magnification
n	refractive index
P	partial dispersion, projection lens diameter
r	radius
S	source or detector linear dimension
SS	secondary spectrum
s	object distance
s'	image distance
t	temperature
u	ray slope
V	Abbe number
y	height above optical axis
α	radiometer field of view, projector field of view
ϕ	component power ($= 1/f$)

*Deceased.

1.2 FIRST-ORDER LAYOUT

First-order layout is the determination of the arrangement of the components of an optical system in order to satisfy the first-order requirements imposed on the system. The term “first-order” means the paraxial image properties: the size of the image, its orientations, its location, and the illumination or brightness of the image. This also implies apertures, f -numbers, fields of view, physical size limitations, and the like. It does not ordinarily include considerations of aberration correction; these are usually third- and higher-order matters, not first-order. However, ordinary chromatic aberration and secondary spectrum are first-order aberrations. Additionally, the first-order layout can have an effect on the Petzval curvature of field, the cost of the optics, the sensitivity to misalignment, and the defocusing effects of temperature changes.

The primary task of first-order layout is to determine the powers and spacings of the system components so that the image is located in the right place and has the right size and orientation. It is not necessary to deal with surface-by-surface ray-tracing here; the concern is with components. “Components” may mean single elements, cemented doublets, or even complex assemblies of many elements. The first-order properties of a component can be described by its Gauss points: the focal points and principal points. For layout purposes, however, the initial work can be done assuming that each component is of zero thickness; then only the component location and its power (or focal length) need be defined.

1.3 RAY-TRACING

The most general way to determine the characteristics of an image is by ray-tracing. As shown in Fig. 1, if an “axial (marginal)” ray is started at the foot (axial intercept) of the object, then an image is located at each place that this ray crosses the axis. The size of the image can be determined by tracing a second, “principal (chief),” ray from the top of the object and passing through the center of the limiting aperture of the system, the “aperture stop;” the intersection height of this ray at the image plane indicates the image size. This size can also be determined from the ratio of the ray slopes of the axial ray at the object and at the image; this yields the magnification $m = u_0/u'_k$; object height times magnification yields the image height.

The ray-tracing equations are

$$y_1 = -l_1 u_1 \quad (1)$$

$$u'_j = u_j - y_j \phi_j \quad (2)$$

$$y_{j+1} = y_j + d_j u'_j \quad (3)$$

$$l'_k = -y_k / u'_k \quad (4)$$

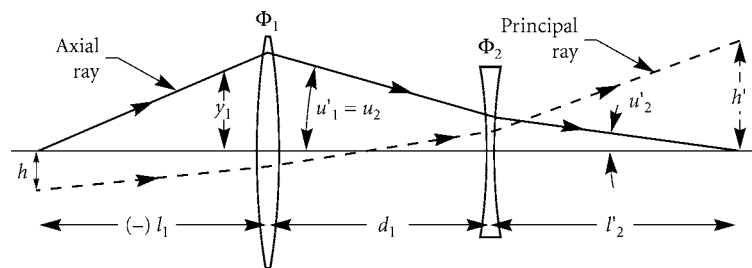


FIGURE 1

where l and l' are the axial intersection distances of the ray before and after refraction by the component, u and u' are the ray slopes before and after refraction, ϕ is the component power ($\phi = 1/f$), y_j is the height at which the ray strikes the j th component, and d_j is the distance from the j th to the $(j + 1)$ th component. Equations (2) and (3) are applied sequentially to the components, from object to image.

These equations can be used in two different ways. When the components and spacings are known, the image characteristics can readily be calculated. In the inverse application, the (unknown) powers and spaces can be represented by symbols, and the ray can be traced symbolically through the postulated number of components. The results of this symbolic ray-tracing can be equated to the required characteristics of the system; these equations can then be solved for the unknowns, which are the component powers and spacings.

As an example, given the starting ray data, y_1 and u_1 , we get

$$\begin{aligned} u'_1 &= u_1 - y_1 \phi_1 \\ y_2 &= y_1 + d_1 u'_1 = y_1 + d_1 (u_1 - y_1 \phi_1) \\ u'_2 &= u'_1 - y_2 \phi_2 \\ &= u_1 - y_1 \phi_1 - [y_1 + d_1 (u_1 - y_1 \phi_1)] \phi_2 \\ y_3 &= y_2 + d_2 u'_2 = \text{etc.} \end{aligned}$$

Obviously the equations can become rather complex in very short order. However, because of the linear characteristics of the paraxial ray equations, they can be simplified by setting either y_1 or u_1 equal to one (1.0) without any loss of generality. But the algebra can still be daunting.

1.4 TWO-COMPONENT SYSTEMS

Many systems are either limited to two components or can be separated into two-component segments. There are relatively simple expressions for solving two-component systems.

Although the figures in this chapter show thick lenses with appropriate principal planes, “thin” lenses (whose thickness is zero and whose principal planes are coincident with the two coincident lens surfaces) may be used.

For systems with infinitely distant objects, as shown in Fig. 2, the following equations for the focal length and focus distance are useful:

$$f_{AB} = f_A f_B / (f_A + f_B - d) \quad (5)$$

$$\phi_{AB} = \phi_A + \phi_B - d \phi_A \phi_B \quad (6)$$

$$B = f_{AB} (f_A - d) / f_A \quad (7)$$

$$F = f_{AB} (f_B - d) / f_B \quad (8)$$

$$h' = f_{AB} \tan u_p \quad (9)$$

where f_{AB} is the focal length of the combination, ϕ_{AB} is its power, f_A and f_B are the focal lengths of the components, ϕ_A and ϕ_B are their powers, d is the spacing between the components, B is the “back

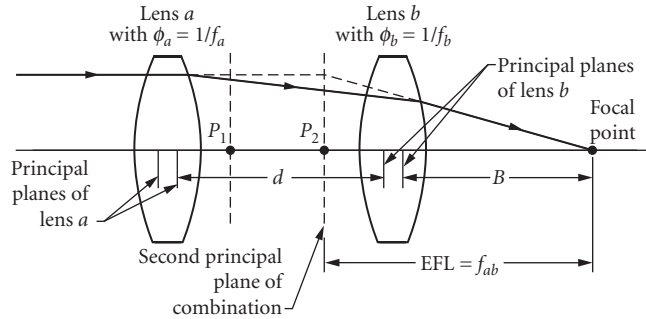


FIGURE 2

focus” distance from the B component, F is the “front focus” distance, u_p is the angle subtended by the object, and h' is the image height.

If f_{AB} , d , and B (or F) are known, the component focal lengths can be found from

$$f_A = df_{AB}/(f_{AB} - B) \tag{10}$$

$$f_B = -dB/(f_{AB} - B - d) \tag{11}$$

These simple expressions are probably the most widely used equations in optical layout work.

If a two-component system operates at *finite* conjugates, as shown in Fig. 3, the following equations can be used to determine the layout. When the required system magnification and the component locations are known, the powers of the components are given by

$$\phi_A = (ms - md - s')/msd \tag{12}$$

$$\phi_B = (d - ms + s')/ds' \tag{13}$$

where $m = h'/h$ is the magnification, s and s' are the object and image distances.

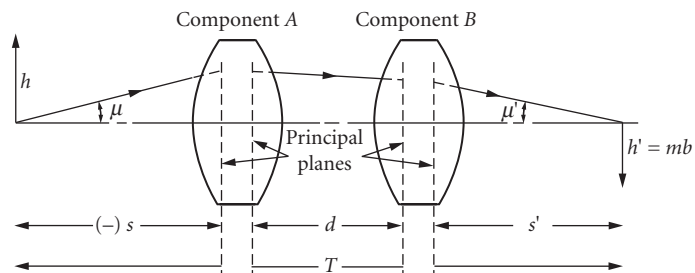


FIGURE 3

In different circumstances, the component powers, the object-to-image distance, and the magnification may be known and the component locations are to be determined. The following quadratic equation [Eq. (14)] in d (the spacing) is solved for d :

$$0 = d^2 - dT + T(f_A + f_B) + (m - 1)^2 f_A f_B / m \quad (14)$$

and then

$$s = [(m - 1)d + T] / [(m - 1) - md\phi_A] \quad (15)$$

$$s' = T + s - d \quad (16)$$

1.5 AFOCAL SYSTEMS

If the system is afocal, then the following relations will apply:

$$MP = -(f_O / f_E) = (u_E / u_O) = (d_O / d_E) \quad (17)$$

and, if the components are “thin,”

$$L = f_O + f_E \quad (18)$$

$$f_O = -L \cdot MP / (1 - MP) \quad (19)$$

$$f_E = L / (1 - MP) \quad (20)$$

where MP is the angular magnification, f_O and f_E are the objective and eyepiece focal lengths, u_E and u_O are the apparent (image) and real (object) angular fields, d_O and d_E are the entrance and exit pupil diameters, and L is the length of the telescope as indicated in Fig. 4.

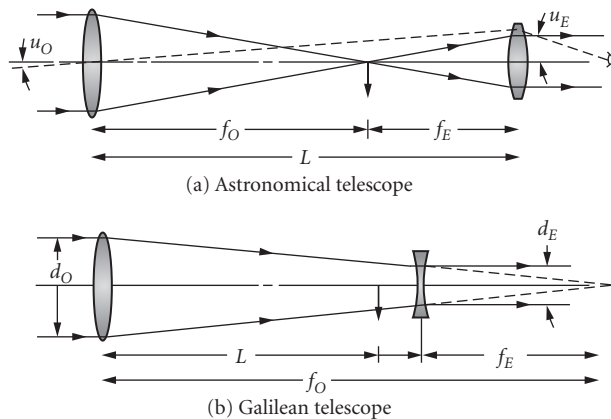


FIGURE 4

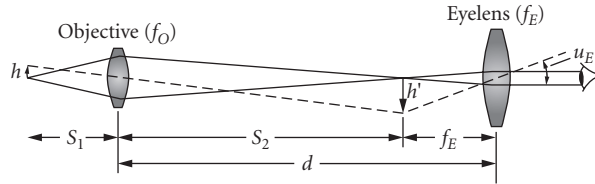


FIGURE 5

1.6 MAGNIFIERS AND MICROSCOPES

The conventional definition of magnifying power for either a magnifier or microscope compares the angular size of the image with the angular size of the object *when the object is viewed from a (conventional) distance of 10 inches*. Thus the magnification can be found from

$$MP = 10''/f \quad (21)$$

for either a simple microscope (i.e., magnifier) or a compound microscope, where f is the focal length of the system. Using the symbols of Fig. 5, we can also write the following for the compound microscope

$$MP = (f_E + f_O - d)10''/f_E f_O \quad (22)$$

$$\begin{aligned} MP &= m_O \times m_E \\ &= (S_2/S_1)(10''/f_E) \end{aligned} \quad (23)$$

1.7 AFOCAL ATTACHMENTS

In addition to functioning as a telescope, beam expander, etc., an afocal system can be used to modify the characteristics of another system. It can change the focal length, power, or field of the “prime” system. Figure 6 shows several examples of an afocal device placed (in these examples) before an imaging system. The combination has a focal length equal to the focal length of the prime system multiplied by the angular magnification of the afocal device. Note that in Fig. 6a and b the same afocal attachment has been reversed to provide two different focal lengths. If the size of the film or detector is kept constant, the angular field is changed by a factor equal to the inverse of the afocal magnification.

1.8 FIELD LENSES

Figure 7 illustrates the function of the field lens in a telescope. It is placed near (but rarely exactly at) an internal image; its power is chosen so that it converges the oblique ray bundle toward the axis sufficiently so that the rays pass through the subsequent component. A field lens is useful to keep the component diameters at reasonable sizes. It acts to relay the pupil image to a more acceptable location.

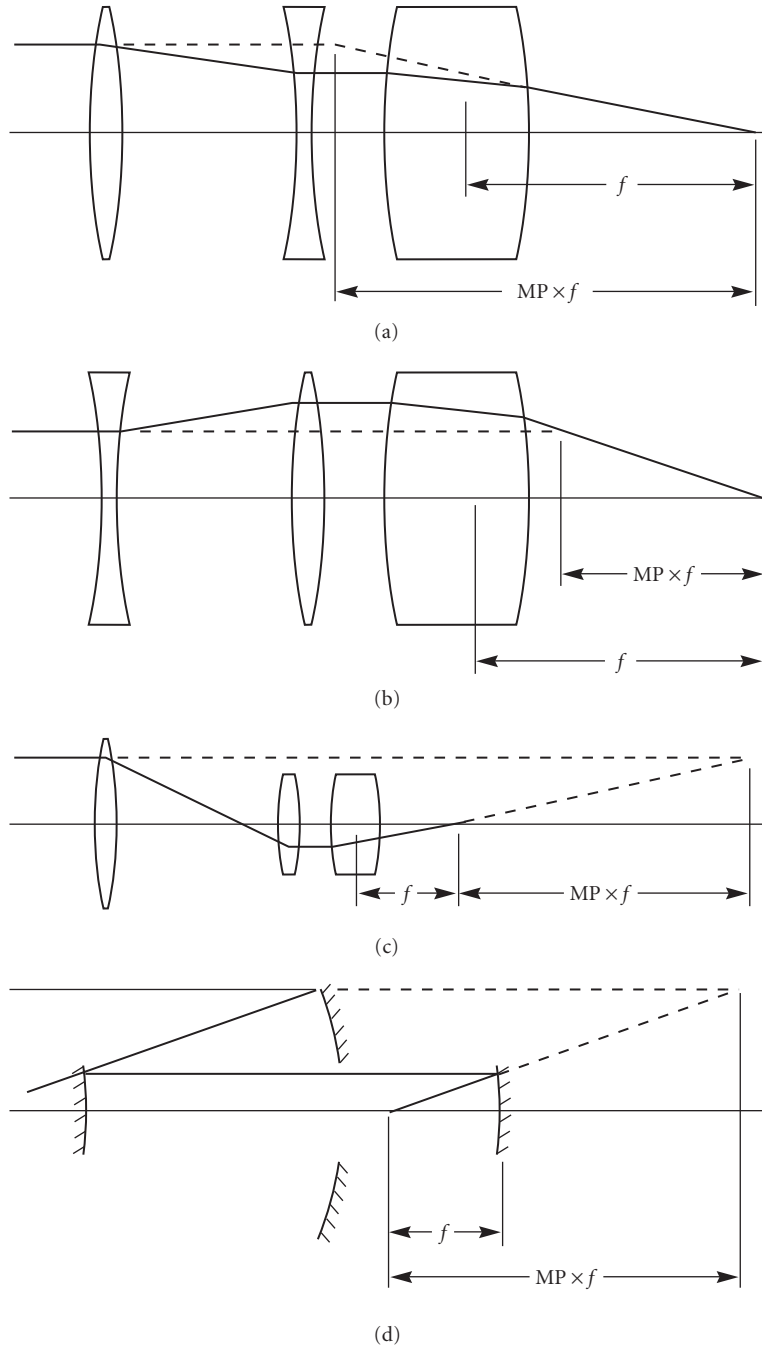


FIGURE 6

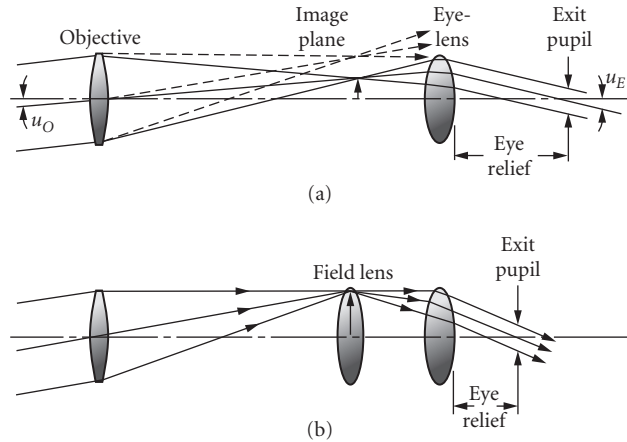


FIGURE 7

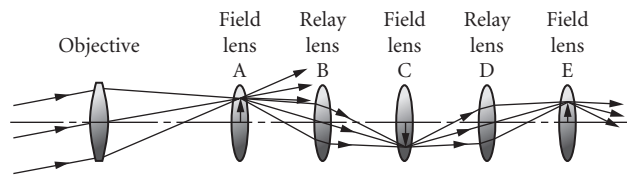


FIGURE 8

The required field lens power is easily determined. In Fig. 7 the most troublesome ray is that from the bottom of the objective aperture; its slope (u) is simply the height that it climbs divided by the distance that it travels. The required slope (u') for the ray after refraction by the field lens is defined by the image height (y), the “eyelens” semidiameter, and the spacing between them. Then Eq. (2) can be solved for the field lens power,

$$\phi = (u - u')/y \quad (24)$$

A periscope is used to carry an image through a long, small-diameter space. As shown in Fig. 8, the elements of a periscope are alternating field lenses and relay lenses. An optimum arrangement occurs when the images at the field lenses and the apertures of the relay lenses are as large as the available space allows. This arrangement has the fewest number of relay stages and the lowest power components. For a space of uniform diameter, both the field lenses and the relay lenses operate at unit magnification.

1.9 CONDENSERS

The projection/illumination condenser and the field lens of a radiation measuring system operate in exactly the same way. The condenser (Fig. 9) forms an image of the light source in the aperture of the projection lens, thereby producing even illumination from a nonuniform source. If the source

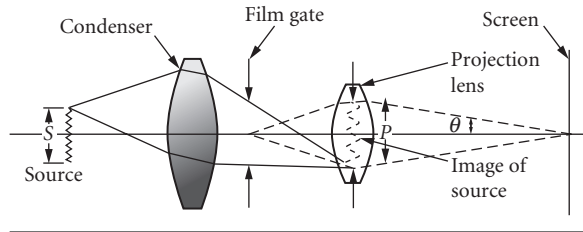


FIGURE 9

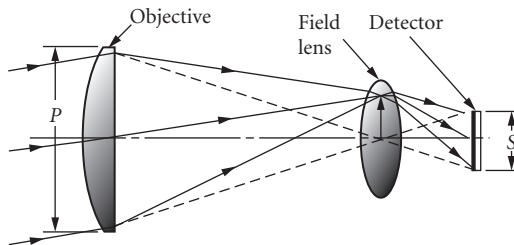


FIGURE 10

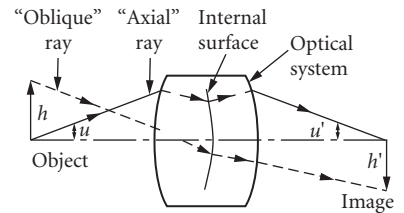


FIGURE 11

image fills the projection lens aperture, this will produce the maximum illumination that the source brightness and the projection lens aperture diameter will allow. This is often called Köhler illumination. In a radiometer type of application (Fig. 10), the field lens images the objective lens aperture on the detector, uniformly illuminating its surface and permitting the use of a smaller detector. Often, the smallest possible source or detector is desired in order to minimize power or maximize signal-to-noise. The smallest possible size is given by

$$S = P\alpha/2n \quad (25)$$

where S is the source or detector size, P is the projection lens or objective aperture diameter, α is the field angle of projection or the radiometer field of view, and n is the index in which the source or detector is immersed. This value for S corresponds to an (impractical) system speed of $F/0.5$. A source or detector size twice as large is a far more realistic limit, corresponding to a speed of $F/1.0$.

The *invariant*, $I = n(y_2u_1 - y_1u_2)$, where y_1 , u_1 , y_2 , and u_2 are the ray heights and slopes of two different rays, is an expression which has the same value everywhere in an optical system. If the two rays used are an axial ray and a principal (or chief) ray as shown in Fig. 11, and if the invariant is evaluated at the object and image surfaces, the result is

$$hnu = h'n'u' \quad (26)$$

1.10 ZOOM OR VARIFOCAL SYSTEMS

If the spacing between two components is changed, the effective focal length and the back focus are changed in accord with Eqs. (5) through (9). If the motions of the two components are arranged so that the image location is constant, this is a mechanically compensated zoom lens, so called because

the component motions are usually effected with a mechanical cam. A zoom system may consist of just the two basic components or it may include one or more additional members. Usually the two basic components have opposite-signed powers.

If a component is working at unit magnification, it can be moved in one direction or the other to increase or decrease the magnification. There are pairs of positions where the magnifications are m and $1/m$ and for which the object-to-image distance is the same. This is the basis of what is called a “bang-bang” zoom; this is a simple way to provide two different focal lengths (or powers, or fields of view, or magnifications) for a system.

1.11 ADDITIONAL RAYS

When the system layout has been determined, an “axial” ray at full aperture and a “principal” ray at full field can be traced through the system. Because of the linearity of the paraxial equations, we can determine the ray-trace data (i.e., y and u) of *any* third ray from the data of these two traced rays by

$$y_3 = Ay_1 + By_2 \quad (27)$$

$$u_3 = Au_1 + Bu_2 \quad (28)$$

where A and B are scaling constants which can be determined from

$$A = (y_3u_1 - u_3y_1)/(u_1y_2 - y_1u_2) \quad (29)$$

$$B = (u_3y_2 - y_3u_2)/(u_1y_2 - y_1u_2) \quad (30)$$

where y_1 , u_1 , y_2 , and u_2 are the ray heights and slopes of the axial and principal rays and y_3 and u_3 are the data of the third ray; these data are determined at any component of the system where the specifications for all three rays are known. These equations can, for example, be used to determine the necessary component diameters to pass a bundle of rays which are A times the diameter of the axial bundle at a field angle B times the full-field angle. In Fig. 12, for the dashed rays $A = +0.5$ and -0.5 and $B = 1.0$. Another application of Eqs. (27) through (30) is to locate either a pupil or an aperture stop when the other is known.

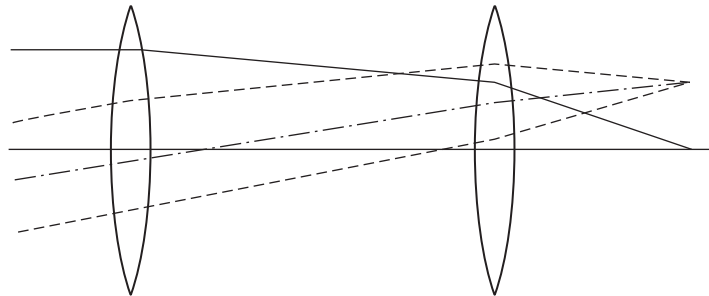


FIGURE 12

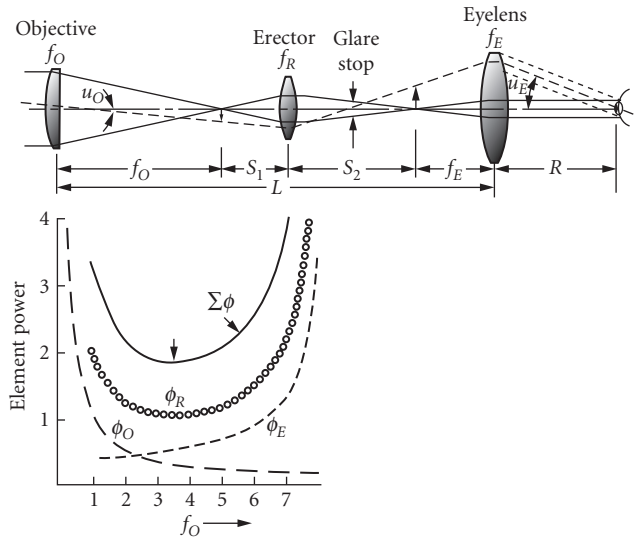


FIGURE 13

1.12 MINIMIZING COMPONENT POWER

The first-order layout may in fact determine the ultimate quality, cost, and manufacturability of the system. The residual aberrations in a system are a function of the component powers, relative apertures, and angular fields. The relationships are complex, but a good choice for a system layout is one which minimizes the sum of the (absolute) component powers, or possibly the sum of the (absolute) $y\phi$ product for all the components.

For example, in Fig. 13 the length, magnification, and the eye relief of the rifle scope are specified. There are five variables: three component powers and two spaces. This is one more variable than is necessary to achieve the specified characteristics. If we take the focal length of the objective component as the free variable, the component powers which satisfy the specifications can be plotted against the objective focal length, as in Fig. 13, and the minimum power arrangement is easily determined.

Minimizing the component powers will strongly tend to minimize the aberrations and also the sensitivity of the system to fabrication errors and misalignments. The *cost* of an optical element will vary with its diameter (or perhaps the square of the diameter) and also with the product of the diameter and the power. Thus, while first-order layout deals only with components, these relationships still apply reasonably well even when applied to components rather than elements. Minimizing the component powers does tend to reduce the cost on these grounds (and also because it tends to reduce the complexity of the components).

1.13 IS IT A REASONABLE LAYOUT?

A simple way to get a feel for the reasonableness of a layout is to make a rough scale drawing showing each component as single element. An element can be drawn as an equiconvex lens with radii which are approximately $r = 2(n - 1)f$; for an element with an index of 1.5 the radii equal the focal length. The elements should be drawn to the diameter necessary to pass the (suitably vignettted) off-axis

bundle of rays as well as the axial bundle. The on-axis and off-axis ray bundles should be sketched in. This will very quickly indicate which elements or components are the difficult ones. If the design is being started from scratch (as opposed to simply combining existing components), each component can be drawn as an achromat. The following section describes achromat layout, but for visual-spectrum systems it is often sufficient to assume that the positive (crown) element has twice the power of the achromat and the (negative) flint element has a power equal to that of the achromat. Thus an achromat may be sketched to the simplified, approximate prescription: $r_1 = -r_2 = f/2$ and $r_3 = \text{plano}$.

Any elements which are too fat must then be divided or “split” until they look “reasonable.” This yields a reasonable estimate of the required complexity of the system, even before the lens design process is begun.

If more or less standard design types are to be utilized for the components, it is useful to tabulate the focal lengths and diameters to get the (infinity) f -number of each component, and also its angular field coverage. The field coverage should be expressed both in terms of the angle that the object and image subtend from the component, and also the angle that the smaller of these two heights subtends as a function of the focal length (rather than as a function of that conjugate distance). This latter angle is useful because the coverage capability of a given design form is usually known in these terms, that is, h/f , rather than in finite conjugate terms. With this information at hand, a reasonable decision can be made as to the design type necessary to perform the function required of the component.

1.14 ACHROMATISM

The powers of the elements of an achromat can be determined from

$$\phi_A = \phi_{AB} V_A / (V_A - V_B) \quad (31)$$

$$\phi_B = \phi_{AB} V_B / (V_B - V_A) \quad (32)$$

$$= \phi_{AB} - \phi_A$$

where ϕ_{AB} is the power of the achromatic doublet and V_A is the Abbe V -value for the element whose power is ϕ_A , etc. For the visible spectral region $V = (n_d - 1)/(n_F - n_C)$; this can be extended to any spectral region by substituting the indices at middle, short, and long wavelengths for n_d , n_F , and n_C .

If the elements are to be spaced apart, and the back focus is B , then the powers and the spacing are given by

$$\phi_A = \phi_{AB} B V_A / (V_A B - V_B / \phi_{AB}) \quad (33)$$

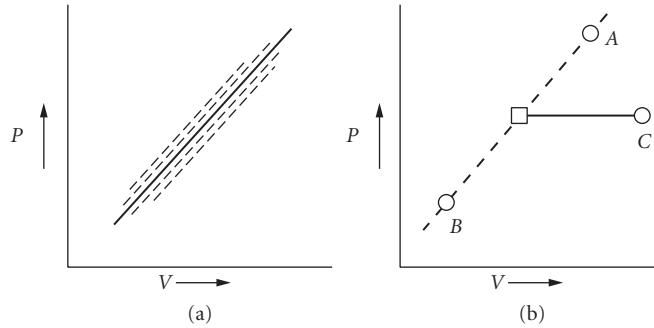
$$\phi_B = -\phi_{AB} V_B / B (V_A B - V_B / \phi_{AB}) \quad (34)$$

$$D = (1 - B \phi_{AB}) / \phi_A \quad (35)$$

For a complete system, the transverse *axial chromatic* aberration is the sum of $y^2 \phi / V u'_k$ for all the elements, where y is the height of the axial ray at the element and u'_k is the ray slope at the image. The *lateral color* is the sum of $y y_p \phi / V u'_k$, where y_p is the principal ray height.

The *secondary spectrum* is the sum of $y^2 \phi P / V u'_k$, where P is the partial dispersion, $P = (n_d - n_c) / (n_F - n_c)$. Summed over two elements, this leads to an expression for the longitudinal secondary spectrum of an achromatic doublet

$$\begin{aligned} \text{SS} &= f(P_B - P_A) / (V_A - V_B) \\ &= -f(\Delta P / \Delta V) \end{aligned} \quad (36)$$


FIGURE 14

This indicates that in order to eliminate secondary spectrum for a doublet, two glasses with identical partial dispersions [so that $(P_A - P_B)$ is zero] are required. A large difference in V -value is desired so that $(V_A - V_B)$ in the denominator of Eqs. (31) and (32) will produce reasonably low element powers. As indicated in the schematic and simplified plot of P versus V in Fig. 14a, most glasses fall into a nearly linear array, and $(\Delta P/\Delta V)$ is nearly a constant for the vast majority of glasses. The few glasses which are away from the "normal" line can be used for apochromats, but the ΔV for glass pairs with a small ΔP tends to be quite small. In order to get an exact match for the partial dispersions so that ΔP is equal to zero, two glasses can be combined to simulate a third, as indicated in Fig. 14b. For a unit power ($\phi = 1$) apochromatic triplet, the element powers can be found from

$$X = [V_A(P_B - P_C) + V_B(P_C - P_A)] / (P_B - P_A) \quad (37)$$

$$\phi_C = V_C / (V_C - X) \quad (38)$$

$$\phi_B = (1 - \phi_C)(P_C - P_A) V_B / [V_B(P_C - P_A) + V_A(P_B - P_C)] \quad (39)$$

$$\phi_A = 1 - \phi_B - \phi_C \quad (40)$$

1.15 ATHERMALIZATION

When the temperature of a lens element is changed, two factors affect its focus or focal length. As the temperature rises, all dimensions of the element are increased; this, by itself, would lengthen the focal length. However, the index of refraction of the lens material also changes with temperature. For many glasses the index rises with temperature; this effect tends to shorten the focal length.

The thermal change in the power of a thin element is given by

$$d\phi/dt = -\phi[a - (dn/dt)/(n-1)] \quad (41)$$

where dn/dt is the differential of index with temperature and a is the thermal expansion coefficient of the lens material. Then for a thin doublet

$$d\phi/dt = \phi_A T_A + \phi_B T_B \quad (42)$$

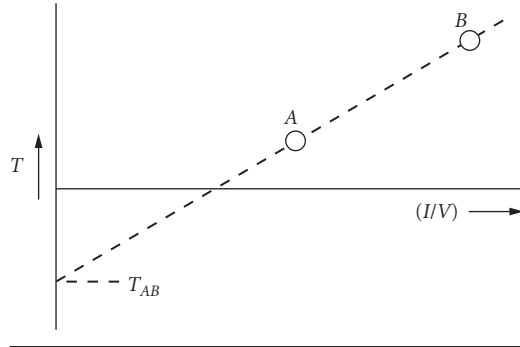


FIGURE 15

where

$$T = [-a + (dn/dt)/(n-1)] \quad (43)$$

and ϕ is the doublet power.

For an athermalized doublet (or one with some desired $d\phi/dt$) the element powers are given by

$$\phi_A = [(d\phi/dt) - \phi T_B] / (T_A - T_B) \quad (44)$$

$$\phi_B = \phi - \phi_A \quad (45)$$

To get an athermalized *achromatic* doublet, a plot of T against $(1/V)$ for all the glasses/materials under consideration is made. A line drawn between two glass points is extended to intersect the T axis as indicated in Fig. 15. Then the value of the $d\phi/dt$ for the achromatic doublet is equal to the doublet power times the value of T at which the line intersects the T axis. A pair of glasses with a large V -value difference and a small or zero T axis intersection is desirable.

An athermal achromatic triplet can be made with three glasses as follows:

$$\phi_A = \phi V_A (T_B V_B - T_C V_C) / D \quad (46)$$

$$\phi_B = \phi V_B (T_C V_C - T_A V_A) / D \quad (47)$$

$$\phi_C = \phi V_C (T_A V_A - T_B V_B) / D \quad (48)$$

$$D = V_A (T_B V_B - T_C V_C) + V_B (T_C V_C - T_A V_A) + V_C (T_A V_A - T_B V_B) \quad (49)$$

See also Chap. 8, "Thermal Compensation Techniques," by Philip J. Rogers and Michael Roberts.

NOTE: Figures 2, 3, 4, 5, 7, 8, 9, 10, 11, and 13 are adapted from W. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, New York, 1990. The remaining figures are adapted from *Critical Reviews of Lens Design*, W. Smith (Ed.), SPIE, vol. CR41, 1992.

ABERRATION CURVES IN LENS DESIGN

Donald C. O'Shea

*Georgia Institute of Technology
School of Physics
Atlanta, Georgia*

Michael E. Harrigan

*Harrigan Optical Design
Victor, New York*

2.1 GLOSSARY

H	ray height
NA	numerical aperture
OPD	optical path difference
P	petzval
S	sagittal
T	tangential
$\tan U$	slope

2.2 INTRODUCTION

Many optical designers use aberration curves to summarize the state of correction of an optical system, primarily because these curves give a designer important details about the relative contributions of individual aberrations to lens performance. Because a certain design technique may affect only one particular aberration type, these curves are more helpful to the lens designer than a single-value merit function. When a design is finished, the aberration curves serve as a summary of the lens performance and a record for future efforts. For applications such as photography, they are most useful because they provide a quick estimate of the effective blur circle diameter.

The aberration curves can be divided into two types: those that are expressed in terms of ray errors and those in terms of the optical path difference (OPD). OPD plots are usually plotted against the relative ray height in the entrance pupil. Ray errors can be displayed in a number of ways. Either the transverse or longitudinal error of a particular ray relative to the chief ray can be plotted as a function of the ray height in the entrance pupil. Depending upon the amount and type of aberration present, it is sometimes more appropriate to plot the longitudinal aberration as a function of field angle.

For example, astigmatism or field curvature is more easily estimated from field plots, described below. Frequently, the curves are also plotted for several wavelengths to characterize chromatic performance. Because ray error plots are the most commonly used format, this entry will concentrate on them.

2.3 TRANSVERSE RAY PLOTS

These curves can take several different forms, depending on the particular application of the optical system. The most common form is the transverse ray aberration curve. It is also called lateral aberration, or ray intercept curve (also referred to by the misleading term “rim ray plots”). These plots are generated by tracing fans of rays from a specific object point for finite object distances (or a specific field angle for an object at infinity) to a linear array of points across the entrance pupil of the lens. The curves are plots of the ray error at an evaluation plane measured from the chief ray as a function of the relative ray height in the entrance pupil (Fig. 1). For afocal systems, one generally plots angular aberrations, the differences between the tangents of exiting rays and their chief ray in image space.

If the evaluation plane is in the image of a perfect image, there would be no ray error and the curve would be a straight line coincident with the abscissa of the plot. If the curve were plotted for a different evaluation plane parallel to the image plane, the curve would remain a straight line but it would be rotated about the origin. Usually the aberration is plotted along the vertical axis, although some designers plot it along the horizontal axis.

The curves in Fig. 1 indicate a lens with substantial undercorrected spherical aberration as evidenced by the characteristic S-shaped curve. Since a change of the evaluation plane serves only to rotate the curve about the origin, a quick estimate of the aberrations of a lens can be made by reading the scale of the ray error axis (y axis) and mentally rotating the plot. For example, the blur spot can be estimated from the extent of a band that would enclose the curve a in Fig. 1, but a similar estimate could be made from the curves b or c , also.

The simplest form of chromatic aberration is axial color. It is shown in Fig. 2 in the presence of spherical aberration. Axial color is the variation of paraxial focus with wavelength and is seen as a difference in slope of the aberration curves at the origin as a function of wavelength. If the slopes of the curves at the origin for the end wavelengths are different, primary axial color is present. If primary axial color is corrected, then the curves for the end wavelengths will have the same slope at the origin. But if that slope differs from the slope of the curve for the center wavelength, then secondary axial color is present.

A more complex chromatic aberration occurs when the aberrations themselves vary with wavelength. Spherochromatism, the change of spherical aberration with wavelength, manifests itself as a difference in the shapes of the curves for different colors. Another curve that provides a measure of lateral color, an off-axis chromatic aberration, is described below.

For a point on the axis of the optical system, all ray fans lie in the meridional plane and only one plot is needed to evaluate the system. For off-axis object points, a second plot is added to evaluate a fan of skew rays traced in a sagittal plane. Because a skew ray fan is symmetrical across the meridional

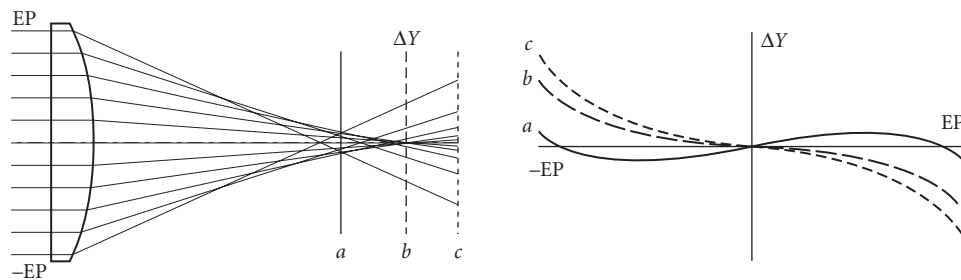


FIGURE 1 (Left) Rays exiting a lens are intercepted at three evaluation planes. (Right) Ray intercept curves plotted for the evaluation planes: (a) at the point of minimum ray error (circle of least confusion); (b) at the paraxial image plane; and (c) outside the paraxial image plane. (See also color insert.)

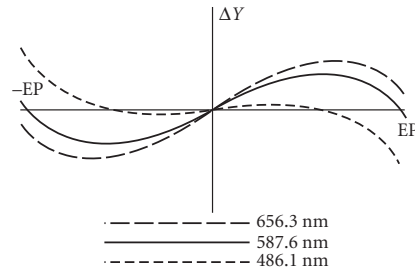


FIGURE 2 Meridional ray intercept curves of a lens with spherical aberration plotted for three colors. (See also color insert.)

plane, only one side of the curve is usually plotted. For all curves the plots are departures from the chief ray location in the evaluation plane (Fig. 3). (In the case of the on-axis point, the chief ray is coincident with the optical axis.) For systems of small-field coverage only two or three object points need to be analyzed, but for wide-angle systems, four or more field points may be necessary.

What can be determined most easily from a comparison between the meridional and sagittal planes is the amount of astigmatism in the image for that field point. When astigmatism is present, the image planes for the tangential and sagittal fans are located at different distances along the chief ray. This is manifested in the ray intercept curves by different slopes at the origin for the tangential and sagittal curves. In Fig. 3 the slopes at the origins of the two curves are different at both 70 percent and full field, indicating astigmatism at both field points. The fact that the difference in the slopes of these two curves has changed sign between the two field points indicates that at some field angle between 70 percent and full field, the slopes are equal and there is no astigmatism there. In addition, the variation of slopes for each curve as a function of field angle is evidence of field curvature.

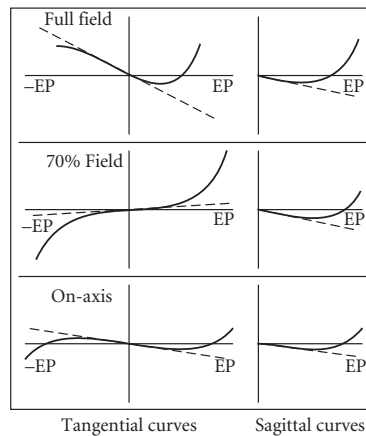


FIGURE 3 Evaluation of a lens on-axis and at two off-axis points. The reduction of the length of the curve with higher field indicates that the lens is vignetting these angles. The differences in slopes (dashed lines) at the origin between the meridional and skew curves indicate that the lens has astigmatism at these field angles. The variation in the slopes with field indicates the presence of field curvature. (See also color insert.)

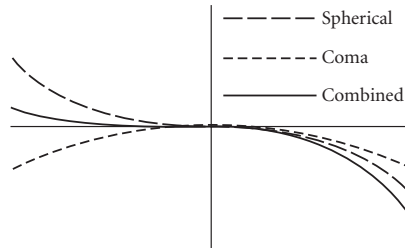


FIGURE 4 Ray intercept curve showing coma combined with spherical aberration. (See also color insert.)

The off-axis aberration of pure primary coma would be evident on these plots as a U-shaped curves for the meridional fan and sagittal fans, the tangential curve being three times larger than the sagittal curve. The “U” will be either upright or upside down depending on the sign of the coma. In almost all cases coma is combined with spherical to produce an S-shaped curve that elongates one of the arms of the “S” and shortens the other (Fig. 4).

The amount of vignetting can be determined from the ray intercept curves also. When it is present, the meridional curves get progressively shorter as the field angle is increased (Fig. 3), since rays at the edges of the entrance pupil are not transmitted. Taken from another perspective, ray intercept curves can also provide the designer with an estimate of how far a system must be stopped down to provide a required degree of correction.

2.4 FIELD PLOTS

The ray intercept curves provide evaluation for a limited number of object points—usually a point on the optical axis and several field points. The field plots present information on certain aberrations across the entire field. In these plots, the independent variable is usually the field angle and is plotted vertically and the aberration is plotted horizontally. The three field plots most often used are: distortion, field curvature, and lateral color. The first of these shows percentage distortion as a function of field angle (Fig. 5).

The second type of plot, field curvature, displays the tangential and sagittal foci as a function of object point or field angle (Fig. 6a). In some plots the Petzval surface, the surface to which the image would

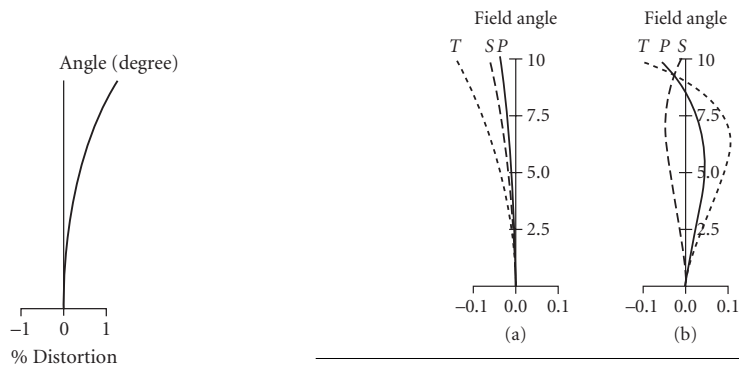


FIGURE 5 Field curve: distortion plot. The percentage distortion is plotted as a function of field angle. Note that the axis of the dependent variable is the horizontal axis. (See also color insert.)

FIGURE 6 Field curve: field curvature plot. The locations of the tangential *T* and sagittal *S* foci are plotted for a full range of field angles. The Petzval surface *P* is also plotted. The tangential surface is always three times farther from the Petzval surface than from the sagittal surface: (a) an uncorrected system and (b) a corrected system. (See also color insert.)

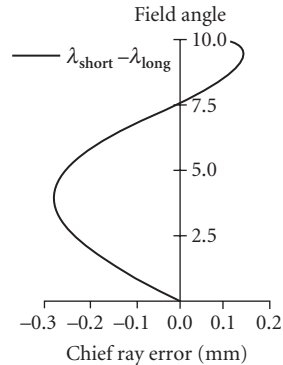


FIGURE 7 Field curve: lateral color plot. A plot of the transverse ray error between red and blue chief ray heights in the image plane for a full range of field angles. Here the distance along the horizontal axis is the color error in the image plane. (See also color insert.)

collapse if there were no astigmatism, is also plotted. This plot shows the amount of curvature in the image plane and amount of astigmatism over the entire field. In cases of corrected field curvature (Fig. 6*b*), this plot provides an estimate of the residual astigmatism between the axis and the corrected zone and an estimate of the maximum field angle at which the image possesses reasonable correction.

The last of the field curves provides information on color error as a function of field angle (Fig. 7). Lateral color, the variation of magnification with wavelength, is plotted as the difference between the chief ray heights at the red and blue wavelengths as a function of field angle. This provides the designer with an estimate of the amount of color separation in the image at various points in the field. In the transverse ray error curves, lateral color is seen as a vertical displacement of the end wavelength curves from the central wavelength curve at the origin.

Although there are other plots that can describe aberrations of optical systems (e.g., plot of longitudinal error as a function of entrance pupil height), the ones described here represent the ensemble that is used in most ray evaluation presentations.

2.5 ADDITIONAL CONSIDERATIONS

In many ray intercept curves the independent variable is the relative entrance pupil coordinate of the ray. However, for systems with high NA or large field of view, where the principal surface cannot be approximated by a plane, it is better to plot the difference between the tangent of the convergence angle of the chosen ray and the tangent of the convergence angle of the chief ray. This is because the curve for a corrected image will remain a straight line in any evaluation plane.¹ When plotted this way, the curves are called *H-tan U* curves.

Shifting the stop of an optical system has no effect on the on-axis curves. However, it causes the origin of the meridional curves of off-axis points to be shifted along the curve. In Fig. 8, the off-axis meridional curves are plotted for three stop positions of a double Gauss lens. The center curve (Fig. 8*b*) is plotted for a symmetrically located stop; the outer curves are plots when the stop is located at lens surfaces before and after the central stop.

It is usually sufficient to make a plot of the aberrations in the meridional and sagittal sections of the beam. The meridional section, defined for an optical system with rotational symmetry, is any plane containing the optical axis. It is sometimes called the tangential section. The sagittal section is a plane perpendicular to the meridional plane containing the chief ray. There are some forms of higher-order

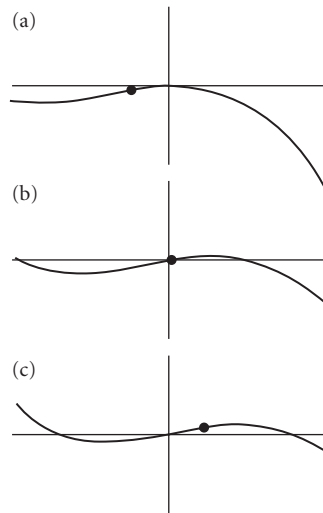


FIGURE 8 The effect of stop shifting on the meridional ray intercept curves of a double Gauss lens. (a) Stop located in front of the normal centrally located stop. (b) Stop at the normal stop position. (c) Stop behind the normal stop position. The dot locates the point on the curve where the origin is located for case (b). (See also color insert.)

coma that do not show in these sections.² In those cases where this aberration is suspected to be a problem, it may be helpful to look at a spot diagram generated from rays in all sections of the bundle.

For a rotationally symmetric system, only objects in a meridional plane need to be analyzed. Also for such systems, only meridional ray errors are possible for purely meridional rays. To observe certain coma types, it is a good idea to plot both the meridional and sagittal ray errors for sagittal rays. It is possible for the meridional section to show no coma and have it show only in the meridional error component of the sagittal fan,² but this aberration is normally small.

In addition to plots of the ray error in an evaluation plane, another aberration plot is one that expresses wavefront aberrations as an optical path difference from a spherical wavefront centered about the image point. These OPD plots are particularly useful for applications where the lens must be close to diffraction-limited.

2.6 SUMMARY

Aberration curves provide experienced designers with the information needed to enable them to correct different types of aberrations. Chromatic effects are much more easily classified from aberration curves also. In comparison to spot diagrams and modulation transfer function curves, the types of aberrations can be more easily seen and quantified. In the case of diffraction-limited systems, modulation transfer functions may provide better estimates of system performance.

2.7 REFERENCES

1. R. Kingslake, *Lens Design Fundamentals*, Academic Press, San Diego, 1978, p. 144.
2. F. D. Cruickshank and G. A. Hills, "Use of Optical Aberration Coefficients in Optical Design," *J. Opt. Soc. Am.* 50:379 (1960).

Douglas C. Sinclair

Sinclair Optics, Inc.

Fairport, New York

3.1 GLOSSARY

a	axial ray
b	chief ray
c	curvature
d, e, f, g	aspheric coefficients
efl	effective focal length
FN	focal ratio
$f()$	function
h	ray height
m	linear, lateral magnification
n	refractive index
PIV	paraxial (Lagrange) invariant
r	entrance pupil radius
t	thickness
u	ray slope
y	coordinate
z	coordinate
α	tilt about x (Euler angles)
β	tilt about y (Euler angles)
γ	tilt about z (Euler angles)
ϵ	displacement of a ray from the chief ray
μ	Buchdahl coefficients
κ	conic constant
ρ	radial coordinate
σ_1	spherical aberration

σ_2	coma
σ_3	astigmatism
σ_4	Petzval blur
σ_5	distortion

3.2 INTRODUCTION

The primary function of optical design software is to produce a mathematical description, or *prescription*, describing the shapes, locations, materials, etc., of an optical system that satisfies a given set of specifications. A typical optical design program contains three principal sections: *data entry*, *evaluation*, and *optimization*. The optical design programs considered here are to be distinguished from *ray-trace* programs, which are mainly concerned with evaluation, and *CAD* programs, which are mainly concerned with drawings. The essence of an optical design program is its optimization section, which takes a starting design and produces a new design that minimizes an *error function* that characterizes the system performance.

The first practical computer software for optical design was developed in the 1950s and 1960s.¹⁻⁴ Several commercially available programs were introduced during the 1970s, and development of these programs has continued through the 1980s to the present time. Although decades have passed since the introduction of optical design software, developments continue in optimization algorithms, evaluation methods, and user interfaces.

This chapter attempts to describe a typical optical design program. It is intended for readers that have a general background in optics, but who are not familiar with the capabilities of optical design software. We present a brief description of some of the most important mathematical concepts, but make no attempt to give a detailed development. We hope that this approach will give readers enough understanding to know whether an optical design program will be a useful tool for their own work.

Of course, many different programs are available, each with its own advantages and disadvantages. Our purpose is not to review or explain specific programs, but to concentrate on the basic capabilities. Some programs work better than others, but we make no quality judgment. In fact, we avoid reference, either explicit or implicit, to any particular program. The features and benefits of particular optical design programs are more than adequately described by software vendors, who are listed in optical industry buyer's guides.⁵

Figure 1 is a flowchart of a typical optical design project. Usually, the designer not only must enter the starting design and initial optimization data, but also must continually monitor the progress of the computer, modifying either the lens data or the optimization data as the design progresses to achieve the best solution. Even when the performance requirements are tightly specified, it is often necessary to change the error function during the design process. This occurs when the error function does not correlate with the desired performance sufficiently well, either because it is ill-conceived, or because the designer has purposefully chosen a simple error function to achieve improved speed.

The fact that there are alternate choices of action to be taken when the design is not good enough has led to two schools of thought concerning the design of an optical design program. The first school tries to make the interface between the designer and the program as smooth as possible, emphasizing the interactive side of the process. The second school tries to make the error function comprehensive, and the iteration procedure powerful, minimizing the need for the designer to intervene.

3.3 LENS ENTRY

In early lens design programs, lens entry was a "phase" in which the lens data for a starting design was read into the computer from a deck of cards. At that time, the numerical aspects of optical design on a computer were so amazing that scant attention was paid to the lens entry process. However, as the use of optical design software became more widespread, it was found that a great deal of a designer's time was spent punching cards and submitting new jobs, often to correct past mistakes. Many times, it turned out that the hardest part of a design job was preparing a "correct" lens deck!

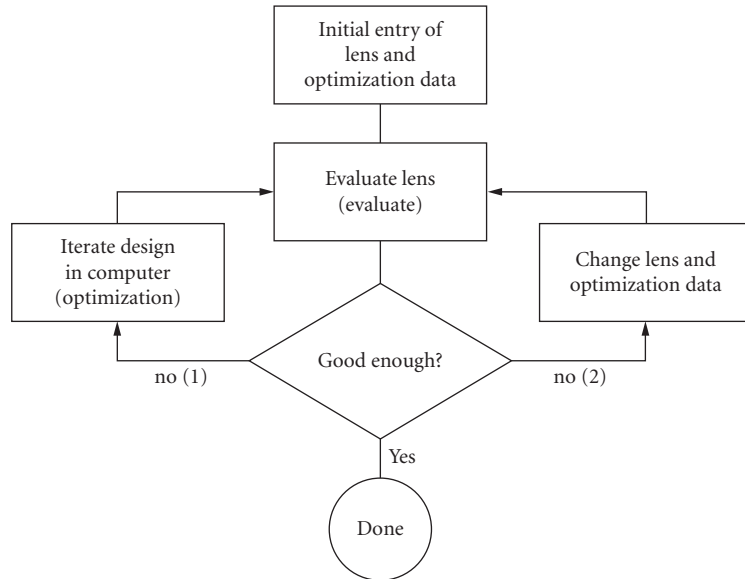


FIGURE 1 Flowchart for the lens design process. The action taken when a design is not satisfactory depends on how bad it is. The designer (or a design program) may change the lens data, or redefine the targets to ones that can be achieved.

Over the years, optical design programs have been expanded to improve the lens entry process, changing the function of this part of the program from simple lens entry to what might be called lens database management. A typical contemporary program provides on-line access to a library of hundreds of lenses, interactive editing, automatic lens drawings, and many features designed to simplify this aspect of optical design.

The lens database contains all items needed to describe the optical system under study, including not only the physical data needed to construct the system (curvatures, thicknesses, etc.), but also data that describe the conditions of use (object and image location, field of view, etc.). Some programs also incorporate optimization data in the lens database, while others provide separate routines for handling such data. In any case, the lens database is often the largest part of an optical design program.

The management of lens data in an optical design program is complicated by two factors. One is that there is a tremendous range of complexity in the types of systems that can be accommodated, so there are many different data items. The other is that the data are often described indirectly. A surface curvature may be specified, for example, by the required slope of a ray that emerges from the surface, rather than the actual curvature itself. Such a specification is called a *solve*, and is based on the fact that paraxial ray tracing is incorporated in the lens entry portion of most optical design programs.

It might seem curious that paraxial ray tracing is still used in contemporary optical design programs that can trace *exact* rays in microseconds. (The term *exact* ray is used in this chapter to mean a real skew ray. Meridional rays are treated as a special case of skew rays in contemporary software; there is not sufficient computational advantage to warrant a separate ray trace for them.) In fact, paraxial rays have some important properties that account for their incorporation in the lens database.

First, paraxial rays provide a linear system model that leads to analysis of optical systems in terms of simple bilinear transforms. Second, paraxial ray tracing does not fail. Exact rays can miss surfaces or undergo total internal reflection. Finally, paraxial rays determine the ideal performance of a lens. In a well-corrected lens, the aberrations are balanced so that the exact rays come to the image points defined by the paraxial rays, not the other way around.

Two basic types of data are used to describe optical systems. The first are the *general* data that are used to describe the system as a whole, and the other are the *surface* data that describe the individual surfaces and their locations. Usually, an optical system is described as an ordered set of surfaces,

beginning with an *object* surface and ending with an *image* surface (where there may or may not be an actual image). It is assumed that the designer knows the order in which rays strike the various surfaces. Systems for which this is not the case are said to contain *nonsequential* surfaces, which are discussed later.

General System Data

The general data used to describe a system include the aperture and field of view, the wavelengths at which the system is to be evaluated, and perhaps other data that specify evaluation modes, vignetting conditions, etc.

Aperture and Field of View The aperture and field of view of a system determine its conditions of use. The aperture is specified by the *axial ray*, which emerges from the vertex of the object surface and passes through the edge of the entrance pupil. The field of view is specified by the *chief ray*, which emerges from the edge of the object and passes through the center of the entrance pupil.

There are various common ways to provide data for the axial and chief rays. If the object is at an infinite distance, the entrance pupil radius and semifield angle form a convenient way to specify the axial and chief rays. For finite conjugates, the numerical aperture in object space and the object height are usually more convenient.

Some programs permit the specification of paraxial ray data by image-space quantities such as the *f*-number and the image height, but such a specification is less desirable from a computational point of view because it requires an iterative process to determine initial ray-aiming data.

Wavelengths It is necessary to specify the wavelengths to be used to evaluate polychromatic systems. Three wavelengths are needed to enable the calculation of primary and secondary chromatic aberrations. More than three wavelengths are required to provide an accurate evaluation of a typical system, and many programs provide additional wavelengths for this reason. There has been little standardization of wavelength specification. Some programs assume that the first wavelength is the central wavelength, while others assume that it is one of the extreme wavelengths; some require wavelengths in micrometers, while others in nanometers.

Other General Data Several other items of general data are needed to furnish a complete lens description, but there is little consistency between programs on how these items are treated, or even what they are. The only one that warrants mention here is the aperture stop. The *aperture stop* is usually defined to be the surface whose aperture limits the angle of the axial ray. Once the aperture stop surface is given, the positions of the paraxial pupils are determined by the imaging properties of the system. Since the aperture and field of view are determined formally by the paraxial pupils, the apertures are not associated with the exact ray behavior.

The “vignetting factor” is used to account for the differences between paraxial and exact off-axis ray heights at apertures. In particular, the vignetting factor provides, in terms of fractional (paraxial) coordinates, the data for an exact ray that grazes the apertures of a system. Typically, there is an upper, lower, and skew vignetting factor. The details of how such factors are defined and handled are program dependent.

Surface Data

Surface Location There are two basic ways to specify the location of surfaces that make up a lens. One is to specify the position of a surface relative to the immediately preceding surface. The other is to specify its position relative to some fixed surface (for example, the first surface). The two ways lead to what are called *local* and *global* coordinates, respectively. For ordinary lenses consisting of a series of rotationally symmetric surfaces centered on an optical axis, local coordinates are more convenient, but for systems that include reflectors, tilted, and/or decentered surfaces, etc., global

coordinates are simpler. Internally, optical design programs convert global surface data to local coordinates for speed in ray tracing.

Most optical design programs use a standard coordinate system and standard sign conventions, although there are exceptions.⁶ Each surface defines a local right-handed coordinate system in which the z axis is the symmetry axis and the yz plane is the meridional plane. The local coordinate system is used to describe the surface under consideration and also the origin of the next coordinate system. Tilted elements are described by an Euler-angle system in which α is a tilt around the x axis, β is a tilt around the y axis, and γ is a tilt around the z axis. Since tilting and decentering operations are not commutative; some data item must be provided to indicate which comes first.

Surface Profile Of the various surfaces used in optical systems, the most common by far is the rotationally symmetric surface, which can be written as⁷

$$z = \frac{cr^2}{1 + \sqrt{1 - c^2(\kappa + 1)r^2}} + dr^4 + er^6 + fr^8 + gr^{10}$$

$$r = \sqrt{x^2 + y^2}$$

c is the curvature of the surface; κ is the conic constant; and d , e , f , and g are aspheric constants. The use of the above equation is almost universal in optical design programs. The description of conic surfaces in terms of a conic constant κ instead of the eccentricity e used in the standard mathematical literature allows spherical surfaces to be specified as those with no conic constant. (The conic constant is minus the square of the eccentricity.)

Although aspheric surfaces include all surfaces that are not spherical, from a design standpoint there is a demarcation between “conic” aspheres and “polynomial” aspheres described using the coefficients d , e , f , and g . Rays can be traced analytically through the former, while the latter require numerical iterative methods.

Many optical design programs can handle surface profiles that are more complicated than the above, including cylinders, torics, splines, and even general aspheres of the form $z = f(x, y)$, where $f(x, y)$ is an arbitrary polynomial. The general operation of an optical design program, however, can be understood by considering only the rotationally symmetric surfaces described here.

As mentioned above, the importance of paraxial rays in optical system design has led to the indirect specification of lens data, using *solves*, as they are called, which permit a designer to incorporate the basic paraxial data describing a lens with the lens itself, rather than having to compute and optimize the paraxial performance of a lens as a separate task. Considering the j th surface of an optical system, let

$$\begin{aligned} y_j &= \text{ray height on surface} \\ u_j &= \text{ray slope on image side} \\ c_j &= \text{curvature of surface} \\ n_j &= \text{refractive index on image side} \\ t_j &= \text{thickness on image side} \end{aligned}$$

The paraxial ray trace equations can then be written as⁸

$$y_j = y_{j-1} + t_{j-1}u_{j-1}$$

$$n_j u_j = n_{j-1}u_{j-1} - y_j c_j (n_j - n_{j-1})$$

These equations can be inverted to give the curvatures and thicknesses in terms of the ray data. We have

$$c_j = \frac{n_{j-1}u_{j-1} - n_j u_j}{y_j(n_j - n_{j-1})}$$

$$t_j = \frac{y_{j+1} - y_j}{u_j}$$

The specification of curvatures and thicknesses by solves is considered to be on an equal basis with the direct specification of these items. The terminology used to specify solves is that the solves used to determine thickness are called *height solves*, and the solves used to determine curvature are called *angle solves*. Often, an axial ray height solve on the last surface is used to automatically locate the paraxial image plane, a chief ray height solve on the same surface to locate the exit pupil, and an axial ray angle solve is used to maintain a given focal length (if the entrance pupil radius is fixed). In some programs, additional types of solves are allowed, such as center of curvature solves, or aperture solves.

Of course, specifying lens data in terms of paraxial ray data means that whenever any lens data is changed, two paraxial rays must be traced through the system to resolve any following data that are determined by a solve. In an optical design program, this function is performed by a *lens setup* routine, which must be efficiently coded, since it is executed thousands of times in even a small design project.

Other functions of the lens setup routine are to precalculate values that are needed for repetitive calculations, such as refractive indices, rotation and translation matrices, etc. Many programs have the capability of specifying certain data items to be equal to (\pm) the value of the corresponding item on a previous surface. These are called *pickups*, and are needed for optimization of systems containing mirrors, as well as maintaining special geometrical relationships. Programs that lack pickups usually have an alternate means for maintaining the required linking between data items. Like solves, pickups are resolved by the lens setup routine, although they do not use paraxial data.

Other Surface Data A variety of other data is required to specify surfaces. Most important are apertures, discussed below, and refractive indices. Refractive indices are usually given by specifying the name of a catalog glass. In the lens setup routine, the actual refractive indices are calculated using an index interpolation formula and coefficient supplied by the glass manufacturer, together with the design wavelengths stored with the lens data. Other surface-related items include phase data for diffractive surfaces, gradient-index data, holographic construction data, and coatings.

Apertures have a somewhat obscure status in many optical design programs. Although apertures have a major role to play in determining the performance of a typical system, they do not usually appear directly in optimization functions. Instead, apertures are usually controlled in optimization by targets on the heights of rays that define their edges. If an aperture is specified directly, it will block rays that pass outside of it and cause typical optimization procedures to become unstable. Accordingly, some programs ignore apertures during optimization. Other programs allow the apertures to be determined by a set of exact “reference rays” that graze their extremities.

Nonsequential Surfaces In some optical systems, it is not possible to specify the order in which a ray will intersect the surfaces as it progresses through the system. The most common examples of such systems are prisms such as the corner-cube reflector, where the ordering of surfaces depends on the entering ray coordinates. Other examples of nonsequential surfaces include light pipes and a variety of nonimaging concentrators. Nonsequential surfaces can be accommodated by many optical design programs, but for the most part they are not “designed” using the program, but rather are included as a subsystem used in conjunction with another part of the system that is the actual

system being designed. Data specification of nonsequential surfaces is more complicated than ordinary systems, and ray tracing is much slower, since several surfaces must be investigated to see which surface is the one actually traversed by a given ray.

Lens Setup

Whenever the lens entry process is completed, the lens must be “set up.” Pickup constraints must be resolved. If the system contains an internal aperture stop, the position of the entrance pupil must be determined. Then paraxial axial and chief rays must be traced through the system so that surface data specified by solves can be computed. Depending on the program, a variety of other data may be precomputed for later use, including aperture radii, refractive indices, and various paraxial constants.

The lens setup routine must be very efficient, since it is the most heavily used code in an optical design program. In addition to running whenever explicit data entry is complete, the code is also executed whenever the lens is modified internally by the program, such as when derivatives are computed during optimization, or when configurations are changed in a multiconfiguration system. Typically, lens setup takes milliseconds (at most), so it is not noticed by the user, other than through its effects.

Programming Considerations

In writing an optical design program, the programmer must make a number of compromises between speed, size, accuracy, and ease of use. These compromises affect the usefulness of a particular program for a particular application. For example, a simple, fast, small program may be well suited to a casual user with a simple problem to solve, but this same program may not be suited for an experienced designer who routinely tackles state-of-the-art problems.

The lens entry portion of an optical design program shows, more than any other part, the difference in programming models that occurred during the 1980s. Before the 1980s, most application programs were of a type called *procedural* programs. When such a program needs data, it requests it, perhaps from a file or by issuing a prompt for keyboard input. The type of data needed is known, and the program is only prepared to accept that kind of data at any given point. Although the program may branch through different paths in response to the data it receives, the program is responsible for its own execution.

With the popularization in the 1980s of computer systems that use a mouse for input the model for an application program changed from the procedural model described above to what is called an *event-driven* model. An event-driven program has a very simple top-level structure consisting of an initialization section followed by an infinite loop usually called something like the *main event loop*. The function of the main event loop is to react to user-initiated interrupts (such as pressing a key, or clicking a mouse button), dispatching such events to appropriate processing functions. In such a program, the user controls the execution, unlike a procedural program, where the execution controls the user.

An event-driven program usually provides a better user interface than a procedural program. Unfortunately, most optical design programs were originally written as procedural programs, and it is difficult to convert a procedural program into an event-driven program by “patching” it. Usually it is easier to start over. In addition, it is harder to write an event-driven program than a procedural program, because the event-driven program must be set up to handle a wide variety of unpredictable requests received at random times. Of course, it is this very fact that makes the user interface better. There is an aphorism sometimes called the “conservation of complexity,” which states that the simpler a program is to use, the more complicated the program itself must be.

The data structures used to define lens data in an optical design program may have a major impact on its capabilities. For example, for various reasons it is usually desirable to represent a lens in a computer program as an array of surfaces. If the maximum size of the array is determined at compile time, then the maximum size lens that can be accommodated is built into the program.

As more data items are added to the surface data, the space required for storage can become unwieldy. For example, it takes about 10 items of real data to specify a holographic surface. If every surface were allowed to be a hologram, then 10 array elements would have to be reserved for each surface's holographic data. On the other hand, in most systems, the elements would never be used, so the data structure would be very inefficient. To avoid this, a more complicated data structure could be implemented in which only one element would be devoted to holograms, and this item would be used as an index into a separate array containing the actual holographic data. Such an array might have a sufficient number of elements to accommodate up to, say, five holograms, the maximum number expected in any one system.

The preceding is a simple example of how the data structure in an optical design program can grow in complexity. In fact, in a large optical design program the data structure may contain all sorts of indices, pointers, flags, etc., used to implement special data types and control their use. Managing this data while maintaining its integrity is a programming task of a magnitude often greater than the numerical processing that takes place during optical design.

Consider, for example, the task of deleting a surface from a lens. To do this, the surface data must of course be deleted, and all of the higher-numbered surfaces renumbered. But, in addition, the surface must be checked to see whether it is a hologram and, if so, the holographic data must also be deleted and that data structure "cleaned up." All other possible special data items must be tested and handled similarly. Then all the renumbered surfaces must be checked to see if any of the "pick up" data from a surface that has been renumbered, and the reference adjusted accordingly. Then other data structures such as the optimization files must be checked to see if they refer to any of the renumbered surfaces, and appropriate adjustments made. There may be several other checks and adjustments that must also be carried out.

Related to the lens entry process is the method used to store lens data on disc. Of course, lens data are originally provided to a program in the form of text (e.g., "TH 1.0"). The program *parses* this data to identify its type (a thickness) and value (1.0). The results of the parsing process (the binary values) are stored in appropriate memory locations (arrays). To store lens data on disc, early optical design programs used the binary data to avoid having to reparse it when it was recovered. However, the use of binary files has decreased markedly as computers have become fast enough that parsing lens input does not take long. The disadvantages of binary files are that they tend to be quite large, and usually have a structure that makes them obsolete when the internal data structure of the program is changed. The alternative is to store lens data as text files, similar in form to ordinary keyboard input files.

3.4 EVALUATION

Paraxial Analysis

Although the lens setup routine contains a paraxial ray trace, a separate paraxial ray trace routine is used to compute data for display to the user. At a minimum, the paraxial ray heights and slopes of the axial and chief ray are shown for each surface, in each color, and in each configuration.

The equations used for paraxial ray tracing were described in the previous section. Although such equations become exact only for "true" paraxial rays that are infinitesimally displaced from the optical axis, it is customary to consider paraxial ray data to describe "formal" paraxial rays that refract at the tangent planes to surfaces, as shown in Fig. 2. Here, the ray ABC is a paraxial ray that provides a first-order approximation to the exact ray ADE. Not only does the paraxial ray refract at the (imaginary) tangent plane BVP, but also it bends a different amount from the exact ray.

In addition to the computation of ray heights and slopes for the axial and chief ray, various paraxial constants that characterize the overall system are computed. The particular values computed depend on whether the system is *focal* (finite image distance) or *afocal* (image at infinity). For focal systems, the quantities of interest are (at a minimum) the focal length f , the f -number FN, the paraxial (Lagrange) invariant PIV, and the transverse magnification m . It is desirable to compute such

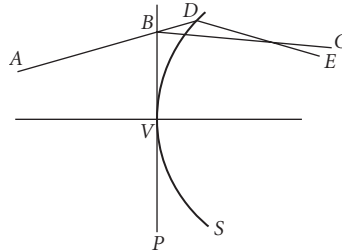


FIGURE 2 Showing the difference between a paraxial ray and a real ray. The paraxial ray propagates along ABC , while the real ray propagates along ADE .

quantities in a way that does not depend on the position of the final image surface. Let the object height be h , the entrance pupil radius be r , the axial ray data in object and image spaces be y , u and y' , u' , the chief ray data be \bar{y} , \bar{u} and \bar{y}' , \bar{u}' , and the refractive indices be n and n' .

The above-mentioned paraxial constants are then given by

$$\text{efl} = \frac{-rh}{\bar{u}'r + u'h}$$

$$\text{FN} = -\frac{n}{2n'u'}$$

$$\text{PIV} = n(\bar{y}u - y\bar{u})$$

$$m = \frac{nu}{n'u'}$$

In addition to the paraxial constants, most programs display the locations of the entrance and exit pupils, which are easily determined using chief-ray data. Surprisingly, most optical design programs do not explicitly show the locations of the principal planes. In addition, although most programs have the capability to display “ $y - \bar{y}$ ” plots, few have integrated this method into the main lens entry routine.

Aberrations

Although most optical designs are based on exact ray data, virtually all programs have the capability to compute and display first-order chromatic aberrations and third-order monochromatic (Seidel) aberrations. Many programs can compute fifth-order aberrations as well. The form in which aberrations are displayed depends on the program and the type of system under study, but as a general rule, for focal systems aberrations are displayed as equivalent ray displacements in the paraxial image plane.

In the case of the chromatic aberrations, the primary and secondary chromatic aberration of the axial and chief rays are computed. In a system for which three wavelengths are defined, the primary aberration is usually taken between the two outer wavelengths, and the secondary aberration between the central and short wavelengths.

The Seidel aberrations are computed according to the usual aberration polynomial. If we let ϵ be the displacement of a ray from the chief ray, then

$$\epsilon_y = \epsilon_{3y} + \epsilon_{5y} + \dots$$

$$\epsilon_x = \epsilon_{3x} + \epsilon_{5x} + \dots$$

For a relative field height h and normalized entrance pupil coordinates r and θ , the third-order terms are

$$\epsilon_{3y} = \sigma_1 \cos \theta r^3 + \sigma_2 (2 + \cos 2\theta) r^2 h + (3\sigma_3 + \sigma_4) \cos \theta r h^2 + \sigma_5 h^3$$

$$\epsilon_{3x} = \sigma_1 \sin \theta r^3 + \sigma_2 \sin 2\theta r^2 h + (\sigma_3 + \sigma_4) \sin \theta r h^2$$

The interpretation of the coefficients is generally as follows, but several optical design programs display tangential coma, rather than the sagittal coma indicated in the table.

σ_1	Spherical aberration
σ_2	Coma
σ_3	Astigmatism
σ_4	Petzval blur
σ_5	Distortion

The fifth-order terms are

$$\epsilon_{5y} = \mu_1 \cos \theta r^5 + (\mu_2 + \mu_3 \cos 2\theta) r^4 h + (\mu_4 + \mu_6 \cos^2 \theta) \cos \theta r^2 h^2$$

$$+ (\mu_7 + \mu_8 \cos 2\theta) r^2 h^3 + \mu_{10} \cos \theta r h^4 + \mu_{12} h^5$$

$$\epsilon_{5x} = \mu_1 \sin \theta r^5 + \mu_3 \sin 2\theta r^4 h + (\mu_5 + \mu_6 \cos^2 \theta) \sin \theta r^3 h^2$$

$$+ \mu_9 \sin 2\theta r^2 h^3 + \mu_{11} \sin \theta r h^4$$

These equations express the fifth-order aberration in terms of the Buchdahl μ coefficients. In systems for which the third-order aberrations are corrected, the following identities exist:

$$\mu_2 = \frac{3}{2} \mu_3$$

$$\mu_4 = \mu_5 + \mu_6$$

$$\mu_7 = \mu_8 + \mu_9$$

μ_1	Spherical aberration
μ_3	Coma
$(\mu_{10} - \mu_{11})/4$	Astigmatism
$(5\mu_{11} - \mu_{10})/4$	Petzval blur

$\mu_4 + \mu_6$	Tangential oblique spherical aberration
μ_5	Sagittal oblique spherical aberration
$\mu_7 + \mu_8$	Tangential elliptical coma
μ_9	Sagittal elliptical coma
μ_{12}	Distortion

Some programs display only the aberrations that have corresponding third-order coefficients, omitting oblique spherical aberration and elliptical coma.

The formulas needed to calculate the chromatic and third-order aberrations are given in the *U.S. Military Handbook of Optical Design*. The formulas for calculating the fifth-order aberrations are given in Buchdahl's book.⁹

Aberration coefficients are useful in optical design because they characterize the system in terms of its symmetries, allow the overall performance to be expressed as a sum of surface contributions, and are calculated quickly. On the negative side, aberration coefficients are not valid for systems that have tilted and decentered elements for systems that cover an appreciable field of view, and the accuracy of aberration coefficients in predicting performance is usually inadequate. Moreover, for systems that include unusual elements like diffractive surfaces and gradient index materials, the computation of aberration coefficients is cumbersome at best.

Ray Tracing

Exact ray tracing is the foundation of an optical design program, serving as a base for both evaluation and optimization. From the programmer's standpoint, the exact ray-trace routines must be accurate and efficient. From the user's viewpoint, the data produced by the ray-trace routines must be accurate and comprehensible. Misunderstanding the meaning of ray-trace results can be the source of costly errors in design.

To trace rays in an optical design program, it is necessary to understand how exact rays are specified. Although the details may vary from one program to the next, many programs define a ray by a two-step process. In the first step, an object point is specified. Once this has been done, all rays are assumed to originate from this point until a new object point is specified. The rays themselves are then specified by aperture coordinates and wavelength.

Exact ray starting data is usually normalized to the object and pupil coordinates specified by the axial and chief rays. That is, the aperture coordinates of a ray are specified as a fractional number, with 0.0 representing a point on the vertex of the entrance pupil, and 1.0 representing the edge of the pupil. Field angles or object heights are similarly described, with 0.0 being a point on the axis, and 1.0 being a point at the edge of the field of view.

Although the above normalization is useful when the object plane is at infinity, it is not so good when the object is at a finite distance and the numerical aperture in object space is appreciable. Then, fractional aperture coordinates should be chosen proportional to the direction cosines of rays leaving an object point. There are two reasons for this. One is that it allows an object point to be considered a point source, so that the amount of energy is proportional to the "area" on the entrance pupil. The other is that for systems without pupil aberrations, the fractional coordinates on the second principal surface should be the same as those on the first principal surface. Notwithstanding these requirements, many optical design programs do not define fractional coordinates proportional to direction cosines.

It is sometimes a point of confusion that the aperture and field of view of a system are specified by paraxial quantities, when the actual performance is determined by exact rays. In fact, the paraxial specifications merely establish a normalization for exact ray data. For example, in a real system the field of view is determined not by the angle of the paraxial chief ray, but by the angle at which exact rays blocked by actual apertures just fail to pass through the system. Using an iterative procedure, it is not too hard to find this angle, but because of the nonlinear behavior of Snell's law, it does not provide a convenient reference point.

There are two types of exact rays: *ordinary* or *lagrangian* rays, and *iterated* or *hamiltonian* rays. The designation of rays as lagrangian or hamiltonian comes from the analogy to the equations of motion of a particle in classical mechanics. Here we use the more common designation as ordinary or iterated rays. An ordinary ray is a ray that starts from a known object point in a known direction. An iterated ray also starts from a known object point, but its direction is not known at the start. Instead, it is known that the ray passes through some known (nonconjugate) point inside the system, and the initial ray direction is determined by an iterative procedure.

Iterated rays have several applications in optical design programs. For example, whenever a new object point is specified, it is common to trace an iterated ray through the center of the aperture stop (or some other point) to serve as a reference ray, or to trace several iterated rays through the edges of limiting apertures to serve as reference rays. In fact, many programs use the term *reference ray* to mean iterated ray (although in others, reference rays are ordinary rays). Iterated rays are traced using differentially displaced rays to compute corrections to the initial ray directions. Because of this, they are traced slower than ordinary rays. On the other hand, they carry more information in the form of the differentials, which is useful for computing ancillary data like field sags.

Reference rays are used as base rays in the interpretation of ordinary ray data. For example, the term *ray displacement* often refers to the difference in coordinates on the image surface of a ray from those of the reference ray. Similarly, the *optical path difference* of a ray may compare its phase length to that of the corresponding reference ray. The qualifications expressed in the preceding sentences indicate that the definitions are not universal. Indeed, although the terms *ray displacement* and *optical path difference* are very commonly used in optical design, they are not precisely defined, nor can they be. Let us consider, for example, the optical path difference.

Imagine a monochromatic wavefront from a specified object point that passes through an optical system. Figure 3 shows the wavefront *PE* emerging in image space, where it is labeled “actual wavefront.” Because of aberrations, an ordinary ray perpendicular to the actual wavefront will not intersect the final image surface at the ideal image point *I*, but at some other point *Q*. The optical path difference (OPD) may be defined as the optical path measured along the actual ray between the actual wavefront and a reference sphere centered on the ideal image point.

Unfortunately, the ideal image point is not precisely defined. In the figure, it is shown as the intersection of the reference ray with the image surface, but the reference ray itself may not be precisely defined.

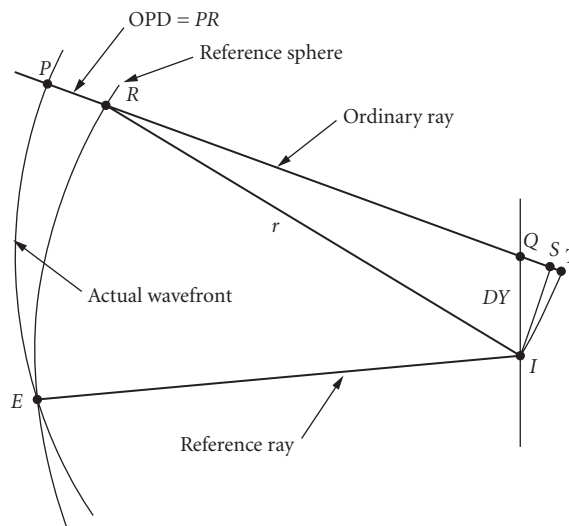


FIGURE 3 The relation between ray trajectories and optical path difference (OPD). See text.

Is it the ray through the center of the aperture stop, or perhaps the ray through the center of the actual vignetted aperture? These two definitions will result in different reference rays, and correspondingly different values for the optical path difference. In fact, in many practical applications neither definition is used, and the actual ideal image point is defined to be the one that minimizes the variance of the optical path difference (and hence maximizes the peak intensity of the diffraction image).

Moreover, the figure shows that even if the ideal image point is precisely defined, the value of the optical path difference depends on the point E where the actual wavefront intersects the reference sphere. For the particular point shown, the optical path difference is the optical length along the ordinary ray from the object point to the point T , less the optical length along the reference ray from the object point to the point I . As the radius of the reference sphere is increased, the point T merges with the point S , where a perpendicular from the ideal image point intersects the ordinary ray.

The above somewhat extended discussion is meant to demonstrate that even “well-known” optical terms are not always precisely defined. Not surprisingly, various optical design programs in common use produce different values for such quantities. There has been little effort to standardize the definitions of many terms, possibly because one cannot legislate physics. In any case, it is important for the user of an optical design program to understand precisely what the program is computing.

Virtually all optical design programs can trace single rays and display the ray heights and direction cosines on each surface. Other data, such as the path length, angles of incidence and refraction, and direction of the normal vector, are also commonly computed. Another type of ray-data display that is nearly universal is the ray-intercept curve, which shows ray displacement on the final image surface versus (fractional) pupil coordinates. A variation plots optical path difference versus pupil coordinates.

In addition to the uncertainty concerning the definition of ray displacement and optical path difference, there are different methods for handling the pupil coordinates. Some programs use entrance pupil coordinates, while others use exit pupil coordinates. In most cases, there is not a significant difference, but in the case of systems containing cylindrical lenses, for example, there are major differences.

Another consideration relating to ray-intercept curves is the way in which vignetting is handled. This is coupled to the way the program handles apertures. As mentioned before, apertures have a special status in many optical design programs. Rays can be blocked by apertures, but this must be handled as a special case by the program, because there is nothing inherent in the ray-trace equations that prevents a blocked ray from being traced, in contrast to a ray that misses a surface or undergoes total internal reflection.

Even though a surface may have a blocking aperture, it may be desirable to let the ray trace proceed anyway. As mentioned before, blocking rays in optimization can produce instabilities that prevent convergence to a solution even though all the rays in the final solution are contained within the allowed apertures. Another situation where blocking can be a problem concerns central obstructions. In such systems, the reference ray may be blocked by an obstruction, so its data are not available to compute the displacement or optical path difference of an ordinary ray (which is not blocked). The programmer must anticipate such situations and build in the proper code to handle them.

In the case of ray-intercept curves, it is not unusual for programs to display data for rays that are actually blocked by apertures. The user is expected to know which rays get through, and ignore the others, a somewhat unreasonable expectation. The justification for allowing it is that the designer can see what would happen to the rays if the apertures were increased.

In addition to ray-intercept curves, optical design programs usually display field sag plots showing the locations of the tangential and sagittal foci as a function of field angle and distortion curves. In the case of distortion, there is the question of what to choose as a reference height. It is generally easiest to refer distortion to the paraxial chief ray height in the final image surface, but in many cases it is more meaningful to refer it to the centroid height of a bundle of exact rays from the same object point. Again, it is important for the user to know what the program is computing.

Spot-Diagram Analysis

Spot diagrams provide the basis for realistic modeling of optical systems in an optical design program. In contrast to simple ray-trace evaluation, which shows data from one or a few rays, spot

diagrams average data from hundreds or thousands of rays to evaluate the image of a point source. Notwithstanding this, it should be understood that the principal purpose of an optical design program is to design a system, not to simulate its performance. It is generally up to the designer to understand whether or not the evaluation model of a system is adequate to characterize its real performance, and the prudent designer will view unexpected results with suspicion.

From a programmer's point of view, the most difficult task in spot-diagram analysis is to accurately locate the aperture of the system. For systems that have rotational symmetry, this is not difficult, but for off-axis systems with vignetted apertures it can be a challenging exercise. However, the results of image evaluation routines are often critically dependent on effects that occur near the edges of apertures, so particular care must be paid to this problem in writing optical design software. Like many other aspects of an optical design program, there is a trade-off between efficiency and accuracy.

A spot diagram is an assemblage of data describing the image-space coordinates of a large number of rays traced from a single object point. The data may be either monochromatic or polychromatic. Each ray is assigned a weight proportional to the fractional energy that it carries. Usually, the data saved for each ray include its xyz coordinates on the image surface, the direction cosines klm , and the optical path length or optical path difference from the reference ray. The ray coordinates are treated statistically to calculate root-mean-square spot sizes. The optical path lengths yield a measure of the wavefront quality, expressed through its variance and peak-to-valley error.

To obtain a spot diagram, the entrance pupil must be divided into cells, usually of equal area. Although for many purposes the arrangement of the cells does not matter, for some computations (e.g., transfer functions) it is advantageous to have the cells arranged on a rectangular grid. To make the computations have the proper symmetry, the grid should be symmetrical about the x and y axis. The size of the grid cells determines the total number of rays in the spot diagram.

In computing spot diagrams, the same considerations concerning the reference point appear as for ray fans. That is, it is possible to define ray displacements with respect to the chief-ray, the paraxial ray height, or the centroid of the spot diagram. However, for spot diagrams it is most common to use the centroid as the reference point, both because many image evaluation computations require this definition, and also because the value for the centroid is readily available from the computed ray data.

$$a_i = \epsilon_{xi} = \text{ray displacement in the } x \text{ direction}$$

$$b_i = \epsilon_{yi} = \text{ray displacement in the } y \text{ direction}$$

$$c_i = k_i/m_i = \text{ray slope in the } x \text{ direction}$$

$$d_i = l_i/m_i = \text{ray slope in the } y \text{ direction}$$

$$w_i = \text{weight assigned to ray}$$

The displacements of rays on a plane shifted in the z direction from the nominal image plane by an amount Δz are given by

$$\delta x_i = a_i + b_i \Delta z$$

$$\delta y_i = c_i + d_i \Delta z$$

If there are n rays, the coordinates of the centroid of the spot diagram are

$$\delta \bar{x} = \frac{1}{W} \sum_{i=1}^n w_i \delta x_i = A + B \Delta z$$

$$\delta \bar{y} = \frac{1}{W} \sum_{i=1}^n w_i \delta y_i = C + D \Delta z$$

where W is a normalizing constant that ensures that the total energy in the image adds up to 100 percent, and

$$A = \frac{1}{W} \sum_{i=1}^n w_i a_i$$

$$B = \frac{1}{W} \sum_{i=1}^n w_i b_i$$

$$C = \frac{1}{W} \sum_{i=1}^n w_i c_i$$

$$D = \frac{1}{W} \sum_{i=1}^n w_i d_i$$

The mean-square spot size can then be written as

$$\text{MSS} = \frac{1}{W} \sum_{i=1}^n w_i \{(\delta x_i - \delta \bar{x})^2 + (\delta y_i - \delta \bar{y})^2\}$$

Usually, the root-mean-square (rms) spot size, which is the square root of this quantity, is reported. Since the MSS has a quadratic form, it can be written explicitly as a function of the focus shift by

$$\text{MSS} = P + 2Q \Delta z + R(\Delta z)^2$$

where

$$P = \frac{1}{W} \sum_{i=1}^n w_i \{(a_i^2 + c_i^2) - (A^2 + C^2)\}$$

$$Q = \frac{1}{W} \sum_{i=1}^n w_i \{(a_i b_i + c_i d_i) - (AB + CD)\}$$

$$R = \frac{1}{W} \sum_{i=1}^n w_i \{(b_i^2 + d_i^2) - (B^2 + D^2)\}$$

Differentiating this expression for the MSS with respect to focus shift, then setting the derivative to zero, determines the focus shift at which the rms spot size has its minimum value:

$$\Delta z_{\text{opt}} = -Q/R$$

Although the above equations determine the rms spot size in two dimensions, similar one-dimensional equations can be written for x and y separately, allowing the ready computation of the tangential and sagittal foci from spot-diagram data. In addition, it is straightforward to carry out the preceding type of analysis using optical path data, which leads to the determination of the center of the reference sphere that minimizes the variance of the wavefront.

Beyond the computation of the statistical rms spot size and the wavefront variance, most optical design programs include a variety of image evaluation routines that are based on spot diagram data. It is useful to characterize them as belonging to geometrical optics or physical optics, according to whether they are based on ray displacements or wavefronts, although, of course, all are based on the results of geometrical ray tracing.

Geometrical Optics Most optical design programs provide routines for computing radial diagrams and knife-edge scans. To compute a radial energy diagram, the spot-diagram data are sorted according to increasing ray displacement from the centroid of the spot. The fractional energy is then plotted as a function of spot radius. The knife-edge scan involves a similar computation, except that the spot-diagram data are sorted according to x or y coordinates, instead of total ray displacement.

Another type of geometrical image evaluation based on spot-diagram data is the so-called *geometrical optical transfer function* (GOTF). This function can be developed as the limiting case, as the wavelength approaches zero, of the actual diffraction MTF, or, alternately, in a more heuristic way as the Fourier transform of a line spread function found directly from spot-diagram ray displacements (see, for example, Smith's book¹⁰). From a programming standpoint, computation of the GOTF involves multiplying the ray displacements by 2π times the spatial frequency under consideration, forming cosine and sine terms, and summing over all the rays in the spot diagram. The computation is quick, flexible, and if there are more than a few waves of aberration, accurate. The results of the GOTF computation are typically shown as either plots of the magnitude of the GOTF as a function of frequency, or alternately in the form of what is called a "through-focus" MTF, in which the GOTF at a chosen frequency is plotted as a function of focus shift from the nominal image surface.

Physical Optics The principal physical optics calculations based on spot-diagram data are the modulation transfer function, sometimes called the "diffraction" MTF, and the point spread function (PSF). Both are based on the wavefront derived from the optical path length data in the spot diagram. There are various ways to compute the MTF and PSF, and not all programs use the same method. The PSF, for example, can be computed from the pupil function using the fast Fourier transform algorithm or, alternately, using direct evaluation of the Fraunhofer diffraction integral. The MTF can be computed either as the Fourier transfer of the PSF or, alternately, using the convolution of the pupil function.¹¹ The decision as to which method to use involves speed, accuracy, flexibility, and ease of coding.

In physical-optics-based image evaluation, accuracy can be a problem of substantial magnitude. In many optical design programs, diffraction-based computations are only accurate for systems in which diffraction plays an important role in limiting the performance. Systems that are limited primarily by geometrical aberrations are difficult to evaluate using physical optics, because the wavefront changes so much across the pupil that it may be difficult to sample it sufficiently using a reasonable number of rays. If the actual wavefront in the exit pupil is compared to a reference sphere, the resultant fringe spacing defines the size required for the spot diagram grid, since there must be several sample points per fringe to obtain accurate diffraction calculations. To obtain a small grid spacing, one can either trace many rays, or trace fewer rays but interpolate the resulting data to obtain intermediate data.

Diffraction calculations are necessarily restricted to one wavelength. To obtain polychromatic diffraction results it is necessary to repeat the calculations in each color, adding the results while keeping track of the phase shifts caused by the chromatic aberration.

3.5 OPTIMIZATION

The function of the optimization part of the program is to take a *starting design* and modify its construction so that it meets a given set of specifications. The starting design may be the result of a previous design task, a lens from the library, or a new design based on general optical principles and the designer's intuition.

The performance of the design must be measured by a single number, often known in optics as the *merit function*, although the term *error function* is more descriptive and will be used here. The error function is the sum of squares of quantities called *operands* that characterize the desired attributes. Examples of typical operands include paraxial constants, aberration coefficients, and exact ray displacements. Sometimes, the operands are broken into two groups: those that must be satisfied exactly, which may be called *constraints*, and others that must be minimized. Examples of constraints might include paraxial conditions such as the focal length or numerical aperture.

The constructional parameters to be adjusted are called *variables*, which include lens curvatures, thicknesses, refractive indices, etc. Often the allowed values of the variables are restricted, either by requirements of physical reality (e.g., positive thickness) or the given specifications (e.g., lens diameters less than a prescribed value). These restrictions are called *boundary conditions*, and represent another form of constraint.

Usually, both the operands and constraints are nonlinear functions of the variables, so optical design involves nonlinear optimization with nonlinear constraints, the most difficult type of problem from a mathematical point of view. A great deal of work has been carried out to develop efficient, general methods to solve such problems. Detailed consideration of these methods is beyond the scope of this chapter, and the reader is referred to a paper by Hayford.¹²

In a typical optical design task, there are more operands than variables. This means that there is, in general, no solution that makes all of the operands equal to their target values. However, there is a well-defined solution called the *least-squares* solution, which is the state of the system for which the operands are collectively as close to their targets as is possible. This is the solution for which the error function is a minimum.

The Damped Least-Squares Method

Most optical design programs utilize some form of the *damped least-squares* (DLS) method, sometimes in combination with other techniques. DLS was introduced to optics in about 1960, so it has a history of 50 years of (usually) successful application. It is an example of what is known as a *downhill* optimizer, meaning that in a system with multiple minima, it is supposed to find the nearest local minimum. In practice, it sometimes suffers from *stagnation*, yielding slow convergence. On the other hand, many designers over the years have learned to manipulate the damping factor to overcome this deficiency, and even in some cases to find solutions beyond the local minimum.

We consider first the case of unconstrained optimization. Let the system have M operands f_i and N variables x_j . The error function ϕ is given by

$$\phi = f_1^2 + f_2^2 + \cdots + f_M^2$$

Define the following:

$$\mathbf{A} = \text{derivative matrix, } A_{ij} \equiv \frac{\partial f_i}{\partial x_j}$$

$$\mathbf{G} = \text{gradient vector, } G_k \equiv \frac{1}{2} \frac{\partial \phi}{\partial x_k}$$

\mathbf{x} = change vector

\mathbf{f} = error vector

With these definitions, we have

$$\mathbf{G} = \mathbf{A}^T \mathbf{f}$$

If we assume that the changes in the operands are linearly proportional to the changes in the variables, we have

$$\mathbf{f} = \mathbf{A}\mathbf{x} + \mathbf{f}_0$$

$$\mathbf{G} = \mathbf{A}^T \mathbf{A}\mathbf{x} + \mathbf{G}_0$$

At the solution point, the gradient vector is zero, since the error function is at a minimum. The change vector is thus

$$\mathbf{x} = -(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{G}_0$$

These are called the least-squares *normal equations*, and are the basis for linear least-squares analysis. When nonlinear effects are involved, repeated use of these equations to iterate to a minimum often leads to a diverging solution. To prevent such divergence, it is common to add another term to the error function, and this limits the magnitude of the change vector \mathbf{x} . In the DLS method, this is accomplished by defining a new error function

$$\varphi = \phi + p\mathbf{x}^T \mathbf{x}$$

A key property of DLS is that the minimum of φ is the same as the minimum of ϕ since, at the minimum, the change vector \mathbf{x} is zero. By differentiating and setting the derivative equal to zero at the minimum, we arrive at the damped least-squares equations

$$\mathbf{x} = -(\mathbf{A}^T \mathbf{A} + p\mathbf{I})^{-1} \mathbf{G}_0$$

which look like the normal equations with terms added along the diagonal. These terms provide the damping, and the factor p is called the *damping factor*. This particular choice of damping is called *additive damping* but, more generally, it is possible to add any terms to the diagonal and still maintain the same minimum. Some optical design programs multiply the diagonal elements of the $\mathbf{A}^T \mathbf{A}$ matrix by a damping factor, while others make them proportional to the second derivative terms. Although theoretical arguments are sometimes advanced to support the choice of a particular method of damping, in practice the choice of damping factor is an ad hoc way to accelerate convergence to a solution by limiting the magnitude (and changing the direction) of the change vector found from the normal equations.

In practical optical design work, it has been found that no single method for choosing the damping factor works best in all cases. In a particular problem, one method may be dramatically better than another, but in a different problem, the situation may be completely reversed. Every optical design program has its unique way of choosing the optimum damping, which makes each program different from the others, and gives it a *raison d'être*.

Although the principal use of the damping factor is to accelerate convergence by limiting the magnitude of the change vector, the damping factor has also been used routinely to increase the magnitude of the change vector to escape a local minimum. During the course of a minimization task, if the solution stagnates, or does not converge to what the designer believes to be an acceptable configuration, it may be possible to force the solution into another region by running one or more iterations with reduced damping in which the error function increases.

Constraints and Boundary Conditions There are two general methods used in optical design programs for handling constraints and boundary conditions. The first is to add a term (called a *penalty function*) to the error function that targets the constraint to its desired value. In the case of boundary violations, “one-sided” terms can be added, or special weighting functions can be constructed that increase in magnitude as a violation goes farther into a forbidden region. The other method augments

the number of equations by the number of constraints and solves the resulting equations using the Lagrange multiplier method. This produces a minimum that satisfies the constraints exactly and minimizes the remaining error function.

The penalty function method is more flexible and faster (since there are fewer equations) than the Lagrange multiplier method. On the other hand, the Lagrange multiplier method gives more precise control over the constraints. Both are commonly used in optical design software.

Other Methods

Although DLS is used in the vast majority of optical design applications, other methods are occasionally used,¹² and two warrant mention. These are *orthonormalization*, which has been used to overcome stagnation in some DLS problems, and *simulated annealing*, which has been used for global optimization.

Orthonormalization The technique of orthonormalization for the solution of optical design problems was introduced by Grey.² Although it solves the same problem as DLS does, it proceeds in a very different fashion. Instead of forming the least-squares normal equations, Grey works directly with the operand equations

$$\mathbf{Ax} = -\mathbf{f}$$

To understand Grey's method, it is best to forget about optics and consider the solution of these equations strictly from a mathematical point of view. The point of view that Grey uses is that \mathbf{f} represents a vector in m -dimensional space. The columns of \mathbf{A} can be regarded as basis vectors in this m -dimensional space. Since there are only n columns, the basis vectors do not span the space. The change vector \mathbf{x} represents a projection of \mathbf{f} on the basis vectors defined by \mathbf{A} . At the solution point, the residual part of \mathbf{f} will be orthogonal to its projection on the basis vectors.

In Grey's orthonormalization method, the solution of the equations is found by a technique similar to Gram-Schmidt orthogonalization, but during the solution process, the actual error function is evaluated several times in an effort to use the best variables to maximum advantage. Because of this, the method is computationally intensive compared to DLS. However, the extra computation is justified by a more accurate solution. The common wisdom is that orthogonalization is superior to DLS near a solution point, and inferior to DLS when the solution is far removed from the starting point.

Simulated Annealing Simulated annealing has been applied to optical design optimization, chiefly in problems where the task is to find a global minimum. The method varies drastically from other techniques. It makes no use of derivative information, and takes random steps to form trial solutions. If a trial solution has a lower error function than the current system, the new system replaces the old. If a trial solution has a higher error function than the current system, it may be accepted, depending on how much worse it is. The probability of acceptance is taken to be $\exp(-\Delta\phi/T)$, where T is an experimentally determined quantity. In general simulated annealing, T is provided by the user. In adaptive simulated annealing, T is reduced automatically according to algorithms that hold the system near statistical equilibrium.

Error Functions

Obviously, the choice of an error function has a major impact on the success of an optical design task. There are a number of requirements that an error function should meet. Most importantly, the error function should accurately characterize the desired properties of the system under design. There is little chance of success if the program is optimizing the wrong thing. Yet this is an area of great difficulty in computer-aided optical design, because it is at odds with efficiency. In order to obtain more accuracy, more extensive computations should be carried out, but this takes time.

There are two schools of thought concerning the implementation of error functions in optical design programs. The first holds that the designer should have complete control over the items included in the error function, while the second holds that the program itself should set up the basic error function, allowing the designer some degree of control through weighting functions. Neither school has demonstrated superiority, but the approach to error function construction taken by various optical design programs accounts for user allegiances that are sometimes remarkably strong.

The different ways that optical design programs handle error functions makes it difficult to discuss the topic here in anything other than broad detail. At one extreme are programs that provide practically no capability for the user to insert operands, displaying only the value of the overall error function, while at the other extreme are programs that make the user enter every operand individually. Regardless of the user interface, however, there are some general concepts that are universally relevant.

Error functions can be based on either aberration coefficients or exact-ray data (or both). In the early stages of design, aberration coefficients are sometimes favored because they provide insight into the nature of the design, and do not suffer ray failures. However, the accuracy of aberration coefficients for evaluating complex systems is not very good, and exact-ray data are used in virtually all final optimization work.

So far as exact-ray error functions are concerned, there is the question of whether to use ray displacements or optical path difference (or both). This is a matter of user (or programmer) preference. The use of ray displacements leads to minimizing geometrical spot sizes, while the use of optical path difference leads to minimizing the wavefront variance.

For exact-ray error functions, a suitable pattern of rays must be set up. This is often called a *ray set*. There are three common methods for setting up a ray set. The first is to allow the designer to specify the coordinates (object, pupil, wavelength, etc.) for a desired set of rays. This gives great flexibility, but demands considerable skill from the user to ensure that the resulting error function accurately characterizes performance.

The other two methods for setting up ray sets are more automatic. The first is to allow the user to specify object points, and have the program define a rectangular grid of rays in the aperture for each point. The second uses a Gaussian integration scheme proposed by Forbes to compute the rms spot size, averaged over field, aperture, and wavelength.¹³ The Forbes method, which is restricted to systems having plane symmetry, leads to dividing the aperture into *rings* and *spokes*. For systems having circular pupils, the Forbes method has both superior accuracy and efficiency, but for vignetted pupils, there is little difference between the two.

Multiconfiguration Optimization

Multiconfiguration optimization refers to a process in which several systems having some common elements are optimized jointly, so that none of the individual systems but the *ensemble* of all of the systems is optimized. The archetype of multiconfiguration systems is the *zoom* system in which the focal length is changed by changing the separation between certain elements. The system is optimized simultaneously at high, medium, and low magnifications to produce the best overall performance.

Most of the larger optical design programs have the capability to carry out multiconfiguration optimization, and this capability is probably used more for non-zoom systems than for zoom systems. A common use of this feature is to optimize a focal system for through-focus performance in order to minimize sensitivity to image plane shifts. In fact, multiconfiguration optimization is used routinely to control tolerances.

Tolerancing

Beyond the task of desensitizing a given design, considerations of manufacturing tolerances become increasingly important as the complexity of optical designs increases. It is quite easy to

design optical systems that cannot be built because the fabrication tolerances are beyond the capability of optical manufacturing technology. In any case, specifying tolerances is an integral part of optical design, and a design project cannot be considered finished until appropriate tolerances are established.

Tolerancing is closely related to optimization. The basic tolerance computation is to calculate how much the error function changes for a small change in a construction parameter, which is the same type of computation carried out when computing a derivative matrix. Even more relevant, however, is the use of *compensators*, which requires reoptimization. A compensator is a construction parameter that can be adjusted to compensate for an error introduced by another construction parameter. For example, a typical compensator would be the image distance, which could be adjusted to compensate for power changes introduced by curvature errors.

There is considerable variation in how different optical design programs handle tolerancing. Some use the reoptimization method described here, while others use Monte Carlo techniques. Some stress interaction with the designer, while others use defaults for more automatic operation.

3.6 OTHER TOPICS

Of course, many other topics would be included in a full discussion of optical design software. Space limitations and our intended purpose prevents any detailed consideration, but a few of the areas where there is still considerable interest are the following.

Simulation

There is increased interest in using optical design programs to simulate the performance of actual systems. The goal is often to be able to calculate radiometric throughput of a system used in conjunction with a real extended source. It is difficult to provide software to do this with much generality, because brute force methods are very inefficient and hard to specify, while elegant methods tend to have restricted scope, and demand good judgment by the person modeling the physical situation. Nevertheless, with the increase in the speed of computers, there is bound to be an increasing use of optical design software for evaluating real systems.

Global Optimization

After several years during which there was little interest in optimization methodology, the tremendous increase in the speed of new computers has spawned a renewal of efforts to find global, rather than local, solutions to optical design problems. Global optimization is a much more difficult problem than local optimization. In the absence of an analytic solution, one never knows whether a global optimum has been achieved. All solution criteria must specify a region of interest and a time limit, and the method cannot depend on the starting point. The simulated annealing method described above is one area of continuing interest. Several methods for what might be called *pseudo-global* optimization have been used in commercial optical design programs, combining DLS with algorithms that allow the solution to move away from the current local minimum.

Computing Environment

Increasingly, optical design programs are used in conjunction with other software. Drawing programs, manufacturing inventory software, and intelligent databases are all relevant to optical design. While the conventional optical design program has been a *stand-alone* application, there is increasing demand for integrating optical design into more general design tasks.

3.7 BUYING OPTICAL DESIGN SOFTWARE

The complexity of the optical design process, together with the breadth of applications of optics, has created an ongoing market for commercial optical design software. For people new to optical design, however, the abundance of advertisements, feature lists, and even technical data sheets doesn't make purchasing decisions easy. The following commentary, adapted from an article, may be helpful in selecting an optical design program.¹⁴ It considers five key factors: hardware, features, user interface, cost, and support.

Hardware

It used to be that the choice of an optical design program was governed by the computer hardware available to the designer. Of course, when the hardware cost was many times higher than the software cost, this made a great deal of sense. Today, however, the software often costs more than the hardware, and many programs can be run on several different computer platforms, so the choice of computer hardware is less important. The hardware currently used for optical design is principally IBM-PC compatible.

To run optical design software, the fastest computer that can be obtained easily is recommended. The iterative nature of optical design makes the process interminable. There is a rule, sometimes called the Hyde maxim, that states that an optical design is finished when the time or money runs out. Notwithstanding this, the speed of computers has ceased to be a significant impediment to ordinary optical design. Even low-cost computers now trace more than 1000 ray-surfaces/s, a speed considered the minimum for ordinary design work, and create the potential for solving new types of problems formerly beyond the range of optical design software.

Before desktop computers, optical design software was usually run on time-shared central computers accessed by terminals, and some programs are still in that mode. There seems to be general agreement, however, that the memory-mapped display found on PCs provide a superior working environment and dedicated desktop computer systems are currently most popular.

Features

If you need a particular feature to carry out your optical design task, then it is obviously important that your optical design program have that feature. But using the number of features as a way to select an optical design program is probably a mistake. There are more important factors, such as cost, ease of use, and scope. Moreover, you might assume that all the features listed for a program work simultaneously, which may not be true. For example, if a vendor states that its program handles holograms and toric surfaces, you might assume that you can work with holographic toroids, but this may not be true.

The continuing growth of optics and the power of desktop computers has put heavy demands on software vendors to keep up with the development of new technology. Moreover, since the customer base is small and most vendors now support the same computer hardware, the market has become highly competitive. These factors have led to a "feature" contest in which software suppliers vie to outdo each other. While this is generally good for the consumer, the introduction of a highly visible new feature can overshadow an equally important but less obvious improvement (for example, fewer bugs or better documentation). In addition, the presence of a number of extra features is no guarantee that the underlying program is structurally sound.

User Interface

There is very little in common between the user interfaces used by various optical design programs. Each seems to have its own personality. The older programs, originally designed to run in batch

mode on a large computer, are usually less interactive than ones that were written specifically for desktop computers. Batch programs tend to be built more around default actions than interactive programs, which require more user input. It would be hard to put any of today's major optical design programs in a box classified as either batch or interactive, but the look and feel of a program has a strong influence on its usefulness.

Many people don't realize that the most important benefit of using an optical design program is often the understanding that it provides the user about how a particular design works. It's often tempting to think that if the computer could just come up with a satisfactory solution, the design would be finished. In practice, it is important to know the trade-offs that are made during a design project. This is where the judgment of the optical engineer comes in, knowing whether to make changes in mechanical or electrical specifications to achieve the optimum balance in the overall system. Lens designers often say that the easiest lens to design is one that has to be diffraction-limited, because it is clear when to stop. If the question of how to fit the optics together with other system components is important, then the ability of the user to work interactively with the design program can be a big help.

Cost

In today's market, there is a wide range of prices for optical design software. This can be very confusing for the first-time buyer, who often can't see much difference in the specifications. The pricing of optical design software is influenced by (at least) three factors.

First, the range of tasks that can be carried out using an optical design program is enormous. The difference in complexity between the job of designing a singlet lens for a simple camera, and that of designing a contemporary objective for a microlithographic masking camera is somewhat akin to the difference between a firecracker and a hydrogen bomb.

Second, all software is governed by the factors originally studied in F. P. Brooks' famous essay *The Mythical Man-Month*.¹⁵ Brooks was director of the group that developed the operating system for the IBM 360, a mainframe computer introduced in the 1960s. Despite its provocative title, Brooks' essay is a serious work that has become a standard reference for software developers. In it, he notes that if the task of developing a program to be used on a single computer by its author has a difficulty of 1, then the overall difficulty of producing integrated software written by a group of people and usable by anyone on a wide range of computers may be as high as 10. In recent years, the scope of the major optical design programs has grown too big for a single programmer to develop and maintain, which raises costs.

Third, there are structural differences in the way optical design software is sold. The original mainframe programs were rented, not sold. If the user did not want to continue monthly payments, the software had to be returned. PC programs, on the other hand, are usually sold with a one-time fee. In the optical design software business, several vendors offer a compromise policy, combining a permanent license with an optional ongoing support fee.

It would be nice if the buyer could feel comfortable that "you get what you pay for," but unfortunately this view is too simplistic. One program may lack essential capabilities, another may contain several unnecessary features when evaluated for a particular installation. Buying on the basis of cost, like features, is probably not a good idea.

Support

Support is an important aspect to consider in selecting an optical design program, and it is often difficult to know what is included in support. Minimal support consists of fixing outright bugs in the program. More commonly, support includes software updates and phone or email assistance in working around problems.

Optical design programs are typically not bug-free. Unlike simple programs like word processors, optimization programs cannot be fully tested, because they generate their own data. One result of this is that software vendors are generally reluctant to offer any warranty beyond a "best-effort"

attempt to fix reported problems. Unfortunately, there is no good way for buyers to know whether and when their particular problems may be fixed; the best approach is probably to assess the track record of the vendor by talking to other users.

Coupled with support is user training. Although it should be possible to use a program by studying the documentation, the major optical design software vendors offer regular seminars, often covering not only the mechanics of using their program, but also general instruction in optical design. For new users, this can be a valuable experience.

3.8 SUMMARY

As stated in the introduction, this chapter is intended as a survey for readers who are not regular users of optical design software. The form of an optical design program described here, consisting of lens entry, evaluation, and optimization sections, is used in many different programs. There has been little standardization in this field, so the “look and feel,” performance features and extent of various programs are quite different. Nonetheless, it is hoped that with a knowledge of the basic features described here, the reader will be in a good position to judge whether an optical design program is of use, and to make an informed decision about whether one particular program is better than another.

3.9 REFERENCES

1. D. P. Feder, “Automatic Optical Design,” *Appl. Opt.* **2**:1209–1226 (1963).
2. D. S. Grey, “Aberration Theories for Semiautomatic Lens Design by Electronic Computers,” *J. Opt. Soc. Am.* **53**:672–680 (1963).
3. G. H. Spencer, “A Flexible Automatic Lens Correction Procedure,” *Appl. Opt.* **2**:1257–1264 (1963).
4. C. G. Wynne and P. Wormell, “Lens Design by Computer,” *Appl. Opt.* **2**:1233–1238 (1963).
5. See, for example, *The Photonic Industry Buyer’s Guide*, Laurin Publishing, Pittsfield, MA 01202.
6. *U.S. Military Handbook for Optical Design*, republished by Sinclair Optics, Fairport, NY 14450 (1987).
7. G. H. Spencer and M. V. R. K. Murty, “Generalized Ray-Tracing Procedure,” *J. Opt. Soc. Am.* **52**:672–678 (1962).
8. D. C. O’Shea, *Elements of Modern Optical Design*, John Wiley & Sons, New York (1985).
9. H. A. Buchdahl, *Optical Aberration Coefficients*, Oxford Press, London (1954).
10. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York (1990).
11. H. H. Hopkins, “Numerical Evaluation of the Frequency Response of Optical Systems,” *Proc. Phys. Soc. B* **70**:1002–1005 (1957).
12. M. J. Hayford, “Optimization Methodology,” *Proc. SPIE* **531**: 68–81 (1985).
13. G. W. Forbes, “Optical System Assessment for Design: Numerical Ray Tracing in the Gaussian Pupil,” *J. Opt. Soc. Am. A* **5**:1943–1956 (1988).
14. D. C. Sinclair, “Optical Design Software: What to Look For in a Program,” *Photonics Spectra*, Nov. 1991.
15. F. P. Brooks, Jr., *The Mythical Man-Month*, Addison-Wesley, Reading, MA (1975).

4

OPTICAL SPECIFICATIONS

Robert R. Shannon*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

4.1 GLOSSARY

ATF	approximate transfer factor
DTF	diffraction transfer function
MTF	modulation transfer function
W	wavefront error in units of wavelengths
W_{rms}	root-mean-square wavefront error
t -number	f -number adjusted for lens transmission
ν	normalized spatial frequency

4.2 INTRODUCTION

Setting the specifications for an optical instrument or system is an essential part of engineering, designing, or purchasing an optical system. Since the optics usually serve as a portion of a larger system, the specifications are frequently set by project managers who do not have specific knowledge in the basics of optical systems. This can at times lead to unrealistic requirements being established; this can profoundly affect the probability of success for the system. Properly drafted specifications can make the entire project successful and cost effective. Poorly written specifications can lead to excess cost and ultimately project failure.

One of the difficulties with setting optical specifications is that the ultimate result of a beam of light passing through a complex assembly of components is affected by each of those components, which in turn need to be specified and tolerances placed upon the fabrication and assembly of those components. In the case of an imaging system, the problem is compounded by the need to describe an optical system which passes many bundles of light across a wide field of view. Even in the case of single beam, optical communications components, indirect issues such as scattered light and environmental stability may prove to be major issues.

*Retired.

In the worst case, the specifications may be set so high that the system is not capable of being manufactured. In most cases, the specifications interact with other devices, such as detector arrays, and matching the quality of the optic to the limits of the sensor is required. In this section some of the principles involved in setting the specifications will be discussed, and guidelines provided for carrying out the process of specification setting. The reader will have to extend these principles to the device or system that is being considered. In this chapter, the stress will be placed on imaging systems.

Specifications for optical systems cover a wide range of needs. Functional specifications of the image quality or other optical characteristics are required for the satisfactory operation of a system. These functional specifications serve as the goal for the design and construction of the optical system. In addition, these specifications are a basis for tolerances placed upon the components of the optical system and lead to detailed component specifications used for procurement of the optical elements of the system. Assembly specifications and detailed specifications of optical parts to be produced by a shop can be written based upon these component specifications. The detail and extent of information required is different at each step. Over- or underspecification can contribute significantly to the cost or feasibility of design of an optical system.

Functional specifications are also used to describe the characteristics that an instrument must demonstrate in order to meet the needs of the user. This may include top-level requirements such as size, weight, image scale, image format, power levels, spectral range, and so on. Component specifications are developed after design of the system and describe the optical components, surface, and materials used in the system to the detail necessary to permit fabrication of the components. Assembly specifications are another derivative of the design and system specifications. These include the statement of tolerances upon location of the components, as well as the procedure to be used in assembling and testing the system.

The development and writing of these specifications is important both for initiating and for tracking the course of development of an optical instrument. In a business or legal sense, specifications are used to establish responsibility for a contractor or subcontractor, as well as to define the basis for bidding on the job. Thus the technical specifications can have business importance as well as engineering significance. "Meeting the customer's specifications" is an essential part of any design and fabrication task. Identifying areas where the specifications could be altered with benefit to all parties is an important business and engineering responsibility.

Specifications are usually communicated as a written document following some logical format. Although there are some international standards that may cover the details of drawings of components, there is no established uniform set of standards for stating the specifications on a system or component. The detailed or component specifications are usually added as explanatory notes to drawings of the components to be fabricated. In modern production facilities, the specifications and tolerances are often part of a digital database that is accessed as part of the production of the components of the system.

The detail and the intent of each of these classes of specifications are different. Optical specifications differ from many mechanical or other sets of specifications in that numbers are applied to surfaces and dimensions that control the cumulative effect of errors imposed on a wavefront passing through the total system. Each of the specifications must be verifiable during fabrication, and the overall result must be testable after completion.

Mechanical versus Optical Specifications

There are two types of specifications that are applied to an optical system or assembly. One set of these includes mechanical tolerances on the shape or location of the components that indirectly affect the optical quality of the image produced by the system. Examples of this include the overall size or weight of the system. The other set consists of specialized descriptions that directly affect the image quality. Examples of this latter type of specification are modulation transfer function (MTF), illumination level, and location of the focal plane relative to the system.

System versus Components Specifications

Some specifications have meaning only with respect to the behavior of the entire optical system. Others apply to the individual components, but may affect the ability of the entire system to function.

An example of a system specification is a set of numbers limiting the range of acceptable values of the MTF that are required for the system. Another system specification is the desired total light transmission of the system.

Examples of component specifications are tolerances upon surface irregularity, sphericity, and scattering. The related component specification based upon the system light transmission specification might provide detailed statements about the nature and properties of the antireflective coatings to be applied to the surface of each element.

Image Specifications

The specifications that are applied to the image usually deal with image quality. Examples are modulation transfer function, fraction of scattered light, resolution, or distortion. In some cases, these specifications can be quite general, referring to the ability of the lens to deliver an image suitable for a given purpose, such as the identification of serial numbers on specific products that are to be read by an automated scanner. In other cases, the requirements will be given in a physically meaningful manner, such as “the MTF will be greater than 40 percent at 50 lines per millimeter throughout the field of view.”

Other criteria may be used for the image specifications. One example is the energy concentration. This approach specifies the concentration of light from a point object on the image surface. For example, the specification might read “75 percent of the light shall fall within a 25- μm -diameter circle on the image.” This quantity is obviously measurable by a photometer with appropriate-size apertures. The function may be computed from the design data by a method of numerical integration similar to that providing the point spread function or modulation transfer function.

Wavefront Specifications

Wavefront specifications describe the extent to which the wavefront leaving the lens or components conforms to the ideal or desired shape. Usually the true requirement for an optical system is the specification of image quality, such as MTF, but there is a relation between the image quality and the wavefront error introduced by the optical system. The wavefront error may be left to be derived from the functional image quality specification, or it may be defined by the intended user of the system.

For example, a wavefront leaving a lens would ideally conform to a sphere centered on the chosen focal location. The departure of the actual wavefront from this ideal, would be expressed either as a matrix or map of departure of the wavefront from the ideal sphere, as a set of functional forms representing the deviation, or as an average [usually root-mean-square (rms)] departure from the ideal surface. By convention, these departures are expressed in units of wavelength, although there is a growing tendency to use micrometers as the unit of measure.

The rms wavefront error is a specific average over the wavefront phase errors in the exit pupil. The basic definition is found by defining the n th power average of the wavefront $W(x, y)$ over the area A of the pupil and then specifically defining W_{rms} or, in words, the rms wavefront error is the square root of the mean square error minus the square of the mean wavefront error:

$$\bar{W}^n = \frac{1}{A} \int W(x, y)^n dx dy$$

$$W_{\text{rms}} = \sqrt{[\bar{W}^2] - [\bar{W}]^2}$$

The ability to conveniently obtain a complete specification of image quality by a single number describing the wavefront shape has proven to be questionable in many cases. Addition of a correlation

length, sometimes expressed as a phase difference between separated points, has become common. In other cases, the relative magnitude of the error when represented by various orders of Zernike polynomials is used.

There is, of course, a specific relationship between the wavefront error produced by a lens and the resulting image quality. In the lens, this is established by the process of diffraction image formation. In establishing specifications, the image quality can be determined by computation of the modulation transfer function from the known wavefront aberrations. This computation is quite detailed and, while rapidly done using present day computer techniques, is quite complex for general specification setting. An approximation which provides an average MTF or guide to acceptable values relating wavefront error and MTF is of great aid.

A perfect lens is one that produces a wavefront with no aberration, or zero rms wavefront error. By convention, any wavefront with less than 0.07 wave, rms, of aberration is considered to be essentially perfect. It is referred to as *diffraction-limited*, since the image produced by such a lens is deemed to be essentially indistinguishable from a perfect image.

The definition of image quality depends upon the intended application for the lens. In general, nearly perfect image quality is produced by lenses with wavefront errors of less than 0.15 wave, rms. Somewhat poorer image quality is found with lenses that have greater than about 0.15 wave of error. The vast majority of imaging systems operate with wavefront errors in the range of 0.1 to 0.25 wave, rms.

There are several different methods that can be used to establish this relationship. The most useful comparison is with the MTF for a lens with varying amounts of aberration. The larger the wavefront error, the lower will be the contrast at specific spatial frequencies. For rms error levels of less than 0.25 or so, the relation is generally monotonic. For larger aberrations, the MTF becomes rather complex, and the relation between rms wavefront error and MTF value can be multiple valued. Nevertheless, an approximate relation between MTF and rms wavefront error would be useful in setting reasonable specifications for a lens.

There are several possible approximate relations, but one useful one is the empirical formula relating root-mean-square wavefront error and MTF given by

$$\text{MTF}(\nu) = \text{DTF}(\nu) \times \text{ATF}(\nu)$$

The functional forms for these values are

$$\text{DTF}(\nu) = \frac{2}{\pi} [\arccos(\nu) - \nu\sqrt{1-\nu^2}]$$

$$\text{ATF}(\nu) = \left[1 - \left(\frac{W_{\text{rms}}}{0.18} \right)^2 \right] (1 - 4(\nu - 0.5)^2)$$

$$\nu = \frac{\text{spatial frequency}}{\text{spatial frequency cutoff}} = \frac{N}{\left(\frac{1}{\lambda f\text{-number}} \right)}$$

These look quite complicated, but are relatively simple, as is shown in Fig. 1. This is an approximation, however, and it becomes progressively less accurate as the amount of the rms wavefront error W_{rms} exceeds about 0.18 wavelength. The approximation remains reasonably valid for lower spatial frequencies, less than about 25 percent of the diffraction limited cutoff frequency. The majority of imaging systems fall into this category.

Figure 1 shows a plot of several values for the MTF of an optical system using this approximate method of computation. The system designer can use this information to determine the appropriate level of residual rms wavefront error that will be acceptable for the system of interest. It is important to note that this is an empirical attempt to provide a link between the wavefront error and the

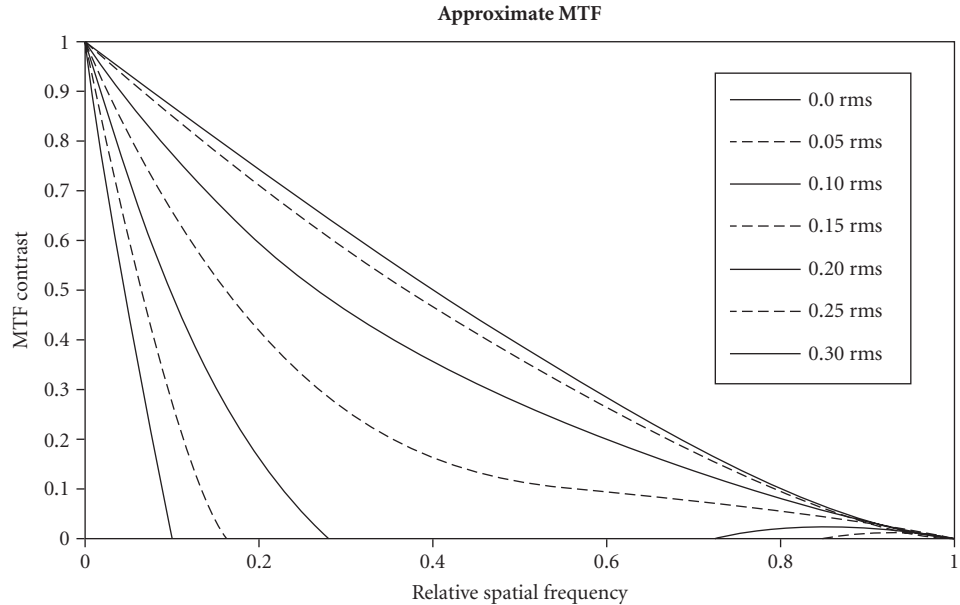


FIGURE 1 Approximate MTF curves from formula.

MTF as a single-value description of the state of correction of a system. Examination of the curves provides a method of communicating the specification to the system designer and fabricator. More detail on applying the rms wavefront error can be found in Chap.5, “Tolerancing Techniques.”

In addition, it must be pointed out that most imaging systems operate over a finite wavelength range. Thus the specification of “wavefront error” can be a bit fuzzy, but is usually meant to mean either the wavefront error at a specified wavelength, or a weighted average over the wavelength band. This should be mentioned when writing the specifications. In either case, the stated wavefront error contains a measure that communicates the extent of perfection required of the optical system performance.

4.3 PREPARATION OF OPTICAL SPECIFICATIONS

Gaussian Parameters

The gaussian parameters determine the basic imaging properties of the lens. They are the starting point for setting the specifications for a lens system. In principle these numbers can be specified precisely as desired. In reality, overly tight specifications can greatly increase the cost of the lens. Some of the important parameters are shown in Table 1.

Table 1 is a sample of reasonable values that may be placed upon a lens. A specific case may vary from these nominal values. The image location, radiometry, and scale are fixed by these numbers. A specific application will require some adjustment of these nominal values. In general, specifications that are tighter than these values will likely result in increased cost and difficulty of manufacture.

There is an interaction between these numbers. For example, the tight specification of magnification and overall conjugate distance will require a very closely held specification upon the focal length. The interaction between these numbers should be considered by the user to avoid accidentally producing an undue difficulty for the fabricator. It may be appropriate to specify a looser tolerance on some of these quantities for the prototype lens, and later design a manufacturing process

TABLE 1 Gaussian Parameters

Parameter	Precision Target	Importance	How Verified
Focal length	1–2%	Determines focal position and image size	Lens bench
<i>f</i> -number	<±5%	Determines irradiance at image plane	Geometrical measurement
Field angle	<±2%	Determines extent of image	Lens bench
Magnification	<±2%	Determines overall conjugate distances	Trial setup of lens
Back focus	±5%	Image location	Lens bench
Wavelength range	As needed; set by detector and source	Describes spectral range covered by lens	Image measurement
Transmission	Usually specified as >0.98 ⁿ for <i>n</i> surfaces	Total energy through lens	Imaging test, radiometric test of lens
Vignetting	Usually by requiring transmission to drop by less than 20% or so at the edge of the field	Uniformity of irradiance in the image	Imaging test, radiometric test of lens

to bring the production values within a smaller tolerance. However, it is appropriate that this be investigated fully at the design stage. The designer should be encouraged to consider the possibility of leaving adjustment possibilities in the lens design, so that a final assembly adjustment can bring the Gaussian parameters into the required tolerance range.

4.4 IMAGE SPECIFICATIONS

Image Quality

The rms wavefront error and the MTF for a lens have been discussed earlier as useful items to specify for a lens. Frequently, the user desires to apply a detection criterion to the image. This is always related to the application for the lens.

The most familiar functional specification that is widely used for system image quality is the resolution of the system. This is usually stated as the number of line pairs per millimeter that need to be visually distinguished or recognized by the user of the system. Since this involves both the physics of image formation as well as the psychophysics of vision, this is an interesting goal, but needs to be specified clearly to be of use to a designer. The reading of the resolution by a human observer is subjective, and the values obtained may differ between observers. Therefore, it is necessary to specify the conditions under which the test is to be carried out.

The type of target and its contrast need to be stated. The default standard in this case is the “standard” U.S. Air Force three-bar target, with high contrast, and a 6:1 ratio of bar length to width. This is usually selected as it will give the highest numerical values, certainly politically desirable. However, studies have shown that there is a better correlation between the resolution produced by a low-contrast target, say 2:1 contrast ratio, or 0.33 modulation contrast and the general acceptability of an image.

The resolution is, of course, related to the value of the MTF in the spatial frequency region of the resolution, as well as the threshold of detection or recognition for the observer viewing the target. If the thresholds are available, the above-described empirical relation between the rms wavefront error and the MTF can be used to estimate the allowable aberration that can be left in the system after design or fabrication.

In the case of a system not intended to produce an image to be viewed by a human, a specific definition of the required image contrast or energy concentration is usually possible. The signal-to-noise ratio of the data transmitted to some electronic device that is to make a decision can be calculated once a model for operation of the detector is assumed. The specification writer can then work backward through the required MTF to establish an acceptable level of image quality. The process is similar to that for the visual system above, except that the threshold is fully calculable.

In some cases, the fractional amount of energy collected by the aperture of the lens from a small angular source, such as a star, falling within the dimensions of a detector of a given size is desired. Such a requirement can be given directly to the designer.

Image Irradiance

The radiometry of the image is usually of importance. With an optical system containing the source, such as a viewer, projector, or printer, the usual specification is of the irradiance of the image in some appropriate units. Specifying the screen irradiance in watts per square centimeter or, more commonly, foot-lamberts, implies a number of optical properties. The radiance of the source, the transmission of the system, and the apertures of the lenses are derived from this requirement.

In the more usual imaging situation, the f -number and the transmission of the lens are specified. If the lens covers a reasonable field, the allowable reduction in image irradiance over the field of view must also be specified. This leads directly to the level of vignetting that can be allowed by the designer in carrying out the setup and design of the lens system.

There is an interaction between these irradiance specifications and the image quality that can be obtained. The requirement of a large numerical aperture leads to a more difficult design problem, as the high-order aberration content is increased in lenses of high numerical aperture.

An attempt to separate the geometrical aperture effects from the transmission of the components of the lens is accomplished through the t -number specification. Since the relative amount of irradiance falling on the focal plane is inversely proportional to the square of the f -number of a lens, the effect of transmission of the lens can be included by dividing the f -number by the transmission of the lens.

$$t\text{-number} = \frac{f\text{-number}}{\sqrt{t}}$$

where t is the transmission factor for the lens. The transmission factor for the lens is the product of the bulk transmission of the glass and the transmission factor for each of the surfaces. When this is specified, the designer must provide a combination of lens transmission and relative aperture that meets or exceeds a stated value.

Depth of Focus

The definition of the depth of focus is usually the result of a tolerance investigation. The allowable focal depth is obtained by determining when an unacceptable level of image quality is obtained. There is an obvious relation between the geometry of the lens numerical aperture and the aberrations that establishes the change of MTF with focal position. This effect can be computed for specific cases, or estimated by recognizing that the relation between rms wavefront error and focus shift is

$$W_{\text{rms}_{\text{def}}} = \frac{\delta l}{8\lambda(f\text{-number})^2}$$

which can be used in the above approximate MTF to provide an estimate of the likely MTF over a focal range.

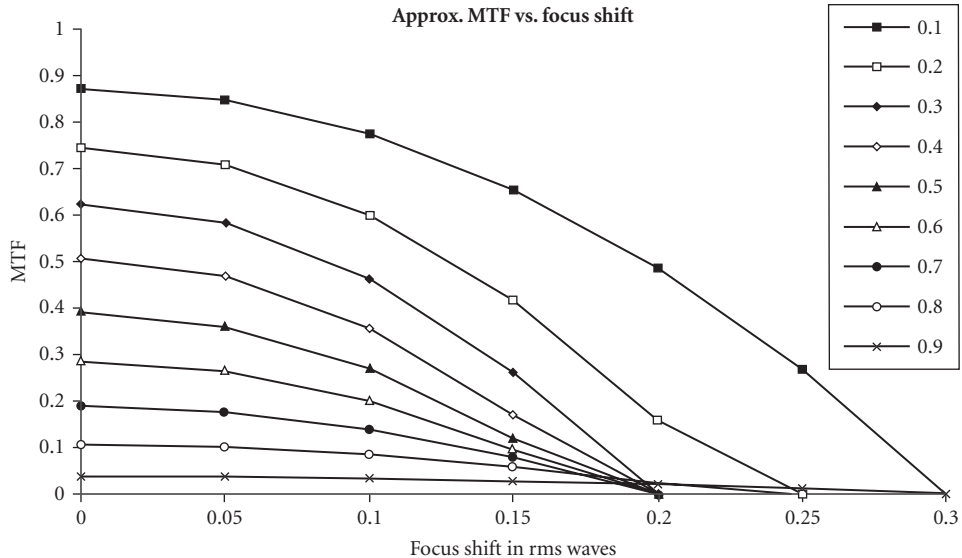


FIGURE 2 Approximate MTF as a function of focal position for various spatial frequencies.

If there is a basic amount of aberration present in the lens, then the approximation is that

$$W_{\text{rms, total}} = \sqrt{W_{\text{rms, def}}^2 + W_{\text{rms, lens}}^2}$$

leading to a calculation of the estimated MTF value for the given spatial frequency and focal position. As an example, Fig. 2 provides a plot of the focal position change of the MTF. The interpretation of this curve is straightforward. Each of the curves provides the MTF versus focus for a spatial frequency that is the stated fraction of the cutoff spatial frequency. The defocus is given in units of rms wavelength error, which can be obtained from reference to the appropriate formula. Using this approximate data, a specification writer can determine whether the requirements for image quality, f -number, and focal depth are realistic.

An additional consideration regarding depth of focus is that the field of the lens must be considered. The approximate model presented here is used at an individual field point. An actual lens must show the expected depth of field across the entire image surface, which places some limits upon the allowable field curvature. In general, it is the responsibility of the specification writer to establish the goal. It is the responsibility of the optical designer to determine whether the goal is realizable, and to design a system to meet the needs. In a sensible project, there will be some discussion between the designer and the engineer writing the specifications in order to avoid an unrealizable set of goals being set.

4.5 ELEMENT DESCRIPTION

Each element of the lens to be fabricated must be described in detail, usually through a drawing. All of the dimensions will require tolerances, or plus and minus values that, if met, lead to a high probability that the specified image quality goals will be met.

Mechanical Dimensions

The mechanical dimensions are specified to ensure that the element will fit into the cell sufficiently closely that the lens elements are held in alignment. This will be a result of tolerance evaluation, and must include allowances for assembly, thermal changes, and so on.

An important item for any lens is the interface specification, which describes the method of mounting the lens to the optical device used with it. For some items, such as cameras and microscopes, there are standard sizes and screw threads that should be used. In other cases, the specification needs to describe a method for coupling or mounting the optics in which there is a strain-free transfer of load between the lens and the mounting.

Optical Parameters

The optical parameters of the lens element relate to the surfaces that are part of the image-forming process. The radius of curvature of the spherical refracting (or reflecting) surfaces needs to be specified, as well as a plus or minus value providing the allowable tolerances. When tested using a test glass or an interferometer, the important radius specification is usually expressed in terms of fringes of spherical departure from the nominal radius. In addition, the shape of the surface is usually specified in terms of the fringes of irregularity that may be permitted.

When specifying a surface that will be measured on an interferometer, adjustment of focus during the test can be made. In this case, the spherical component of the surface, that is, the fringes of radius error, can be specified independently of the irregularity fringes that are applicable to the surface. When test glasses are to be used, the spherical component must be fabricated to within a small level of error to permit accurate reading of the irregularity component of the surface.

The cosmetic characteristics of the surface also need to be stated. The specification for this is as yet a bit imperfect, with the use of a scratch-and-dig number. This is actually intended to be a comparison of surface scratches with a visual standard, but is generally accepted to be in terms of a ratio, such as 20:10, which means, more or less, scratches of less than 2- μm width and digs of less than 100- μm diameter. This specification is described in MIL-O-13830, and is referred to a set of standards that are used for visual comparison to the defects on the surface. There have been several attempts to quantify this specification in detail, but no generally accepted standard has yet been achieved. A broader description of these specifications is found in International Standards 10110 and 9211, discussed later in this chapter.

Material Specifications

The usual material for a lens is optical glass, although plastics are becoming more commonly used in optics. The specification of a material requires identification of the type, as, for example, BK7 glass from Schott. Additional data upon the homogeneity class and the birefringence needs to be stated in ordering the glass. The homogeneity is usually specified by class, currently P1 through P4 with the higher number representing the highest homogeneity, or lowest variation of index of refraction throughout the glass. The method of specifying glass varies with the manufacturer, and with the catalog date. It is necessary to refer to a current catalog to ensure that the correct specification is being used.

Similar data should be provided regarding plastics. Additional data about transmission is usually not necessary, as the type of material is selected from a catalog which provides the physical description of the material. Usually, the manufacturer of the plastic will be noted to ensure that the proper material is obtained.

Materials for reflective components similarly have catalog data describing the class and properties of the material. In specifying such materials, it is usually necessary to add a description of the form and the final shape required for the blank from which the components will be made.

Coating Specifications

The thin film coatings that are applied to the optical surfaces require some careful specification writing. In general, the spectral characteristics need to be spelled out, such as passband and maximum reflectivity for an antireflection coating. Requirements for the environmental stability also need to be described, with reference to tests for film adhesion and durability. Generally, the coating supplier will have a set of “in-house” specifications that will guarantee a specific result that can be used as the basis for the coating specification.

4.6 ENVIRONMENTAL SPECIFICATIONS

Temperature and Humidity

Specification setting should also include a description of the temperature range that will be experienced in use or storage. This greatly affects the choice of materials that can be used. The humidity and such militarily favorite specifications as salt spray tests are very important in material selection and design.

Shock and Vibration

The ruggedness of an instrument is determined by the extent to which it survives bad handling. A requirement that the lens shall survive some specified drop test can be used. In other cases, stating the audio frequency power spectrum that is likely to be encountered by the lens is a method of specifying ruggedness in environments such as spacecraft and aircraft. In most cases, the delivery and storage environment is far more stressing than the usage environment. Any specification written in this respect should be careful to state the limits under which the instrument is actually supposed to operate, and the range over which it is merely meant to survive storage.

4.7 PRESENTATION OF SPECIFICATIONS

Published Standards

There are published standards from various sources. The most frequently referred to are those from the U.S. Department of Defense, but a number of standards are being proposed by the International Standards Organization.

Format for Specifications

The format used in conveying specifications for an optical system is sometimes constrained by the governmental or industrial policy of the purchaser. Most often, there is no specific format for expressing the specifications.

The best approach is to precede the specifications with a brief statement as to the goals for the use of the instrument being specified. Following this, the most important optical parameters, such as focal length, f -number, and field size (object and image) should be stated. In some cases, magnification and overall object-to-image distance along with object dimension will be the defining quantities.

Following this, the wavelength range, detector specifications, and a statement regarding the required image quality should be given. The transmission of the lens is also important at this stage.

Following the optical specifications, the mechanical and environmental requirements should be stated. The temperature and humidity relations under which the optical system needs to operate as

well as a statement of storage environment are needed. Descriptions of the mechanical environment, such as shock and vibration, are also important, even if expressed generally.

Other important pieces of information, such as a desired cost target, can then be included. Any special conditions, such as the need to be exposed to rapid temperature changes or a radiation environment, should be clearly stated. Finally, some statement of the finish quality for the optical system should be given.

In many cases, a list of applicable governmental specifications will be listed. In each case it is appropriate to ensure that these referenced documents are actually available to the individual who has to respond to the specification.

Use of International Standards

An important tool in writing specifications is the ability to refer to an established set of standards that may be applicable to the system being designed. In some cases, the development of specifications is simplified by the specification writer being able to refer to a set of codified statements about the environment or other characteristics the system must meet. In other cases, the established standards can be used as a reference to interpret parts of the specifications being written. For example, there may be a set of standards regarding interpretation of items included in a drawing.

Standards are an aid, not an end to specification. If the instrument must be interchangeable with parts from other sources, then the standards must be adhered to carefully. In other cases, the standards can serve as an indicator of accepted good practice in design or fabrication. It is the responsibility of the specification writer to ensure that the standard is applicable and meaningful in any particular situation.

At the present time there is growing activity in the preparation of standards for drawings, interfaces and dimensions, MTF, and other properties of optical and electro-optical systems. The efforts in this direction are coordinated by the International Standards Organization, but there are a number of individual standards published by national standards organizations in Germany, England, Japan, and the United States. The first major publications are ISO 10110, detailing preparation of drawings for optical elements and systems and ISO 9211, on optical coatings. Other standards on optical testing and environmental requirements are in draft form.

The ISO standards are expected to provide significant detail on various standards issues, and should become the principal guiding documents. At present, the standards documents that are most used in the United States are the various military specifications, or MIL-SPEC documents, that cover many different aspects of optical systems.

Information on published standards is available from the American National Standards Institute (ANSI), 11 West 42d Street, New York, NY 11036, or may be downloaded at www.webstore.ansi.org. A recent (2008) review of this website showed over 350 individual documents dealing with these issues, the majority of which deal with issues regarding fiber optical systems. Information on U.S. Department of Defense standards can be searched for through the National Search Engine for Standards at www.nssn.org. Additional information about worldwide standards is available at www.worldwidestandards.com.

4.8 PROBLEMS WITH SPECIFICATION WRITING

Underspecification

Failure to specify all of the conditions leaves the user vulnerable to having an instrument that will not operate properly in the real world. In many cases, the designer may not be aware of situations that may arise in operation that may affect the proper choice of design methods. Therefore, the design may not meet the actual needs.

The engineer developing a specification should examine all aspects of the problem to be solved, and carefully set the boundaries for the requirements on an optical system to meet the needs. All of

the pertinent information about the image quality, environment, and relation to other systems that may interact with the lens being specified should be considered. The specifying engineer should also review the physical limits on the image quality and ensure that these are translated into realistic values.

Overspecification

Specifying image quality and focal position requirements too tightly can lead to problems. Overspecification would seem to ensure that the needs will be met, but difficulties in meeting these requirements can lead to designs that are difficult and expensive to build. Achieving the goals can be costly and may fail. In such cases, the penalty for not quite meeting very tight specifications can be serious, both economically and technically.

Boundary-Limit Specification

In most cases, the statement of goals or boundaries within which a lens must operate is better than stating specific values. This leaves the designer with some room to maneuver to find an economic solution to the design. Obviously, some fixed values are needed, such as focal length, f -number, and field angle. However, too-tight specifications upon such items as weight, space, and materials can force the design engineer into a corner where a less desirable solution is achieved.

Negotiation of Specifications

Finally it is important to note that unless there is an existing closely defined set of established specifications for an specific optical device (such as a fiber optics coupler, for example) each specification is the product of a single individual or group and reflects the experience and understanding of that individual. The procuring official should be prepared in some cases to act as a negotiator between the engineer and the supplier to ensure that a reasonable and successful set of verifiable specifications has been stated.

4.9 REFERENCES

There are many useful references on optical specifications that deal with specific topics not directly covered by the general discussion in this chapter. The most useful suggestion is that the users having the task of setting specifications on a specific product or system use the massive capabilities of Internet search engines to look for specific data applicable to that task. A general Google search on "Optical Specifications" provided 360,000 hits, of which probably less than 10 will be applicable to any specific problem.

TOLERANCING TECHNIQUES

Robert R. Shannon*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

5.1 GLOSSARY

a	relative tolerance error
BK7, SF2	types of optical glass
C to F	spectral region 0.486 to 0.656 μm
f -number	relative aperture as in $F/2.8$
n	refractive index
r_1, r_2, r_3, r_4	radii of curvature of surfaces
t_2	airspace thickness in sample lens
V	Abbe number or reciprocal dispersion
W	wavefront or pupil function
x	change factor
Δ	finite change in a parameter
δ	small change in a parameter
λ	wavelength

5.2 INTRODUCTION

Determination of the tolerances on an optical system is one of the most important parts of carrying out an optical design. No component can be made perfectly; thus, stating a reasonable acceptable range for the dimensions or characteristics is important to ensure that an economical, functioning instrument results. The tolerances attached to the dimensions describing the parts of the lens system are an important communication by the designer to the fabrication shop of the precision required in making the components and assembling them into a final lens.

*Retired.

The tolerances are related to but are not the same as the specifications. Setting specifications is discussed in Chap. 4. The tolerances are responsive to the requested system specifications and are intended to ensure that the final, assembled instrument meets the requested performance. The specifications placed on the individual lens elements or components are derived from the tolerances. Thus there is an interactive relation between the tolerancing activity and the setting of specifications. The system specifications drive the tolerances that need to be determined, and the tolerances are used in setting the specifications for the components of the system. The reality is that neither of these processes can be done fully independently.

At this point it is important to note that most optical design programs include a tolerancing utility that can be used to generate and distribute tolerances automatically once a few questions about goals have been answered by the designer. This appears to be a seductively simple process that is usually quite useful, but can be very disastrous if used uncritically by anyone who does not understand the basics of the process being carried out.

There are three principal issues in optical tolerancing. The first is the setting of an appropriate goal for the image quality or transmitted wavefront to be expected from the system. The second is the translation of this goal into allowable changes introduced by errors occurring on each component of the system. The third is the distribution of these allowable errors against all of the components of the system, in which some components of the optical system may partially or completely compensate for errors introduced by other components.

In this chapter, some basic approaches to distributing tolerances within an optical assembly are discussed. The examples will deal with tolerancing to meet a specified wave-front error and level of image quality. Similar principles apply to nonimaging optical systems, once the procedures necessary for relating errors in components or alignment to the specified operating requirements are established. The user will obviously have to adapt these approaches to the specific system being tolerated.

Optical versus Mechanical Tolerances

The tolerances on mechanical parts, in which a dimension may be stated as a specific value, plus or minus some allowable error, are familiar to any engineer. For example, the diameter of a rotating shaft may be expressed as 20.00 mm + 0.01/−0.02 mm. This dimension ensures that the shaft will fit into another component, such as the inner part of a bearing, and that fabricating the shaft to within the specified range will ensure that proper operational fit occurs. These tolerances may include the effect of environmental effects, such as operating temperature or lubrication needs, on the mechanical assembly.

Optical tolerances are more complicated, as they are generally stated as a mechanical error in a dimension, but the allowable error is determined by the effect upon an entire set of wavefronts passing through the lens. For example, the radius of curvature of a surface may be specified as 27.00 mm ± 0.05 mm. The interpretation of this is that the shape of the optical surface should conform to a specific spherical form, but remain within a range of allowable curvatures. Meeting this criterion indicates that the surface will perform properly in producing a focused wavefront, along with other surfaces in the optical system. Verification that the specific component tolerance is met is usually carried out by an optical test, such as examining the fit to a test plate. Verification that the entire system operates properly is accomplished by an assembled system test in which a specified image quality criterion is measured.

Basis for Tolerances

The process involved in setting tolerances begins with setting of the minimum level of acceptable image quality. This is usually expressed as the desired level of contrast at a specific spatial frequency as expressed by the modulation transfer function. Each parameter of the system, such as a radius of curvature of a surface, is individually varied to determine how large an error in each component is allowed before the contrast is reduced to the specified level. This differential change is then used to set the allowable range of error in each component.

In most cases, direct computation of the change in the contrast is a lengthy procedure, so that a more direct function, such as the rms wavefront error, is used as the quality-defining criterion. In other cases, the quantity of importance will be the focal length, image position, or distortion.

In nonimaging systems, the beam divergence or the uniformity of illumination after passage through the system may be the criterion of interest.

Relating the computed individual errors in the system to the tolerances to be specified is not always a simple matter. If there are several components, some errors may compensate other errors. Thus, it would be easy to assign too tight a tolerance for each surface unless these compensating effects, as well as the probability of a specific distribution of errors, are used in assigning the final tolerances.

Tolerance Budgeting

The method of incorporating compensation of one error by another, as well as the likelihood of obtaining a certain level of error in a defined fabrication process, is called *tolerance budgeting*. As an example, in a lens system it may be found that maintaining the thickness of a component may be easier than keeping the surfaces of the component at the right spherical form. The designer may choose to allow a looser tolerance for the thickness and use some of the distributed error to tighten the tolerance on the radius of curvature. In other cases, the shop carrying out the fabrication may be known to be able to measure surfaces well, but has difficulty with the centering of the lenses. The designer may choose to trade a tight tolerance on the irregularity of the surfaces for a looser tolerance on the wedge in the lens components.

Finally, the effect of a plus error on one surface may be partially compensated by a minus error on another surface. If the probability of errors is considered, the designer may choose to budget a looser tolerance to both surfaces.

This budgeting of tolerances is one of the most difficult parts of a tolerancing process, since judgment, rather than hard numbers, is very much a part of the budget decisions. It is advisable for the engineer carrying out the tolerance budgeting to do some modeling of the system performance using trial sets of parameter variations based on the tolerances that have been obtained. This verification serves as a method of ensuring that the tolerances are indeed reasonable and justified.

Tolerance Verification

Simply stating a set of allowable errors does not complete the integrated process of design and fabrication. The errors must be measurable. Measurement of length can be gauged, but has to be within the capabilities of the shop fabricating the optics. Measurement of error in radius of curvature requires the use of an interferometer or test plates to determine the shape of the surface. Measurement of the nonspherical component of the surface, or the irregularity, requires either an estimate from the test plate, or a computation of the lack of fit to a spherical surface based on measured fringes.

Finally, the quality of a completed lens must be measurable. Use of a criterion that cannot be measured or controlled by the shop or by the user is not acceptable. The contrast mentioned is not always measurable by the optical shop. The surface errors as measured by an interferometer or by test plates are common.

As shown in the Chap. 4 on specifications, the average or rms wavefront error can be related to the level of contrast, or modulation transfer function (MTF), that can be expected in the image. In addition, measurement of the final wavefront from an assembled optical system is most frequently obtained by an optical shop in a summary method by using an interferometer. For this reason, wavefront tolerancing methods have become the most commonly used methods of defining and verifying tolerances.

5.3 WAVEFRONT TOLERANCES

The rms wavefront error tolerancing method will now be used as an example of the approach to evaluating the tolerances required to fabricate a lens. An example which discusses the axial image tolerances for a doublet will be used to provide insight into the tolerancing of a relatively simple system. Most optical tolerancing problems are far more complex, but this example provides an insight into the methods applied. The specific example selected for this chapter is an airspaced achromatic

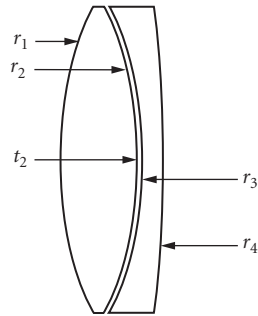


FIGURE 1 Drawing of the doublet lens used for tolerancing.

doublet using BK7 and SF2 glasses, $F/2.8$, 100-mm focal length. The lens design is nominally of moderately good quality, and is optimized over the usual C to F spectral range, with balanced spherical aberration, and is corrected for coma.

Figure 1 is a drawing of the sample doublet used in this chapter. The locations of the four radii of curvature and the airspace to be toleranced are indicated. The number of possible errors that actually can occur in such a simple lens is surprisingly large. There are four curvatures, two thicknesses, one airspace, and two materials that may have refractive index or dispersion errors. In addition to these seven quantities, there are two element wedge angles, four possible tilts, and four decenter possibilities, plus the irregularity on four surfaces and the homogeneity of two materials. So far, there are 21 possible tolerances that are required in order to completely define the lens. For interfacing to the lens mount, the element diameters, roundness after edging, and cone angle on the edges must be considered as well. More complex systems have far more possible sources of error. For the example here, only the four radii and the airspace separation will be considered.

Parameter Error Quality Definitions

The starting point for a tolerance calculation is the definition of a set of levels that may be used to define the initial allowable range of variation of the parameters in the lens. The magnitude of these classes of errors is determined by experience, and usually depends upon the type of fabrication facility being used. Table 1 presents some realistic values for different levels of shop capability.

These values are based on the type of work that can be expected from different shops, and serves as a guide for initiating the tolerancing process. It is obvious that the degree of difficulty in meeting the quality goals becomes more expensive as the required image quality increases.

Computation of Individual Tolerances

The individual tolerances to be applied to the parameters are obtained by computing the effect of some arbitrary but reasonable parameter changes upon the image-quality function. For the example doublet, if made perfectly with no errors in the individual components or assembly, the nominal amount of rms wavefront error at the central wavelength is 0.116 waves, rms. It is determined by the user from consideration of the needs for the application that the maximum amount of error that is acceptable is 0.15 waves, rms. Thus a distribution in allowable errors that results in no greater than about a 0.15 wavelength rms wavefront error would produce an acceptable system. The tolerancing task is to specify the tolerances on the radii of curvature and the separation between the component surfaces such that the goal is met.

TABLE 1 Reasonable Starting Points for Tolerancing a Lens System

Parameter	Commercial	Precision	High Precision
Wavefront residual	0.25 wave rms 2-wave peak	0.1 wave rms 0.5 wave peak	<0.07 wave rms <0.25 wave peak
Thickness	0.1 mm	0.01 mm	0.001 mm
Radius	1.0%	0.1%	0.01%
Index	0.001	0.0001	0.00001
V-number	1.0%	0.1%	0.01%
Homogeneity	0.0001	0.00001	0.000002
Decenter	0.1 mm	0.01 mm	0.001 mm
Tilt	1 arc min	10 arc sec	1 arc sec
Sphericity	2 rings	1 ring	0.25 ring
Irregularity	1 ring	0.25 ring	<0.1 ring

The reason for the choice of 0.15 wave rms is indicated by Fig. 2. The designer experimented with several choices of focus position to obtain a set of plots of the MTF for different amounts of error. This is not a completely general conclusion since the source of error produces an rms error which may not be the same for every source of error. But the samples permit the intelligent selection of an upper bound to the required error. In a lens with more sources of error, and with larger amounts of aberration in the basic design, setting up an example such as this is extremely important to avoid an error in the goal for the final image quality.

The first computation of the effect of nominal changes in the parameters on the rms wavefront error leads to the results in Table 2. (Only the radii and thickness are considered for this example.) The two columns for rms wavefront effect are first, for the aberration if no adjustment for best focus is made, and second, permitting the establishment of best focus after assembly of the lens.

Table 2 shows that the effect of a change in the parameter will have an effect proportional to the change, but that the factor relating the change to the resulting wavefront error is different for the various parameters. The amount of change permitted if the parameter is not compensated by allowing an adjustment for best focus is quite small. Any one of the parameters would have to be maintained within a range far less than the delta used in computation. The allowance of a compensating focal shift does greatly loosen the tolerance.

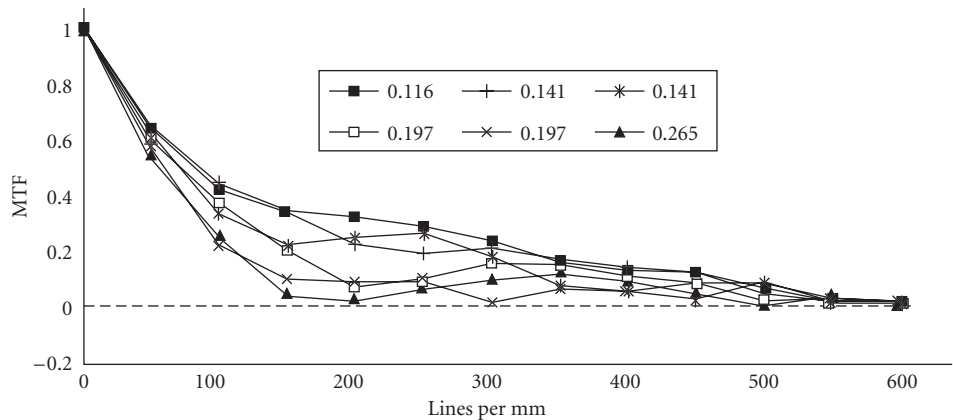
**FIGURE 2** Some examples of the effect of various rms errors on the MTF of the sample doublet. (The rms error is stated in wavelengths.)

TABLE 2 Finite Differentials for Computing Tolerances

Parameter	Delta	Rms Uncompensated	Rms Compensated
$r1$	0.1%	0.740	0.117
$r2$	0.1%	1.187	0.171
$r3$	0.1%	1.456	0.157
$r4$	0.1%	0.346	0.110
$t2$	0.025 mm	1.155	0.152

Since the acceptable goal is 0.15 waves, rms, the amount of change of an individual parameter to attain the acceptable level is about 50 times that of $r1$ and $r4$, but the change of 0.1 percent would be excessive for $r2$ and $r3$. The allowable change for $t2$ alone would be just about the delta of 0.025 mm.

Combination of Tolerances

No parameter in a lens lives alone. The effect upon the image will be the result of combining the effect of all of the errors. If the errors are uncorrelated, then the usual statistical summing of errors can be used. This states that the total amount of aberration produced by the errors can be found by using

$$W_{\text{rms}} = \sqrt{\sum_i W_i^2} = \sqrt{\sum_i \left(a_i x_i \frac{\partial W_{\text{rms}}}{\partial x_i} \right)^2}$$

where the sum is taken over the i parameters of interest. The x factors are the amount of change used in computing the change of wavefront error and the a factors are the relative amount of tolerance error allotted to each parameter in units of the delta used in the computation.

There are implicit assumptions in the application of this method to distributing tolerances. The principal assumption is that the fabrication errors will follow normal gaussian statistics. For many fabrication processes, this is not true, and modification of the approach is required.

For the example of the doublet, Table 3 can be generated to evaluate the different possibilities in assigning tolerances. The allowable change in rms wavefront error is 0.033 waves; thus the root sum square of all of the contributors must not exceed that amount.

In Table 3 the first column identifies the parameter, the second states the delta used in the computation, the third states the amount of wavefront error caused by a delta amount. The final two columns show different budgeting of the allowable error. Distribution 1 loosens the outer radii and the thickness at the cost of maintaining the inner radii very tightly. Distribution 2 tightens the outer radii and spacing tolerances, but loosens the inner radii tolerances. Depending upon the capabilities of the shop selected to make the optics, one of these may be preferable.

The interpretation of these statistical summations is that they are the sum of a number of different random processes. Thus, if the interpretation of each of the values given is the width of a normal distribution, which implies that 67 percent of the samples lie within that value, then 67 percent of the resulting combinations will lie within that range. If the interpretation is a two- or three-sigma value, the interpretation of the result follows similarly.

TABLE 3 Two Possible Tolerance Distributions for the Doublet

Parameter	Delta	Coefficient	Distbn. 1	Distbn. 2
$r1$	0.1%	0.00071	10.0	0.1
$r2$	0.1%	0.05440	0.25	0.4
$r3$	0.1%	0.04069	0.25	0.4
$r4$	0.1%	0.00069	10.0	0.1
$t2$	0.025 mm	0.03589	0.75	0.5
		rms change	0.033	0.033

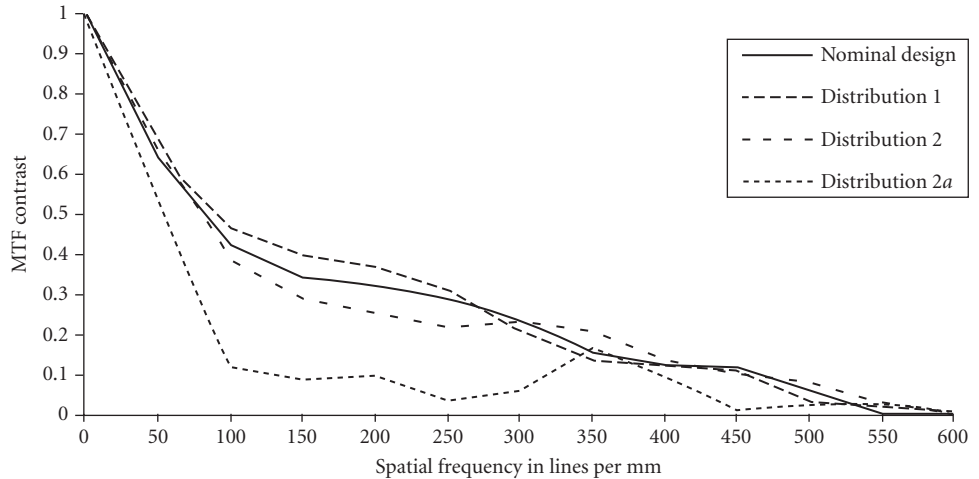


FIGURE 3 The resulting effect of different tolerance budgets for the sample doublet.

Figure 3 shows the effect of the various choices of allowable error distribution on the modulation transfer function of the lens. Either of the distributions stated in Table 3 provides an acceptable lens. For comparison, the allowable tolerances were doubled to provide distribution 2a, which is clearly not an acceptable lens. A spot check of some sample distributions is always relevant when doing tolerancing just to ensure that a reasonable relation between the tolerated system and acceptable image quality exists.

In any manufacturing process, the individual statistics will not necessarily follow a random rule. The interpretation is somewhat modified, but the principle still remains. In some cases, parameters may not be independent. For example, in the doublet there is a linking between the values of r_2 and r_3 that would loosen the tolerances if both are in error in the same direction. This could be used to advantage if the manufacturing process is carefully defined.

Use of Compensators

The use of compensators to loosen the tolerances was indicated above. An example of compensation for aberrations can be seen from a single lens element. If the first curvature is varied, both the element power and the spherical aberration from the element will change. However, a specific change of the second radius can restore the focal position and reduce the change of spherical aberration. Thus the tolerance allowed to the first curvature needs to take into account the possibility of a correlated or deliberate change in the second curvature. It is evident that the proper use of compensators can greatly loosen the tolerances applied to a surface.

A compensation that is frequently employed is the establishment of the correct focal position after assembly of the lens. If this procedure is followed, the individual tolerances on the surface of the elements can be loosened. It is obvious that the tolerancing and the development of a plan for fabrication and assembly must be coordinated.

5.4 OTHER TOLERANCES

Often, a particular optical parameter for the lens must be specified and maintained. Sometimes, for example, the focal length or back focus must be obtained within some tolerance. The computation of these paraxial constants for the lens can be made in the usual manner, and tolerances

obtained by using differentials relating each of the parameters to the quantity, such as the focal length, and then distributing the tolerances in a manner similar to that shown above for the doublet in Table 3.

Boresight

The pointing direction, or boresight, for a lens is sometimes of interest. Errors in boresight are usually due to asymmetric fabrication or mounting errors for the lens. In the simplest manner, one needs simply to trace an axial ray through the lens, and evaluate the direction of this ray as a result of introducing tilts and decenters of the surfaces, or of entire components. Tolerances on the lens parameters can be obtained by the procedure described for the doublet above, substituting the boresight error for the wavefront error.

Distortion

Distortion is the failure of the lens to provide a constant mapping from object space to image space. There are alternate interpretations for this error, which can have radial components due to symmetric errors in the lens, as well as tangential components from tilt and decenter of the lens components. This can be toleranced in the usual manner, but may be related to some general properties of the lens, such as the overall glass thickness of the components. In the simplest case, the tolerances upon distortion may be obtained from direct aberration computation. In complex cases, it may be necessary to compute the actual location of the centroid of the image as a function of image position and in the best image location.

Assembly

Assembly tolerances are related to the tolerances on image quality. The elements must be located and held in position so that the resulting image-quality goals are met. There are additional questions of allowing sufficient clearance between the elements and the lens barrel so that the elements can be inserted into the barrel without breaking or being strained by the mountings. These must be considered in stating the allowable dimensional range in the diameter, wedge, and concentricity of the edge of the lens.

5.5 STARTING POINTS

Shop Practice

Table 3, given as part of the sample tolerancing of the doublet, provides some estimates of the accuracies to which an optical shop may operate. These generic levels of error convey what is likely to be possible. The designer carrying out a tolerance evaluation should consult with probable fabrication shops for modification to this table. The tolerances that are ultimately assigned relate errors in the system to acceptable errors in the image. However, an understanding of shop practice is of great assistance in intelligent budgeting of tolerances.

Measurement Practice

Contemporary practice in optical fabrication and testing is to use interferometry to define wavefronts and surfaces. A convention that has become common in recent years is the use of polynomials

fitted to the wavefront as a method of describing the wavefront. There are several representations used, the most frequent of which is a limited set of Zernike polynomials. These ideally serve as an orthonormal set describing the wavefront or surface up to a specific order or symmetry. In tolerancing, the principal use for the coefficients of the Zernike set has been the easy computation of the rms wavefront error fitted to a given order. Thus the residual error in the system can be described after removal of low-order error such as focus or coma, which can sometimes be attributed to properties of the test setup.

5.6 MATERIAL PROPERTIES

The most important material is optical glass. Specification of the material usually includes some expected level of error in index of refraction or dispersion. In addition, glass is offered having several different levels of homogeneity of index of refraction. The usual range allows for grades of glass having index of refraction inhomogeneity ranging from ± 0.00001 to less than ± 0.0000001 within a single glass blank. It is usually assumed that this variation will be random, but the process of glass manufacturing does not guarantee this.

To place a tolerance upon the required glass homogeneity variation, the concept of wavefront tolerancing can be used. In general, the amount of wavefront error that can be expected along a glass path of length t through the glass is

$$\delta W = \frac{\delta N \times t}{\lambda}$$

For example, if a lens has a glass path of 5 cm but an error of 0.01 wavelengths is assigned to glass homogeneity, then the allowable glass homogeneity is about 0.0000013 within the glass. Thus precision-quality glass is needed for this application. For less glass path or looser tolerance assignment to glass homogeneity error, the required glass precision can be loosened. For a prism, the light path may be folded within the glass, so that an effective longer glass path occurs.

5.7 TOLERANCING PROCEDURES

The example of the doublet serves to illustrate the basic principles involved in determining the tolerances on a lens system. Most lens design programs contain routines that carry out tolerancing to various degrees of sophistication. Some programs are capable of presenting a set of tolerances automatically with only limited input from the designer. The output is a neat table of parameters and allowable ranges in the parameters that can be handed to the shop. This appears to be a quite painless method of carrying out a complex procedure, but it must be remembered that the process is based on application of a set of principles defined by the program writer, and the result is limited by the algorithms and specific logic used. In most cases, some trials of samples of the suggested tolerance distribution will suggest changes that can be made to simplify production of the lens system.

Direct Calculation

The preceding discussion describes methods used in calculating the tolerance distribution for a lens system. Frequently, tolerance determinations for special optical systems are required that either do not require the formal calculation described above, or may be sufficiently unusual that the use of a lens-design program is not possible. In that case, application of the principles is best accomplished directly.

The procedure is first to decide on a meaningful measure of the image quality required. In fact, the term "image quality" may require some broader interpretation. For example, the problem may be to optimize the amount of energy that is collected by a sensor in an optical communication system;

or the goal may be to scan a specific pattern with specific goals on the straightness of the projected spot or line during the scan.

The next step is to express the desired image quality in a numeric form. Usually, the rms wavefront error is the useful quantity. In some cases, other values such as the focus location of a beam waist or the size of the beam waist may be pertinent. In radiometric cases, the amount of flux within a specified area on the image surface (or within a specified angular diameter when projected to the object space) may be the pertinent value.

Once this is accomplished, the third step is to determine the relation between small changes in the parameters of the optical system and changes in the desired image quality function. This is usually accomplished by making small changes in each of the parameters, and computing the value of each differential as

$$\frac{dW}{dx_j} = \frac{\Delta W}{\Delta x_j}$$

where the right side is a finite differential. On occasion, the magnitude of this relation is nonlinear, and may require verification by using different magnitudes of change in the parameter.

This computation provides relations for independent, individual changes of the parameters. The possibility of compensation by joint changes in two or more parameters also has to be investigated. The best approach is to compute changes in the image-quality parameter in which specific coupling of parameters is included. For example, the differentials for variation of the curvatures of the two surfaces of a lens independently will be significantly different than the coupled changes of both surfaces simultaneously, either in the same or different directions.

Spreadsheet Calculation

The differentials must be combined in some manner to provide insight into the tolerance distribution. The best way to accomplish this is to develop a spreadsheet which allows simultaneous evaluation of the combination of errors using the equation for rms summation stated earlier. The use of spreadsheets for calculation is so common today that details need not be covered in this chapter.

Lens-Design Programs

The use of lens-design programs for tolerance calculation has become very widespread because of the proliferation of programs for use on the PC-level computer. The status of lens-design programs changes rapidly, so that any specific comments regarding the use of any program is sure to be out of date by the time this book appears in print. Suffice it to say that all of the principal programs have sections devoted to establishing tolerances. Usually the approach follows the procedures illustrated earlier in this chapter, with finite changes, or sometimes true computed analytical derivatives, used to establish a change table relating parameters to changes in the state of correction of the lens. In this case, the tolerances would be a listing of the allowed changes in the parameter to remain within some specified distance from the design values in aberration space. Some programs use a more complex approach where the allowable change in such quantities as the contrast value of the modulation transfer function at specified spatial frequencies is computed.

The distribution of tolerances is usually established according to the statistical addition rules given above. Some programs permit the user to specify the type of distribution of errors to be expected for various types of parameters.

As recommended above, it is strongly suggested that the user or designer not accept blindly the results of any tolerancing run but, rather, do some spot checking to verify the validity of the range of numbers computed. It is frequently found that alterations in the specified tolerances will occur as a result of such an investigation.

5.8 PROBLEMS IN TOLERANCING

Finally, it is useful to recite some of the problems remaining in establishing tolerances for a lens system. Even though the computational approach has reached a high level of sophistication for some lens-design programs, there are aspects of tolerancing that more closely approach an art than a science. The judgment of the designer or user of a tolerance program is of importance in obtaining a successful conclusion to a project.

Use of Computer Techniques

The use of a computer program mandates the application of rules that have been established by the writer of the program. These rules, of necessity, are general and designed to cover as many cases as possible. As such, they are not likely to be optimum for any specific problem. User modifications of the weighting, aberration goals, and tolerance image-quality requirements are almost always necessary.

Overtightening

The safe thing for a designer to do is to require very tight tolerances. This overtightening may ensure that the fabricated system comes close to the designed system, but the cost of production will likely be significantly higher. In some cases, the added cost generated by the overtight tolerances can raise the cost of the lens to the point where the entire project is abandoned.

The designer should consult with the fabricators of the optical system to develop an approach to assembly and testing that will allow the use of more compensating spacings or alignments to permit loosening of some of the tight tolerances.

Overloosening

A similar set of comments can be made about too-generous tolerances. In many schemes for production, these loose tolerances are justified by inclusion of an alignment step that corrects or compensates for cumulative system error. Too casual an approach to developing tolerances that require specific assembly processes, which are not fully communicated to the project, can result in a lens which is initially inexpensive to build, but becomes expensive after significant rework required to correct the errors.

Judgment factors

The preceding two sections really state that judgment is required. There is no completely “cookbook” approach to tolerancing any but the very simplest cases. The principles stated in this chapter need to be applied with a full knowledge of the relation between a change in a system parameter and the effect upon the image quality. In some cases, a completely novel relationship needs to be developed, which may include, for example, the connection between the alignment of a laser cavity and a nonlinear component included within the cavity. Finite difference calculations to obtain the output level can be developed using whatever computation techniques are appropriate. These values can be combined in a spreadsheet to examine the consequence of various distributions of the allowable errors.

5.9 REFERENCES

There are many useful references on optical tolerances that deal with specific topics not directly covered by the general discussion in this chapter. The most useful suggestion is that the users having the task of setting specifications on a specific product or system use the massive capabilities of internet search engines to look for specific data applicable to that task. A general Google search on “Optical Tolerances” provided 1,600,000 hits, of which probably less than 10 will be applicable to any specific problem.

This page intentionally left blank.

6

MOUNTING OPTICAL COMPONENTS

Paul R. Yoder, Jr.

*Consultant in Optical Engineering
Norwalk, Connecticut*

6.1 GLOSSARY

a_G	acceleration factor
D_G	diameter of optic
E	Young's modulus
ID	internal diameter
K	constant factor
m	mass
OD	outer diameter
P	total preload
S_Y	yield stress
SPDT	single point diamond turning
t	thickness
Δ	deflection of spring or flange
ν	Poisson's ratio

6.2 INTRODUCTION AND SUMMARY

This chapter summarizes the techniques most commonly used to mount lenses, windows, small mirrors, and similar optical components as well as moderate-sized mirrors, and prisms within their mechanical surrounds to form optical instruments. Because of space limitations, mountings suitable for large (i.e., >85-cm diameter) optics are not discussed here. Two basic approaches for mounting optical components are considered: those in which the optic is held firmly against mechanical reference surfaces by applied forces (hard mounting) or those supported by benign means that do not inherently apply force (soft mounting). Descriptions of hard mountings include ones using threaded retaining rings, flanges, or springs while descriptions of soft

mountings include ones using flexures, elastomeric encapsulation, or bonding to mechanical pads. With either type of mounting, the required location and orientation of the optic relative to other portions of the instrument, that is, its alignment, is established during assembly in order to maximize performance. An important aspect of mounting design considered here is how the adverse influences of shock, vibration, temperature change, and moisture on alignment and system performance can be minimized. References cited here provide equations for designing and analyzing a large variety of mountings. Although we speak here of optics as if they are always made of glass and to mechanical parts (housings, cells, spacers, retainers, etc.) as if they are always made of metal, it should be understood that many of these mounting considerations also apply to other materials such as crystals, plastics, and composites.

6.3 MOUNTING INDIVIDUAL ROTATIONALLY SYMMETRIC OPTICS

Hard Mounting Techniques

In order to constrain an optic and preserve its alignment relative to other critical components of an optical instrument, hard mountings apply compressive forces to the glass at discrete locations or along line contacts. These forces, called preloads, are established during assembly and generally are of sufficient magnitude to hold the optic against appropriately located mechanical reference surfaces in the mount under all environmental conditions, including shock, vibration, and temperature changes. The magnitude of the preload P in N applied along any axis should be at least $9.81ma_G$, where m is the mass of the optic and any related components to be held by a single constraining means and a_G is the worst case acceleration expected to be encountered by the subassembly. This term a_G is understood to be a multiple of ambient gravity. If the optic is rotationally symmetric, glass-to-metal contact can be provided at the optic's cylindrical rim, at its ground bevels, or at its polished surfaces. Five degrees of freedom (three translations and two tilts) must be controlled. The sixth degree of freedom (rotation about the optic axis) is also adjusted and controlled in some cases to improve performance in the presence of residual optical wedge or if nonsymmetrical aspheric surfaces are involved. All six degrees of freedom must be constrained for noncircular optics, such as prisms.

The forces applied at the interfaces as well as those from gravity or imposed accelerations may distort the optical surfaces (thereby affecting performance) and introduce stress into the glass. Stress is known to cause birefringence, or, in extreme cases, damage to the optic—especially at low temperatures where shrinkage of the metal exerts maximum force on the glass. To minimize these adverse effects, forces must be kept within acceptable limits. Very few closed-form equations are available for predicting refracting or reflecting surface deformations due to applied forces. Finite element analyses are most frequently used for this purpose.^{1,2} Explanation of how this is done is beyond the scope of this presentation.

Relatively simple analytical means for estimating compressive and tensile stresses introduced by mounting forces are detailed in the literature.^{3,4} Most of these techniques are based on adaptations of standard formulations by Roark⁵ and Timoshenko and Goodier.⁶ The magnitude of the stress generated by a force depends not only upon the magnitude of that force, but also on the shapes of the surfaces in contact and Young's modulus and Poisson's ratio values for the glass and metal involved.

Statistical analyses backed by experimentation indicate that an optical component made by conventional high-quality grinding and polishing methods can usually withstand tensile stress as large as ~ 6.9 MPa without failure. This value is generally accepted as a "rule-of-thumb" tolerance for survival of the optic under stress.⁷ Optics made by "controlled grinding" techniques,⁸ polished and assembled with great care, and not scratched or otherwise damaged during use, might well survive long-term stress about 1.7 times greater.⁹

Under the more benign conditions of the operating environment (wherein the instrument must perform to specifications), survival is not a concern, but distortions of optical surfaces

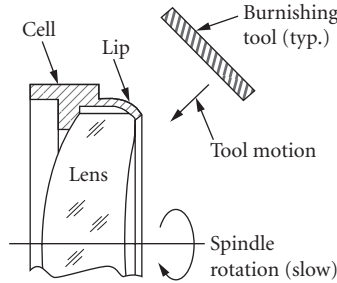


FIGURE 1 A small lens burnished into a cell made of malleable metal.

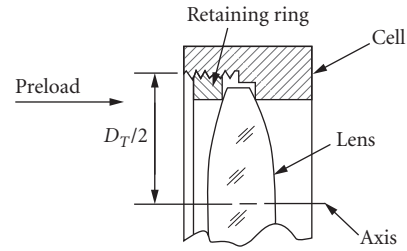


FIGURE 2 A lens preloaded in its cell with a threaded retaining ring.

due to mounting forces may degrade performance. High-performance optical systems and those using polarized light may also be especially sensitive to stress-induced birefringence. A commonly applied tolerance for stress in the glass in such cases is ~ 3.4 MPa. Analytical methods outlined in this chapter allow mounting stresses to be estimated to predict the potential success of a given design.

Burnished Mountings Figure 1 shows a very simple way to mount a lens element in a tubular cell. This cell has an internal shoulder against which the lens is to be held. The cell is mounted on a spindle, the lens is inserted and held in place gently, and the assembly is slowly rotated. The cell has a lip that extends beyond the rim of the lens. That lip is burnished with one or more hardened rod-shaped tool(s) over the edge of the lens as indicated in the right-hand view. The cell material must be malleable so it can be bent easily. Brass or annealed aluminum are common choices. The magnitude of the force, if any, introduced by the mounting cannot be quantified in this case because the bent metal tends to spring slightly away from the glass once tool pressure is removed. This type mounting is most suitable for use with small elements used in some endoscopes, simple microscope objectives, or low-cost cameras.

Mounts Using Threaded Retaining Rings This mounting, shown schematically in Fig. 2, is the type most frequently used to secure a lens element in its mount. Torque Q in N-mm applied to the ring with a wrench creates axial preload P to hold the lens against the shoulder very approximately as $5Q/D_T$, where D_T is the pitch diameter of the thread in millimeters.

The fit of the mating threads in the cell and retainer should be loose enough for the retainer to align itself to the centered lens surface; otherwise, lens alignment may be altered when the retainer is tightened. Such a fit may be specified as Class-1 or -2 per ANSI/ASME B1.1-2003.* During assembly, the lens should first be aligned in the cell and then held in place as the retainer is tightened to the required torque.

Mounts Using Annular Flanges Figure 3 shows a lens element preloaded against a shoulder by an annular flange that is deflected axially by a distance Δ from its nominal flat shape. Adapting an equation from Roark,⁵ the deflection Δ required to produce a given preload P in N equals $(K_A - K_B)P/t^3$ where t is the flange thickness in millimeters and the constants K_A and K_B are determined by the material properties and the dimensions a and b indicated in the figure.

For a given design, the required deflection may be obtained by customizing the thickness of the spacer located under the flange and should be at least 10 times larger than the resolution capability of the device to be used to measure the flange deflection at the time of assembly. As the flange is bent, stress is developed within that component. To prevent damage to the flange, its thickness should equal $K_C Pf_S/S_Y$ where the constant K_C depends upon the dimensions a and b and the flange material

*Unified Inch Screw Threads (UN and UNR Thread Form).

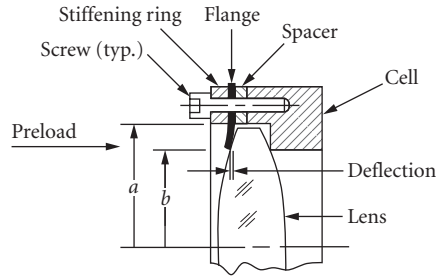


FIGURE 3 A lens preloaded in its cell with a deflected continuous ring flange.

properties. The quantity S_y is the yield stress of the flange material and f_s is the desired safety factor. The stiffening ring shown next to the flange maintains uniform flange deflection between the attaching screws.

A significant advantage of the annular flange as compared to the threaded retaining ring is that the flange can be calibrated before installation by measuring the actual preload developed as a function of deflection. Then, one can be quite confident that the preload on the lens is as stated by the above relationship when that flange is deflected by the specific distance Δ . This level of confidence cannot be achieved with a threaded ring.

Soft Mounting Techniques

Elastomeric Mountings A convenient technique for mounting a lens in a cell is to inject a continuous annular ring of an elastomeric material such as room temperature vulcanizing (RTV) sealing compound between the lens rim and the inside surface of the cell (see Fig. 4). This is sometimes called an “elastomeric ring mounting.” The thickness t_e equals $K_D D_G$, where D_G is the lens diameter and the constant K_D is determined from the material properties using a relationship attributed by Herbert¹⁰ to R. Vanbezooijen. The lens is then virtually free of radial stress at all temperatures. This is because the elastomer expands or contracts with temperature changes just enough to always fill the radial gap between the glass and the metal.

Some designs using the elastomeric ring approach also benefit from the fact that a continuous ring of this material effectively seals the lens to its mount so, if this subassembly forms part of the exterior skin of an optical instrument, leakage of gases and moisture through that interface is prevented.

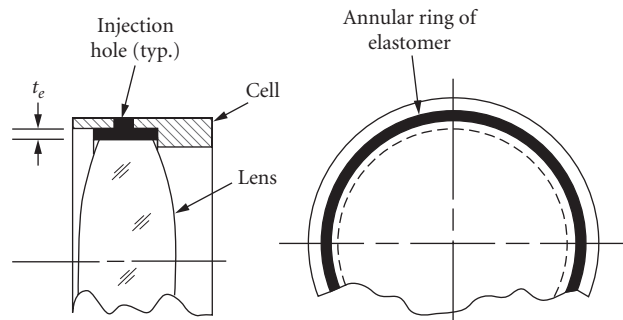


FIGURE 4 A lens supported in its cell by a continuous annular ring of elastomer.

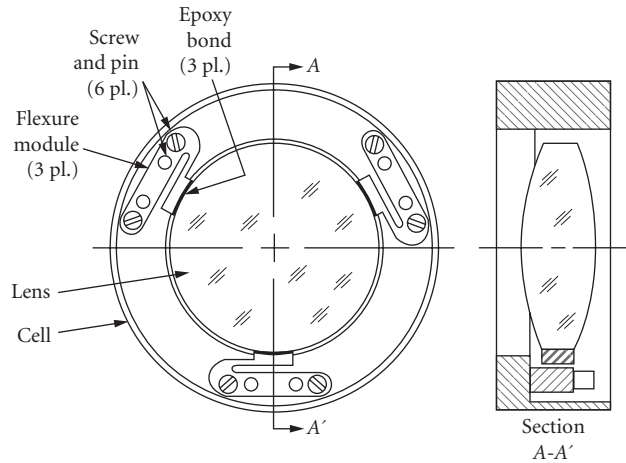


FIGURE 5 A lens bonded to three flexure modules attached to a cell.

Another type of elastomeric mounting for lenses uses discrete pads of elastomer located between the lens rim and the cell's inside surface. At least three such pads are needed to fully constrain the optic. They should be symmetrically distributed around the lens. The radial thicknesses of these pads can be sized as described above for the continuous ring.

Flexure Mountings High-performance optical systems require the optical axes of their lenses to be precisely centered mechanically with respect to some mechanical reference and to remain in that condition when the temperature changes. Because metal and glass components expand or contract with temperature at different rates, the optics may become decentered, tilted, or stressed when temperature changes occur in the above-described hard mountings. A properly designed support configuration using three or more symmetrically located identical flexures between the mount and lens rim will ensure that the lens stays as originally aligned and free of stress in spite of such changes.

Figure 5 shows a concept for a simple flexure mount design suggested by Ahmad and Huse.¹¹ Three identical flexure modules are made with narrow slots cut into them (by an electric discharge machining method) to form cantilevered flexure blades. Each blade has, at its free end, a curved pad shaped to interface with the lens rim. These modules are attached to the lens cell with screws passing through slightly oversized holes. In an alignment fixture, the optical axis of the lens is centered with respect to the axis of the cell. The modules are then adjusted to provide specific gaps between the pads and the lens rim and pinned in place. Epoxy is injected into those gaps and cured. Because the flexures are separate from the cell, they can be made from a material (such as titanium) with a higher yield stress than the cell. The cell is typically made of less expensive yet dimensionally stable material (such as stainless steel). More complex flexure designs and ones featuring a larger number of radial flexures have also been described.¹²⁻¹⁶ Because of their inherent flexibility, flexure mountings should be analyzed to determine their responses to externally imposed shock and vibration.

6.4 MULTICOMPONENT LENS ASSEMBLIES

Groups of lenses used in optomechanical assemblies typically are individually mounted and constrained in seats machined into a housing or are separated axially by spacer rings within a common cylindrical bore in the housing. It is important for those lenses to have a common optical axis and the correct axial airspaces within allowable tolerances. We here consider several ways in which such assemblies can be designed.

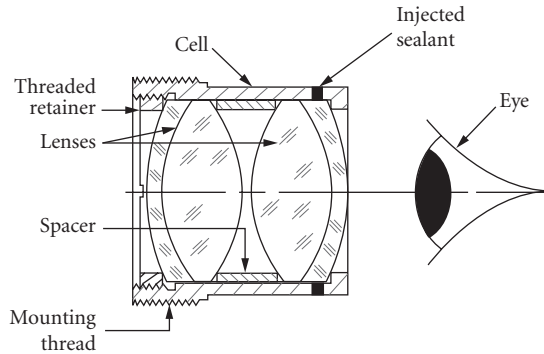


FIGURE 6 A telescope eyepiece with two lenses and a spacer assembled by the “drop-in” method.

Drop-In Assembly

If alignment requirements are not too demanding, lenses, housings, and spacers can be machined to reasonable tolerances and simply assembled without further machining or adjustment other than tightening a retainer to provide prescribed preload. Figure 6 shows an eyepiece for a low-power telescope assembled in this way. Radial clearances are typically ~ 0.075 mm, so individual lenses can easily be inserted into their seats. The eyepiece is configured to thread into a cylindrical opening in the telescope housing and to be focused by rotating the entire eyepiece. In some cases, the lenses are sealed to the cell and the threaded joint with the telescope also is sealed.

Many all-plastic lens assemblies used in consumer products are designed for swift drop-in assembly. The example shown in Fig. 7 is the objective for a rear-projection television system. Flexible tabs molded into the inside walls of both halves of the plastic housing project inward to form pockets for insertion of the three injection-molded plastic lens elements. Grooves (not shown) molded into the inside walls of the housings reduce stray light that otherwise could reduce contrast of the image. The housings are fastened together by self-tapping screws passing through flanges along each side, as indicated in the end view. Optical alignment relies on accuracy of the molding processes and is adequate for the application.

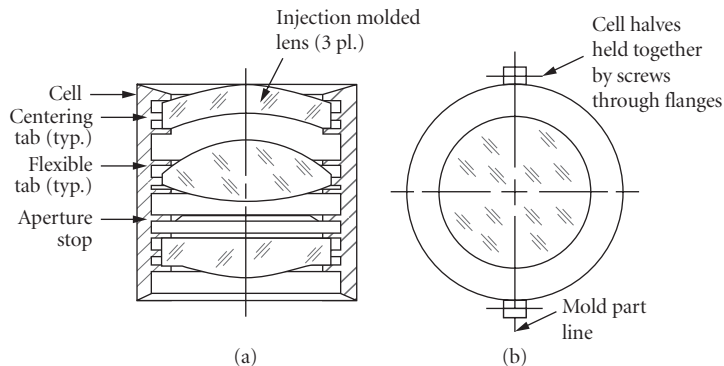


FIGURE 7 An all-plastic projection lens assembled by the “drop-in” method.

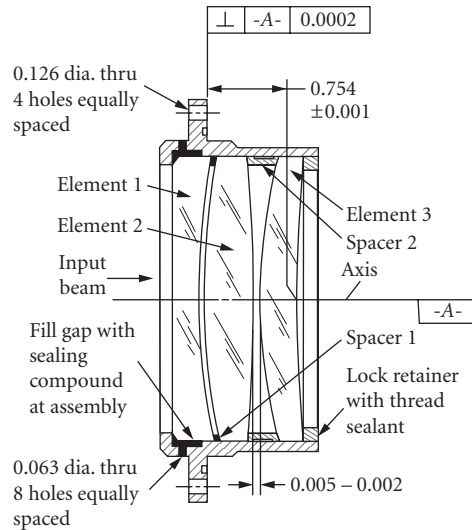


FIGURE 8 A telescope objective assembly with alignment resulting from tightly tolerated dimensions. Dimensions are in inches. (From Yoder.¹⁷)

Tightly Toleranced Assemblies

When higher performance is required, the dimensional tolerances are tightened and better optical alignment is achieved. Figure 8 shows the objective lens assembly for a military telescope.¹⁷ The three lenses are edged to fit the cell inside diameter with nominal radial clearances of 0.012 mm. All metal parts are made of stainless steel. The first spacer is made of sheet metal stock 0.025 ± 0.005 mm thick. It conforms to the spherical shapes of the adjacent lens surfaces under preload. The axial thicknesses of the lenses are tolerated to ± 0.005 mm. Residual optical wedge tolerances for the lenses are 12 arcsec. The beam deviation from these wedges is minimized at assembly by rotating (i.e., clocking) two lenses about their axes relative to the third lens to obtain maximum symmetry of the image of an on-axis artificial star. This image is observed with a microscope during alignment.

Lathe Assembly

A technique that is frequently used to obtain lens centration by minimizing radial clearances between lens ODs and cell IDs is called “lathe assembly” because it is done on a machinist’s lathe. The diameters and thicknesses of a selected set of lenses are measured and recorded. The required central air spaces and their tolerances are obtained from the optical system design. Actual lens surface radii are obtained from interferometric measurements made during lens manufacture. This data accompanies the lenses to the machine shop where a partially machined cell or housing is customized to provide conical or toroidal interfaces with the polished surfaces of that particular set of lenses and to provide all other required dimensions for the optomechanical assembly within the required tolerances. Radial clearances of ~ 0.005 mm can be achieved by this method. This clearance is adequate for careful assembly of the lenses into the cell.

Figure 9¹⁷ shows an air-spaced doublet lens subassembly created by this process. The individual lens seats are finish-machined at the time of assembly to fit those lenses. The length of the spacer (dimension E) and the location of the mounting flange (Datum B) relative to the front lens vertex

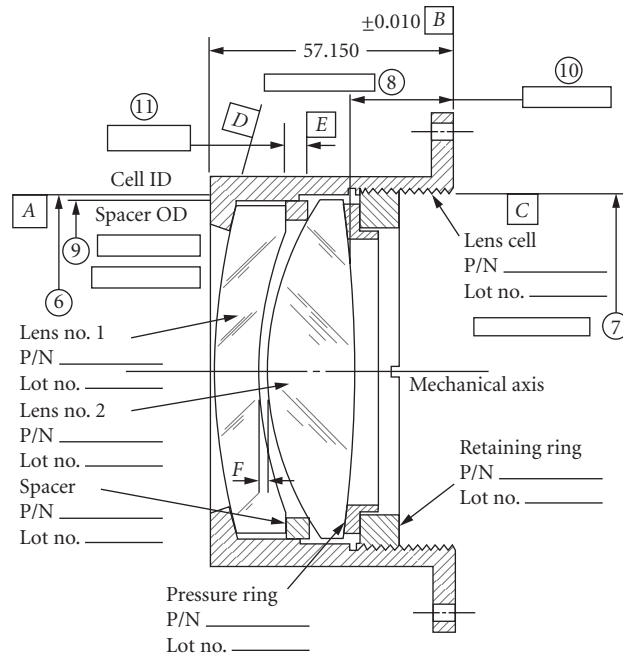


FIGURE 9 An air-spaced doublet assembled by the “lathe assembly” process. (From Yoder.¹⁷)

are also machined to produce the proper air space and overall length. Both lenses are secured by the retainer. The actual values of the numbered dimensions are recorded in the boxes. If required, a complete pedigree of that particular assembly can be established for future reference from the measured data and inspection reports. This type of construction is especially suited for applications involving high accelerations. Bayar described an aerial camera lens assembled by this method.¹⁸

“Poker Chip” Assembly

Figure 10 is a partial section view through a lens assembly that features seven lenses: four doublets and three singlets. Each lens, except the largest, was centered interferometrically to the mechanical axis of its cell OD and held in place in that cell with annular rings of epoxy nominally 0.381 mm thick. After the adhesive was cured, the axial thicknesses of the cells were final machined so all axial air spaces would be within design tolerances. The cell subassemblies Numbers 6 through 2, which had been machined to the same ODs within tight tolerances, were then inserted into the stainless steel housing and secured by Cell No. 1 that was threaded into the housing to act as a retainer. The largest lens (No. 12) was held directly in the housing by its own retainer. Accuracy of internal alignment was built into the assembly by the fabrication process.¹⁹ This type of construction is frequently referred to as “poker chip” assembly because the individual lens/cell subassemblies are stacked on top of each other inside the housing.

Lenses Adjusted at Assembly

Many complex lens assemblies to be used in very high-performance applications such as micro-lithographic projection systems need positional adjustment of a few carefully selected elements at the final stage of assembly. This is because application of the best possible optical and mechanical

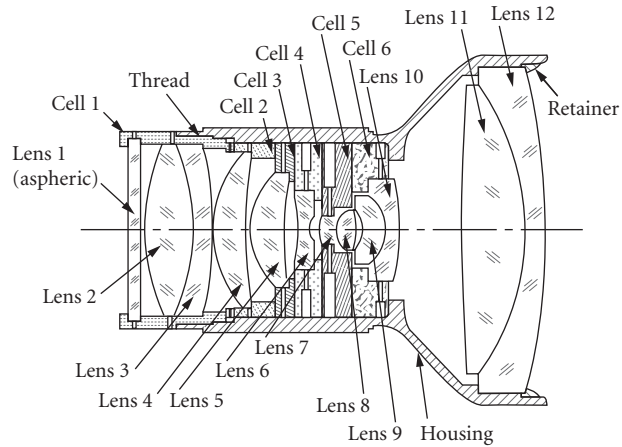


FIGURE 10 A projection lens comprising a stack of “poker chip” lens/cell subassemblies inserted into the bore of the mount. (From Fischer.¹⁹)

manufacturing processes and extremely tight dimensional tolerances cannot make the lenses and mechanical parts accurately enough to obtain the full level of performance required by the application. An example is shown schematically in Fig. 11.²⁰

This optomechanical system comprises twelve air-spaced “poker chip” subassemblies, stacked on top of each other with custom-made spacers placed between lapped coplanar pads on the cell faces

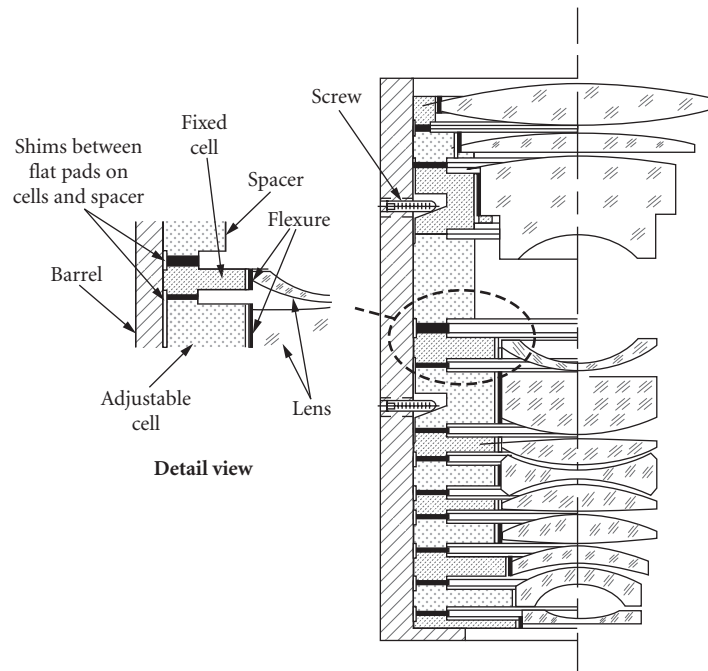


FIGURE 11 Partial section view of a “poker chip” lens assembly with two lenses adjusted after assembly to optimize performance. (From Yoder.⁴)

to control axial air spaces. The lenses are mounted on flexures machined directly into the interior surfaces the cells. Ten of the subassemblies fit closely inside the stainless steel barrel. Two subassemblies are adjustable laterally in orthogonal directions from the outside. Optical performance of the system is measured interferometrically in near real time, while the adjustable lenses are moved very slightly until optimal performance is achieved. The adjustment mechanisms are then locked and the lens is installed into the microlithography system.

Determination of which lens elements to move in a given optical system to correct residual aberrations is a job for the lens designer working with mechanical engineers and metrology experts who help decide how to incorporate the needed mechanisms and to conduct the necessary tests. The sensitivities of spherical, coma, astigmatism, and distortion aberration contributions from each lens to lateral and axial displacements are determined by raytracing. The ideal candidates for correcting each aberration are lens shifts that modify that aberration significantly, but do not excessively affect the other aberrations. The results are reviewed to determine which lens movements are best to minimize each aberration. Williamson²⁰ outlined a procedure in which the aberration contributions of the optical system shown in Fig. 12a for specific axial and lateral displacements of each element are plotted as shown in Fig. 12b and c, respectively. Using phase-measuring ultraviolet interferometry

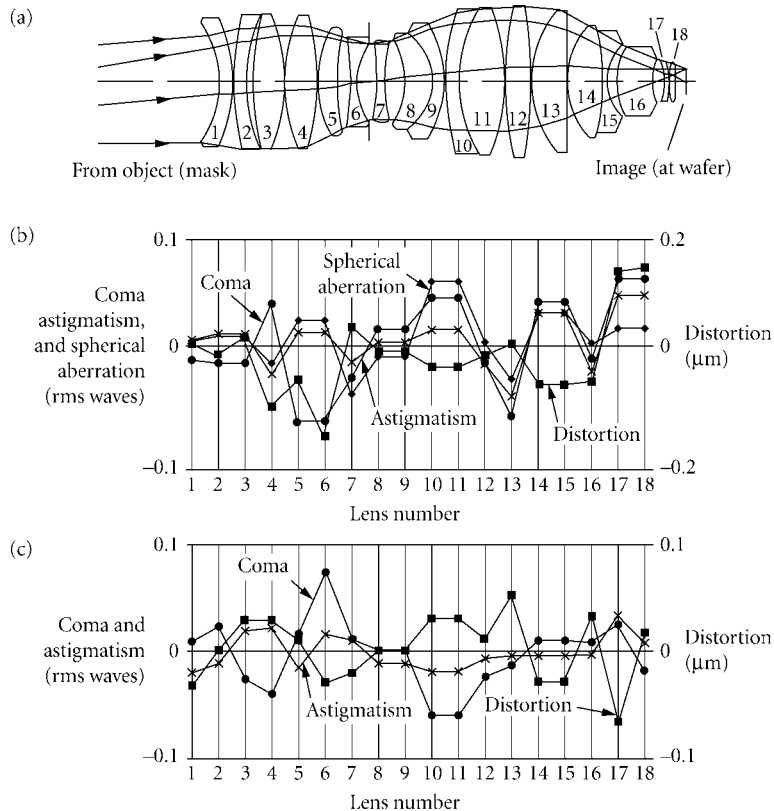


FIGURE 12 (a) Optical schematic of an 18-element lithographic projection lens. (b) The effects on aberrations of individually displacing each element axially by $25\ \mu\text{m}$. (c) Similar effects of displacing each element laterally by $5\ \mu\text{m}$. The best lenses to adjust for optimum system performance can be determined. (From Williamson.²⁰)

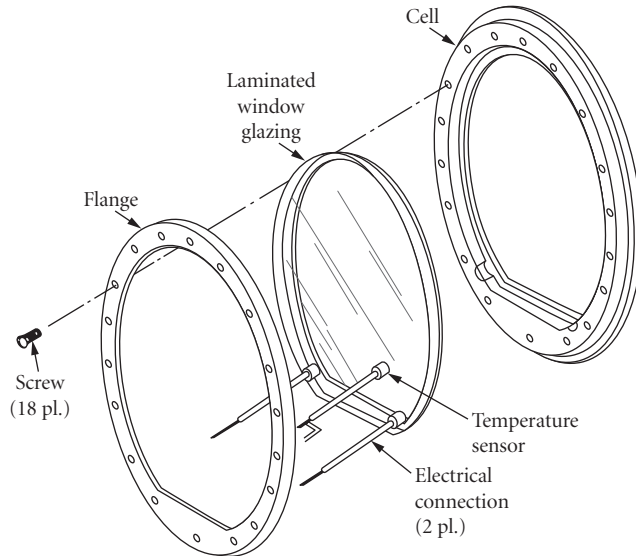


FIGURE 13 Exploded view of a heated window assembly used in a military application. (Courtesy of Goodrich Corporation, Danbury, CT)

as the image quality monitor, all these aberrations can be significantly reduced by iteration of lens movements and the final system performance of production lens assemblies greatly enhanced.

6.5 MOUNTING WINDOWS AND DOMES

Small circular windows usually are secured in a mount with a threaded retainer or by an elastomeric ring. Noncircular ones are best held in place with elastomer. The continuous flange-mounting method can be used to advantage with larger windows. The one shown in Fig. 13 has an elliptical aperture of 20.32×30.5 cm.⁴ The electrical connections shown provide current to a conductive coating on a buried surface, which keeps the window free of fog in high-humidity situations. The flange preloads the window into an aluminum cell. An elastomeric sealant is injected into a groove around the window's rim to seal it to the cell. The cell is sealed to the instrument housing with a gasket or an O-ring.

Figure 14 shows typical mountings for deeply curved spherical windows, called *shells* or *domes*. That in Fig. 14a is sealed and secured with a Neoprene gasket clamped in place by a flange²¹ while that in Fig. 14b is secured and sealed with a continuous ring of elastomer.⁴ In some more elaborate designs, an elliptically shaped sapphire dome is brazed with special metallic alloys to a titanium mount.²²

6.6 MOUNTING SMALL MIRRORS AND PRISMS

General Considerations

The appropriateness of designs for mechanical mountings for small mirrors and prisms depends upon a variety of factors including: tolerable rigid body movement of the optic and distortion of its reflecting and/or refracting surface(s); the magnitudes, application locations, and directions

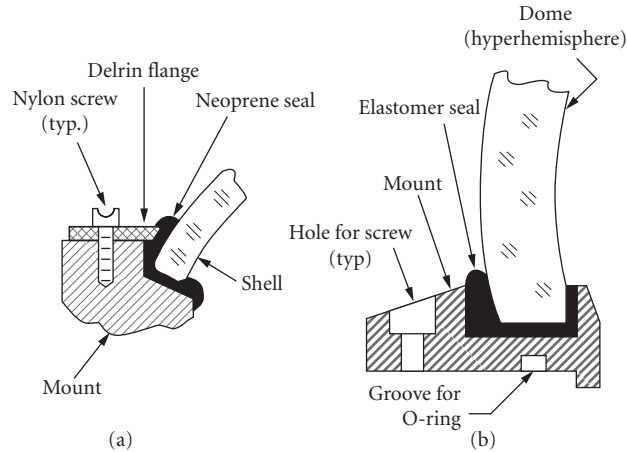


FIGURE 14 Typical mountings for (a) a thin shell and (b) a hyperhemispherical dome. [(a) From Vukobratovich.²¹ (b) from Yoder.⁴]

of forces tending to move the optic with respect to its mount; steady-state and transient thermal effects (including gradients); the sizes and kinematic compatibility of interfacing optomechanical surfaces; and the rigidity and long-term stability of the structure supporting the optic. In addition, the designs must be compatible with assembly, maintenance, package size, weight, and configuration constraints, as well as being cost effective. The representative mounting designs described in the following sections illustrate proven mounting techniques.

Mechanically Clamped Mountings

Figure 15 shows a simple means for attaching a first surface flat mirror to a mechanical bracket.²³ Three cantilevered springs press the reflecting surface against three pads that have been lapped coplanar. The contacts between the springs and the mirror's back face are directly opposite the pads to minimize bending moments. This design constrains one translation and two tilts. Translations in the plane of the reflecting surface can most easily be constrained by dimensioning the spacers supporting the springs so as to just clear the rim of the mirror at minimum temperature. Rotation in that plane usually does not need to be constrained. Given the number of springs N , the spring material's Young's modulus E_M , its yield stress S_Y , and its Poisson's ratio ν_M , the spring lengths L and widths b , and an appropriate safety factor f_S , the spring thickness t that will provide a total preload P to the mirror is determined as $[K_{S1}PLf_S/(bS_YN)]^{1/2}$. The length of the spacer located under each spring is chosen to cause that spring to be deflected from its flat condition by Δ equal to $(K_{S2}L^3)/(1 - \nu_M^2)(E_Mbt^3N)$. In these relationships, K_{S1} is 4 and K_{S2} is 0.75.

Elastomeric Mountings for Mirrors

Small mirrors can often be mounted in the manner illustrated by Fig. 4 for a lens. In applications where the optic does not need to be sealed in place with a continuous ring of elastomer, three or more discrete pads located between the lens rim and the cell ID can support it. Vanbezooijen's equation is again used to determine the pad thicknesses.¹⁰ The lateral dimensions of the pads have, in some designs, been determined by finite element analysis that predicts the dynamic response of the subassembly to vibration inputs from the environment.²⁴

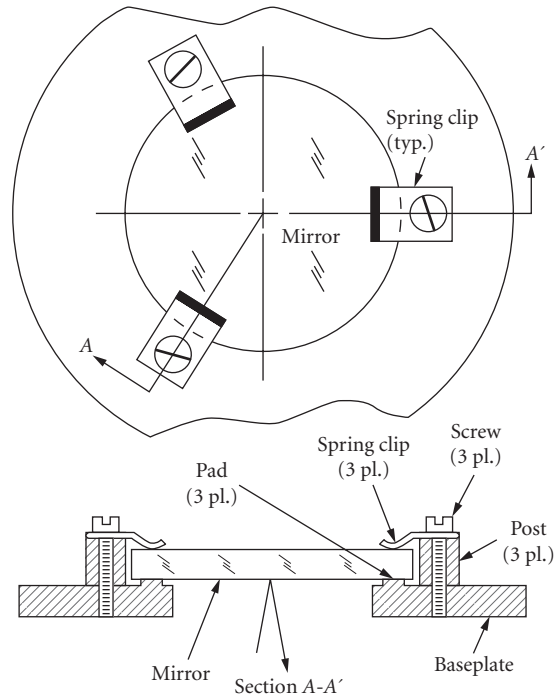


FIGURE 15 A simple mounting for a flat mirror preloaded with deflected cantilevered springs. (From Yoder.²³)

Spring-Loaded Mountings for Prisms

A spring-loaded mounting for a prism is illustrated in Fig. 16.²⁵ Here, a penta prism is preloaded against three coplanar pads on the baseplate by three cantilevered springs supported by posts with spacers machined to produce the necessary spring deflection and resultant preload, as described earlier for a mirror mounting. Constraint in the plane parallel to the pad surfaces is provided by a single spring (called a straddling spring) that is supported at each end and presses against the end of the prism. The dimensions and deflection of this spring are chosen to preload the prism against three locating pins that are pressed or threaded into strategically located holes in the baseplate. The relationships for t and Δ given in Sec. “Mechanically Clamped Mountings” apply also to the straddling spring, but K_{s1} equals 0.75, $N = 1$, and $K_{s2} = 0.0625$. How the direction of the force exerted by the straddling spring can be optimized to nearly equalize the stresses created in the prism at the interfaces with the pins has been explained in the literature.²⁶

Bonded Mountings for Small Mirrors and Prisms

A widely used and successful technique for mounting small mirrors and prisms is to bond them directly to a plate or bracket with an adhesive such as epoxy. Any alignment adjustments that are needed should be built into the mount rather than into the glass-to-metal joint. Figure 17 shows a typical mirror mount of this type.²³ It has proven satisfactory for cases where the diameter-to-thickness ratio for the mirror substrate is at least 6:1. The mirror should then be stiff enough not to be excessively

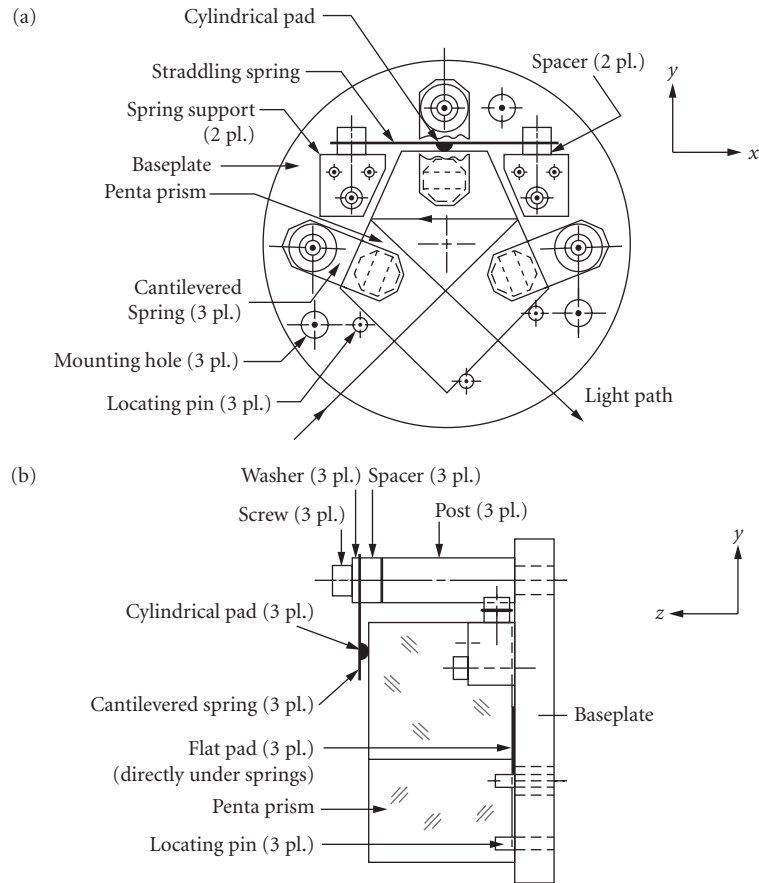


FIGURE 16 A penta prism preloaded against lapped pads on a baseplate with cantilevered and straddling springs. (From Yoder.²⁵)

distorted by shrinkage of the adhesive as it cures. Most prisms are thick enough that they are not distorted by this shrinkage.

All adhesive bonds need to have sufficient area for the joint to be strong enough to support the optic under all anticipated levels of acceleration. The minimum bond area Q_{MIN} is $9.81ma_g f_s/J$, where J is the strength of the cured adhesive joint and all other terms are as previously defined.²⁷ Bonding should be done on a fine ground surface of the optic for maximum joint strength. A typical value for J for a two-part epoxy such as 3M 2216B/A with bond thickness of 0.100 ± 0.025 mm is ~ 17.2 MPa. For conservative design, the factor f_s should be ~ 4 . Successful bonding requires careful cleaning of the surfaces to be bonded and adequate curing time. The adhesive manufacturer's recommendations should be followed unless tests indicate otherwise for a specific application.

The 29-mm aperture roof penta prism in Fig. 18 is bonded in cantilevered fashion to a bracket nominally oriented vertically. The circular bond area is adequate to withstand a severe military shock and vibration environment. Some designs work better if the prism is supported from both sides. Figure 19 shows one way to do this. It was adapted from Beckmann.²⁸ The mount is designed with two arms, one of which has a hole bored through it. The prism is supported by a fixture in the proper location and orientation relative to the mount and epoxy bonded to the flat pad on the left arm. A plug made of the same

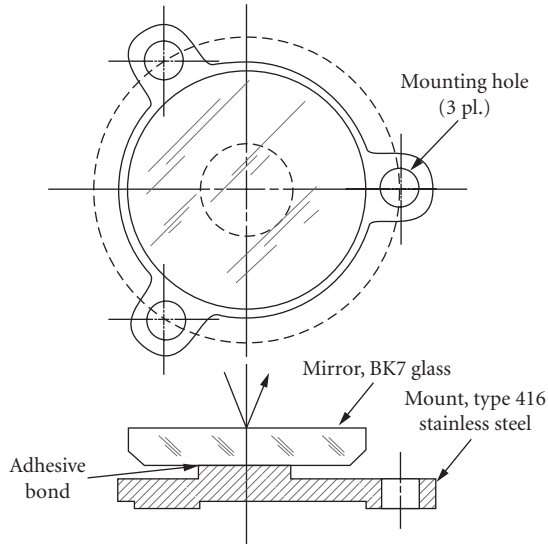


FIGURE 17 A flat mirror bonded on its back to a pad on its mount. (From Yoder.²³)

metal as the mount is then centered in the hole in the right arm and bonded to the right surface of the prism. When those bonds have cured, the plug is bonded into the right arm.

Flexure Mountings for Small Mirrors and Prisms

Circular mirrors as large as ~15-cm diameter have been successfully mounted on flexures in the general manner shown in Fig. 5 for a lens. Usually, these are image-forming mirrors, perhaps aspheric, that need to have constant centration relative to a system axis.

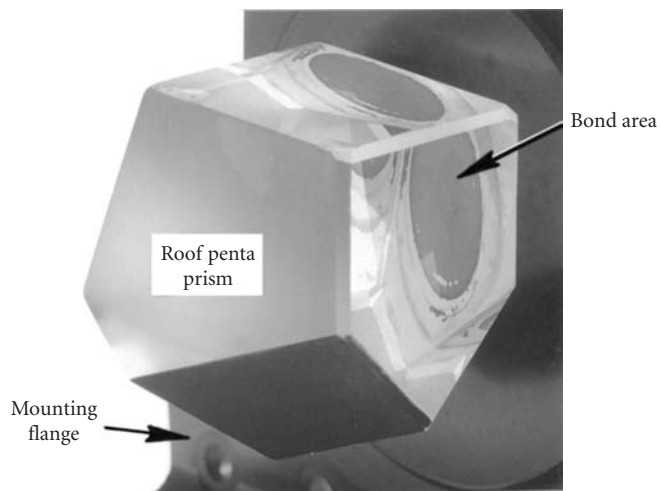


FIGURE 18 A roof penta prism bonded to a pad on a mounting flange.

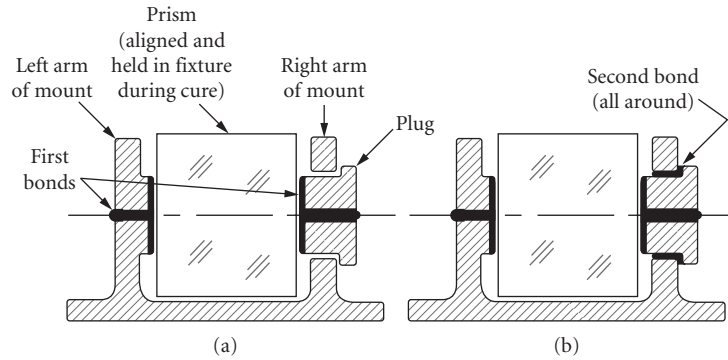


FIGURE 19 Concept for supporting a prism from both sides. (a) Prism bonded to left arm of mount and plug bonded to prism. (b) Plug bonded into right arm of mount. (Adapted from Beckmann.²⁸)

Prisms intended for use in relatively benign environments can be mounted on flexures. One way to do this is by attaching the prism to instrument structure through three posts with integral flexures at each end. Figure 20 shows a large multiple component Zerodur prism mounted in this manner. It has two wing prisms optically contacted to a third (base) prism, to which the flexures are bonded. The wing prism surfaces are perpendicular to each other and form a 15.2-cm-wide roof mirror that is inclined at 45° to the vertical. The reflected image is inverted horizontally as the optic turns the incident beam axis 90° . The orientations of three of the flexure joints are

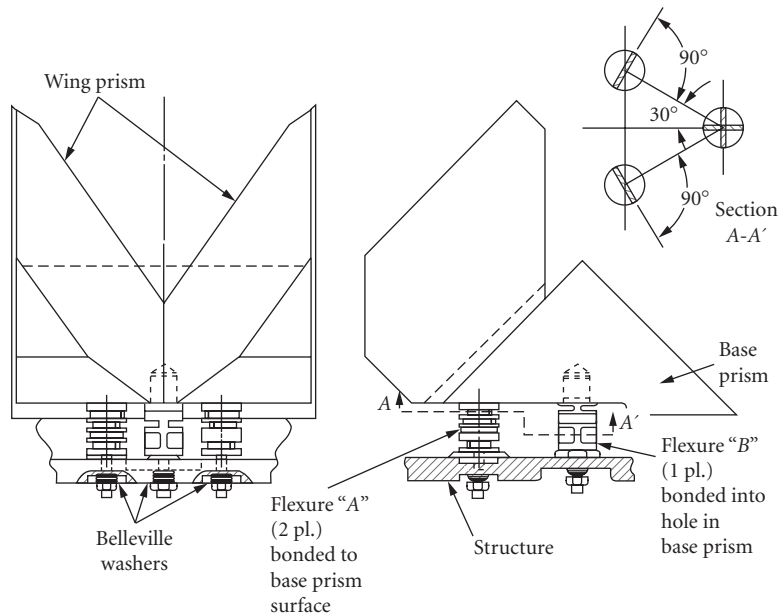


FIGURE 20 Optomechanical configuration of a large prism assembly with three flexure mounting posts to isolate the optic from dimensional changes under temperature changes. (Courtesy of ASML Lithography, Wilton, CT.)

as indicated in the section view A-A'. If attached to a structure that expands or contracts more than the prism as the temperature changes, the flexures simply bend very slightly and prevent the introduction of mounting forces that could distort the reflecting surfaces and interfere with performance of the optical system.

6.7 MOUNTING MODERATE-SIZED MIRRORS

General Considerations

The simple mirror mountings described earlier are not satisfactory for mirrors larger than about 15-cm diameter because they are too flexible to be treated as rigid bodies. The important criteria for selecting a suitable mounting are orientation with respect to gravity, performance level required, substrate material stiffness, and weight limitations. The mounting for a mirror to be used in a fixed horizontal- or vertical-axis orientation can be figured during polishing to compensate for gravity effects. Variable orientation applications require mounts that change their force distribution with inclination to keep surface deflections within tolerance. Both axial and radial supports are required. Mirrors to be used in space have the added requirement of release of gravitational force after being fabricated, tested, and installed into the instrument in a normal gravity environment. Choice of mounting depends strongly on the substrate configuration. Weight constraints generally lead to solid substrates with shaped back surfaces or ones built-up from multiple parts that are attached together. We here describe a few typical ways to support mirrors of various shapes as large as ~85 cm. Designs appropriate to both nonmetallic and metallic mirrors are considered.

Substrate Configurations

Figures 21*b* through *e* shows half-section views of four first-surface mirror solid substrates of the same diameter and material with concave surfaces of the same radius. Their back surface shapes differ and reduce the mirror weight as compared to a flat-back baseline design (Fig. 21*a*). All these mirrors, except one, can be supported within the telescope housing on a hub passing through the mirror's central perforation. For example, see Fig. 22. Here, the hub has a toroidal-shaped located land that supports the 41-cm-diameter meniscus-shaped mirror radially and a shoulder that locates it axially. The radial support lies in the mirror's neutral plane where fore and aft bending moments are balanced. A threaded retaining ring provides axial preload. To focus, the locating ring is moved on the hub and secured with the clamping ring. The substrate configuration from Fig. 21 that cannot be hub mounted is the double arch configuration (Fig. 21*e*). It is best supported on flexures at three or more points spaced equally around the zone of greatest thickness.

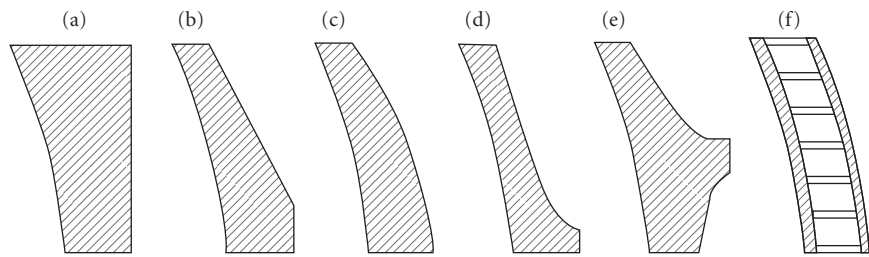


FIGURE 21 Sectional views of baseline concave-plane (a) and lightweighted mirror substrates (b) through (e) with contoured backs. Figure (f) shows a built-up substrate configuration.

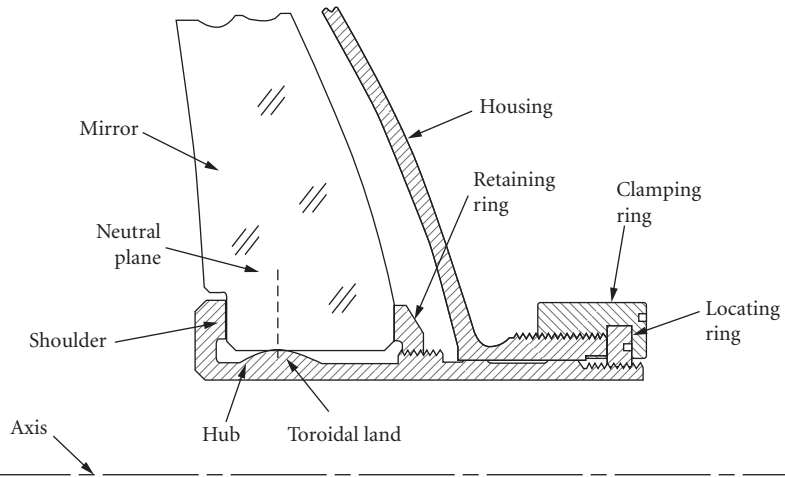


FIGURE 22 Hub mounting for a meniscus-shaped telescope mirror. Focus adjustment means is illustrated.

Lighter-weight mirror constructions typically employ built-up substrates such as that shown in Fig. 21*f*. A very successful type is the monolithic meniscus construction illustrated by Fig. 23. Such mirrors are usually made of Corning ULE. Strips of the material form the webs of a core to which front and back facesheets are fused. All joints in the core also are fused together. The spacing of the webs is large except at locations where axial and radial supports attach to the substrate. There, the spacing is considerably smaller to increase strength.

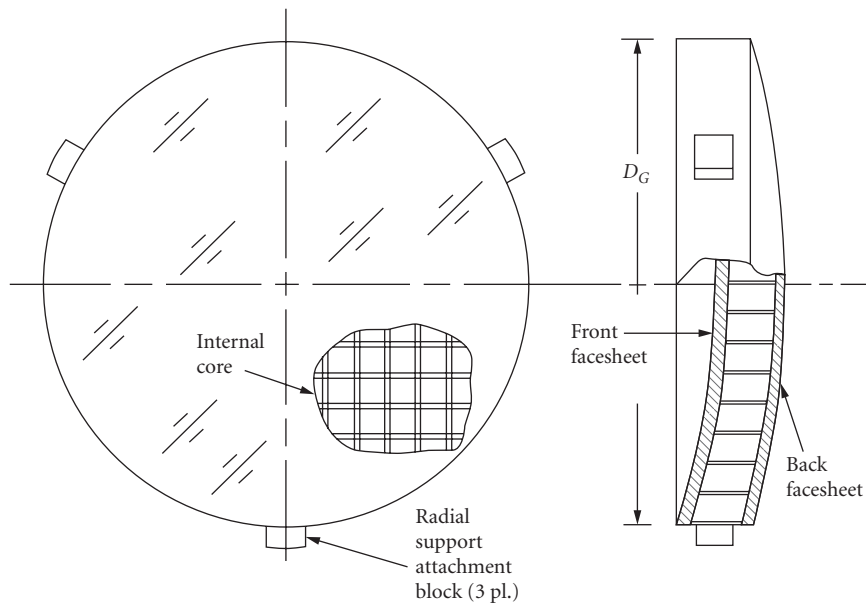


FIGURE 23 A completely fused (monolithic) built-up lightweight mirror substrate.

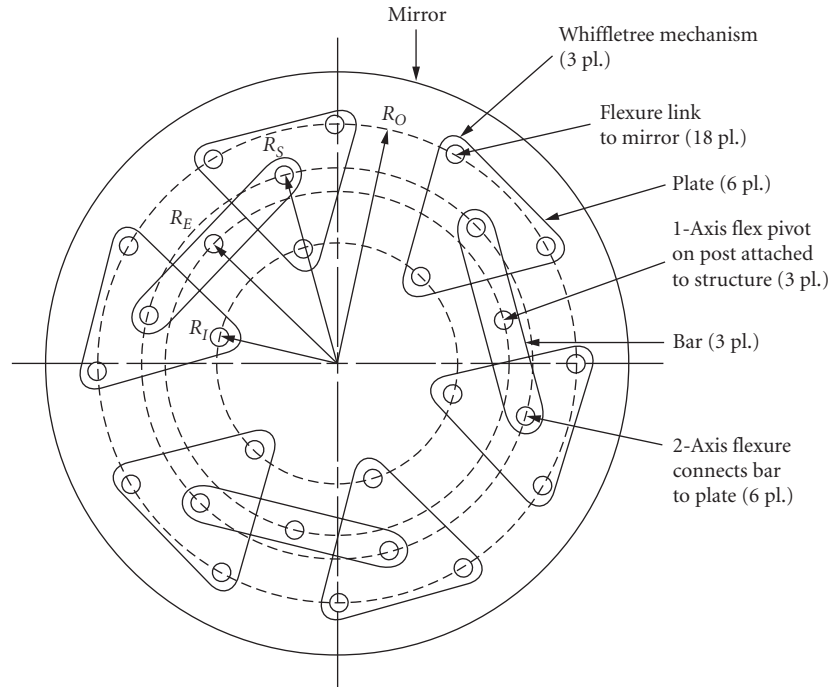


FIGURE 24 An 18-support Hindle-type mirror mount supporting the optic at multiple points on rings of radii R_O and R_I from three posts attached to structure at radius R_E . The whiffletree plates are centered on the ring of radius R_S .

Lever Mechanism Mountings

Because of the flexibility of lightweighted mirrors, axial support is frequently provided at many points on the back of the substrate. Hindle mounts²⁹ using multiple lever mechanisms (called “whiffletrees”) are commonly used. Figure 24 shows such a mount with 18 supports for the mirror. The number of supports needed is the minimum number that keeps the gravitational sag of the reflecting surface between support points smaller than the deflection tolerance when the mirror axis is vertical.⁴ To avoid friction effects, flexures (sometimes called “Flex-Pivots”) are typically used as single-axis bearings in these mounts. Dual-axis bearings are usually necked-down posts that serve as flexures.

A mirror on a Hindle mount also needs radial support if it is to be used in any orientation other than axis vertical. This might be in the form of three or more mechanical links with universal-joint flexures at each end that are oriented tangent to the rim of the mirror and connect the mirror rim to the surrounding structure. Provision for such a support is shown in the mirror of Fig. 23. Multiple-point whiffletree radial supports have also been used for this purpose.²¹

Mountings for Metallic Mirrors

Metallic mirrors are generally easier to support than nonmetallic ones because attachments can be made directly to the substrate through, for example, threaded holes for screws. The metallic substrate may also be stiffer than the glass counterpart. An example is the aluminum mirror shown

by section and back views in Fig. 25a.³⁰ Here, a single-point diamond-turning (SPDT) method is employed to machine the optical surface and the axial and radial mounting interface surfaces on the mirror's back. In this method, many extremely fine cuts are made with a precision diamond tool as the substrate rotates about a common mechanical and optical axis. The tool moves on a prescribed path under interferometric control. This results in very accurate surface shapes and surface interrelationships, as well as smooth surfaces and very low residual stresses in the parts. The mating surfaces on the mount are also created by diamond turning. The mirror is shown installed in its mount in Fig. 25b. Optical surface distortions due to mounting forces are minimal because the contacting surfaces on the optic and its mount are parallel when drawn together.³¹ When the mirror and its mount are made of the same material, the effects of temperature changes are minimized.

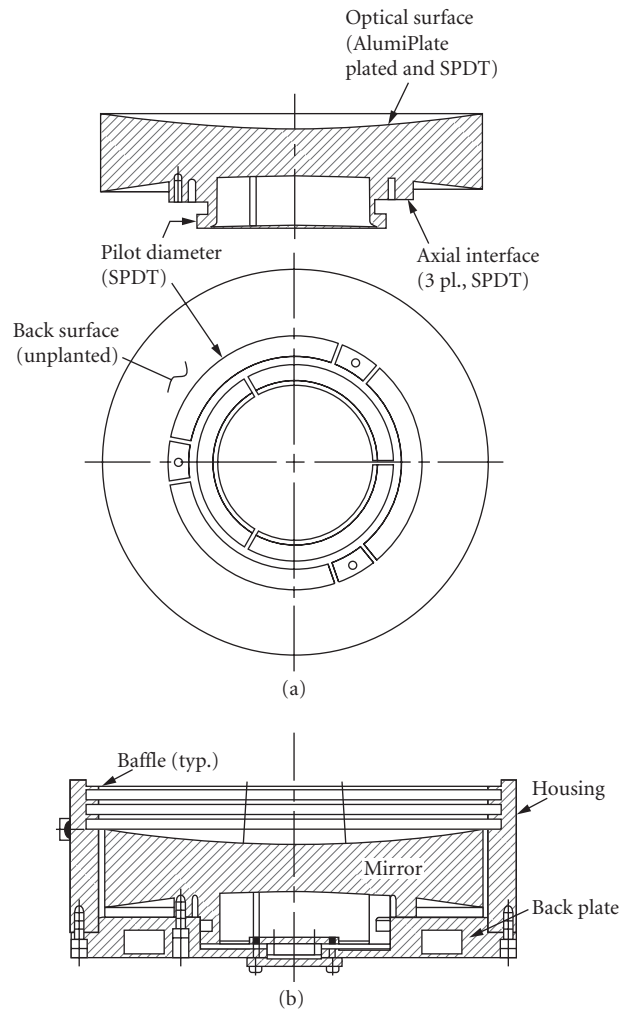


FIGURE 25 Optomechanical configuration (b) of an aluminum mirror (a) with optical and mounting surfaces machined by SPDT methods. The radial and axial interfaces are shown. (From Vukobratovich et al.³⁰)

6.8 CONTACT STRESSES IN OPTICS

The shape of the mechanical surface touching a lens, mirror, or prism surface is typically spherical, cylindrical, conical, or flat. Point contacts occur at spherical pads attached to the ends of springs while short line contacts occur if cylindrical pads are used on springs or if pins are used to locate the optic. Lenses, windows, and mirrors preloaded against mechanical constraints by a retaining ring or flange typically have circular line contacts with the metal around the edges of their apertures. The metal surface typically is conical for a convex glass surface and toroidal for a concave glass surface. The preloads applied through all of these interfaces cause elastic deformations of the glass and metal parts. Associated with these deformations are compressive and tensile stresses in those materials. Up-to-date analytical methods for estimation of these stresses have been presented in detail elsewhere.⁴ Space constraints preclude discussion of those methods here. Once the tensile stress to be expected in a given optomechanical design has been quantified, it can be compared to the aforementioned rule-of-thumb tolerance to predict success or failure of the optomechanical design. Should the stress appear to be too large, certain design changes that can be made to reduce it are suggested in the referenced publication.

6.9 TEMPERATURE EFFECTS ON MOUNTED OPTICS

General Considerations

Because the temperature environment of any optical instrument is seldom constant, we should anticipate changes in dimensions of all parts, in refractive indices, and in material parameters [such as coefficient of thermal expansion (CTE) and Young's modulus] to occur throughout the lifetime of the device. These changes may defocus the system, change aberration balance, or degrade alignment. Athermalized designs are created in a manner to reduce the magnitudes of these effects to tolerable levels.

Prevention of Axial Gaps

Differential expansions and contractions of all types of materials with temperature changes may change the axial and/or radial relationships, that is, alignment between optics and their mechanical reference surfaces. Optomechanical assemblies that are adequately preloaded at assembly will tend to maintain optic-to-mount contact, but this preload will change as the temperature changes. It may disappear completely at elevated temperatures. Then the optics may be free to move if externally disturbed, as by vibration or shock. These component shifts may become permanent if the optic is decentered or tilted when the temperature drops and the mount reapplies forces to the optics.

To reduce this effect, each optical assembly might be designed to compensate for axial dimensional changes so axial preload changes are reduced to insignificance.³² For example, the air-spaced triplet assembly of Fig. 26a is constructed of three optical glasses, an aluminum cell, and two aluminum spacers. The scale of the figure is as indicated. At maximum temperature, the physical separation of the interfacing points *A* and *B* in this particular assembly changes by 0.015 mm if computed for a path through the lenses and spacers, but changes by 0.030 mm if computed for a path through the cell. One or more axial air gaps totaling 0.015 mm would then exist somewhere within the assembly and the lenses might move or tilt within that space. If the design were to be modified by changing the metals in the cell and in one spacer, lengthening that spacer, and providing space for the larger spacer by adding a step bevel to the second lens—as indicated in Fig. 26b—the chosen materials and component dimensions would make the *A*-to-*B* separation remain equal for both paths for all temperature changes. Preload would then remain unchanged and misalignment would not occur.

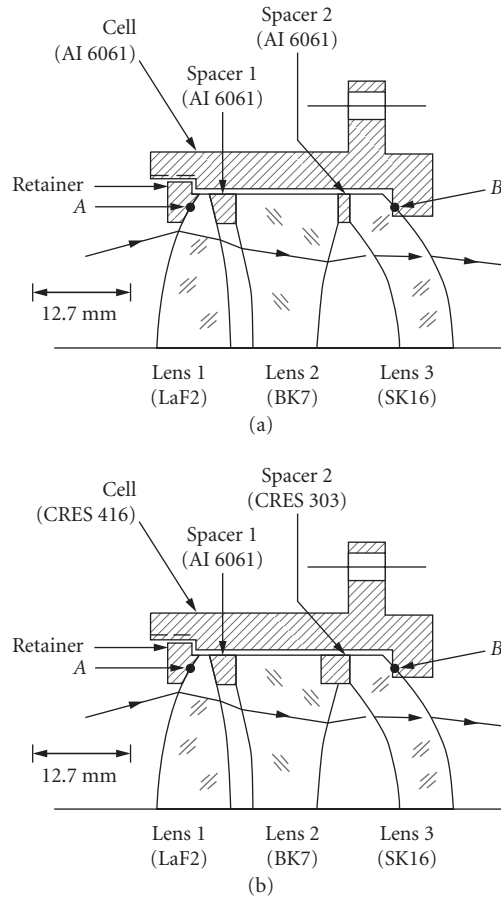


FIGURE 26 An air-spaced triplet lens assembly (a) in which an axial gap between glass and metal parts exists at high temperature, possibly allowing the lenses to become misaligned. Modified design (b) is athermalized to maintain registry of the optics in the mount. (Adapted from Yoder and Hatheway.³²)

Focus Athermalization Techniques

Single Material Designs Figure 27 shows a reflecting telescope made of a single material, in this case, aluminum.³³ All dimensions change, but the assembly remains in alignment and the optical performance is unchanged (other than a small change in image scale) as the temperature changes. This telescope is an example of the use of single-point diamond-machining methods as all optical and mounting surfaces are precisely made in the proper geometric relationships so alignment accuracy is built-in.

Passive Athermalization Figure 28a illustrates the use of materials with dissimilar CTEs and carefully chosen axial dimensions so the axial distance between optical components (in this case, the two mirrors) remains constant when the temperature changes.³⁴ This keeps the optical performance within required limits. Control of the mirror separation of this telescope is modeled schematically in Fig. 28b. Positive signs associated with lengths of low and high CTE materials indicate how the

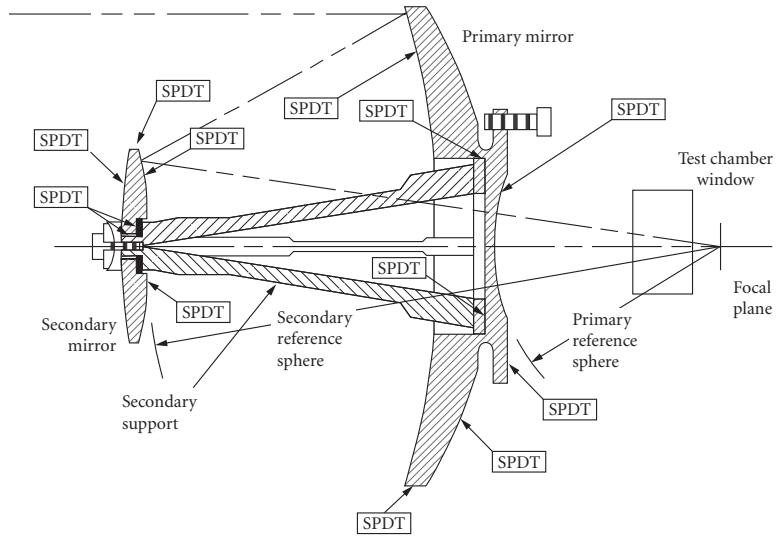


FIGURE 27 Schematic of an all-aluminum (athermalized) telescope objective with optical and mechanical interface surfaces finished by SPDT for ease of assembly without alignment. (From Erickson et al.³³)

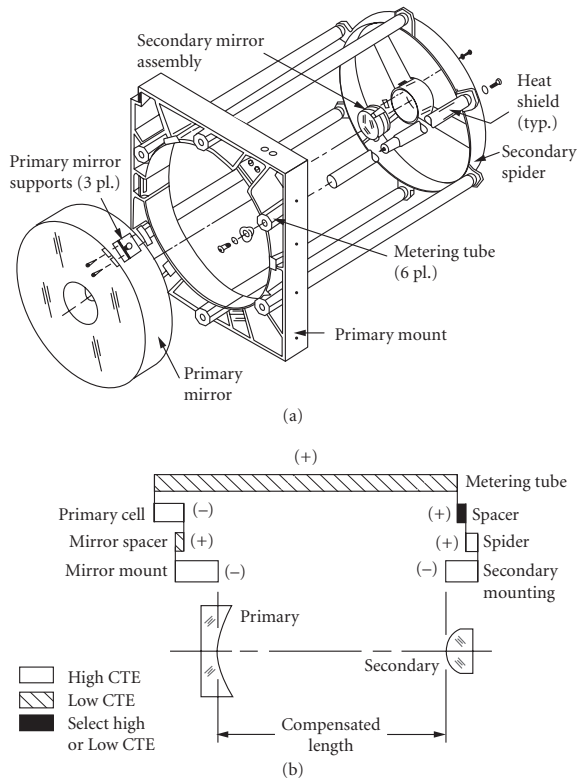


FIGURE 28 A passively athermalized telescope structure using Invar metering tubes to connect the primary and secondary mirror mounts. (a) Exploded view of the telescope. (b) Model of the compensation system. (From Zurmehly and Hookman³⁴)

mirror separation changes as the temperature rises. Proper choices of materials for their coefficients of thermal expansion and dimensions make the mirror separation remain constant as the temperature changes.

Active Athermalization When a source of power is available, components in an optical system can be physically moved to compensate for the effects of temperature changes. For example, Fig. 29a shows a concept for a zoom lens system in which locations of the moveable lenses are varied by motors as commanded by an internal microprocessor that monitors the temperature of the system.³⁵ As indicated in Fig. 29b, desired magnification inputs from the operator are automatically converted into the lens shifts required to focus properly on the object at the measured temperature.

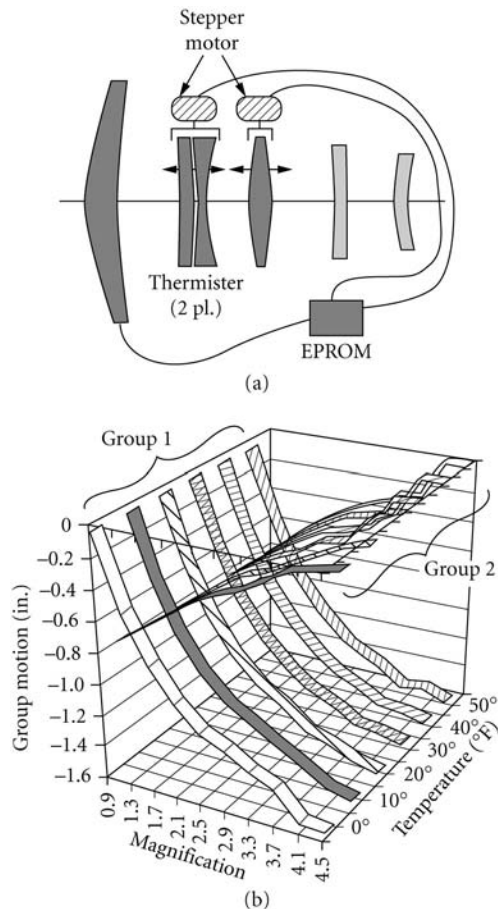


FIGURE 29 An actively athermalized zoom lens system that drives two lens groups to maintain focus at selectable magnification settings in spite of temperature changes. (From Fischer and Kampe.³⁵)

6.10 REFERENCES

1. V. L. Genberg, "Structural Analysis of Optics," Chapter 8 in *Handbook of Optomechanical Engineering*, A. Ahmad, (ed.), CRC Press, Boca Raton, 1997.
2. K. B. Doyle, V. L. Genberg, and G. J. Michels, *Integrated Optomechanical Analysis*, SPIE Press, Bellingham, 2002.
3. P. R. Yoder, Jr., *Opto-Mechanical Systems Design*, 3rd ed., CRC Press, Boca Raton, 2005.
4. P. R. Yoder, Jr., *Mounting Optics in Optical Instruments*, 2nd ed., SPIE Press, Bellingham, 2008.
5. R. J. Roark, *Formulas for Stress and Strain*, 3rd ed., McGraw-Hill, New York, 1954.
6. S. P. Timoshenko and J. N. Goodier, *Theory of Elasticity*, 3rd ed., McGraw-Hill, New York, 1970.
7. K. B. Doyle and M. Kahan, "Design Strength of Optical Glass," *Proc. SPIE*, 5176, 14, 2003.
8. R. Stoll, P. F. Forman, and J. Edelman, "The Effect of Different Grinding Procedures on the Strength of Scratched Fused Silica," *Proc. of the Symposium on the Strength of Glass and Ways to Improve It*, Union Scientifique du Verre, Florence, 1961.
9. K. B. Doyle, private communication, 2008.
10. J. J. Herbert, "Techniques for Deriving Optimal Bondlines for Athermal Bonded Mounts," *Proc. SPIE* **6288OJ-1**, 2006.
11. A. Ahmad and R. L. Huse, "Mounting for High Resolution Projection Lenses," *U.S. Patent 4,929,054*, 1990.
12. J. J. Bacich, "Precision Lens Mounting," *U.S. Patent 4,733,945*, 1988.
13. J. H. Bruning, F. A. DeWitt, and K. E. Hanford, "Decoupled Mount for Optical Element and Stacked Annuli Assembly," *U.S. Patent 5,428,482*, 1995.
14. E. T. Kvamme, D. Trevias, R. Simonson, and L. Sokolsky, "A Low Stress Cryogenic Mount for Space-Borne Lithium Fluoride Optics," *Proc. of SPIE* **58770T**, 2005.
15. E. T. Kvamme and Michael Jacoby, "A Second Generation Low Stress Cryogenic Mount for Space-Borne Lithium Fluoride Optics," *Proc. SPIE* **66920I**, 2007.
16. E. T. Kvamme and M. Jacoby, "Opto-Mechanical Testing Results for the Near Infra-red Camera on the James Webb Space Telescope," *Proc. SPIE* **7010**, 2008.
17. P. R. Yoder, Jr., "Lens Mounting Techniques," *Proc. SPIE* **389**: 2, 1983.
18. M. Bayar, "Lens Barrel Optomechanical Design Principles," *Opt. Eng.* **20**:181, 1981.
19. R. E. Fischer, "Case Study of Elastomeric Lens Mounts," *Proc. SPIE* **1533**: 27, 1991.
20. D. M. Williamson, "Compensator Selection in the Tolerancing of a Microlithography Lens," *Proc. SPIE* **1049**: 178, 1989.
21. D. Vukobratovich, *Introduction to Opto-Mechanical Design*, SPIE Short Course SC014, 2003.
22. W. Sunne, "Dome Attachment with Brazing for Increased Aperture and Strength," *Proc. SPIE* **5078**: 121, 2003.
23. P. R. Yoder, Jr., "Non-Image-Forming Optical Components," *Proc. SPIE* **531**: 206, 1985.
24. P. Mammini, B. Holmes, A. Nordt, and D. Stubbs, "Sensitivity Evaluation of Mounting Optics Using Elastomer and Bipod Flexures," *Proc. SPIE* **5176**: 26, 2003.
25. P. R. Yoder, Jr., "Mounting-Induced Contact Stresses in Prisms," *Proc. SPIE* **3429**: 71, 1998.
26. P. R. Yoder, Jr., "Improved Semikinematic Mounting for Prisms," *Proc. SPIE* **4771**:173, 2002.
27. P. R. Yoder, Jr., "Design Guidelines for Bonding Prisms to Mounts," *Proc. SPIE* **1013**: 112, 1988.
28. L. H. J. F. Beckmann, private communication, 1990.
29. J. H. Hindle, "Mechanical Floatation of Mirrors," *Amateur Telescope Making, Advanced*, A.G., Ingalls, ed., Scientific American, New York, 1945: 229. (Reprinted 1996 as Chapter B.8 in *Amateur Telescope Making, 2*, William-Bell, Inc., Richmond.)
30. D. Vukobratovich, A. Gerzoff, and M. K. Cho, "Therm-Optic Analysis of Bi-Metallic Mirrors," *Proc. SPIE* **3132**: 12, 1997.
31. R. L. Rhorer and C. J. Evans, "Fabrication of Optics by Diamond Turning," Chapter 41 in *Optical Society of America Handbook of Optics*, 2nd ed., Vol. I, Bass, M., Van Stryland, E. W., and Wolfe, W. L., eds., McGraw-Hill, New York, 1995.

32. P. R. Yoder, Jr. and A. E. Hatheway, "Further Considerations of Axial Preload Variations with Temperature and the Resultant Effects on Contact Stresses in Simple Lens Mountings," *Proc. SPIE* **587705**, 2005.
33. D. J. Erickson, R. A. Johnston, and A. B. Hull, "Optimization of the Opto-Mechanical Interface Employing Diamond Machining in a Concurrent Engineering Environment," *Proc. SPIE* **CR43**: 329, 1992.
34. G. E. Zurmehly and R. Hookman, "Thermal/Optical Test Setup for the Geostationary Operational Environmental Satellite Telescope," *Proc. SPIE* **1167**: 360, 1989.
35. R. E. Fischer and T. U. Kampe, "Actively Controlled 5:1 Afocal Zoom Attachment for Common Module FLIR," *Proc. SPIE* **1690**: 137, 1992.

CONTROL OF STRAY LIGHT

Robert P. Breault

*Breault Research Organization
Tucson, Arizona*

7.1 GLOSSARY

A	area
BRDF	bidirectional reflectance distribution function
GCF	geometric configuration factor
L	radiance
R	distance
θ, ϕ	angles
Φ	power
Ω	solid angle

7.2 INTRODUCTION

The analysis of stray light suppression is the study of all unwanted sources that reduce contrast or image quality. The control of stray light encompasses several very specialized fields of both experimental and theoretical research. Its basic input must consider (1) the optical design of the system; (2) the mechanical design, size, and shape of the objects in the system; (3) the thermal emittance characteristics for some systems; and (4) the scattering and reflectance characteristics of each surface for all input and output angles. It may also include spectral characteristics, spatial distribution, and polarization. Each of these areas may be concentrated on individually, but ultimately the analysis culminates in the merging of the various inputs.

Developments in detector technology, optical design software, diffraction-limited optical designs, fabrication techniques, and metrology testing have created a demand for sensors with lower levels of stray radiation. Ways to control stray light to meet these demands must be considered during the “preliminary” conceptual design. Decisions made at this time are, more often than not, irrevocable. This is because parallel studies based upon the initially accepted starting design are often very expensive. The task of minimizing the stray radiation that reaches the detector after the system has been

designed by “adding on” a suppression system is very difficult. Therefore, every effort should be made to start off with a sound stray light design. To ensure a sound design, some stray light analysis should be incorporated in the earliest stages of a preliminary design study.

This chapter presents some basic concepts, tools, and methods that you, the optical or mechanical designer, can consider when creating a sensor system. You do not need to be very experienced in stray light suppression to design basic features into the system, or to consider alternative designs that may significantly enhance the sensor’s performance. The concepts are applicable to all sizes of optical instrumentation and to virtually all wavelengths. In some cases, you can use the concepts to rescue a design when experimental test results indicate a major design flaw.

7.3 CONCEPTS

This section outlines some concepts that you can use to reduce stray radiation in any optical system. The section also contains some experimental and computer-calculated data as examples that should give you some idea of the magnitude of the enhancement that is possible.

The power on a collector depends on the following:

1. The power from the stray light source.
2. The surface scatter characteristics of the source; these characteristics are defined by the bidirectional scatter distribution function (BSDF).
3. The geometrical relationship between the source and collector. This relationship is called the geometrical configuration factor (discussed later in this section).

To reduce the power on the detector, we can try to reduce the contributions from these elements:

$$\Phi_{\text{collector power}} = \Phi_{\text{source power}} \times \text{BRDF}_{\text{source}} \times \text{GCF}_{\text{source collector}} \times \pi \quad (1)$$

Ways to reduce each of these factors are discussed below. The creative use of aperture stops, Lyot stops, and field stops is an important part of any attempt to reduce the GCF term of the power transfer equation.

For the discussion that follows, examples from a two-mirror Cassegrain design, with the aperture stop at the primary, will be used to stimulate thoughts about stray light reduction possibilities for other sensors.¹ The system is shown in Fig. 1.

Critical Objects

The most fundamental concept is to start the stray light analysis from the detector plane of the proposed designs. The most critical surfaces in a system are those that can be seen from the detector position or focal surface. These structures are the only ones that contribute power to the detector. For this reason, direct your initial attention toward minimizing their power contributions by removing them from the field of view of the detector.

The basic idea is to visualize what would be seen if you were to look out of the system from the image plane. Unlike most users of optical instruments, the stray light designer’s primary concern is seldom the object field, but rather all the interior surfaces that scatter light. It is necessary to look beyond the radii of the imaging apertures to find the sources of unwanted energy. Removing these sources from the field of the detector is a real possibility, and will result in a significant improvement in the system.

Real-Space Critical Objects

I will start out by identifying a particular critical real-space object that can be seen by the detector in our example; it is the inside of the secondary baffle. The direct view discussed here is different than the image of the same baffle reflected by the secondary which is discussed in the next section.

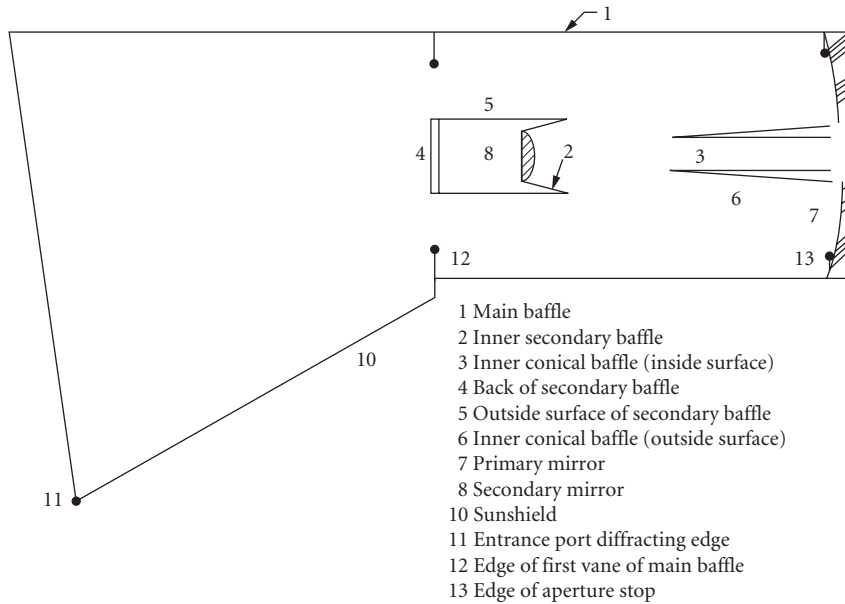


FIGURE 1 Typical Cassegrain design with the aperture stop at the primary. (Ref. 2, p. 52.)

Many Cassegrain secondary baffles have been designed to be cone-shaped (Fig. 2), usually approximating the converging cone of light from the primary. From the detector, portions of this secondary cone are seen directly as a critical surface. Since most of the unwanted energy is incident on this baffle from nearly the same direction as this surface is seen from the detector, the addition of vane structures would be of little help, assuming an optimum coating is used on the simple baffle. If the cone is made *more* cylindrical, the amount of critical cone area is reduced, and the angle at which the surface is seen gives a smaller projected area (Fig. 3).

Avoid making the baffle cylindrical because the outside of it would be seen. Since the detector is of a finite size there is a fan of rays off the primary representing the field of view of the telescope.

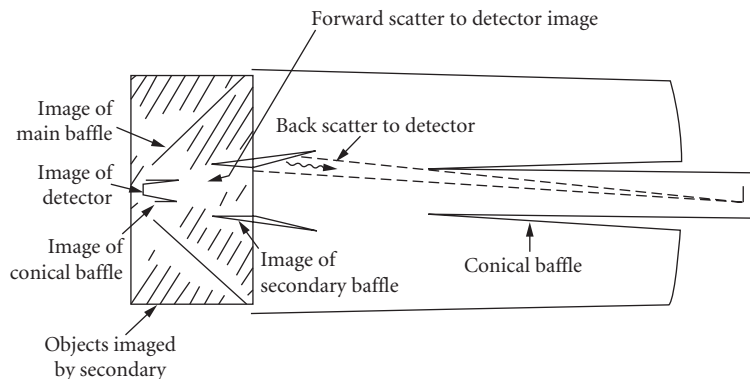


FIGURE 2 Direct and reflected scatter from the cone-shaped secondary baffle. (Ref. 1, p. 4.)

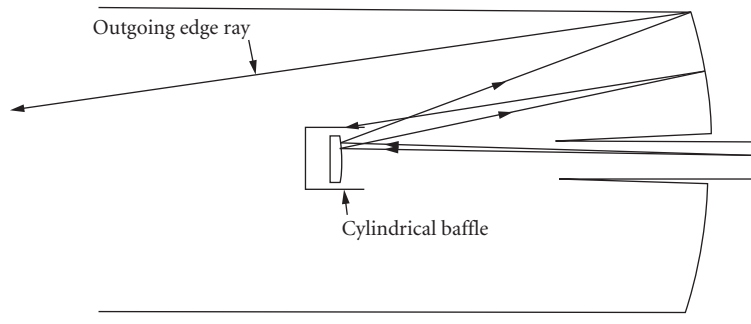


FIGURE 3 Reduced scatter from an almost cylindrical-shaped secondary baffle.

Although collimated for any point on the detector, any point not on axis would have its ray bundle at some angle to the optical axis, hence a cylindrical secondary baffle would be seen from off-axis positions on the detector.

Imaged Critical Objects

Imaged objects are often critical objects. They too can be seen from the detector. Determining which of the imaged objects are critical requires a bit of imagination and usually some calculations; stray light software can help you make the calculations. The Y-Y bar diagram can help you to conceptualize and determine the relative image distances and sizes with a minimum number of calculations.³ The same could be done with other first-order imaging techniques (see Chap. 1 in this volume). Using the Cassegrain example again (Fig. 2), you can see that reflected off the secondary mirror are the images of the detector and the inside of the inner conical baffle (object 3 in Fig. 1). In some designs the outside of the conical baffle will be seen in reflection. These are *imaged critical objects*. If you wish, you can eliminate some of these images with a central obscuration on the secondary, or for the conical baffle, with a spherical mirror concentric about the image plane. The direct view to the inner conical baffle will remain, but the path from the image of it is removed.

The cone-shaped secondary baffle is also seen in reflection (Fig. 2). For the incident angles of radiation on this surface, the near specular (forward-scattering) characteristics will often be one of the most important stray light paths because the image of the detector is in that direction. This is an extension of *starting from the detector*. There is an image of the detector at the prime focus of the primary mirror. Often, as in this case, one location may be easier than another for you to determine what could be seen. By making the baffle more cylindrical, part of the *image* of the baffle is removed from the detector's view; as a result, the power that can scatter to the detector is reduced. Furthermore, it is sometimes possible to baffle most of this power from the field of view with one or two vanes (Fig. 4).

Continue the process of removing critical surfaces until all the critical surfaces have been considered for all points in the image plane. The power contributions from these surfaces will either go to zero, or at least be lessened after you reduce the area of the sections seen.

There is still more that can be done, since only the GCF term in the power transfer equation to the detector has been reduced. It is also possible to minimize the power onto the critical sections, which will become the source of power, Φ , at the next level of scatter. This approach can be very similar to the approach used for minimizing the power scattered to the image plane. The viewing is now forward from the critical surfaces instead of the image plane. By minimizing the BRDF and GCF factors of the surface scattering to the critical sections, the power incident on the critical surfaces will be reduced. Hence, the power to the detector is also reduced.

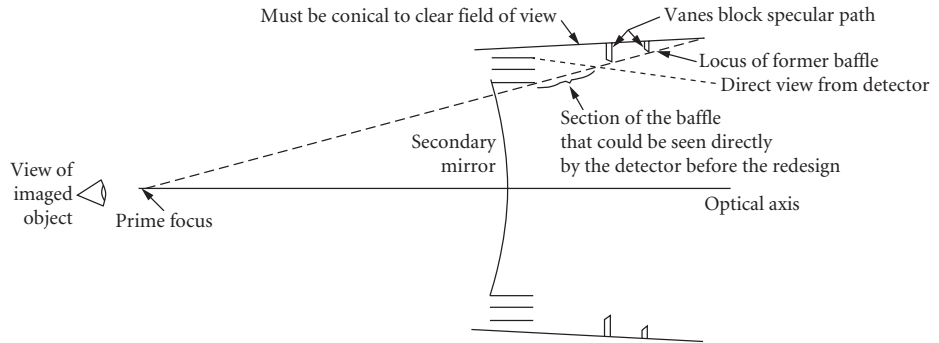


FIGURE 4 A cylindrical secondary baffle can be seen from off-axis positions on the detector.

Illuminated Objects

Minimizing the GCFs and BRDFs for the specific input and output angle is sometimes easier if you look into the system from the position of the stray light source in object space. By doing this, you can identify the surfaces that directly receive the unwanted energy. I will call these the *illuminated objects*. If any of these illuminated surfaces contain sections that the detector can see, then you should direct your initial efforts toward eliminating these paths. These paths will usually dominate all other stray light paths because there is only a single scatter before the stray light reaches the detector. An example of such a path that is often encountered is from the source onto the inner conical baffle of multimirror systems (Fig. 1). Some of the ways that the direct radiation can be eliminated is by extending the main baffle tube, increasing the obscuration ratio by increasing the diameter of the secondary baffle (Fig. 5), or by narrowing the field of view, which will allow you to extend the secondary baffle and the inner conical baffle toward each other.

The effect of eliminating this path is shown in a composite Point Source Transmittance (PST) plot in Fig. 6.⁴ The PST plot is defined as the reference plane (detector plane in most cases) irradiance divided by the input irradiance along the line of sight. (See the section called “Point Source Transmittance Definitions” for a more detailed definition of PST.) For the case shown, the unwanted irradiance on the detector is reduced by over an order of magnitude.

Aperture Placement

I will now focus on the optical design aspects of a stray light suppression system, and give a qualitative discussion of some general aspects that you might consider. All optical systems will have at least

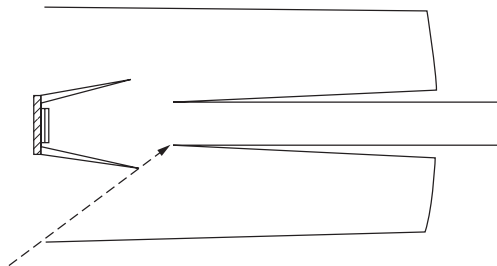


FIGURE 5 Increased obscuration ratio blocks direct path to inner conical baffle. (Ref. 1, p. 6.)

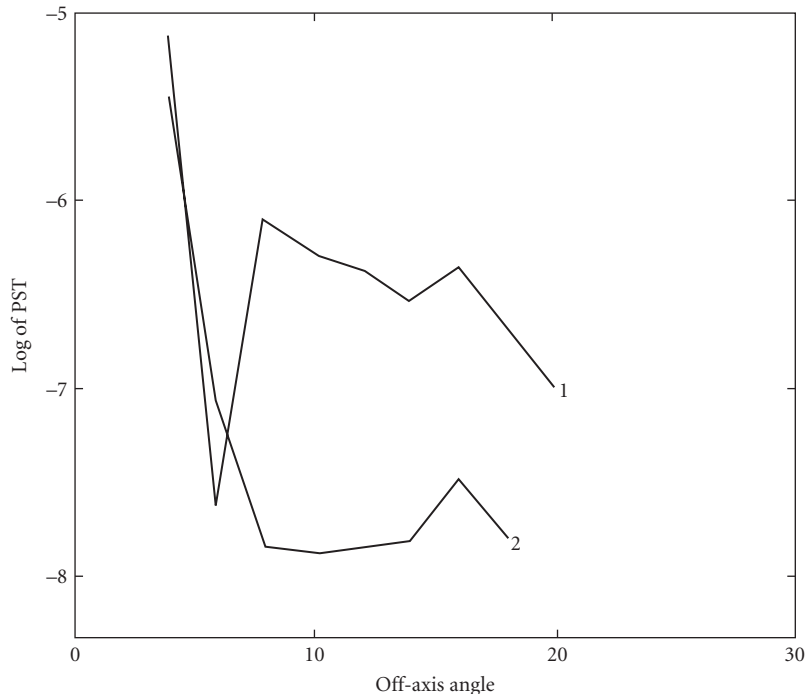


FIGURE 6 Point source transmittance with obscurations of (1) 0.333 and (2) 0.4. The 0.4 obscuration removed the direct path from the source to the inner conical baffle. (Ref. 1, p. 6.)

one aperture, called an *aperture stop*, that limits the size of the bundle of the incoming signal rays. Some systems will have field stops and/or Lyot stops. Each type of stop has a clearly defined role in stray radiation suppression, which is discussed in the following sections.

In many cases stop placement will have a much more noticeable effect on system performance than any vane structure, coating, or baffle redesign. Probably the only factor with more effect on the PST curve is the off-axis position of the source. Therefore, the benefits of any of the stops cannot be overemphasized.

Aperture Stops The aperture stop is the aperture that limits the size of the cone of radiation that will reach a point on the image plane. Sometimes shifting this stop allows the optical designer to better balance the aberrations. In a stray radiation suppression design, it plays a similar important role. All objects in the spaces preceding the stop in the optical path will not be seen unless they are imaging elements, central obscurations, or objects that vignette the field of view. Only a limited number of critical objects is possible before the aperture stop. In the intervening spaces from the stop to the image plane it is likely that many of the baffle surfaces will be seen. Figure 7 represents a two-mirror design, and Fig. 8 represents a three-element refracting system; both have the stop at the first element. In both cases the second element is oversized to accommodate the field of view from a point in the field stop; the amount depends on the full field of view of the design. Because the elements are oversized, the main baffle following the first element will be seen. This baffle will be a critical object, a direct path of unwanted energy. The “overviewing” is characteristic of all of the optical elements past the aperture stop.

If you move the stop along the optical path toward the detector plane, its performance as a stray radiation baffle will improve. If you shift the stop to the second element, the intermediate baffle will

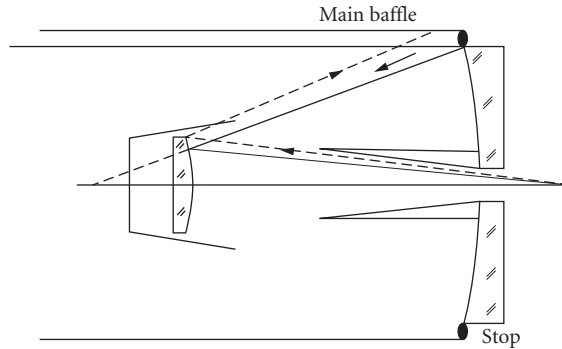


FIGURE 7 The oversized secondary allows the main baffle to be seen in reflection. (*Ref. 1, p. 8.*)

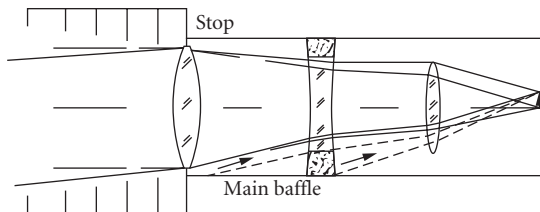


FIGURE 8 The main baffle is seen through the oversized secondary and tertiary. (*Ref. 1, p. 8.*)

not be seen. It is removed from the field of view of the detector, since the stop now eliminates direct paths from baffles in all spaces that precede it. Figure 9 shows the improvement in the PST curve for a two-mirror system. By moving the stop you have reduced the PST by a factor of 10. This is a desirable feature to consider for stray radiation reduction.

Direct paths from central obscurations can be blocked by a central disk located at some location deeper into the system; however, because of the parallax involved between the central obscuration disk and this central disk, the central disk obscuration will usually be a larger obstruction to imaging rays. In a reimaging design it is often possible to locate a central disk conjugate to the actual central obscuration.

Field Stops An aperture can be placed at intermediate *images* in a system to limit the field of view. Such an aperture will usually prevent any stray light from outside of the field of view from being directly imaged into the system beyond this field stop aperture. In a sense, its operation is just opposite that of an aperture stop. Baffle surfaces following a field stop cannot be seen from outside the field of view in the object plane, unless they are central obscurations. Note that with just a field stop, succeeding optical elements may allow out-of-field critical sections to be seen *through* the field stop, from within the field in the image plane (Fig. 11). Aperture stops are necessary to block such paths. Figures 10 and 11 show two such cases. Although for some designs the field stop is not 100 percent effective because of optical aberrations, its small size limits most of the unwanted stray light. Field stops therefore do not remove critical sections, but rather limit the propagation of power to illuminated objects. In reflecting systems, take care that the object side of the field stops does not become a critical area, which can be seen directly or in reflection from the image plane because unwanted energy is being focused onto them.⁵

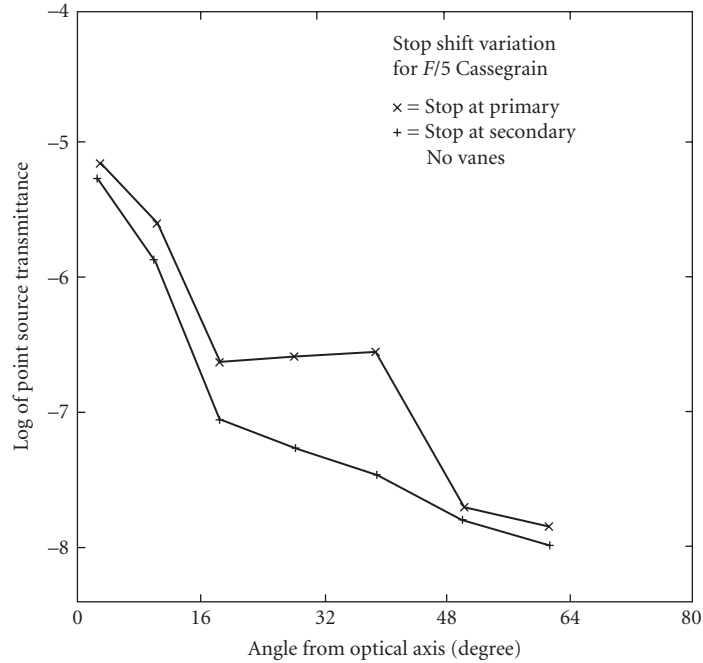


FIGURE 9 PST improvement with stop shift for the two-mirror system. (Ref. 1, p. 8.)

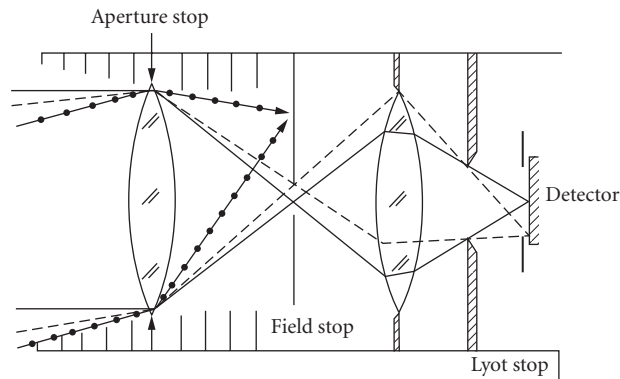


FIGURE 10 The addition of a Lyot stop prevents preceding baffles from being seen from the image.

Lyot Stops A limiting aperture placed at the location of the image of an aperture stop, sometimes called a glare stop or Lyot stop, has the same property as described for aperture stops. It should be slightly smaller than the image of the aperture stop. It limits the critical sections which are out of the field of view to those objects in succeeding spaces only. Since Lyot stops are by definition further along the optical path to the detector, the number of critical surfaces seen by the detector will be reduced. Usually, these stops are incorporated into the design to block the diffracted energy from an aperture stop and field stop pair, so that only secondary or tertiary diffracted energy reaches the image. Nevertheless, both diffracted and scattered energy are removed from the direct

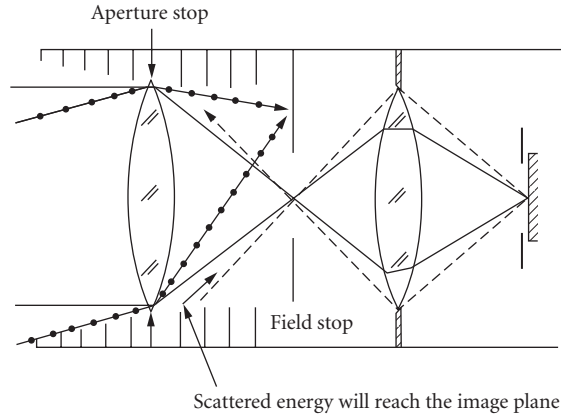


FIGURE 11 Out-of-field energy in object plane will not be imaged beyond field stop. Out-of-field energy elsewhere may be seen. (Ref. 1, p. 9.)

view of the image, re imaging the largest optical element as the stop takes full advantage of both the light-gathering power of the optics and the stray radiation suppression features provided by the stop. Figure 10 shows a system with a Lyot stop.

On space-based telescopes the image plane is often *shared* by one or more instruments. Each instrument reimages the telescope's image through some optical train, and eventually onto the detector. In that optical train there could be a logical place to use a Lyot stop to improve the stray light performance of the viewing instrument well beyond that of the telescope.

It is the combination of these different stops or apertures that helps minimize the propagation of unwanted energy by limiting the number of critical objects seen by the detector, and the objects illuminated by the stray light source.

When all direct paths have been eliminated, the next step is to determine the relationship between the sections that received power (illuminated objects) and the critical surfaces. This relationship takes the form of scattering *paths*; that is, stray light can scatter from the illuminated objects to the critical objects. To start reducing the stray light contributions from these paths, you can start at the critical surfaces as described above. But now you have more knowledge about where the direct incident power is being distributed throughout the system, since you can also look into the system from the source side to find the surfaces receiving direct power. With this information you can identify the possible paths between the illuminated and critical surfaces.

Design considerations for extended baffle shields (Figs. 12 and 13) provide a good example of starting from the source side to identify possible paths. In the examples, object 2 (an optical surface) is the largest contributor of scattered radiation and is the best superpolished mirror available

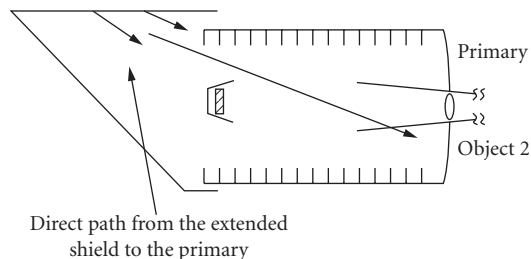


FIGURE 12 There is a direct path from the baffle to the primary mirror. (Ref. 1, p. 7.)

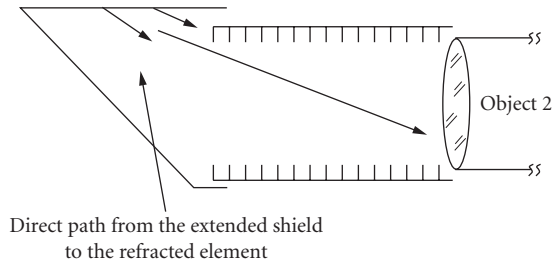


FIGURE 13 There is a direct path from the baffle to the refracting element. (Ref. 1, p. 7.)

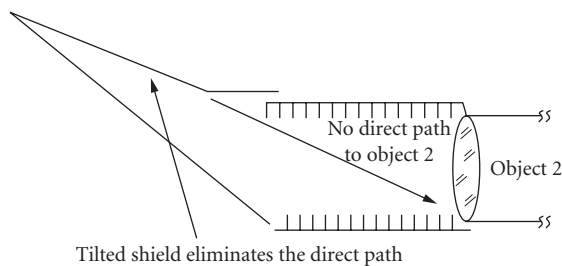


FIGURE 14 Direct paths are removed by properly tilting the shield. (Ref. 1, p. 7.)

(Fig. 12), or if it is a lens as in Fig. 13, it has the lowest possible scattering characteristics. It cannot be removed from the view of the image plane. If the initial power incident on object 2 is only from the extended shield, then by tilting the shield (Fig. 14), the power on the shield must first scatter to the main tube and then to the optical element. The combination is then referred to as a two-stage baffle. If vanes are added to the main baffle, the scattered radiation incident on the optical element will be reduced by many orders of magnitude when all the scattering solid angles and the number of absorbing surfaces are considered. Note that without the tilt to the shield, vanes on the main baffle are worthless because there is a direct path to the objective.

Figure 15 is an abstract representation of the process of reducing stray light in a sensor system. Start at the detector, then work from its conjugate image locations. Starting from the detector simplifies the analysis and directs your attention to the most productive solutions, because you can identify all the possible sources of stray light to the detector. You can then work at decreasing their number by slightly redesigning the baffles and stops. Next, identify which objects are illuminated. Discover how energy may propagate between them and you have identified the paths of stray light propagation. From then on the process of moving objects or blocking paths is quite simple, although the quantitative calculations might get difficult and may require some analysis software.

Baffles and Vanes

A few definitions are required to define baffles and vanes. Other authors have used their own different definitions. In this section the term *baffle* is used to describe conical structures (including cylindrical) that can also be described as tubelike structures. Their function is to shade, or occult, stray light from the source to one or more system components. The main baffle shields the primary

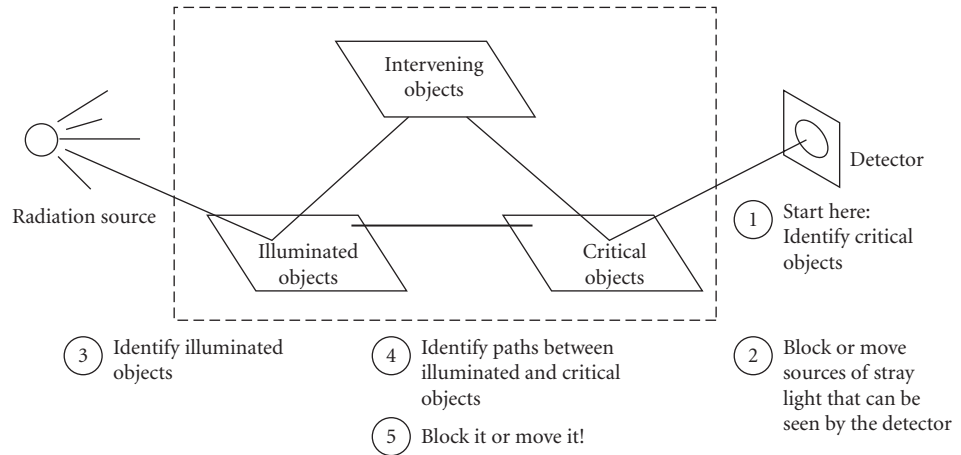


FIGURE 15 The first step in a stray light analysis begins from the detector plane, not from the source.

mirror from direct radiation at the larger off-axis angles. *Vanes* are structures put *on* baffles to affect the scatter characteristics of the surface. Other authors have used the term “baffles,” or “glare stops,” to describe these vanes.

Baffles In a well-designed system vanes play an important role only at large off-axis angles. For example, when one-tenth of the stray light falls on the primary of the Cassegrain design, then the main baffle receives the remaining 90 percent of the stray light. When the main baffle has properly designed vanes on it, light that falls on the baffles is attenuated by five orders of magnitude before it reaches the primary mirror. The resultant power on the primary mirror is then about 9.0×10^{-5} compared to the direct 10 percent that fell on the mirror. This results in less than 0.1 percent of the total on the primary. In addition, most of the subsequent scatter off the primary will be at much higher scatter angles. This will cause the scattered energy to have much lower BRDFs off the primary mirror when scattered in the direction of the detector, further reducing the scatter contribution from the baffle.

Only when no power illuminates the objective will the baffles play a significant role in the propagation paths of the stray light. Usually the system’s performance merit function is then very good. Only if the stray light source has a tremendous amount of energy, like the sun, does the stray light become measurable.

Vanes The depth, separation, angle, and bevel of vanes are variables that need to be evaluated for every design. In the following paragraphs stray light analysis results are presented for both a centrally obscured system (Cassegrain) and an unobscured eccentric pupil design (Z-system).⁶ Profiles of these systems are given in Figs. 2 and 16. Of the two designs, only the eccentric pupil design has a reimager that would allow for the placement of an intermediate field stop and an accessible Lyot stop, as discussed above.

The APART stray light analysis program was used to analyze the two designs. The APART program was a substantial software package that performed deterministic calculations of stray light propagation in optical systems.^{2,7,8}

As an example of vane design considerations, the design of vanes on a main baffle tube will be explained. With minor differences, the design steps are the same for the Cassegrain and the eccentric pupil designs. In a reimaging system, vane structure deep in the system is usually not necessary, but there are exceptions. Figure 17 shows a collecting optical element that has some small field of view (FOV). The optical element could represent a primary mirror or a refractive element.

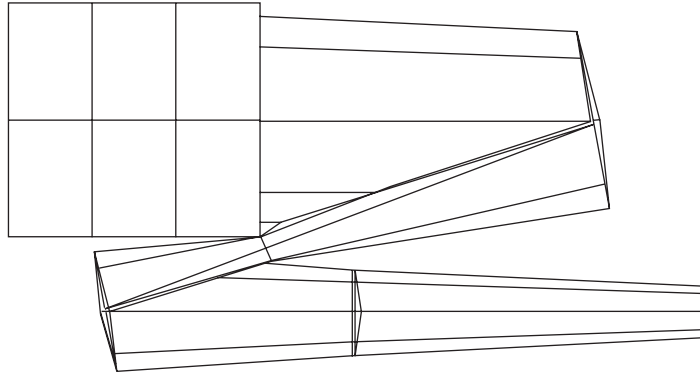


FIGURE 16 Confocal mirror system, eccentric pupil, no obscuration, low-scatter system. (Ref. 6, p. 91.)

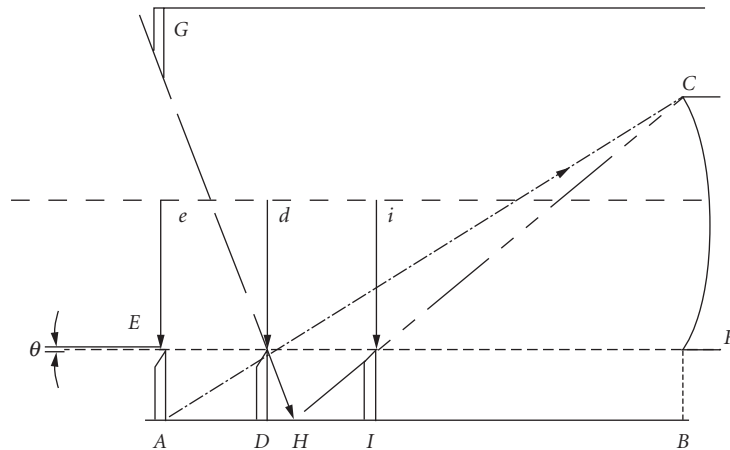


FIGURE 17 Vane placement design, lowercase letters are radii (measured from the optical axis), uppercase are z locations. (Ref. 6, p. 94.)

The placement of a straight, diffusely coated cylindrical tube would block the direct radiation from an external source, such as the sun, from reaching the optical element for a certain range of off-axis angles. If it were at a large off-axis angle, the forward scatter off the inside of the tube would be so high that it would normally not be acceptable. The solution is to add vanes to block this path.

Figures 18 and 19 depict the two cases that could represent the scatter from a baffle. In one case there are no vanes; in the other case there are vanes. This example shows how a propagation path is blocked by vanes. Vanes are useful, but a better approach is to make the solid angle (Ω_c) from the baffle (not the vanes) to the collectors of the scattered light go to zero, so that there is no path from the baffle and vane structure to the collecting object. By moving the baffle out of the field of view of the collector, the baffle's contribution goes to zero. There is no edge scatter, and no edge diffraction effects. That topic is in the realm of baffle design, which has already been touched on, and is well covered in the literature.^{9,10}

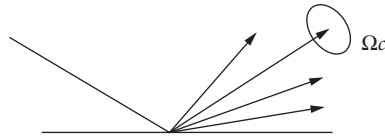


FIGURE 18 High forward scatter path. (Ref. 6, p. 92.)

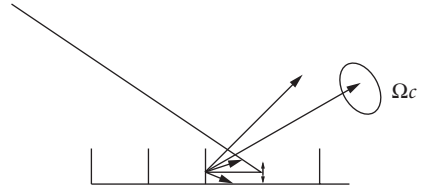


FIGURE 19 Forward scatter path highly attenuated by the vane structure. (Ref. 6, p. 92.)

Designing the Vane Spacing and Depth¹¹

A first vane is most often placed at the entrance of the baffle and an external ray is brought in from object space at a maximum off-axis position. If there is no forebaffle the angle is 90° off axis. The depth of the vane cavity is normally dictated by space and weight requirements. Too little depth will dictate the requirement for many vanes. Then vane edge scatter eventually becomes the major source of scatter instead of multiple vane scatter.

The initial ray will strike the side wall at the base of the first vane (point *A* in Fig. 17). From this point, a design line is drawn/calculated (*AC*) from the wall to the edge of the optical element on the opposite side. This line (*AC*) intersects the edge ray (*EF*), at *z* position *D*. At this point a vane could be placed. Mathematically, this assures that any point below *C*, including those on the optical element, would not see any directly illuminated side wall. However, practicality dictates that some offset of point *D* to a point *D'* (not shown) is required to allow for tolerance errors in fabrication of the vane, thermal effects, assembly errors, and for stray light edge scatter and diffraction effects. The tolerance allowance is company-, material-, and design-dependent. Acceptable numbers are often about 0.125 mm for fabrication and assembly tolerances. For the rest of this analysis, assume that this is accounted for.

Continue the design process by constructing another line from the edge of the entrance aperture to the tip of the second vane to the wall (line *GH*). Draw a new *HC* line to the area near the objective and determine the placement of the third vane (at *I*); once tolerances are considered, iterate the process to reach a final design. In some cases you may have to consider more than just the scatter path to the objective. In the Cassegrain design you may also have to consider the inner conical baffle opening. It is beyond the present scope to go into further detail.¹²

Bevel Placement on Vanes In this short discussion on baffle-vane design and placement, I did not mention the placement of a slanted surface, or bevel, sometimes placed on a vane edge as shown in Fig. 20. Which side should the bevel go on? The answer is usually dictated by first-order scatter principles.^{13,14} Near the front of the tube, direct radiation from a source at large off-axis angles will strike this bevel. If it is placed on the right side (Fig. 21), then the illuminated bevel will scatter its radiation all the way down into the tube to some optical surface. If placed on the left side, as depicted in Fig. 20, then it will go only 16° deeper into the system to the opposing vanes, a much better solution. For vanes deeper into the system, the bevel is placed on the right side. The point at which this is done is determined by the angle of the bevel and the diameter of the baffle tube. At some point, external radiation will not be able to directly strike the beveled edge if it is on the right-hand side of the vane. Only the nonbeveled, straight side will be illuminated. Therefore, the vane can rescatter only in the left side of the hemisphere, which is in the direction of going out of the system. If the bevel is placed on the left side, it can scatter 16° (in the example) deeper into the sensor; this is usually a needless design shortcoming that could be a significant error.

Vane Angle Considerations Another variation on the design feature of vanes that has sometimes been incorporated onto baffles in an optical system is angled vane structure. These vanes are non-planar objects. This makes them quite tedious to cut out of sheet metal, fabricate, and install.

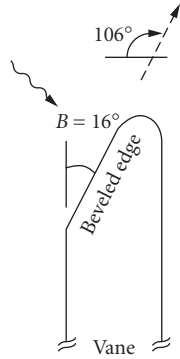


FIGURE 20 Placement of the bevel on the left side of the vane structure. (Ref. 6, p. 96.)

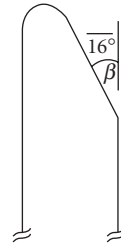


FIGURE 21 Placement of the bevel on the right side of the vane structure. (Ref. 6, p. 96.)

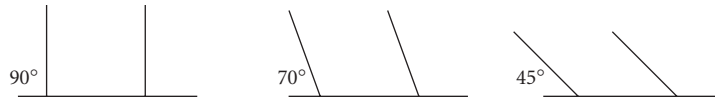


FIGURE 22 Vane structure angled at 90, 70, and 45°, respectively. (Ref. 6, p. 97.)

The next few paragraphs will present computer analysis results from two designs to show the effect of vanes on the propagated stray light. The vane angles used were 90, 70, and 45°, as depicted in Fig. 22.

The comparative stray light results for the Cassegrain system (Fig. 1) with a Martin Black coating on the vanes are shown in Fig. 23; in this system the vanes are on the main baffle, but not on the sunshade. There is no difference in the performance as the vane angle is varied from 45° to 90° (all three curves lie one on top of the other).

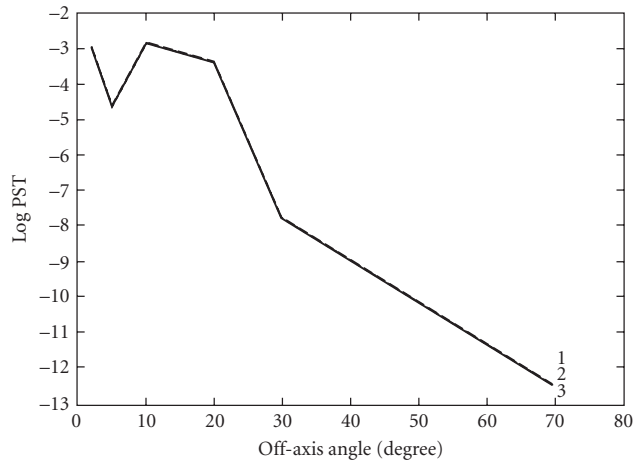


FIGURE 23 Cassegrain with Martin Black. Vane angles 1 = 90°, 2 = 70°, 3 = 45°. Log PST = detector irradiance divided by input irradiance. (Ref. 6, p. 97.)

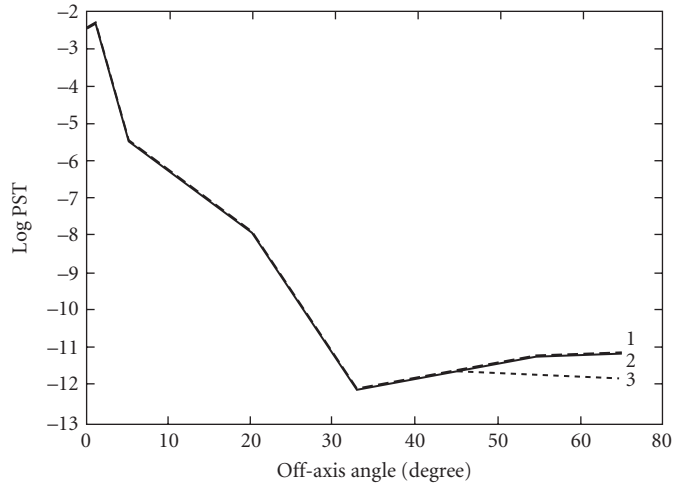


FIGURE 24 Z-system with Martin Black. Vane angles 1 = 45°, 2 = 70°, 3 = 90°. Vanes are on the sunshield. (Ref. 6, p. 114.)

The comparative results for the Z-system (Fig. 16) with vanes on the sunshield are shown in Fig. 24. The results differ from the Cassegrain results for source located at angles greater than 45° off-axis. This is because the primary side of the baffle is illuminated and scatters light directly to the primary mirror. The 70° baffles would fail for sources beyond 70° off-axis. The Cassegrain system has vanes on the main baffle (not the sunshade) and the sunshade occulted the direct illumination of the primary side of the 45° vanes. This accounts for the subtle but important difference in the results.

Usually the first-order scattering properties of the vane structures are more important than whether the vanes are angled or not. The results presented above confirm this statement. There are occasions where angled vanes would be beneficial, but to fully understand those cases a much longer explanation of diffuse vane baffle scatter is necessary. These results are detailed elsewhere.^{15,16}

There are special situations where angled vanes will have a significant advantage over annular vanes. One example is a bright source at a fixed offset angle. I have seen such a feature on a spaceborne telescope on a platform where there was nearby a brightly sunlit rocket-thruster casing at a fixed angle outside the field of view. Figure 25 shows the design where the vanes were aimed at the thruster at an angle where the primary mirror side (right side in Fig. 25) could not be directly illuminated by the sunlight scattered off the thruster. Under those circumstances most of the stray light

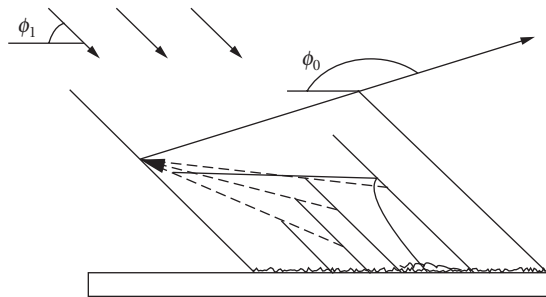


FIGURE 25 Angle-staggered vanes for fixed input angle. (Ref. 6, p. 104.)

had to make three scatters before exiting the vane cavity. In general, as soon as the position of the bright object is moved over a range of angles, the advantage of the angled vanes is lost. Nevertheless, there are many occasions within a sensor where the relative positions of a scattering source and a collecting object are fixed along a major stray light path. The front parts of the main-barrel baffle and the opening of the inner conical baffle in the Cassegrain design is an example. Many more examples could be cited. But the point is that you, as a designer, should first consider the first-order, single scatter paths off the baffle wall, each side of the vanes, and the bevel, for the full range of input values. Based on that information you can make the decision to use planar or angular vanes.

Vane Depth Considerations By varying the vane depth in the example analysis we can evaluate how the vane spacing-to-depth ratio affects system performance. Figure 26 gives the results of an analysis of the Cassegrain system with varying vane depths on the main baffle of 0.2, 0.4, and 0.8 inches. Figure 27 gives similar output from the Z-system analysis results. The performance of the system gets better as the vane depth increases from 0.2 to 0.4 inches, but there is little performance difference between the 0.4- and 0.8-inch baffle depths. The latter is the normal case. The 0.2-inch vane depth allows for a single path from the walls of the baffle tube, which increases the stray light propagation. Once that path is blocked by a greater vane depth, no further improvement should be expected due to further increases in vane depth.

The intent of presenting the two different optical designs was not to trade off one optical design against another. It needs to be made clear that the two optical sensors being used as examples are intentionally not equivalent from stray light design considerations. This is why the changes in performance are design-dependent. The nominal design of the eccentric pupil has a reimager, and the Cassegrain does not. The Cassegrain could have a reimager, in which case the stray light performance of both could be made essentially equal. It would depend on the optical design characteristics, $F/\#$, field of view, obscuration ratio, etc. The Cassegrain design has a specular sunshield and the Z-system has a vaned diffuse baffle structure. Which would perform better could only be determined after all of these features are considered.

To summarize, the general points being made in this section are

1. Usually, angled vane structure has little, if any, additional benefit over straight, annular vanes, and the annular vanes are much easier to fabricate and assemble.
2. Once the depth of a diffuse black vane structure is deep enough to block the single scatter path, further increases will not improve performance.

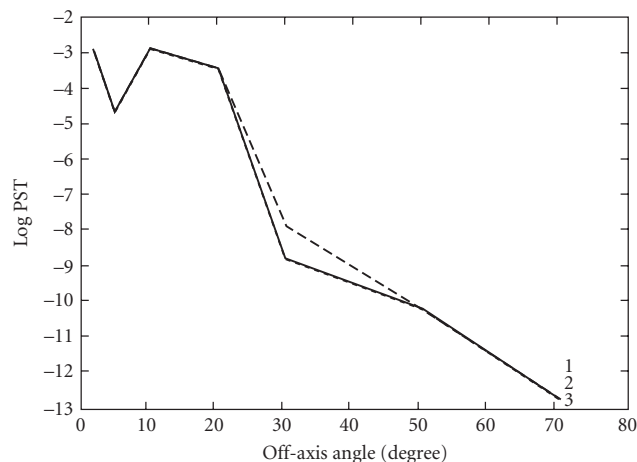


FIGURE 26 Cassegrain 90° baffles, coated with Martin Black, at varying depths; 1 = 0.2-inch, 2 = 0.4-inch, 3 = 0.8-inch depth. (Scatter is dominated by baffles.) (Ref. 6, p. 98.)

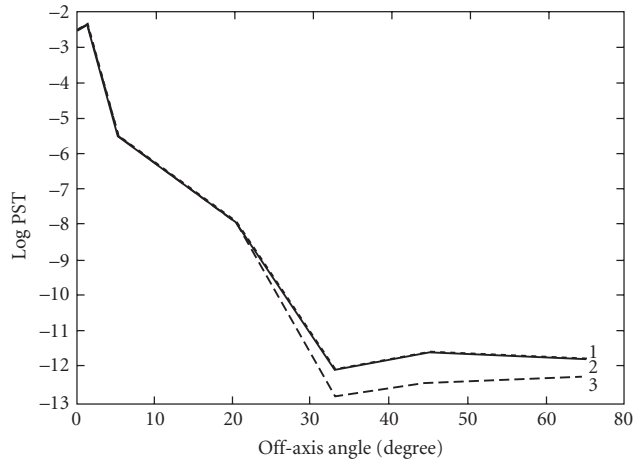


FIGURE 27 Z-system with varying vane depths. 1 = 0.2-inch, 2 = 0.4-inch, and 3 = 0.8-inch depth. (Ref. 6, p. 114.)

Specular Vanes Another aspect about vane structure that has been explored, but only in a limited way, is the specular vane cavity. Previous studies indicated that specular vanes have a problem with the aberrated rays and near specular angle scatter; this problem is severe enough to degrade the performance significantly.^{17,18} In another study by Freniere this was not always true.¹⁹ The ASAP²⁰ stray light software was used to evaluate the Z-system (Fig. 28) with (1) no vane structure, but with the main barrel baffle coated with Martin Black; (2) with Martin-Black-coated vanes; and (3) with a specular vane structure. The results show a dramatic degradation in the stray light performance without the coating on the main baffle tube. A subsequent specular baffle design developed by Nick Stavroudis has been shown to be a major improvement over previous concepts.²¹

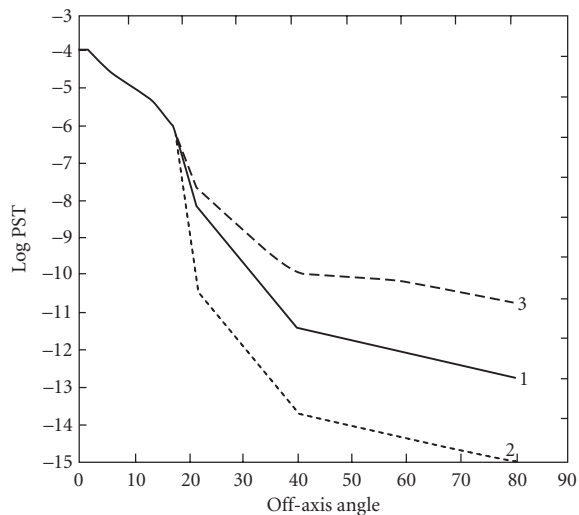


FIGURE 28 PST for unobscured pupil design without vane structure, with diffuse vane structure, and with specular vanes. 1 (solid) = no vanes, diffuse black coating; 2 (dotted) = diffuse vanes on main tube; 3 (dashed) = specular vanes on main tube. (Ref. 6, p. 115.)

Contamination Levels

Light scattered from a particulate-contaminated surface can have a pronounced effect on the stray light performance of a system.²²

I will now relate the performance of both designs (the centrally obscured Cassegrain, and the unobscured eccentric pupil) as a function of the level of scatter, per MIL-STD 1246A.²³ This analysis evaluates the sensor for different amounts of contamination on the optics only. The levels of contamination as defined in IEST-STD-CC1246D are for a distribution of particles with a specified range in particle sizes.

Ray Young used Mie scattering theory to predict the BRDF of a mirror covered with such MIL-STD distributions.²⁴ Table 1 was generated from Young's work for the 10 μm radiation. This table shows the base BRDF value and the BRDF slope that would be used in a typical stray light analysis program for input. The base value is the BRDF value at $(\beta - \beta_0) = 0.01$ and the second term is the slope of the BRDF in a (log-log) plot of BRDF versus $(\beta - \beta_0)$. β_0 is the sine of the angle of incidence and β is the sine of the observation angle.²⁵ The terms work equally well for out-of-plane values, but the above definitions, for simplicity, assume in-plane scattering data. See also the works of Spyak.²²

Spyak and Wolfe²⁶ did a series of experiments and calculations that relate BRDF to particulate contamination. They counted and sized particles on a mirror surface, and then measured the contribution of these particles to the mirror's BRDF at both visible (633 nm) and infrared (10.6 μm) wavelengths. They also performed Mie theory calculations and compared their calculations with the measured BRDFs. At both visible and infrared wavelengths, Mie theory calculations were a reasonable estimate for contribution of particulates to a mirror's BRDF. In most cases agreement between Mie calculations and their measurements were within a factor of two. Spyak and Wolfe also published Mie calculations of the BRDF expected from the MIL-STD-1246B (now IEST-STD-CC1246D) standard at 633 nm and 10.6 μm .*

Michael Dittman²⁷ has published a series of Mie calculations at five wavelengths. His calculations were done for the IEST 1246D distributions, but he considered an additional distribution in which the "particle slope" on the distribution was reduced from 0.926 to 0.383. The latter slope results in more large particles, and is commonly thought to be a more realistic distribution on surfaces that are exposed in a cleanroom environment.*

There is a problem with specifying the optics with this standard because it is difficult to reliably relate a level of contamination by particles to a BRDF performance. Two equal sizes and distributions of particulates may not give the same BRDF, because the index of refraction, the reflectivity, and the roughness of the particulates enter the calculations. In general, few people go to the trouble to determine these other factors. These factors will vary from one distribution to another. BRDF is the most usable value when performing a stray light analysis, so it should be the stray light specification. For manufacturing specifications, other parameters may be more appropriate, but they are not as good as BRDF for a stray light specification.

The level of scatter is also given in Table 1 along with the BRDF. The BRDF data from particulate scatter for the 5-, 10-, and 20- μm wavelengths for the 100, 300, and 500 contamination level have

TABLE 1 Mirror Scatter Relationships [Wavelength = 10 μm , BRDF Slope in Log $(\beta - \beta_0)$]

BRDF at $(\beta - \beta_0) = 0.01$	BRDF Slope	Cleanliness Level
0.02	-1.17	500
0.01	-1.17	454
0.001	-1.17	300
0.0001	-1.17	204
0.00001	-1.17	100

*Personal communication on partial contamination paragraphs, contributed by Dr. Gary Peterson, Brealut Research Organization, Inc.

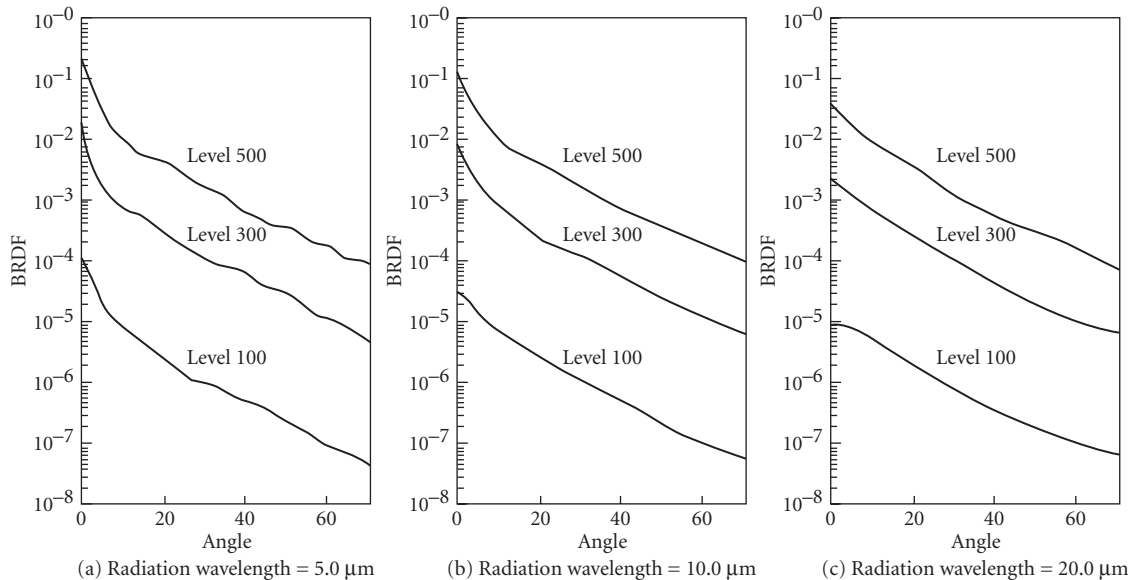


FIGURE 29 Predicted BRDFs on particles deposited on low-scatter mirror for cleanliness levels of 100, 300, and 500 at radiation wavelengths of 5.0, 10.0, and 20.0 μm .²⁴ (Ref. 6, p. 105.)

been plotted in Fig. 29. Consensus, not factually documented, indicates that the current state of the art of contamination control is at the cleanliness level of 300 to 500 for the 10- μm -wavelength region. Measured BRDFs below level 200 are achievable in the lab for short periods of time. A stray light analyst is strongly advised not to predict a system's performance with values below $1.0\text{E-}3 \text{ sr}^{-1}$ in the 10- μm region. Based on historical performance, mirrors in the IR (10- μm region) consistently degrade to this value, usually because of particulate scatter. Research work performed under Rome Air Development Center contract for the detection, prevention, and removal of contamination from the ground and in space could greatly reduce the degradation currently experienced by IR sensors.²⁸

Hal Bennett presents the significance of particulate scatter, as shown in Fig. 30.²⁹ This figure shows an agreement between measured data and theoretical data, and illustrates why IR sensors are usually more sensitive to particulate scatter than RMS scatter; the opposite is true in the visible. Figure 30 also indicates why the wavelength scaling law does not usually relate visible BRDF measurements to BRDF measurements in the IR. The physical process is different.

Figures 31 and 32 are the representative point source transmittances (defined as the irradiance on the detector divided by the incident irradiance) for the cleanliness levels of 100 through 500 for each design. The Cassegrain is much less affected by changes in contamination level, because the scatter from the black-coated surfaces dominates all other scatters. If the system had a reimager its performance would be better because these black surfaces would be blocked from the field of view of the detector, and the stray light performance would be due to the cleanliness level of the optics. The eccentric pupil design is sensitive to changes in the mirror coatings because it does have a reimager, and the major source of scatter is from the mirror surfaces.

In summary, the impact of particle contamination on the performance of a system will depend on how well the system is designed to suppress stray light. The goal is to be limited by a single optical element, such as the collecting lens or mirror, which is the objective of the system. The eccentric pupil design (Z-system) has this design feature. The better the optical design from a stray light point of view, the more the system's performance will be degraded by particle contamination. The more the system performance is determined by the black coatings, the more it will be sensitive to degradations in the coatings on the baffles.

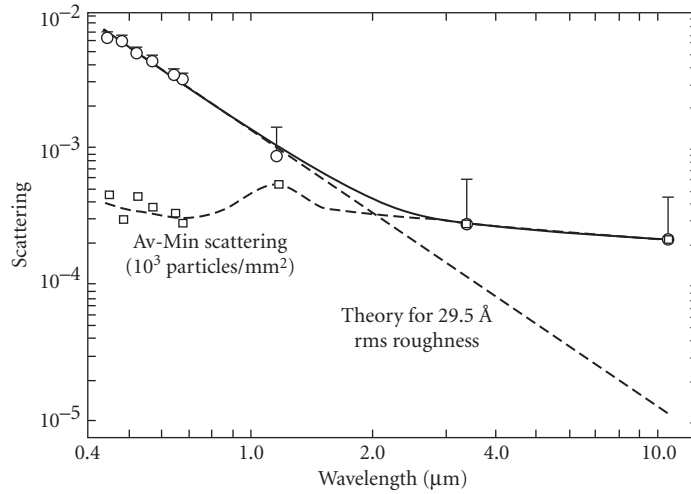


FIGURE 30 Scattering from polished dense flint glass. The diagonal line gives the contribution predicted for microirregularity scattering by a 29.5 Å rough surface. Circles indicate the minimum scattering observed, and the bars and squares between the average and minimum scattering observed at several points on the surface. This difference may be related to particulate scattering. (Ref. 29, p. 32.)

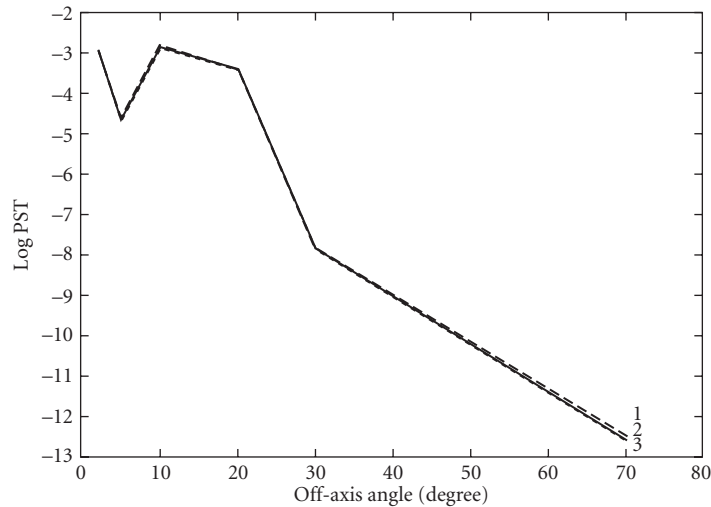


FIGURE 31 Cassegrain system with mirrors at all five contamination levels. 1 = 100, 2 = 204, 3 = 300, 4 = 454, 5 = 500. (Ref. 5, p. 113.)

Strut Design

In a centrally obscured system the central obscuration must be supported. In some designs (Schmidt-Cassegrains) the obscuration can be supported by a refractive element, but in most designs some form of struts are used. The most common error in strut design is to specify manufacture

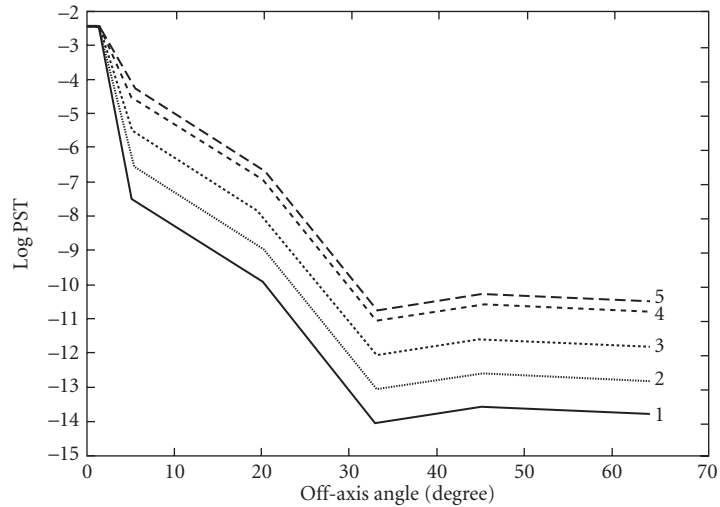


FIGURE 32 Z-system with mirrors at all five contamination levels. 1 = 100, 2 = 204, 3 = 300, 4 = 454, 5 = 500. (Ref. 6, p. 113.)

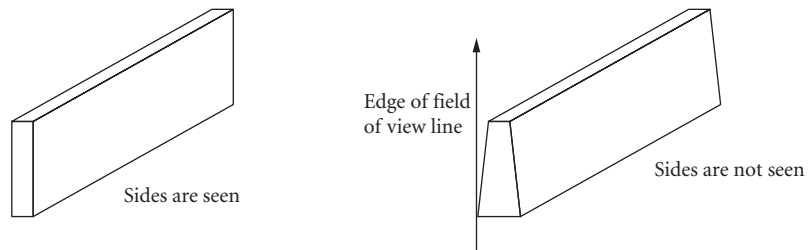


FIGURE 33 Angled strut design does not allow the detector to see the sides of the strut.

from a slab or plate of coated metal. Because all detectors have some finite field of view, the scatter from the sides of the struts can be seen from the image plane. Usually the struts are out “front” and exposed to more stray light sources than the objects deeper into the system. The near off-axis angles of incidence of scattered light off of the strut sides make for very high scattering toward the detector.

The proper strut design will preclude this path by making the object end of the strut narrower than the side nearest the objective (primary). This shape, shown in profile in Fig. 33, does not allow the detector to see the sides of the struts. The angle of the taper depends upon the object space field of view of the detector. It requires only a small change in design to remove this stray light path.

Basic Equation of Radiation Transfer

This section briefly discusses the most fundamental equation needed to perform the quantitative calculations of a stray light analysis. It reinforces the concept of first identifying what the detector can see and working on the geometry of the system to limit the stray light propagation, and not the BRDF term.

The fundamental equation relating to power transfer from one section to another is:

$$d\Phi_c = L_s(\theta_c, \phi_c) dA_s \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{R_{sc}^2} \quad (2)$$

where $d\Phi_c$ is the differential power transferred, $L_s(\theta_c, \phi_c)$ is the radiance of the source section, dA_s and dA_c are the elemental areas of the source and collector, and ϕ_s and ϕ_c are the angles that the line of sight from the source to the collector makes with their respective normals. This equation can be rewritten as three factors that help clarify the reduction of scattered radiation.

$$d\Phi_c = \frac{L_s(\theta_c, \phi_c)}{E(\theta_i, \phi_i)} E(\theta_i, \phi_i) dA_s \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{R_{sc}^2} \quad (3)$$

$$d\Phi_c = \text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c) d\Phi_s(\theta_i, \phi_i) d\Omega_{sc} \cos(\phi_s) \quad (4)$$

$$d\Phi_c = \text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c) d\Phi_s(\theta_i, \phi_i) \text{GCF}_{sc} \pi \quad (5)$$

$E(\theta_i, \phi_i)$ is the incident irradiance on the source section dA_s . GCF_{sc} is the projected solid angle from the source to the collector divided by π .

The GCF is independent of the first two terms and solely determined by the geometry of the system, including obscurations. The first term, $\text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c)$, is the bidirectional reflectance distribution function. It is usually considered independent of the second term, the incident power, and is therefore a function of the surface characteristics only. When reducing stray radiation propagation, one or more of these terms must be reduced. If any one of these terms is reduced to zero, no power will be transferred between the source and collector.

Stray Radiation Paths

Since the third term (GCF) in Eq. (4) is the *only* term that can be reduced to *zero*, it should receive attention first. This is a crucial point in a stray light analysis. Therefore, the logical starting place for stray light reduction is with the critical objects, since it is the GCF terms for these transfers which can be reduced to zero. Most novice analysts make the mistake of working on the BRDF term first.

$$\text{GCF} = \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{\pi R_{sc}^2}$$

The apparent possibilities for decreasing the GCF are to increase R_{sc} , ϕ_s , ϕ_c or to reduce the area dA_c . Not readily apparent is that the GCF is limited by apertures and obstructions. These features will, in some cases, block out the entire view of the source section from the collector so that there is no direct path. This is the mathematical basis for the logical approach, discussed at the beginning of the chapter. First block off as many direct paths of unwanted energy to the detector as possible, and then minimize the GCF for the remaining paths.

Point Source Transmittance Definitions

There are five common ways to define the merit function of the stray light in an optical sensor. The most common and preferred method is to define it as the output irradiance divided by the input irradiance, in terms of the *normalized detector irradiance* (NDI),³⁰ or in terms of the *point*

source normalized irradiance transmittance (PSNIT).³¹ This merit function is appropriate because it describes an irradiance transmittance, and it is relatively independent of the detector size.

A term often used in the past was the *off-axis rejection* (OAR), defined as the detector power divided by the input power from the same source *on-axis*. The term *rejection* is a misnomer because by definition the term describes a power transmittance, which can have little correlation with the rejected stray light. The second objection is that as a merit function it varies significantly with the detector size. If you double the area of the detector, the OAR will increase by about the same factor even though the system hasn't performed significantly worse in any way.

Another term commonly used is the system's stray light *point source power transmittance* (PSPT), or its reciprocal, the *attenuation* of the system. The PSPT is the detector power divided by the input power into the sensor from the specified *off-axis angle*. Again, this term varies with the detector size. Sometimes there is no well-defined entrance port so the denominator is impossible to define. Note that the magnitude of attenuation would normally be expressed in terms of a positive exponential. Beware that attenuations are often incorrectly called out with negative exponents.

A final PST definition that is sometimes specified is the *point source irradiance transmittance* (PSIT), defined as the output irradiance divided by the entrance port input irradiance. This definition becomes inappropriate when there is no clearly defined entrance port.

Surface Scattering Characteristics

Of the three potentially important factors in scattered radiation analysis cited above (the radiance of the undesirable source or sources, the geometry of the scattered radiation paths (GCF), and the surface scattering characteristics, (BRDF)), usually the first possibility considered is to improve the surface coatings or the addition of vane structure. In concept it *appears* to be the right place to start and that it is straightforward. Neither is the case; the BRDF never goes to zero as does the GCF, and the BRDF varies with input and output angles. However, with accurate *bidirectional reflectance distribution function* (BRDF) data and knowledge about the variations with applications, time, wavelength, and other factors, BRDF problems can be dealt with. The scattering characteristics of surfaces are discussed by Church, and the scattering characteristics of black coatings by Pompea and Breault elsewhere in this *Handbook*. The addition of vane sections on baffles can usually be considered as a specialized "coating" with its own specialized BRDF.

BRDF Characteristics

Usually, BRDF data that are presented represent only one profile of the BRDF, and many such profiles for various angles of incidence are necessary for understanding the scattering characteristics. However, studies have shown that a single profile of a mirror's surface scattering characteristics can be used, with some approximations, to define the BRDF for all angles of incidence.³² This is a significant achievement. It reduces the amount of data that must be taken, and it makes it easier to calculate, or estimate, the BRDF value for any set of input and output angles. The BRDF can also be reconstructed for cases where only a single profile of the function has been presented, which has been the usual practice.

The approximation has its limitations, as clearly detailed by Stover.³³ The approximation is quite good for nominal angles of incidence (see Fig. 34).³⁴ However, it breaks down for very high θ_i and high observation angles θ_o .

It is important to understand qualitatively the scattering characteristics of diffuse black coatings. Figure 35 shows the BRDF profile of Martin Black at 10.6 μm for several angles of incidence.³⁵ At near-normal angles of incidence the BRDF values are bowl-shaped; the values increase at large observation angles from the normal. At high angles of incidence the BRDF values in the near specular direction have increased by 2.5 orders of magnitude. There is a good discussion of the qualitative characteristics of diffuse black surfaces by Pompea and Breault elsewhere in this *Handbook*.

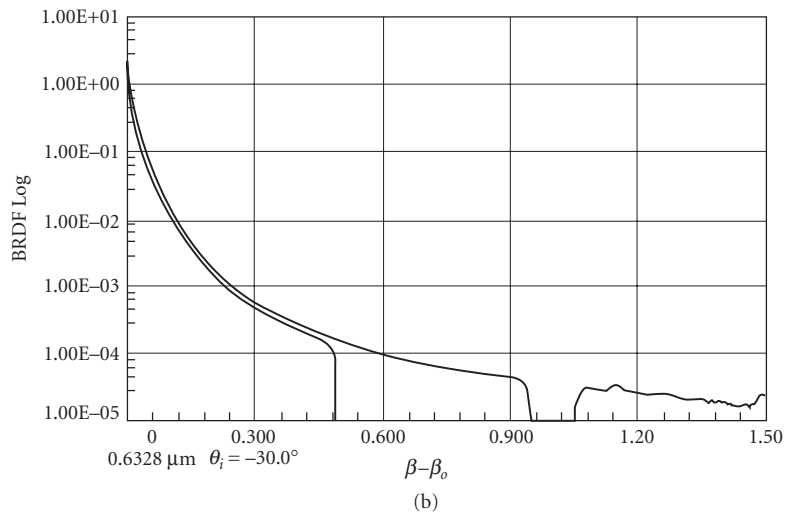
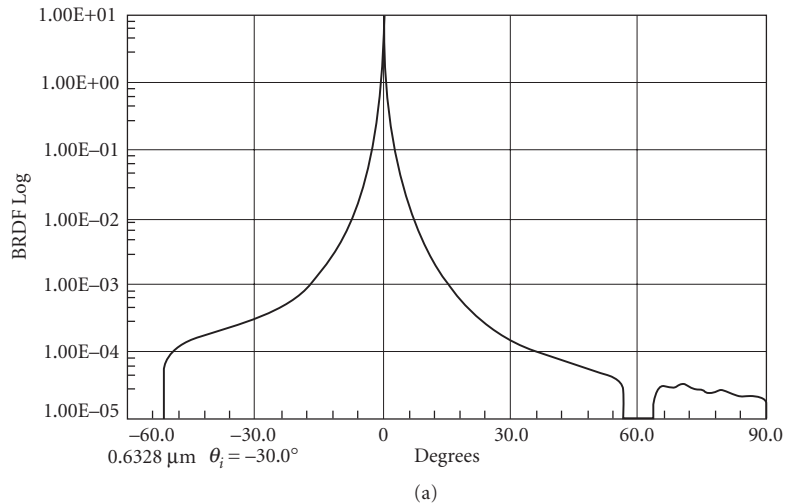


FIGURE 34 (a) The BRDF is asymmetrical when plotted against $\theta_s - \theta_i$, (b) The data in (a) exhibits near symmetry when plotted against $|\beta - \beta_o|$. The slight deviation from symmetry is due to the factor $(\cos \theta_s Q)$, where Q is a polarization factor. (Ref. 34, p. 69, reprinted with permission.)

7.4 OPTICAL SOFTWARE FOR STRAY LIGHT ANALYSIS

There is a small bevy of commercial optical software programs on the market that perform stray light analyses or some aspects of the stray light problems that are typically encountered. One must be knowledgeable of what each package can accomplish so the best thing to do is to ask for a demonstration of what each optical software manufacturer has to offer that is relevant to your design challenge. Here's a summary of a few commercially available optical software codes for stray light analysis. Note also that the programs' capabilities are always in a state of flux.

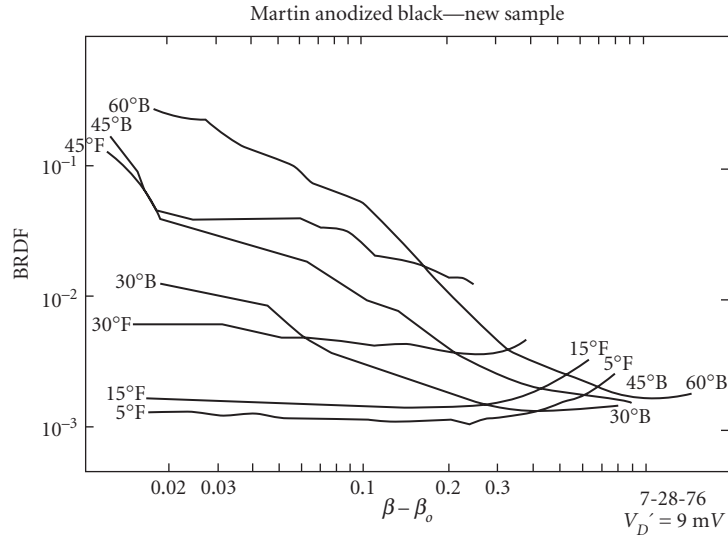


FIGURE 35 BRDF profile of Martin Black at $10 \mu\text{m}$. (F. O. Bartell et al., "A Study Leading to Improvements in Radiation Focusing and Control in Infrared Sensors," Final Report, Army Materials and Mechanics Research Center, December 1976.)

ASAP, by Breault Research Organization, Inc.

ASAP Optical Software was developed to meet design and analysis criteria of imaging and illumination systems and the unique challenges of stray light analysis with CAD interoperability. ASAP is powered by the ASAP nonsequential ray-tracing engine—known throughout the optics industry for its accuracy and efficiency. Rays can encounter surfaces in any order and any number of times, with automatic ray splitting. Optimized for speed, ASAP will trace millions of rays in minutes. The standard edition of ASAP includes a license of the SolidWorks Parts Only 3D Modeling Engine—an intuitive 3D-design environment optimized for use with ASAP. The user can write ASAP geometry files from within SolidWorks, import XML files, or use BRO's proprietary smart-IGES system to import system models from any CAD package while maintaining fast, efficient ray trace speed.

Use ASAP to model complex imaging systems, illumination systems, and light-concentrating devices. Create highly accurate source models using source images, point sources, ray grids, and fans. Model incandescent bulbs, LEDs, CCFLs, and HID arc lamps, or import from the BRO Light Source Library. Perform the analyses necessary to validate your designs without experimental prototyping.

ASAP includes a distributed-processing capability allowing the user to complete big design and analysis jobs effectively in short-time span—spawn up to 5 additional ASAP sessions on other local area network (LAN), without leaving your desk. Web site: www.breault.com

FRED, by Photon Engineering

FRED is an optical engineering software package that uses a statistical ray sampling approach to analyzing incoherent stray light mechanisms in optical systems. The user can assign one or more of many different BSDF scatter functions to a surface. When a ray is incident on the surface, a specified number of scatter rays are generated in random directions into the hemisphere according to a uniformly random angular distribution, although a Monte-Carlo technique can be used

to generate scatter rays with directional ray density proportional to the BSDF function. The power assigned to any particular scatter ray is proportional to the incident ray power and value of the BSDF function evaluated in the direction of the scatter ray. Uniform sampling of angle space has the advantage that lower power scatter paths will be realized with higher probability. Ray-density directional sampling has the advantage that more rays are directed in the higher power directions which decreases the statistical noise in those directions.

When analyzing scatter from a surface, the user is often interested in only a subset of the entire hemispherical angular range. In these cases a technique called importance sampling can often be used to great advantage. The user can specify that the same number of scatter rays be generated into the importance sample direction as was originally directed into the hemisphere to decrease the statistical noise. An alternative is to decrease the number of rays but still achieve the same noise level. The desired angular range can be specified in a variety of ways and FRED will then direct scatter rays randomly into the specified range. The power assigned to each scatter ray is adjusted to account for the fact that it can be directed only into a subset of the entire hemisphere. Web site: www.photonengr.com

LightTools and CodeV, by Optical Research Associates

LightTools is a complete illumination design and analysis software product. It combines full optical accuracy, powerful optical and illumination analysis, and an intuitive graphical user interface in a 3D solid modeling environment where models interact with rays to produce virtual prototypes of manufacturable systems. A fully integrated illumination optimization capability automatically improves model performance. *LightTools'* Monte Carlo ray tracing facilitates accurate spectral modeling of the illuminance, luminance, intensity distributions, and CIE colorimetric data anywhere in the optomechanical model.

In addition to illumination system design, *LightTools* supports many other applications, from packaging studies to stray light analysis. For example, its ray path visualization collects and displays information about ray-surface interactions to identify system elements that are contributing to light loss, scatter, unintentional reflections, or ghost images. A unique point-and-shoot ray tracing capability allows rapid, interactive evaluation of optical behavior.

CODE V is used for the optimization, analysis, and tolerancing of image-forming optical systems and free-space photonic devices. Its many capabilities include powerful local and global optimization for optics, fast wavefront differential tolerancing that allows as-built considerations to be evaluated throughout the design process, and highly accurate diffraction beam propagation analysis. For stray light applications, CODE V can be used to analyze ghosts in imaging systems due to Fresnel reflections. Web site: www.opticalres.com

ZEMAX, by ZEMAX Development Corporation

The ZEMAX program, from ZEMAX Development Corporation, has two modes of use. Its primary use is as a sequential optical design (optimization) program. In this mode it has tools to help identify location of ghost pupils and images resulting from Fresnel reflections. A separate nonsequential mode has many capabilities necessary for stray light analysis, including scatter modeling with importance sampling.

ZEMAX's Nonsequential ray-tracing capabilities can further be extended to finding rays which have specific characteristics or properties. For example, imagine you are studying the stray light in a telescope:

How significant are rays which "ghost" reflected off of various surfaces (both mechanical and optical)?

Rays which are experience multiple reflections may be important, but how significant are those which experience more than four reflections?

How effective is a strategically placed baffle in terms of limiting the amount of stray light on the detector?

Website: www.zemax.com

TracePRO, by Lambda Research Corporation

TracePro, from Lambda Research Corporation, is a 3D Computer Aided Design (CAD) program for simulating the performance of illumination and optical systems. TracePro can model the propagation of light in imaging and nonimaging optomechanical systems. Models are created by combining imported lens designs, imported CAD geometry (IGES, STEP, SolidWorks, Pro/E, CATIA, or Inventor files), and geometrical objects created using TracePro's user interface. Optical properties are then assigned to each solid and surface using the TracePro interface or through the TracePro Bridge for SolidWorks. Source models are added by specifying grids, surface emitters, ray file data or by using the surface source utility. Rays are ray traced through the model, while keeping track of absorption, specular reflection and transmission, fluorescence and scatter at each intersected surface or volume scatter site.

From TracePro models, the user may ray trace and analyze:

- Light distributions in illumination and imaging systems
- Stray light, scattered light, and aperture diffraction
- Throughput, loss, or system transmittance
- Flux or power absorbed by surfaces and bulk media
- Light scattering in biological tissue
- Polarization effects
- Fluorescence effects
- Birefringence effects
- Lit Appearance

Website: www.lambdare.com

SPEOS, by Optis

OPTIS' simulation software family, SPEOS and OptisWorks. It manages and optimizes many of the optical aspects of a broad range of sensors: reflection, refraction, scatter from surfaces, diffraction, absorption, polarization, and Gaussian beam propagation. It calculates stray light, illumination, and realistic optical simulations. Any product that needs to manage interactions of light and surfaces is calculated. It deals with the various types of light sources also. The simulations limit the need for costly prototyping of systems.

OPTIS simulation software allows the designer to "see" and realistically render products to depict what the final performance of the illuminator will look like in its applied application with stunning similarity. Its software produces a unique and accurate physiological human vision model of the final lit product for comfort, safety, and performance.

7.5 METHODS

There are two distinct methods that have been used to evaluate a system for stray radiation. You can either build the system and test it, or you can model the system and try to predict its performance. Both methods have advantages and disadvantages. Taken *together* the two methods provide the means to ensure that the system will perform as desired.

Build-and-Test Approach

A common approach is to make the system and either use it or test it for stray radiation rejection. Certainly if the system consistently performs satisfactorily *in its operational environment*, it has passed the ultimate test. But what if it does not meet the desired or expected level of performance? Making more systems to test becomes expensive rapidly. In fact, for very large systems, usually only modifications (“fixes”) can be contemplated because of the high cost. This is not the only argument against the build-and-test approach. The tests are rarely designed to determine *how* the scattered radiation is propagating through the system and which surfaces contribute most of the undesired radiation.

It is this information, and a thorough knowledge of the surface scattering characteristics, that is necessary to make measurable improvements to the system. Such a test, when determining the propagation paths, should yield information about how the system is reacting to its *test* environment, including the test equipment. Unless the tests are being conducted in the environment for which the system was designed, it is imperative to determine that the *test* environment is not giving erroneous results (either better or worse). Without analyzing the test configuration, you should expect that the environment *will* affect the system stray light measurements. It is also incorrect to assume that the test environment can only add to the stray light background. It is sometimes assumed that if the system passes the stray light tests in the lab it will only perform better in space or wherever its intended environment is. This is not necessarily true.

Now that several points have been made about the difficulty of making valid experimental tests, it must be stated that valid tests can and should be made. The measurement costs need not be prohibitive. Even relatively large optical systems have been fabricated and then modestly redesigned. Changes to the system can be made until the desired information and stray radiation rejection is attained. In some cases it will be less expensive to test an existing system and modify it if necessary than to analyze the system with computer software.

The system-level test need not be extensive; it is not necessary to have an all-encompassing measurement from on-axis to 90° off-axis. Indeed, few facilities are capable of making such tests when the attenuation gets even modestly high. An important point to recognize is that the most important paths to check are those at the nearer off-axis angles where the attenuation is not so high. These can usually be measured reliably.

At small off-axis angles the stray light noise is more often much higher than the detector background noise, while at the higher off-axis angles the stray light noise is well below the electronic/detector noise. From a performance point of view, at the higher off-axis angles there is usually only one additional scattering object (scatter from the main baffle) before these same near off-axis angle paths are encountered or are reinvolved. The validation of the analysis will only be susceptible to the scatter from this one object that can't be fully tested at the system level, but most of the scatter paths, and usually all the most important ones, will have been validated by the near off-axis measurements.

This one additional surface scatter most often (especially on space-based sensors) involves the vanes on the main baffle that shields the primary objective. It will normally reduce the optical noise incident upon it by four to five orders of magnitude. That's why the optical noise goes dramatically below the electronic noise. Its most important role is to occult the direct illumination of the objective which is usually part of the most significant direct scatter path. The performance of this baffle and its vane structure could be analyzed separately and then measured independently to confirm that it too will perform as predicted.

NOTE: Contrary to some published papers you cannot, in general, multiply the stray light transmittance of two parts of a sensor and determine the system's overall performance. Although the main baffle system can be analyzed (or measured) independently from the rest of the system, it is not correct to take its performance and multiply it times the stray light performance of the rest of the system. The stray light propagation paths are far more important than the magnitudes of the two parts. In the above analysis where it was proposed that the main baffle could be measured independently it was to confirm its performance alone. A full-system stray light analysis was assumed.

Computer Analysis

As with the experimental tests, computerized analyses are also subject to errors. The three most significant ones are software limitations, scatter data of samples (not the real system), and user error. No software is capable of putting in every detail of a complex design, yet the computer model must faithfully represent the actual performance of the system. On the other hand, the software can put in “parts” with far greater mathematical precision than these parts can actually be assembled. Unless special studies are made the analyst does not usually account for assembly errors that might affect the actual system. The scatter characteristics of the surfaces, usually defined in terms of the bi-directional reflectance distribution function (BRDF), are usually measured on sample substrates, and controls must be exercised to ensure that the samples tested represent the sensor’s actual coatings, and that they do not change with time. The stray light analysis programs are also subject to errors in determining the significant paths. The experimental test is for the actual design, with real coatings, and will include any extraneous unintentional paths due to misalignment or other causes.

On the positive side, a software program can point out many flaws in the system that contribute stray radiation by considering the input BRDF characteristics of the coatings. A program can also do trade-off studies, parametric analysis, and in many other ways aid in the study of alternate designs. The analysis of the paths of scatter will suggest meaningful modifications and help to discard impossible designs. These analyses allow designers to test designs and make modifications before the design goes into production. This is very useful, since rejecting a sensor design is much easier when it is on paper than after it has already been built. It is usually much more cost-efficient than cutting new hardware, redesigning the system, or making fixes on the built system.

If you are in a field related to the optical design of a sensor, be it at the design level or the system level, you know that it would be preposterous to perform the optical design analysis and then put the system together without testing it for its image quality. Yet that is how far the pendulum has swung in favor of performing a stray light analysis over making a system-level stray light test. It reflects a major change in attitude since the early 1970s. It has been stated by stray light analysts that the reliability of a stray light analysis is now much higher than experimental test results, so some people avoid the latter. While there is a degree of truth in this statement, it is wrong to omit the stray light test at the system level.

The advantages and disadvantages of the two methods are summarized in Fig. 36. The disadvantages of the build-and-test approach are the strengths of the analysis method, whereas the strengths of the build-and-test approach cover the weaknesses of an analysis. Taken together these two methods give the greatest amount of reliable information which you can use to create the optimal system and have confidence in its performance. Jointly, they indicate the reliability of the analysis and test results.

Strength	Inexpensive	High	Easy	Real performance complete with manufacturing error	No missed paths	Real BRDF
Weakness	Expensive	Limited	Hard	Models only	Programmer error	Sampled BRDF measurements
	Cost	Information level	Changes	Real-world performance	Error	BRDF

Build and test
 Analysis

FIGURE 36 Build-and-test and analysis methods complement each other.

7.6 CONCLUSION

In summary, the issues involved in designing a system with stray light suppression in mind are

- I. System design concepts
 - A. Critical objects seen by the detector
 - B. Illuminated objects
 - C. Lyot stops
 - D. Field stops
 - E. Optical designs
- II. Baffle and vane design
 - A. Diffuse and specular vane cavities
 - B. Vane edge scatter
- III. Diffraction
- IV. Strut design
- V. Scattering theory
- VI. BRDF data
 - A. Log BRDF versus θ
 - B. Log BRDF versus $\log(\beta - \beta_0)$
 - C. Polar plots
 - D. Isometric projections (3-D characteristics)
- VII. Coatings
 - A. Paints and anodized surfaces
 - B. AR coatings and other thin films
 - C. Mirror coatings
- VIII. Thermal emission
- IX. Ghost images
- X. Software
- XI. Detection, prevention, and removal of contamination

A step by step procedure that can help you to improve your system is:

- I. Start from the detector and identify what objects, called “critical objects,” can be seen from various positions on the detector. Be sure to include a point near the edge of the detector.
- II. Work to remove the number of critical objects that the detector can see.
- III. Determine what objects the source of unwanted radiation can see, called the “illuminated objects.”
- IV. If possible, reduce the number of illuminated objects seen.
- V. If there are illuminated objects that are also critical objects, work very hard on these paths. Orders of magnitude in improvement will be your reward.
- VI. If task V is not possible, then the computations are quite easy.
 - A. Calculate the power incident on the illuminated/critical objects.
 - B. Use Eq. (1) to calculate the transfer of power from the critical objects to the detector. Remember to properly account for the input and output angles when calculating the BRDF. *Do not* use a straight lambertian scatter distribution; there is no such distribution in reality.
- VII. Find all the paths connecting the illuminated objects to the critical objects.
- VIII. Evaluate the corresponding input and output angles at the illuminated and critical objects.

- IX. Determine if vane structure will help, or if some other redesign will effectively block these paths.
- X. For the calculated input and output angles, evaluate which coating would be lowest.
- XI. Perform the stray light calculation using Eq. (1) in an iterative fashion. This should determine the most significant stray light path and quantify the amount of stray light on the detector
- XII. Perform the above tasks for a series of off-axis positions of the point source.

7.7 SOURCES OF INFORMATION ON STRAY LIGHT AND SCATTERED LIGHT

- J. D. Lytle and H. Morrow (eds.), "Stray Light Problems in Optical Systems," *Proc. SPIE*, vol. 107, April 18–21, 1977 (22 papers).
- M. Kahan (ed.), "Optics in Adverse Environments," *Proc. SPIE*, vol. 216, Feb. 4–5, 1980 (30 papers).
- G. H. Hunt (ed.), "Radiation Scattering in Optical Systems," *Proc. SPIE*, vol. 257, Sept. 30–Oct. 1, 1980 (28 papers).
- S. Musikant (ed.), "Scattering in Optical Materials," *Proc. SPIE*, vol. 362, Aug. 25–27, 1982 (28 papers).
- R. P. Breault (ed.), "Generation, Measurement, and Control of Stray Radiation III," *Proc. SPIE*, vol. 384, Jan. 18–19, 1983 (15 papers).
- R. P. Breault (ed.), "Stray Radiation IV," *Proc. SPIE*, vol. 511, Aug. 23, 1984 (14 papers).
- R. P. Breault (ed.), "Stray Radiation V," *Proc. SPIE*, vol. 675, Aug. 18, 1986 (46 papers).
- R. P. Breault (ed.), "Stray Light and Contamination in Optical Systems," *Proc. SPIE*, vol. 967, Aug. 17–19, 1988 (33 papers).
- J. C. Stover (ed.), "Scatter from Optical Components," *Proc. SPIE*, vol. 1165, Aug. 8–10, 1989 (42 papers).
- R. P. Breault (ed.), "Stray Radiation in Optical Systems," *Proc. SPIE*, vol. 1331, July 12–13, 1990 (29 papers).
- J. C. Stover (ed.), "Optical Scatter: Applications, Measurement, and Theory," *Proc. SPIE*, vol. 1530, July 21–27, 1991.
- R. P. Breault (ed.), "Stray Light and Contamination in Optical Systems II," *Proc. SPIE*, vol. 1753, July 21–23, 1992.
- F. O. Bartell et al., "A Study Leading to Improvements in Radiation Focusing and Control in Infrared Sensors," *Final Report Prepared for Army Materials and Mechanics Research Center*, December 1976.
- J. A. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers," M.S. thesis, University of Arizona, 1977.
- P. J. Peters, "Stray Light Control, Evaluation, and Suppression," *Proc. SPIE*, vol. 531, January 1985.
- T. W. Stuhlinger, "Bidirectional Reflectance Distribution Function (BRDF) of Gold-Plated Sandpaper," M.S. thesis, University of Arizona, 1981.
- A. W. Greynolds, "Computer-Assisted Design of Well-Baffled Axially Symmetric Optical Systems," M.S. thesis, University of Arizona, 1981.
- J. W. Figoski, "Interferometric Technique for the Reduction of Scattered Light," M.S. thesis, University of Arizona, 1977.
- J. S. Fender, "An Investigation of Computer-Assisted Stray Radiation Analysis Programs," Ph.D. dissertation, University of Arizona, 1981.
- D. A. Thomas, "Light Scattering from Reflecting Optical Surfaces," Ph.D. dissertation, University of Arizona, 1980.
- R. P. Breault, "Suppression of Scattered Light," Ph.D. dissertation, University of Arizona, 1979.
- P. R. Spyak, "A Cryogenic Scatterometer and Scatter from Particulate Contaminants on Mirrors," Ph.D. dissertation, University of Arizona, 1990.
- F. O. Bartell, "Blackbody Simulator Cavity Radiation Theory," Ph.D. dissertation, University of Arizona, 1978.

- L. D. Brooks, "Microprocessor-based Instrumentation for BSDF Measurements from Visible to FIR," Ph.D. dissertation, University of Arizona, 1982.
- A. G. Lusk, "Measurements of the Light Scattering Profile of Small Size Parameter Fibers," M.S. thesis, University of Arizona, 1987.
- K. Nahm, "Light Scattering by Polystyrene Spheres on a Conducting Plane," Ph.D. dissertation, University of Arizona, 1985.
- G. W. Videen, "Light Scattering from a Sphere on or Near an Interface," Ph.D. dissertation, University of Arizona, 1992.
- Y. Wang, "Comparisons of BRDF Theories with Experiment," Ph.D. dissertation, University of Arizona, 1983.
- S. J. Wein, "Small-Angle Scatter Measurement," Ph.D. dissertation, University of Arizona, 1989.
- J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, Washington, D.C., 1989.
- J. C. Stover, *Optical Scattering Measurement and Analysis*, McGraw-Hill, Inc., New York, 1990.

7.8 REFERENCES

1. R. P. Breault, "Problems and Techniques in Stray Radiation Suppression," *Stray Light Problems in Optical Systems*, J. D. Lytle and Howard Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 2–23.
2. R. P. Breault, A. W. Greynolds, and S. R. Lange, "APART/PADE Version 7: A Deterministic Computer Program Used to Calculate Scattered and Diffracted Energy," *Radiation Scattering in Optical Systems*, G. Hunt (ed.), *Proc. SPIE* **257**, 1980, pp. 50–63.
3. R. V. Shack, "Analytic System Design with a Pencil and Ruler—The Advantages of the $y-\bar{y}$ Diagram," *Applications of Geometrical Optics*, *Proc. SPIE* **39**, 1973.
4. S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program," *Stray-Light Problems in Optical Systems*, J. D. Lytle and H. Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 89–97.
5. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
6. R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis," *Stray Light and Contamination in Optical Systems*, *Proc. SPIE* **967**, 1988.
7. Ibid., S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program."
8. Ibid., R. P. Breault, A. W. Greynolds, and M. A. Gauvin, "Stray Light Analysis with APART/PADE, Version 8.7."
9. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
10. Ibid., S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program."
11. This section is a part of a more complete mathematical description of the process which can be found in R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis." (Ibid.)
12. R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis" (Ibid.) contains a more detailed description of the process.
13. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
14. R. P. Breault, "Suppression of Scattered Light," Ph.D. dissertation, University of Arizona, 1979.
15. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
16. J. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers," M. S. thesis, University of Arizona, 1977.
17. Ibid., R. P. Breault, A. W. Greynolds, and S. R. Lange, "APART/PADE Version 7: A Deterministic Computer Program Used to Calculate Scattered and Diffracted Energy."
18. Ibid., R. P. Breault, A. W. Greynolds, and M. A. Gauvin, "Stray Light Analysis with APART/PADE, Version 8.7."
19. E. R. Friere and D. L. Skelton, "Use of Specular Black Coatings in Well-Baffled Optical Systems," *Stray Radiation V*, Robert Breault (ed.), *Proc. SPIE* **675**, 1986, pp. 126–132.

20. A. W. Greynolds, computer code ASAP, Breault Research Organization, Inc., 1992.
21. G. L. Peterson, S. C. Johnston, and J. Thomas, "Specular Baffles," *Stray Radiation in Optical Systems II*, Robert P. Breault (ed.), *Proc. SPIE* **1753**, 1992.
22. P. R. Spyak, "A Cryogenic Scatterometer and Scatter from Particulate Contaminants on Mirrors," Ph.D. dissertation, University of Arizona, 1990.
23. "Product Cleanliness Levels and Contamination Control Program," *MIL-STD-1246A*, Dept. of Defense, Global Engineering Documents, Santa Ana, Calif., 18 Aug. 1967.
24. R. Young, "Low-Scatter Mirror Degradation by Particle Contamination," *Optical Engineering* **15**, no. 6, Nov.–Dec. 1976.
25. J. E. Harvey, "Light-Scattering Characteristics of Optical Surfaces," Ph.D. dissertation, University of Arizona, 1976.
26. P. R. Spyak and W. L. Wolfe, "Scatter from Particulate Contaminated Mirrors," *Optical Engineering* **31**, no. 8, Aug. 1992, pp. 1746–1784.
27. M. G. Dittman, "Contamination Scatter Functions for Stray Light Analysis," *Optical System Contamination: Effects, Measurements, and Control VII*, P. T. Chen and O. M. Uy (eds.), *Proc. SPIE*, **4774**, 2002, pp. 99–110.
28. Contact Dierdre Dykeman at Rome Air Development Center for the best way to access the results of this work.
29. H. E. Bennett, "Reduction of Stray Light from Optical Components," *Stray Light Problems in Optical Systems*, J. D. Lytle and H. Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 24–33.
30. Name coined by D. Rock, Hughes Aircraft Co., El Segundo, Calif.
31. Name coined by R. P. Breault; this is the name most often used in an APART stray light analysis.
32. *Ibid.*, J. E. Harvey, "Light-Scattering Characteristics of Optical Surfaces."
33. J. C. Stover, *Optical Scattering: Measurement and Analysis*, New York: McGraw-Hill, Inc., New York, 1990.
34. *Ibid.*, J. C. Stover, *Optical Scattering: Measurement and Analysis*, p. 69.
35. *Ibid.*, J. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers."

This page intentionally left blank.

THERMAL COMPENSATION TECHNIQUES

Philip J. Rogers and Michael Roberts

*Pilkington Optronics
Wales, United Kingdom*

8.1 GLOSSARY

c	surface curvature of an optical element
D	diameter
FN	F-number or focal ratio
f	paraxial focal length
G	thermo-optical constant (normalized thermal change of OPD)
h	(subscript) signifies “pertaining to the optic housing”
i	(subscript) number of a specific optical element
j	signifies a number of optical elements
K	Kelvin
k	thermal conductivity
n	refractive index
OPD	optical path difference
T	temperature
t	thickness
V	Abbe number of a refracting optical material
ν	spatial frequency
α	linear coefficient of thermal expansion
γ	thermal glass constant (normalized thermal change of optical power)
Δ	small, finite change
λ	wavelength
δ	infinitely small change of a parameter
ϕ	optical power (reciprocal of focal length)

8.2 INTRODUCTION

In the following, the thermal effects for which compensation is required are taken to be those that affect the focus and image scale of an optical system. Methods for quantifying and offsetting these effects were described some time ago,¹ similar information being provided by several other authorities.^{2,3,4} The thermal compensation techniques described in this chapter, with the exception of intrinsic athermalization, involve either mechanical movement of one or more parts of the optical system, or compensation achieved solely by choice of optical materials. Except in Sec. 8.5 titled “Effect of Thermal Gradients,” a homogeneous temperature change of all parts of the optical system is assumed.

Most optical materials undergo a change of refractive index n with temperature T , conveniently quoted as a rate of change $\delta n/\delta T$. The usual values of n and $\delta n/\delta T$ given for a material (and assumed in this chapter unless stated otherwise) are those relative to the surrounding air rather than the absolute values with respect to vacuum. Air has a $\delta n/\delta T$ of -1×10^{-6} at $T = 288$ K and 1 atmosphere air pressure for wavelengths between 0.25 and 20 μm .⁵ allowance for this must be made when a lens operates in vacuum or in an enclosed space where the number of air molecules per unit volume does not change with temperature. The absolute $\delta n/\delta T$ of an optical material can be found from

$$\left(\frac{\delta n}{\delta T}\right)_{\text{abs}} = n_{\text{air}} \left(\frac{\delta n}{\delta T}\right) + n \left(\frac{\delta n}{\delta T}\right)_{\text{air}} \quad (1)$$

where the value of n_{air} is approximately 1.0.

8.3 HOMOGENEOUS THERMAL EFFECTS

Thermal Focus Shift of a Simple Lens

The rate of change of the power ϕ (reciprocal of the focal length f) of an optical element with temperature T can be obtained by differentiating the thin lens power equation $\phi = c(n - 1)$, where c is the total surface curvature of the element. For a linear thermal expansion coefficient α of the material from which the element is formed this gives

$$f = \frac{1}{\phi} \quad \frac{\delta f}{\delta T} = -\frac{1}{\phi^2} \frac{\delta \phi}{\delta T}$$

$$\frac{\delta \phi}{\delta T} = +\phi \left(\frac{\delta n/\delta T}{n-1} - \alpha \right) \quad (2)$$

Therefore

$$\frac{\delta f}{\delta T} = -f \left(\frac{\delta n/\delta T}{n-1} - \alpha \right) \quad (3)$$

The material-dependent factor inside the parenthesis in Eqs. (2) and (3) is known as the thermal “glass” constant (γ) and represents the thermal power change due to an optical material normalized to unit ϕ and unit change of T . Tables 1 to 3 give γ values for a selected number of visual and

TABLE 1 Optical and Thermal Data for a Number of Visual Waveband Materials

Schott Glass	Optical Plastic Type	Refractive Index, n_e^*	Abbe Number, V_e^\dagger	Thermal Constant, $\gamma(\times 10^6)^\ddagger$	Thermo-Optical Constant, $G(\times 10^6)^\ddagger$	Thermal Conductivity, $k(W \cdot m^{-1} \cdot K^{-1})$
	FK52	1.487	81.4	-27	+1	0.9
	FK5	1.489	70.2	-11	+4	0.9
	BK7	1.519	64.0	-1	+7	1.1
	PSK53A	1.622	63.2	-13	+4	—
	SK5	1.591	61.0	+1	+7	1.0
	BaLKN3	1.521	60.0	-3	+7	1.0
	BaK2	1.542	59.4	-5	+6	—
	SK4	1.615	58.4	-2	+7	0.9
	LaK9	1.694	54.5	-1	+8	0.9
	KzFSN4	1.617	44.1	+4	+8	0.8
	LF5	1.585	40.6	-6	+7	0.9
	BaSF51	1.728	37.9	+8	+14	0.7
	LaFN7	1.755	34.7	+6	+12	0.8
	SF5	1.678	32.0	0	+11	—
	SFN64	1.711	30.1	-4	+9	—
	SF6	1.813	25.2	+6	+18	0.7
	Acrylic [§]	1.497	57	-279	-71	0.2
	Polycarbonate [§]	1.590	30	-247	-68	0.2

*At $\lambda = 546$ nm.†Defined as $(n_{546} - 1)/(n_{480} - n_{644})$.‡At $\lambda = 546$ nm and $T = 20$ °C.§Values (except conductivity) from Waxier et al. *Appl. Opt.* 18:102 (1979).

infrared materials along with the relevant V value (Abbe number) and other data. The much higher level of γ for infrared as opposed to glass optical materials indicates that thermal defocus (focus shift) is generally a much more serious problem in the infrared wavebands. The actual value of γ varies with both wavelength and temperature due to variations in the value of $\delta n/\delta T$ and α . In general, this is unlikely to cause major problems unless a wide wavelength or temperature range is

TABLE 2 Optical and Thermal Data for Selected 3- to 5- μ m Waveband Infrared Materials

Optical Material	Refractive Index, $n_{4\mu}$	Abbe Number, $V_{3-5\mu}$	Thermal "Glass" Constant, γ	Thermal Thermo-Optical Constant, G	Conductivity, $k(W \cdot m^{-1} \cdot K^{-1})$
Silicon	3.43	2.4×10^2	$+6.3 \times 10^{-5}$	$+1.7 \times 10^{-4}$	1.5×10^2
KRS5*	2.38	2.3×10^2	-2.3×10^{-4}	-1.5×10^{-4}	5.0×10^{-1}
AMTIR1†	2.51	1.9×10^2	$+3.9 \times 10^{-5}$	$+9.5 \times 10^{-5}$	3.0×10^{-1}
Zinc selenide	2.43	1.8×10^2	$+3.6 \times 10^{-5}$	$+7.3 \times 10^{-5}$	1.8×10^1
Arsenic trisulfide	2.41	1.6×10^2	-1.9×10^{-5}	$+3.4 \times 10^{-5}$	1.7×10^{-1}
Zinc sulfide	2.25	1.1×10^2	$+2.6 \times 10^{-5}$	$+5.2 \times 10^{-5}$	1.7×10^1
Germanium	4.02	1.0×10^2	$+1.3 \times 10^{-4}$	$+4.2 \times 10^{-4}$	5.9×10^1
Calcium fluoride	1.41	2.2×10^1	-5.1×10^{-5}	-1×10^{-6}	9
Magnesium oxide	1.67	1.2×10^1	$+1.9 \times 10^{-5}$	$+2.6 \times 10^{-5}$	4.4×10^1

*Thallium bromo-iodide.

†Ge/As/Se chalcogenide from Amorphous Materials Inc.

TABLE 3 Optical and Thermal Data for Selected 8- to 12- μm Waveband Infrared Materials

Optical Material	Refractive Index, $n_{10\mu}$	Abbe Number, $V_{8-12\mu}$	Thermal "Glass" Constant, γ	Thermo-Optical Constant, G	Thermal Conductivity, $k(\text{W} \cdot \text{m}^{-1} \cdot \text{K}^{-1})$
Germanium	4.00	8.6×10^2	$+1.2 \times 10^{-4}$	$+4.1 \times 10^{-4}$	5.9×10^1
Cesium iodide	1.74	2.3×10^2	-1.7×10^{-4}	-5.3×10^{-5}	1
Cadmium telluride	2.68	1.7×10^2	$+5.3 \times 10^{-5}$	$+1.1 \times 10^{-4}$	6
KRS5	2.37	1.7×10^2	-2.3×10^{-4}	-1.6×10^{-4}	5.0×10^{-1}
AMTIR1	2.50	1.1×10^2	$+3.6 \times 10^{-5}$	$+9.0 \times 10^{-5}$	3.0×10^{-1}
Gallium arsenide	3.28	1.1×10^2	$+7.6 \times 10^{-5}$	$+2.0 \times 10^{-4}$	4.8×10^1
Zinc selenide	2.41	5.8×10^1	$+3.6 \times 10^{-5}$	$+7.2 \times 10^{-5}$	1.8×10^1
Zinc sulfide	2.20	2.3×10^1	$+2.6 \times 10^{-5}$	$+5.0 \times 10^{-5}$	1.7×10^1
Sodium chloride	1.49	1.9×10^1	-9.5×10^{-5}	-3×10^{-6}	6

being considered.⁶ Thermal defocus results not only from a change of optical power but also from the thermal expansion coefficient α_h of the housing. Equation (3) can be modified to allow for the effect of the latter:

$$\text{Single thin lens:} \quad \Delta f = -f(\gamma + \alpha_h)\Delta T \quad (4)$$

$$j \text{ thin lenses in contact:} \quad \Delta f = -f \left[f \sum_{i=1}^j (\gamma_i \phi_i) + \alpha_h \right] \Delta T \quad (5)$$

Thermal Defocus of a Compound Optical Construction

Consider a homogeneous temperature change in an optical system that comprises two thin-lens groups separated from each other, the normalized thermal power change being the same in each lens group. Taking the thermal defocus calculated from Eq. (4) as unity, then that due to a compound optic comprising two separated components and of the same overall power can be estimated from Fig. 1.⁷ The latter shows scaling of thermal defocus with respect to a simple thin lens, relative to front lens/image plane distance (overall length) for three different positions of the second lens group. The graph is divided into three basic lens constructions distinguished from each other by overall optical length and/or the sign of the power of the front lens group.

Figure 1 assumes germanium optics in an aluminum housing but change of either material, while altering values slightly, has no effect on the following two conclusions:

1. Telephoto/inverted telephoto constructions always give more (and Petzval lenses always give less) thermal defocus than an equivalent simple lens.
2. Thermal defocus reduces as the second lens group is moved toward the image plane, irrespective of lens construction: the efficacy of this procedure is limited, however, by the increased imbalance of optical powers between the groups.

The thermal defocus scaling technique could be extended to cover an optic comprising more than two lens groups. This extension has been carried out for the Cooke triplet construction⁸ and for a series of separated thin lenses,⁶ although only for the case of a zero-expansion housing.

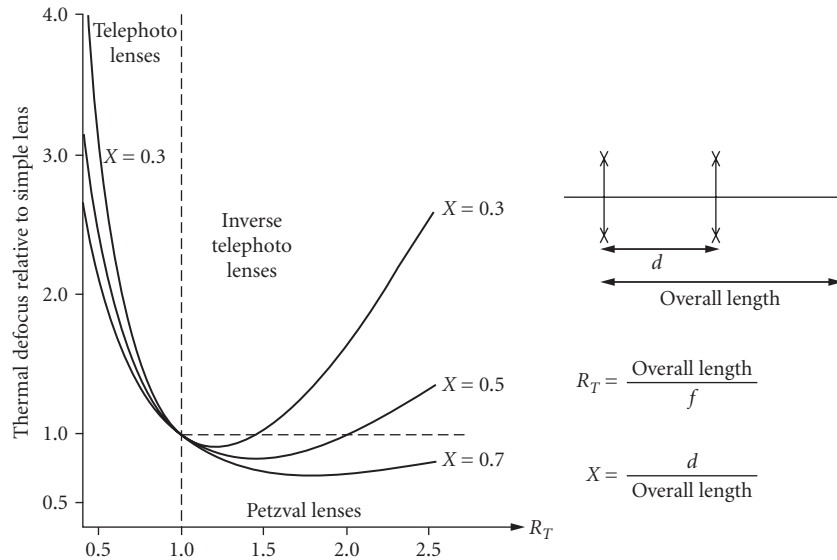


FIGURE 1 Effect of compound lens construction on thermal defocus. (From Rogers.⁷)

8.4 TOLERABLE HOMOGENEOUS TEMPERATURE CHANGE (NO COMPENSATION)

Diffraction-Limited Optic

Equation (5) can be used to establish the temperature change ΔT that will result in a quarter-wave of longitudinal thermal defocus, a reasonable limit for a simple optic that is nominally diffraction-limited. Given an optic of diameter D and focal ratio FN imaging at a mean wavelength of λ :

$$\text{Diffraction-based depth of focus:} \quad \Delta f = \pm 2\lambda(\text{FN})^2 \quad (6)$$

$$\text{Combining Eqs. (5) and (6):} \quad \Delta T = \pm \frac{2\lambda(\text{FN})}{D \left[f \sum_{i=1}^i (\gamma_i \phi_i) + \alpha_h \right]} \quad (7)$$

Figure 2 gives ΔT against D results for a simple 8- to 12- μm bandwidth germanium optic in an aluminum housing;⁹ the curves illustrate the small temperature change that can be tolerated in germanium optics before focus compensation is required. Partial avoidance of this particular problem may be achieved by the replacement of germanium by other infrared optical materials having lower values of γ : this may also be desirable to reduce high-temperature absorption but generally leads to much greater optical complexity.

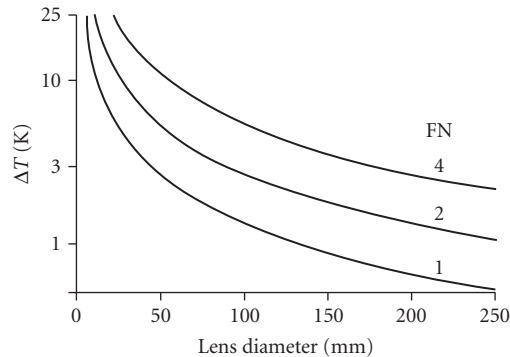


FIGURE 2 Tolerable temperature change for a simple germanium infrared lens. (From Rogers.⁹)

Nondiffraction-Limited Optic

The depth of focus of an optic having a nominal performance far from the diffraction limit is a function of the residual aberration level and balance in the optic as well as its first-order parameters. An estimate related to a cutoff spatial frequency ν that gives a reasonable approximation in many cases can be obtained¹⁰ from

$$\text{Approximate depth of focus:} \quad \Delta f = \pm \frac{(\text{FN})}{\nu} \quad (8)$$

$$\text{Combining Eqs. (5) and (8):} \quad \Delta T = \pm \left\{ D\nu \left[f \sum_{i=1}^j (\gamma_i \phi_i) + \alpha_h \right] \right\}^{-1} \quad (9)$$

Notice that, given the approximation of this method, the value of ν can be determined by extending a straight line MTF from 1.0 response at zero spatial frequency, through the MTF point of interest, to the intersection of the line with the spatial frequency axis.

8.5 EFFECT OF THERMAL GRADIENTS

The previous sections assume a homogeneous temperature change in all parts of the optical system: in situations where steady-state or transient temperature gradients exist, further consideration is required.¹

Allowance for the effect of steady-state longitudinal gradients can be made by applying a different value of T to each lens group and an average local temperature to each portion of the housing that separates two adjacent lens groups. Transient longitudinal gradients are a more difficult problem and, if severe, may require individual athermalization of each lens group in its own housing domain.

Steady-state or transient radial thermal gradients cause at least a shift of focus position, with the possible addition of a change of aberration correction. A localized radial temperature difference of

ΔT through the thickness t of a plane-parallel plate will cause a deviation of a ray of light¹¹ that can be quantified as an optical path difference (OPD):

$$\text{OPD} = t[\alpha(n-1) + \delta n / \delta T] \Delta T \quad (10)$$

The expression in the square bracket is often referred to as the thermo-optical constant G and is an approximate measure of the sensitivity of an optical material to radial gradients. More thorough analysis of the effects produced by radial thermal gradients includes computation of thermally induced stress and consequent anisotropic change of refractive index: in some cases, this may be a significant factor in image degradation.¹²⁻¹⁴

Tables 1 to 3 give values of G for the selected optical materials. Also tabulated is the thermal conductivity k , as in many cases G/k is a more appropriate measure of sensitivity given the greater ability of high-conductivity materials to achieve thermal equilibrium.

8.6 INTRINSIC ATHERMALIZATION

The need for athermalization can be avoided or minimized for some applications by employing optical power and mounting techniques that are inherently insensitive to temperature change. A concave spherical mirror fabricated from the same material that separates the mirror from its focal plane (e.g., an aluminum mirror in an aluminum housing) is in effect “self-athermalized” for a homogeneous distribution of temperature. The optical performance of a single spherical mirror is limited, but the above principle applies for more complex all-reflective optical constructions employing conic or other aspheric surface forms. A glass spherical mirror, although not thermally matched to an aluminum mounting, may be used as part of a self-athermalized catadioptric afocal in the infrared, a germanium Mangin being used in this case as a secondary mirror lens.¹⁵ The high thermal power change of the negatively powered lens in the germanium Mangin, used in double-pass, compensates for the thermal defocus due to the glass primary, the housing, and the remaining germanium optics in the afocal—Fig. 3.¹⁶

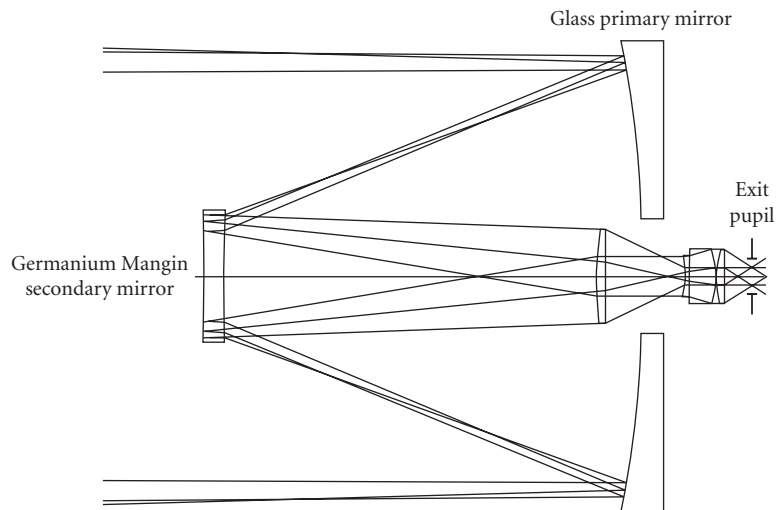


FIGURE 3 High magnification self-athermalized catadioptric afocal. (From Norrie.¹⁶)

An alternative approach to the above is to use glass-ceramic mirrors within a nickel-iron alloy housing, since they can have thermal expansion coefficients approaching zero. A major advantage of this approach is its insensitivity to thermal gradients.

8.7 MECHANICAL ATHERMALIZATION

General

Mechanical athermalization essentially involves some agency moving one or more lens elements by an amount that compensates for thermal defocus—a simple manual option being to use an existing focus mechanism. Automatic methods are, however, preferable in many cases and can be divided into passive or active. Passive athermalization employs an agency, often involving materials (including liquids) with abnormal thermal expansion coefficients, to maintain focus without any powered drive mechanism being required. Automatic active athermalization involves the computation of focus compensation algorithms that are stored (usually electronically) and implemented by a powered device such as an electric motor. The following sections refer to a number of passive and active athermalization methods, although the list is by no means exhaustive.

Passive Mechanical Athermalization

The principal advantages of passive thermal compensation methods are their relative simplicity and potential reliability. Disadvantages are their inadequate response to transient temperature gradients and, generally, lack of adjustment to allow for errors or unforeseen circumstances. Passive methods are ideal in glass optics¹⁷ where thermal effects are low, although here it is not too difficult to achieve optical athermalization (see later under “Optical Athermalization”) except where very low secondary spectrum is required. In the infrared wavebands, where thermal effects are much greater due to the nature of the optical materials, it is difficult to achieve simple passive mechanical athermalization due to the large refocusing movement required, typically 1.5×10^{-4} per unit focal length per Kelvin for an aluminum-housed germanium optic. An exception to the above is the combination of silicon and germanium in 3- to 5- μm optics, where thermal defocus results largely from the expansion of the housing. In this case, use of more than one nonmetallic housing material can result in an athermalized optic, even one having two fields of view¹⁸—Fig. 4.

For infrared cases other than the above, the options are either to provide a mechanism that modifies mechanical expansion effects or to reduce the required refocusing movement by optical means.

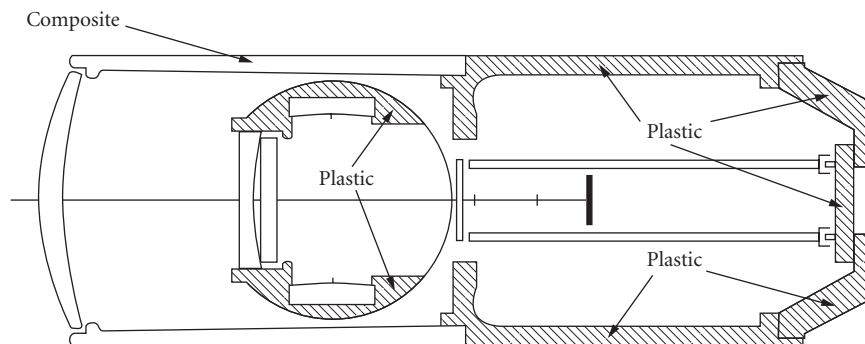


FIGURE 4 Part composite/part plastic mounting structure used for athermalization of a 3- to 5- μm IR optic. (From Garcia-Nuñez and Michika.¹⁸)

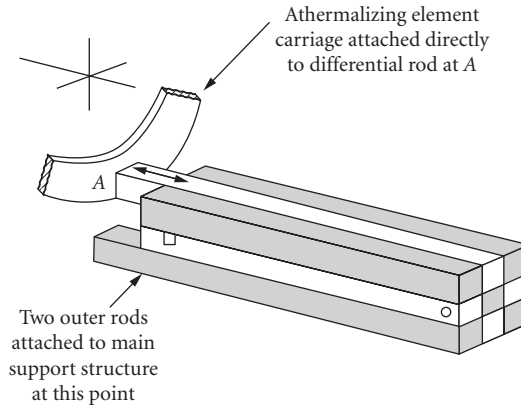


FIGURE 5 Passive mechanical thermal compensation using differential expansion rods. (From Povey.¹⁹)

Examples of the former include¹⁹ a series of linked rods of alternatively high and low expansion coefficient—Fig. 5—and a hydraulic method where the fluid contained in a large-volume reservoir expands into a small-bore cylinder—Fig. 6.

An interesting alternative employs shape-memory metal²⁰ to provide a large movement over a relatively small temperature range.¹⁹ Another alternative is to employ a geodetic arrangement: in this method¹⁹ an athermalizing adjustment of, for example, the separation between primary and secondary mirrors in a catadioptric, is produced by differing expansion coefficients of the primary mirror

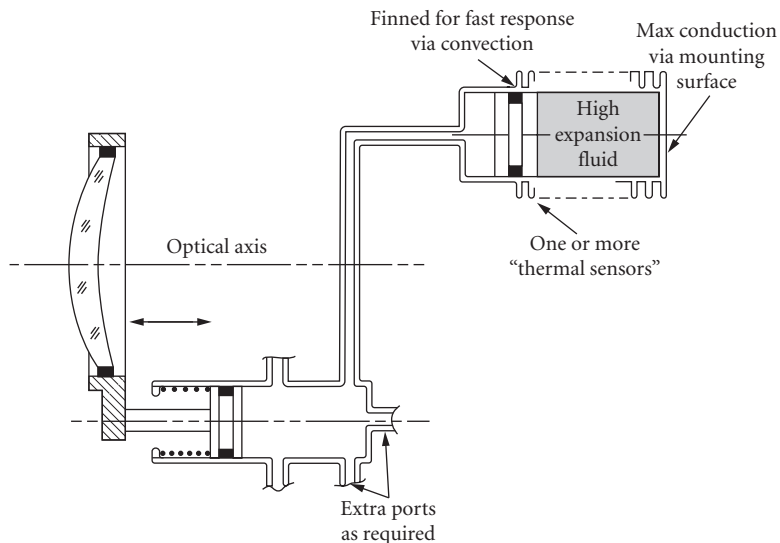


FIGURE 6 Passive mechanical thermal compensation using high-expansion-fluid thermal sensors. (From Povey.¹⁹)

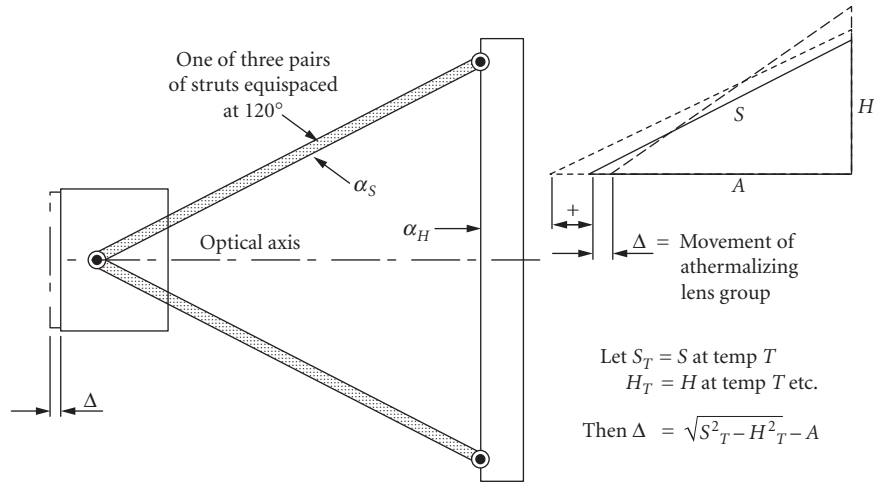


FIGURE 7 Geodetic support structure for positive or negative thermal compensation movement. (From Povey.¹⁹)

mount and the secondary mirror struts—Fig. 7. Where none of the above methods are desirable, the option to reduce the necessary movement may be the only alternative. This may be achieved by an optical layout⁹ configured such that the required athermalizing movement is reduced typically by a factor of four, but at the expense of somewhat greater optical complexity—Fig. 8.

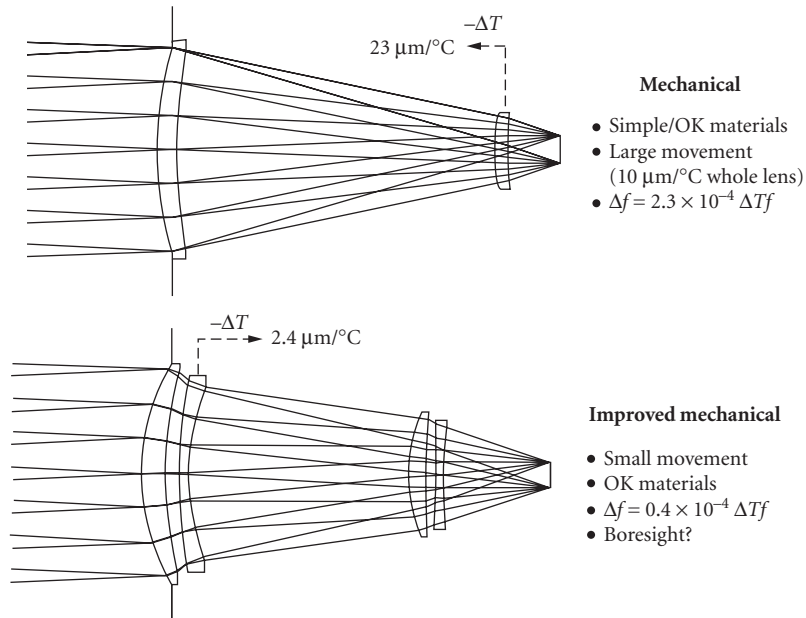


FIGURE 8 Alternative optical configurations for mechanically athermalized forward looking infrared (FLIR) systems. (From Rogers.⁹)

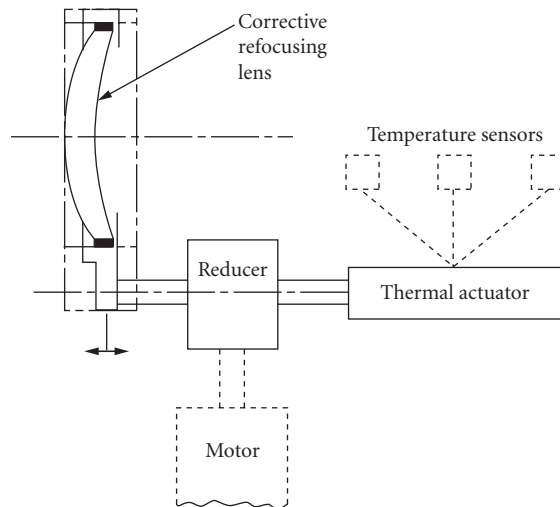


FIGURE 9 Active electromechanical athermalization—schematic.

Active Mechanical Athermalization

Active mechanical athermalization in its simplest form can be manual adjustment of a lens element or group for refocusing. For more complex optics, such as multi-field-of-view, a procedure can be specified for manual (or motorized) adjustment of several lens elements to maintain focus over a range of magnifications and temperatures.^{21,22} Where automatic athermalization is required, a method can be employed that uses a combination of electronics and mechanics—Fig. 9. One or more temperature sensors located along the body of the optic feed their signals into an algorithm that calculates the required movement of a compensating lens and then initiates the motion. For simplicity, the compensating lens may be that which already provides close-distance focusing, thus requiring only an increase in the range of movement for athermalization. The location of sensors is especially important for infrared optics and should be dependent on the thermal sensitivity variations within the optical system.

Active electromechanical thermal compensation is particularly suitable where transient longitudinal temperature gradients are expected and for multi-field-of-view optics where thermal defocus is dependent on field-of-view setting. The algorithm required for elimination of the effects of a combination of both of the above is complex, but compensation may be accomplished by a single mechanical motion.²³

A single motion does not, however, guarantee athermalization of image scale, in which case more than one compensatory movement may be required. Two motion athermalization in a zoom or dual-field-of-view infrared telescope can take advantage of the existing mechanisms required for field-of-view change. Also, by utilizing internal lens elements, problems associated with hermetically sealing an external focusing lens element can be avoided.^{24–28} In order to maintain stability of aberration correction in infrared zoom telescopes, particularly those having a large zoom range, three-motion athermalization has been proposed.^{29,30}

Active/Passive Athermalization

An improvement over simple manual active athermalization is to include partial passive athermalization. This is best suited to systems that already contain axially moving lens components, for example, a dual-field-of-view infrared telescope³¹—Fig. 10.³² Here the majority of the athermalization is provided

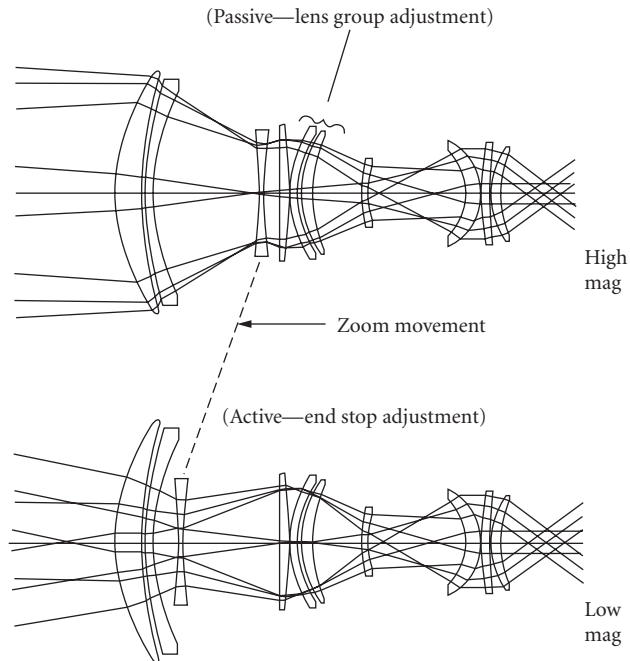


FIGURE 10 Part active, part passive mechanical athermalization. (From Roberts.³²)

by a mechanically passive device that adjusts the position of the rear lens group in the objective. The residual focus error is then corrected by small manual adjustments to the magnification change element. This technique can minimize the change of image scale and aberrations with temperature. A potential problem, however, is the subjective nature of best-focus determination.

Athermalization by Image Processing

Athermalization by image processing is suitable for some applications. A range of automatic focusing techniques exists but, while this approach has the advantage of not requiring temperature sensors, it does suffer the potential disadvantage of misinterpretation of image information.

8.8 OPTICAL ATHERMALIZATION

General

Athermalization of the focus position of an optical system by choice of refractive materials has been described extensively in the literature.^{33–42} The requirements of overall optical power, achromatism, and athermalism demand that three conditions be satisfied for j thin lens elements in contact:

Power:
$$\sum_{i=1}^j \phi_i = \phi \quad (11a)$$

TABLE 4 Unity Focal Length Athermal Two-Material Achromatic Combinations

Material Type	Material Combination	Total Curvatures	Secondary Spectrum [†]	Petzval Sum	Normalized Mass
Optical glasses	BaLKN3 + KzFSN4	+7.24/−4.49	3.6×10^{-4}	0.77	2.1
	BaK2 + LaFN7	+4.43/−1.86	4.8×10^{-4}	0.76	1.4
	FK5 + LF5	+4.85/−2.34	4.8×10^{-4}	0.73	1.3
	PSK53A + SFN64	+3.06/−1.28	5.0×10^{-4}	0.65	1.0
	BaLKN3 + BaSF51	+5.21/−2.35	5.2×10^{-4}	0.79	1.5
Stabilized optical glasses	SK4 + KzFSN4	+7.30/−5.64	1.8×10^{-4}	0.62	2.0 [‡]
	SK5 + SF5*	+3.97/−1.99	10.6×10^{-4}	0.67	1.0 [‡]
3- to 5- μm Materials	As ₂ S ₃ + MgO	+0.77/−0.12	8.6×10^{-4}	0.40	0.8 [§]

*Thermally invariant housing, all others aluminum.

[†]Over wavebands of 480 to 644 nm, 546 to 852 nm, and 3 to 5 μm respectively.

[‡]Relative to SK5/SF5 solution.

[§]Relative to lowest value in Table 5.

Source: Rogers.⁴³

$$\text{Achromatism:} \quad \sum_{i=1}^j \frac{\phi_i}{V_i} = 0 \quad (11b)$$

$$\text{Athermalism:} \quad \sum_{i=1}^j (\gamma_i \phi_i) + \phi \alpha_h = 0 \quad (11c)$$

The presence of three conditions implies the need for three different materials in order to obtain an exact solution. It is possible, however, to find achromatic combinations of two materials that are also athermal, provided that a simple condition is satisfied:⁴¹

$$V_1(\gamma_1 + \alpha_h) = V_2(\gamma_2 + \alpha_h) \quad (12)$$

Suitable combinations for thin-lens athermal achromats can be found by plotting a range of materials on a graph of γV against V ; the slope of the line joining a chosen pair representing the required thermal expansion coefficient of the housing.⁶

A number of approximately athermal optical glass achromats exist of which those listed in Table 4⁴³—with the exception of the last entry—represent examples with low to moderate secondary spectrum over the visible to near infrared waveband. The data given for these achromats are lens element total curvatures for unity focal length; secondary spectrum (second-order color); thin-lens Petzval sum; and an approximate indication of mass, normalized to the lowest value. The pairing of radiation-stabilized versions of SK5 and SF5, both of which have a low value of γ makes a good choice for athermalized space optics in a temperature-invariant mount.⁶

In the infrared wavebands the options are far more limited: at least one 3- to 5- μm waveband two-material athermal combination exists, namely, arsenic trisulfide and magnesium oxide, but there is currently no realistic pairing of materials in the 8- to 12- μm band.

Athermal Laser Beam Expanders

Many more two-material athermal combinations exist if the requirement for achromatism [Eq. (11b)] is removed. This is the situation that occurs with a (preferably) galilean laser beam expander, although here the two lens materials are separated.⁴⁴ From Eq. (4), making the thermal defocus Δf values equal

and opposite for the two lenses leads to a value for magnification at which two given materials in a specific housing material will provide an athermal beam expander (for a homogeneous temperature distribution):

$$\text{Magnification} = \frac{\gamma_1 + \alpha_h}{\gamma_2 + \alpha_h} \quad (13)$$

Three-Material Athermal Solutions

Graphical methods have been described that allow investigation of preferred three-material athermalized achromatic solutions.⁴¹ An alternative method is the systematic evaluation of all possible combinations of three materials selected from a short list, each combination being allocated a risk factor dependent on material characteristics and solution sensitivity.⁹ The optical powers of the three in-contact thin-lens elements are determined by solving Eq. (11a to c) which give for a unity focal length:

$$a = \frac{V_1 V_2 - V_2 V_3}{V_1 V_3 - V_2 V_3}, \phi_3 = \frac{(1-b)\gamma_1 + b\gamma_2 + \alpha_h}{(1-a)\gamma_1 + a\gamma_2 + \gamma_3} \quad (14a)$$

$$b = \frac{V_2}{V_2 - V_1}, \phi_2 = b - a\phi_3 \quad (14b)$$

$$\phi_1 = 1 - (\phi_2 + \phi_3) \quad (14c)$$

Tables 5 and 6 give a selection of lower-risk three-material solutions, in approximate order of increasing risk, for 3- to 5- and 8 to 12- μm infrared combinations, respectively. The data given are similar to those in Table 4, but the housing is assumed to be aluminum in all cases. Note that these tables are intended as a guide only and are based on currently available material data.

Athermalization of Separated Components

In many ways, thermal defocus and thermal change of focal length are analogous to longitudinal and lateral chromatic aberration, having the same first-order dependencies. For this reason it has been suggested that a thermal Abbe number, defined as $\gamma^{-1.3}$ be used to replace the chromatic

TABLE 5 Unity Focal Length Athermal Three-Material Achromatic Combinations for the 3- to 5- μm Waveband

Material Combination	Total Curvatures	Petzval Sum	Normalized Mass
Si + Ge + ZnS	+0.72/-0.36/+0.27	0.39	1.3
ZnSe + Ge + MgO	+1.16/-0.21/-0.06	0.51	1.8
Si + Ge + KRS5	+0.69/-0.26/+0.08	0.34	1.0
[ZnS + MgO + Ge]	+1.28/-0.17/-0.16	0.52	1.5
AMTIR1 + Ge + Si	+0.56/-0.32/+0.46	0.42	1.4
Si + MgO + KRS5	+0.31/-0.08/+0.22	0.31	1.1
ZnSe + ZnS + Ge	+1.80/-0.69/-0.23	0.50	2.8
Si + CaF ₂ + KRS5	+0.32/-0.25/+0.24	0.29	1.1

[] Low residual high-order chromatic aberration.

TABLE 6 Unity Focal Length Athermal Three-Material Achromatic Combinations for the 8- to 12- μm Waveband

Material Combination	Total Curvatures	Petzval Sum	Normalized Mass
KRS5 + ZnSe + Ge	+0.34/−0.15/+0.25	0.30	1.0
ZnSe + ZnS + Ge	+2.05/−0.92/−0.26	0.50	2.5
GaAs + ZnS + KRS5	+0.38/−0.20/+0.26	0.31	1.0
AMTIR1 + ZnS + Ge	+1.42/−0.35/−0.24	0.48	1.3
{CdTe + ZnSe + KRS5}	+0.72/−0.37/+0.22	0.37	1.5
GaAs + ZnSe + KRS5	+0.68/−0.71/+0.33	0.25	1.8
[AMTIR1 + ZnSe + KRS5]	+1.19/−0.72/+0.16	0.39	1.9
[CsI + NaCl + GaAs]	+0.68/−0.32/+0.29	0.38	1.1

{ } Very high transmission.

[] Low residual high-order chromatic aberration.

Abbe number (V value) in the usual chromatic aberration equations. Thermal expansion of the housing—obviously not present in chromatic calculations—does, however, complicate the situation a little.

In equations thus far, Δf has meant both thermal defocus and focal length change, as numerically these are the same for a thin lens. For separated components, rules similar to those for chromatic aberration apply, for example, two separated thin-lens groups—such as those described by Fig. 1—must be individually athermalized if both types of thermal “aberration” are to be corrected simultaneously. More complex optics (for example, multistage) may have transfer of thermal aberration between constituent lens groups but may still be corrected simultaneously for thermal focus shift and focal length change as a whole. This procedure can, however, lead to one lens group requiring excessive optical powers in order to achieve full overall correction—transient longitudinal thermal gradients may also cause problems.

Use of Diffractive Optics in Optical Athermalization

The term “hybrid optic” is generally used to signify a combination of refractive and diffractive means in an optical element. The diffractive part of the hybrid is usually a transmission hologram which for high efficiency would be of surface relief form, the surface structure being machined or etched onto the refractive surface.⁴⁵ The diffractive surface acts as a powered diffraction grating, producing large amounts of chromatic aberration which could be employed in an optic where a lightweight optically athermalized combination of two materials could be chosen without regard to achromatism: residual chromatic aberration could then be corrected by the hologram.⁴⁶

8.9 REFERENCES

1. J. W. Perry, *Proc. Phys. Soc.*, vol. 55, 1943, pp. 257–285.
2. J. Johnson and J. H. Jeffree, U.K. Patent No. 561 503, 1942 U.K. priority.
3. D. S. Grey, *J. Opt. Soc. Am.*, vol. 38, 1948, pp. 542–546.
4. D. S. Volosov, *Opt. Spectrosc.*, U.S.S.R., vol. 4, pp. 663–669 and pp. 772–778, vol. 5, 1958, pp. 191–199.
5. R. Penndorf, *J. Opt. Soc. Am.*, vol. 47, 1957, pp. 176–182.
6. H. Köhler and F. Strähle, *Space Optics*, B. J. Thompson and R. R. Shannon (eds.), National Academy of Sciences, 1974, pp. 116–153.
7. P. J. Rogers, *SPIE Critical Reviews*, vol. CR38, 1991, pp. 69–94.
8. L. R. Estelle, *SPIE*, vol. 237, 1980, pp. 392–401.

9. P. J. Rogers, *SPIE*, vol. 1354, 1990, pp. 742–751.
10. M. Laikin, *Lens Design*, Marcel Dekker, New York, 1991, p. 28.
11. G. G. Slyusarev, *Opt. Spectrosc.*, U.S.S.R., vol. 6, 1959, pp. 134–138.
12. W. H. Turner, *Optical Sciences Center*, vol. 4, University of Arizona, Tucson, Arizona, 1970, pp. 123–125.
13. V. M. Mit'kin and O. S. Shchavlev, *Sov. J. Opt. Tech.*, vol. 40, 1973, pp. 558–561.
14. F. Reitmayer and H. Schroeder, *Appl. Opt.*, vol. 14, 1975, pp. 716–720.
15. P. J. Rogers, *SPIE*, vol. 147, 1978, pp. 141–148; and U.K. Patent No. 2,030,315.
16. D. G. Norrie, *Opt. Eng.*, vol. 25, 1986, pp. 319–322.
17. J. Angénieux et al., *SPIE*, vol. 399, 1983, pp. 446–448.
18. D. S. Garcia-Núñez and D. Michika, *SPIE*, vol. 1049, 1989, pp. 82–85.
19. V. Povey, *SPIE*, vol. 655, 1986, pp. 142–153.
20. A. D. Michael and W. B. Hart, *Metallurgist and Materials Technologist*, August 1980, pp. 434–440.
21. G. V. Thompson, U.S. Patent No. 4,148,548, 1976 U.K. priority.
22. P. J. Rogers and G. N. Andrews, *SPIE*, vol. 99, 1976, pp. 163–175.
23. P. M. Parr-Burman and A. Gardam, *SPIE*, vol. 590, 1985, pp. 11–17.
24. I. A. Neil and W. McCreath, U.K. Patent No. 2,141,260, 1983 U.K. priority.
25. I. A. Neil, U.S. Patent No. 4,659,171, 1984 U.K. priority.
26. I. A. Neil and M. Y. Turnbull, *SPIE*, vol. 590, 1985, pp. 18–29.
27. P. Nory, *SPIE*, vol. 590, 1985, pp. 30–39.
28. P. M. Parr-Burman and P. Madgwick, *SPIE*, vol. 1013, 1988, pp. 92–99.
29. R. C. Simmons and P. A. Blaine, *SPIE*, vol. 916, 1988, pp. 19–26.
30. M. Shechterman, *SPIE*, vol. 1442, 1990, pp. 276–285.
31. M. Roberts and P. R. Crew, U.K. Patent No. 2,201,011, 1987 U.K. priority.
32. M. Roberts, *SPIE*, vol. 1049, 1989, pp. 72–81.
33. C. B. Estes, U.S. Patent No. 3,205,774, 1961 U.S. priority.
34. R. C. Gibbons, *ERIM Report No. 120200-I-X*, 1976, p. 71.
35. K. Straw, *SPIE*, vol. 237, 1980, pp. 386–391.
36. M. O. Lidwell, U.S. Patent No. 4,494,819, 1980 U.K. priority.
37. T. H. Jamieson, *Opt. Eng.*, vol. 20, 1981, pp. 156–160.
38. I. A. Neil, U.S. Patent No. 4,505,535, 1982 U.K. priority.
39. M. Roberts, U.S. Patent No. 4,679,891, 1984 U.K. priority.
40. P. J. Rogers, Thermal Imaging course notes, Institute of Optics, Summer School on Lens Design, Rochester, 1988 (unpublished).
41. J. L. Rayces and L. Lebich, *SPIE*, vol. 1354, 1990, pp. 752–759.
42. M. Yatsu et al., *SPIE*, vol. 1354, 1990, pp. 663–668.
43. P. J. Rogers, *SPIE*, vols. 1780 and 1781, 1992, pp. 36–48.
44. J. M. Palmer, U.K. Patent No. 2,194,072, 1986 U.K. priority.
45. G. J. Swanson and W. B. Veldkamp, *Opt. Eng.*, vol. 28, 1989, pp. 605–608.
46. P. J. Rogers, *SPIE*, vol. 1573, 1992, pp. 13–18.

PART

2

FABRICATION

This page intentionally left blank.

Michael P. Mandina

*Brandon Light
Optimax Systems, Inc.
Ontario, New York*

9.1 INTRODUCTION

The novel creations of optical designers have been limited by the fabricator's ability to manufacture and measure the elements of the optical prescription. A solution to a design criteria often existed only on paper as the required elements were not physically realizable. Optics manufacturing technology innovations continually expand the possibilities for optical components. Increasingly, manufacturing is tethered to metrology. Creation of optics metrology instruments with accuracy equal to that of optics manufacturing equipment and vice versa has driven process development. It is this developmental symbiosis that has brought determinism to the art of precision optics manufacturing. Metrology and machine innovations offer optics of higher quality and complexity in predictable timeframes. The requirement for skilled technicians is still vital in the manufacturing process; however, the skill set is increasingly one of craft in combination with science. Artisan opticians of yesteryear still provide value; however, the future of optics manufacturing is in the hands of the 21st century optics technicians.

The methods described below are the most common for typical optical components used in industrial, aerospace, and defense applications. For the spherical lens section, a brief overview of the traditional process is described first and then the latest methods. The remaining sections will provide general overviews. The focus will be exclusively on brittle materials. For our purposes, brittle materials are defined as those where the removal process is achieved by applying mechanical forces that fracture the surface, releasing fragmented particles in a controlled manner. Much has been documented on fine finishing of brittle materials such as optical glasses, ceramics, and crystals. Works by Preston¹, Silvernail², Izumitani³, Buijs⁴, Bach⁵, Kaller⁶, Lambropoulos⁷, Golini⁸, Cook⁹, Jacobs¹⁰ and DeGroot¹¹ have contributed greatly to the understanding of optics finishing processes.

9.2 MATERIAL FORMS OF SUPPLY

Optical glass is available in boule, slab, and gob forms. Boules are formed in special disposable pots that yield a batch or glass melt of a specific glass type, such as the borosilicate glass BK7, but whose detailed characteristics are unique to that batch. Slab is yielded from a continuous flow process.

Materials are homogeneously mixed and heated, and a continuous ribbon of glass is produced. These ribbons are cut into slabs that are generally 250 mm wide, 25 mm thick, and 350 mm long, although sizes vary significantly based on supplier and material. Gobs are also yielded from a continuous flow process; however, the molten glass flows through an orifice and is sliced like cookie dough at a predetermined frequency that ensures the desired volume for the application. Gobs are almost always made to customer specifications for glass type and volume. They are used as the preblank to produce near net shape molded blanks for high-volume lens systems. Many glass suppliers also provide polished preforms, usually balls. These are used for glass molding finished optics components.

Manufacturers will order the form that best suits their purpose. The closer to final form the material, the less waste and time consumed in bulk removal operations. When rough shaping material blanks from boule or slab forms, most manufacturers use diamond impregnated saw blades and core drills to yield a part appropriate to yield the finished optic, generally called disks. Typically blanks or disks are several millimeters oversized from the final part's critical dimensions.

9.3 BASIC STEPS IN SPHERICAL OPTICS FABRICATION

Generating

This is a bulk material removal operation that starts with a near net shape molded blank or a disk.

Generating—Traditional The removal is accomplished through the application of diamonds embedded in a matrix on the cutting surface of a cup-shaped ring tool. The material is ground away as the diamonds create cracks in the surface, sweeping away glass particles where the fractures intersect.¹² The machine accuracy is generally akin to manually set, mechanically based control production equipment used in the machine tool industry. The operator continually monitors results and modifies machine settings as the cupped ring tools wear. Machine precision is adequate to control thicknesses to ± 0.025 mm and radius to ± 0.010 -mm sagittal height, but the machine is only capable of coarse removal. Subsequent lapping operations are required in order to reduce subsurface damage¹³ to a level where polishing is possible.

Generating—Modern The advent of deterministic microgrinding processes spearheaded by the work at the University of Rochester's Center for Optics Manufacturing¹⁴ in the 1990s shifted the paradigm for finishing expectations from the generating operation. As a result, machines used in the generating operation have evolved to precision machine tool status. Generators manufactured by mainstream optics manufacturing equipment providers such as OptiPro,¹⁵ Satisloh,¹⁶ Schneider¹⁷ and others, have created grinding solutions that enable the generating operation to predictably yield surfaces that are ready for polishing operations. This modern equipment makes use of CNC (computer numeric control) systems, robust motion and motion control systems such as precision linear ball slides, advanced machine base materials, structure design, and improved positioning feedback through optical encoders and other submicron feedback systems. Additionally, most of the machine builders provide in situ metrology options that enable operator assisted or completely automated parameter adjustment optimization. This is an important feature as the tool consumes itself during the process.

Complementing the advent of advanced generating machine tools for optics generating has been the increased understanding of fixed abrasive grinding mechanisms. Deterministic microgrinding is typically preferred to loose abrasive lapping when fabricators have a choice. The residual damage from microgrinding can be estimated based on glass properties.¹⁸ This aids in determining prepolish finish requirements so the overall process time for the optics can be optimized. Even with recent advancements in understanding the microgrinding process, the industry is far from offering a directory of ring tools optimized for the array of optical materials. Therefore the industry continues to rely heavily on empirical results to determine optimal setups.

Lapping

This process reduces subsurface damage left from generating to a manageable level in preparation for polishing.

Lapping—Traditional Lapping is the application of loose abrasive particles applied as slurry and pressed into the work surface by nominally constant applied pressure.¹⁹ The process typically consists of applying the abrasive slurry between a cast-iron-rotating lapping tool and the optic. Both surfaces abrade away as they remain in random dynamic contact. The fabricator controls the material removal so the operation yields the desired surface radius and smoothness. The abrasive material, often aluminum oxide, is typically between 30- and 5- μm particle size. The operator steps through particle sizes, using progressively smaller abrasives. Removal amounts account for the subsurface damage from the prior generating or lapping operation, ideally completely removing it.

Lapping—Modern The use of diamond particles embedded in a resin or metal matrix have been popular for some time. Initially, these matrices were fabricated in pellets, fastened as desired on metal backing plates, and used as laps. Abrasive work is done by the diamonds and coolant serves as lubricant and carries the glass particulate away. Unlike loose abrasive lapping, the slurry is not the abrasive. Diamond tool manufacturers also make diamond-sheet material for ready application to tools, and more recent products such as resin-bonded sheet materials from abrasive manufacturers such as 3M²⁰ can be used the same way.

Polishing

Polishing converts the finely fractured surface from the lapping or deterministic grinding operation [typical roughness of about 1- μm rms (root mean square)] into a specular surface of a surface roughness typically 1 to 3 nm rms. Polishing is a chemical-mechanical process. Water attacks the surface creating a chemically softened layer, and then the mechanical action of the abrasive in the polishing slurry, usually ceria based for optical glasses, removes the chemically softened outer layer of glass.³

Polishing—Traditional The polishing process is expected to remove the damage left from preceding operations, typically 5 to 20 μm of material. The intimate contact between the polishing tool and the optic, working with the slurry, slowly enhances the surface finish. The process is feedback based, and the fabricator works the part for a while and checks the outcome. Reacting to the results, the experienced fabricator controls various parameters to yield the desired form and finish of the polished surface.

The most basic polishing tool is a pitch polisher. Optical polishing pitch is a viscoelastic material. To form a pitch polisher, a metal tool of proper radius is coated with a 4 to 5-mm layer of polishing pitch. The pitch is warmed and formed to the optic. Once cool the brittle pitch will be cut to allow irrigation grooves for slurry access. When performed by artisans, pitch polishing can yield form errors equal to fractional wavelength of visible light routinely. Less capable fabricators may be limited to commercial quality, multiple wavelength form error outcomes.

By the 1980s, high-speed polishing had become very popular. One of the key innovations was the use of polyurethane polishing pads as a replacement for pitch. Polyurethane pads are a viscoelastic thermoplastic material with a higher viscosity than pitch. Polyurethane remains a polishing material staple and is the polishing material of choice for the fast removal seen in high-volume optics manufacturing.

Polishing—Modern Advances in deterministic polishing are dramatically changing the demands placed on optics manufacturers. Deterministic polishing is a feed-forward process, where the outcome is reasonably certain. Industry leaders in deterministic polishing development are QED/Schneider consisting of Magnetorheological Finishing (MRF)²¹ and Zeeko/SATISLOH²² who promote

a precessions polishing and air bladder solution. Each has created opportunities for manufacturers to produce optics at more predictable cost. Their application of CNC machining systems to the polishing process is revolutionizing the precision optics industry.²³

All these new solutions rely on subaperture small “pad” polishing with a known removal rate, where the “pad” may be in the form of variable stiffness polishing fluid or compliant tool made from a variety of materials and consistencies. Originally used to finish large astronomical telescope optics, small pad methods have advanced in recent years to scale cost and size down so these technologies are available to the broader population of precision optics manufactures.

Deterministic subaperture polishing solutions combine a tool’s known removal rate with an error map of the optic to produce a removal schedule. This feed-forward process relies completely on the accuracy of actual surface form information. In most cases this information is acquired from a variety of instruments such as coordinate measurement machines (CMM) or surface profilometers like the Taylor-Hobson Form TalySurf.²⁴ These instruments themselves or the software of the polishing tool convert points of data into an error map for a continuous surface. The removal profile dictates the dwell time for the small aperture polishing pad, and in general form error decreases by a factor of five per iteration. For example, if the form error is 1 wave, it is reasonable to expect that after a deterministic polishing iteration the form error will be $\sim 1/5$ wave, and after another iteration would be $\sim 1/25$ wave.

Newer technology that is also under development at a number of equipment manufacturers including QED and ZEEKO²⁵ incorporates fluid jet technology. Surfaces are corrected by directing a jet of abrasive/fluid mixture at a surface, the flow generates sufficient surface shear stress that chemical-mechanical polishing occurs.^{26,27} The jet-polishing technology is especially promising for difficult to reach areas seen in asphere and conformal optical surfacing.

Edging

Most applications of lenses require mounting into a lens housing. Lens system performance is maximized when the centers of curvature reside on the cylindrical axis of the housing. The edging operation simultaneously creates a precise (± 0.025 mm or less) diameter for mounting and aligns the centers of curvature on the mechanical centerline of the lens.

Edging—Traditional²⁸ Earlier pitch-based methods consisted of using a precision spindle where a brass cup was trued using a cutting tool. This was basically a lathe-type operation and required a skilled combination of heat, pitch, spindle velocity, timing, and consistent axial force applied by a skilled artisan in order to set the lens in a “trued” position. Once the lens was blocked, a diamond wheel ground the diameter to final dimension. Lenses are typically brought to final polished state with the diameter of the lens 1- to 3-mm oversized.

This pitch method was almost entirely replaced with mechanical bell-clamping edging machines. Bell clamping employs two opposed coaxial synchronized precision spindles and is a pitch-free process. Each spindle is affixed with a precision cup of the appropriate size to capture the lens and allow auto alignment by virtue of the mechanical forces on the variably sloped surfaces. Once the lens is “clamped” into true position, an operator mechanically defines and initiates an automated grinding sequence.

Edging—Modern In recent years, the use of CNC edging equipment is enabling a single setup for multiple grinding operations. For example, it is fairly routine to process the diameter, sagitta with a step, bevel and a fiduciary flat, all in one operation. The CNC controller interface shows a series of cross-sections, and the operator fills in inputs for what is the starting point and what features are needed in the end. Simultaneous creation ensures the features will all run true relative to one another. In addition to facilitating grinding of more complex features, optional features such as micropositioning air blasts for automated alignment optimization and measurement enable precision placement of the optic. Lenses are still mechanically bell clamped.

9.4 PLANO OPTICS FABRICATION

A plano surface has a radius equal to infinity. Typically plano form specification does not differentiate between spherical power and irregularity, specifying lump sum reflected errors as flatness. Therefore, maintaining perfect flatness is critical during plano surface finishing. The process steps for plano surfaces are exactly the same as for spheres. Planos have the advantage of fixed radius, so often, companies, departments within companies, personnel and equipment will be plano specific. This specificity allows economies of scale and development of plano-specific solutions. An example of this are continuous polishers (CPs), in which a large (40–60 inches in diameter) annular lap is “conditioned” to maintain lap flatness independent of the workpiece size. The lap is forced by a large glass (or similar material) “conditioner” to stay flat. This persuasion by the conditioner imprints onto the work piece and maintains flatness as a result.

Double-sided CPs polish both sides of a window simultaneously. Much of the recent technology used in the precision plano window manufacturing has been taken from semiconductor industry’s work-optimizing silicon wafer processes.

9.5 ASPHERE OPTICS FABRICATION

Aspheric lenses contain at least one optically active surface of nonconstant curvature. This is the primary differentiator from a spherical lens. Rotationally symmetric aspheric lenses are solids of revolution, where a general equation describes the cross section to be revolved (Fig. 1). Lenses of this style are capable of higher aberration order correction than spherical lenses. While the forms and their promise have been known to optical designers for centuries, for most of that time only the mildest forms have been physically realizable. The methods, machinery, and metrology are specific to asphere manufacturing.¹⁴

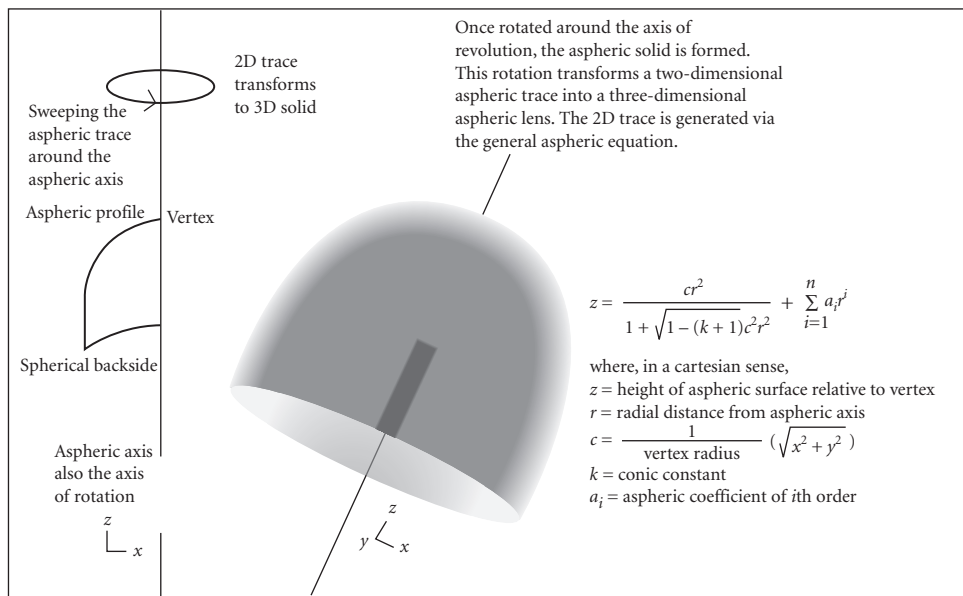


FIGURE 1 Sample general asphere form. (Brandon Light, Optimax Systems, Inc.)

Traditional full-aperture fabrication methods are not capable of manufacturing aspheric surfaces due to their nonconstant curvature. By changing the amount of contact from full to a region where change in local curvature approaches zero, some portions of traditional spherical lens-manufacturing techniques can be applied. Brittle removal by high-speed diamond grinding followed by ductile removal using a polishing slurry (ceria, alumina, etc.) can be used to prepare aspheric surfaces. Instead of full contact, curvature-insensitive local contact is used in grinding and polishing.

In aspheric grinding, a peripheral diamond wheel on a CNC platform traces the surface to generate the aspheric profile. In grinding, machine accuracy determines profile accuracy. A more accurate ground profile makes a more accurate polished profile more likely, since there's less correction needed. Particular attention must be paid to wheel wear, wheel balance, positional accuracy, and overall stiffness of the grinding platform. Imperfection in any of these grinding parameters will leave signatures in the ground surface.

The surface is then polished by working only a small area at a time. All of this must be done while maintaining location of the aspheric axis, the axis around which the solid of revolution was formed. Each iteration has an error inducement associated with it, so making as few correction runs as possible is a primary focus. Typically, asphere polishing is a feed-forward, deterministic process. While the local curvature may be constant, globally it is not. Polishing requires an adaptive tool and knowledge of what's ahead. The polishing tool needs to change to suit local curvature at a suitable rate of change. This requires knowledge of how the tool will evolve and how much removal is needed in which region. Deterministic processes provided by Zeeko/Satisloh and QED machinery, discussed earlier, are examples of such tools. These processes characterize the removal rate as a function of curvature for a given tool and combine that with an error map of the surface to be worked. The resulting removal schedule accommodates for volumes to be removed and tool performance at that local curvature.

Conventional interferometric techniques do not translate to aspheric manufacturing either. Since local curvature is nonconstant, interferometric techniques for aspheres are lock and key. The setup and equipment can be unique for a given aspheric form, so time and money demands are large. For example, form-specific computer-generated holograms (CGHs) may be required to provide feedback to the closed-loop deterministic polishing process. For a more cost-effective solution profilometry is the main two-dimensional compromise, and it is the current industry standard. Although more generalized interferometric solutions are beginning to be offered by QED, Zygo, and others.

Errors in centration are unrecoverable. In centering a spherical lens, errors can be removed. With sufficient diameter overage both centers of curvature could be positioned on the same axis and that axis could be made concurrent and coincidental with the mechanical axis. Since an aspheric surface is centered about an axis and not a point such realignment is not possible. Therefore, centration must be conserved throughout processing.

9.6 CRYSTALLINE OPTICS

As more optical work occurs outside the visible spectrum, use of optics made from nonglass brittle materials will grow. Single crystalline and polycrystalline materials are transparent far outside the usual spectral transmission range of glass. In many cases, the surfaces of these materials have differing hardness values depending on the orientation of the crystal boundaries. Soft laps tend to accentuate the grain boundaries of these materials, and that can lead to wavefront errors, mottled surfaces, and scattering. The traditional optical fabrication process can be adapted to crystal materials if some substitutions are made. The lapping process may substitute finely graded diamond for alumina and tin or zinc laps in place of the typical cast iron. Similarly, diamond suspensions are often used in polishing in place of ceria. Polishing laps may consist of synthetic materials like polyurethane or beeswax instead of optical pitches.^{29,30}

9.7 PURCHASING OPTICS

There are a number of companies who offer lines of standard optical components. These suppliers can provide off-the-shelf optics in a variety of sizes, shapes, and quality levels. Most optics providers have areas of specialization, and the informed optics buyer will select vendors that match their specific optics requirements. When custom optics are required, it is best to understand the capabilities of prospective suppliers. Most optics companies promote a broad range of capabilities, but many tend to specialize in some manner. Professionals who are engaged in optics purchasing on a regular basis learn where to go for specific optics requirements. Often this education is paid for by awarding of numerous contracts across a broad array of parts and suppliers and experiencing the consequences of the decisions. Much is learned in the contract's postmortem review.

Optics purchasing is further complicated with the predominance of the internet as a research tool. Web sites and promotional materials often do not reflect a supplier's true capability and know-how.

Whether buying off-the-shelf or custom optics, it is always best to engage potential suppliers in dialog, preferably addressing tolerances and other manufacturing cost drivers. The buyer should be satisfied the supplier has the ability to meet and measure all critical criteria. For optics that approach a manufacturer's limits, it is especially important to understand the test and acceptance process, as there can be quite a divergence of metrology equipment and methods available for testing various parameters.³¹ This is especially true for aspheres, where full format phase measuring interferometry or transmitted wavefront testing is not always within a supplier's capability.

9.8 CONCLUSION

Optics fabrication requires serial application of relatively simple steps. In the past, these steps were carefully carried out by artisans using traditional techniques. Modern approaches incorporate scientific research into the manufacturing process. Artisan skills integrate with the scientific know-how yielding a new breed of technology workers of the twenty-first century. Nevertheless, the basic process steps of grinding followed by polishing have remained. Introduction of new optical materials, more complex shapes and more narrow tolerance budgets will enable designers to develop improved solutions to old problems over an expanded spectrum, and modern manufacturing methods can make the optics physically realizable. Finally, the tendency for specialization among optics supplier requires open dialog between supplier and designer as a means to optimize successful relationships.

9.9 REFERENCES

1. F. W. Preston, "Structure of Abraded Glass Surfaces," *Trans. Opt. Soc.* **23**:141 (1922).
2. W. W. Silvernail, "Role of Cerium Oxide in Glass Polishing," in *The Science of Polishing*, Duncan Moore (ed.), OSA: Washington, D.C., 1984.
3. T. S. Izumitani, *Optical Glass*, American Institute of Physics: New York, 1986.
4. M. Buijs and K. Korpel-van Houten, "A Model for Lapping of Glass," *J. Mater. Sci.* **28**(11):3014–3020 (1993).
5. H. Bach, "Analysis of Subsurface Layers and Spots and the Reactivity of Glass Components," in *The Science of Polishing*, Duncan Moore (ed.), OSA: Washington, D.C., 1984.
6. A. Kaller, "Properties of Polishing Media for Polishing Optics," *Glastechnische Berichte—Glass Sci. Technol.* **71**(6):174–183 (1998).
7. J. C. Lambropoulos, S. Xu, and T. Fang, "Loose Abrasive Lapping Hardness of Optical Glasses and Its Interpretation," *Appl. Opt.* **36**(7):1501–1516 (1997).
8. D. Golini and S. D. Jacobs, "The Physics of Loose Abrasive Microgrinding," *Appl. Opt.* **30**:2761–2777 (1991).
9. L. M. Cook, "Chemical Processes in Glass Polishing," *J. Non-Cryst. Solids* **120**:152–171 (1990).

10. S. D. Jacobs, D. Golini, Y. Hsu, et al., "Magnetorheological Finishing: A Deterministic Process for Optics Manufacturing," *Opt. Fabrication and Testing*, T. Kasai (ed.), *SPIE* 2576:372–383 (1995).
11. J. E. DeGroote, A. E. Marino, J. P. Wilson, et al., "Removal Rate Model for Magnetorheological Finishing of Glass," *Appl. Opt.* 46:7927–7941 (2007).
12. R. E. Parks, "Optical Fabrication," in *Handbook of Optics*, 2d ed., M. Bass, E. W. Van Stryland, D. R. Williams, and W. L. Wolfe (eds.), McGraw-Hill: New York, 1995, Vol. 1. Chap. 41.
13. P. Hed and D. F. Edwards, "Relationship between Subsurface Damage Depth and Surface Roughness during Grinding of Optical Glass with Diamond Tools," *Appl. Opt.* 26(13):2491 (1987).
14. H. Pollicove and D. Golini, "Computer Numerically Controlled Fabrication," Chap. 7: *International Trends in Applied Optics*, A. H. Guenther (ed.), *SPIE*: Bellingham, Wash., Vol. PM119 (2002).
15. www.Optipro.com, Optipro Systems, Optical Fabrication Equipment, April 21, 2009.
16. www.Satisloh.com, Optical Manufacturing Solutions, Products, Precision Optics, April 21, 2009.
17. www.schneider-om.com/home.html, Schneider GmbH & Co. KG—Fascination for Innovation: Products, April 21, 2009.
18. J. C. Lambropoulos, "Surface Microroughness of Optical Glasses under Deterministic Microgrinding," *Appl. Opt.* 35:4448–4462.
19. J. C. Lambropoulos, "Using the Grinding Merit Function (GMF): What Quality of Grind Can You Expect in the Shop?" *Convergence, Newsletter of the Center for Optics Manufacturing*, Sept./Oct. 1998.
20. www.3M.com, April 21, 2009.
21. D. Golini, G. Schneider, P. Flug, M. Demarco, "Magnetorheological Finishing," *Optics and Photonics News*, October 2001, pp. 20–24.
22. D. D. Walker, D. Brooks, A. King, et al., "The 'Precessions' Tooling for Polishing and Figuring Flat, Spherical and Aspheric Surfaces," *OSA*, 21 April, 2003; *Opt. Exp.* 11(8):958–964.
23. S. D. Jacobs, "Innovations in Polishing of Precision Optics," *Convergence, Newsletter of the Center for Optics Manufacturing*, 1st/2nd qtr. 2003.
24. www.taylor-hobson.com, Taylor Hobson—Surface Profilers, April 21, 2009.
25. www.zeeko.co.uk, Zeeko Ltd. Ultra-Precision Polishing Solutions for Optics and Other Complex Surfaces, April 21, 2009.
26. W. I. Kordonsk, A. Shorey, and M. Tricard, "Magnetorheological Jet (MRJet™) Finishing Technology," *ASME*, 128:20–26, Jan. 2006.
27. S. M. Booij, O. W. Föhnle, and J. J. M. Braat, "Shaping with Fluid Jet Polishing by Footprint Optimization," *Appl. Opt.* 43:67–69.
28. D. F. Horne, *Optical Production Technology*, Crane, Russack & Co.: New York, 1988.
29. G. W. Fynn and W. J. A. Powell, *Cutting and Polishing Optical and Electronic Materials*, 2d ed., Adam Hilger: Philadelphia, Pa., 1988.
30. R. Sumner, "Polishing of IR Materials," in *The Infrared Handbook*, W. Wolfe and G. Zissis (eds.), ERIM: Ann Arbor, Mich., 1978.
31. Y. A. Carts, "How to Buy Custom Optics That Meet Your Specifications," *Laser Focus World*, Aug. 1992, pp. 91–100.

FABRICATION OF OPTICS BY DIAMOND TURNING*†

Richard L. Rhorer

*National Institute of Standards and Technology
Gaithersburg, Maryland*

Chris J. Evans

*Zygo Corporation
Middlefield, Connecticut*

10.1 GLOSSARY

f	feed rate
h	peak-to-valley height
R	tip radius of diamond tool

10.2 INTRODUCTION

The use of special machine tools with single-crystal diamond-cutting tools to produce optical surfaces on some metals and a limited range of other materials is called *diamond turning*. Over the last 50 years or so, diamond turning has matured to become the method of choice for producing some optical surfaces; in other applications, diamond turning provides a critical process step with radically different characteristics from most other optical fabrication methods.

In terms of geometry and motions required, the diamond-turning process is much like the step of “generating the optical surface” in traditional optical fabrication. However, the diamond-turning machine is a more sophisticated piece of equipment that produces the final surface, which frequently does not need the traditional polishing operation. The surface quality produced by the “best” diamond turning does not yet match the best produced by conventional polishing practice. The limits of diamond turning for both figure and surface-finish accuracy have not yet been reached—and diamond turning can be combined with postpolishing to improve surface finish and reduce scatter.¹ Also subaperture processing with small polishing tools or magnetorheological finishing (MRF) can be used to improve figure.

There are several important advantages of using diamond turning, including the ability to produce good optical surfaces to the edge of the element, to fabricate soft ductile materials that are difficult to polish, to eliminate alignment adjustments in some systems, and to fabricate shapes difficult to produce by other methods. The latest generation of diamond-turning machines incorporates up to five axes

*Official contribution of the National Institute of Standards and Technology; not subject to copyright in the United States.

†Certain commercial equipment, instruments, or materials are identified in this chapter. Such identification does not imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products identified are necessarily the best available for the purpose.

of computer-controlled motion, allowing for production of anamorphic optics. Use of form tools on multiaxis machines enables production of “structured” optical surfaces² ranging from subwavelength structures through diffractive/refractive infrared (IR) elements to optical component molds.

If the advantages of diamond turning suggest this fabrication method, then it is important to determine early in the design phase of a project whether the material specified is appropriate for diamond turning and whether slideway travels and linear and rotary axis controls are available on the diamond-turning machine to support fabrication of complex structures.

Sections in this chapter highlight the following:

- The diamond-turning process
- The advantages of diamond turning
- Diamond-turnable materials
- Comparison of diamond turning and traditional optical fabrication
- Machine tools for diamond turning
- Basic steps in diamond turning
- Surface finish of diamond-turned optics
- Metrology of diamond-turned optics
- Conclusions

10.3 THE DIAMOND-TURNING PROCESS

The diamond-turning process produces finished surfaces by very accurately cutting away a thin chip or layer of the surface. Thus, it is generally applicable to ductile materials that machine well rather than to hard brittle materials traditionally used for optical elements. However, by using a grinding head on a diamond-turning machine in place of the tool, hard brittle materials can be finished. At very small effective depths of cut, brittle materials behave in an apparently ductile manner. This attribute allows fracture-free grinding of glasses and ceramics as well as diamond turning of optical surfaces on materials such as germanium, zinc selenide, and potassium dihydrogen phosphate (KDP).

In diamond turning, both the figure and surface finish are largely determined by the machine tool and the cutting process. Note, however, that material characteristics such as grain size and inclusion size limit the ultimate surface finish achievable. The tool has to be very accurately moved with respect to the optical element to generate a good optical surface, and the edge of the diamond tool has to be extremely sharp and free of defects.

10.4 THE ADVANTAGES OF DIAMOND TURNING

Diamond turning fits within a broad spectrum of optics fabrication processes. When compared with traditional optical fabrication methods of lapping and polishing (see, for example, Chap. 9, “Optical Fabrication,” by Michael P. Mandina) diamond turning has several advantages.

- It can produce good optical surfaces clear to the edge of the optical element. This is important, for example, in making scanners, polygons, special shaped flats, and when producing parts with interrupted cuts.
- It can produce optical surfaces on soft ductile materials that are extremely difficult to polish.
- It can easily produce off-axis parabolas and other difficult-to-lap aspherical shapes.
- It can produce optical elements with a significant cost advantage over conventional lapping and polishing where the relationship of the mounting surface—or other feature—to the optical surface is very critical. Expressed differently, this feature of diamond turning offers the opportunity to eliminate alignment adjustments in some systems.

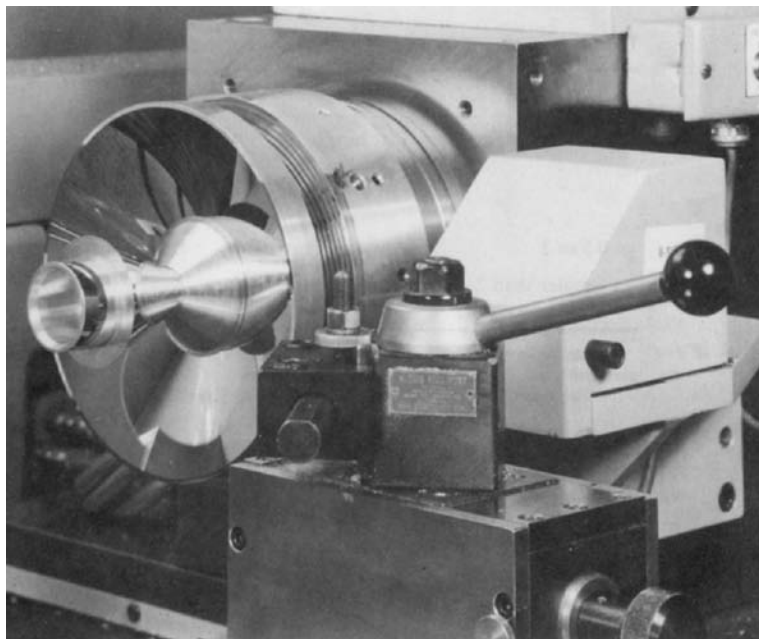


FIGURE 1 An axicon optical element being diamond turned. (Courtesy of Rank Taylor Hobson, Keene, New Hampshire.)

- It can fabricate optical shapes such as axicons, faceted optics, and grazing incidence X-ray optics that would be extremely difficult to fabricate by methods other than diamond turning (see Fig. 1).

Conflicts between optical requirements and diamond turnability on the one hand, and mechanical considerations on the other, often lead to the use of platings. Plating deficiencies, however, can cause as much trouble as poor bulk materials. For example, small changes in the composition of plated electroless nickel may cause dramatic changes in tool wear.³

Residual stress in the mirror blank, whether plated or not, can lead to changes in mirror shape with time. It is essential to pay careful attention to stress-relief prior to final diamond turning.

A decision to diamond turn an optical element, rather than fabricate it by the conventional polishing techniques, might be based on several different considerations such as type of element, size, and material. A general guide to different considerations in selecting diamond turning as a fabrication technique is presented in Table 1.

TABLE 1 General Guide to Optical Fabrication Methods

Size, m	Shape	Material	Preferred Method
Less than 0.5	Flat or sphere	Glass/ceramic	Polish
		Ductile metal	Diamond turn
	Asphere	Glass/ceramic	Grind/polish*
		Ductile metal	Diamond turn
0.5 to 2.0	Any axisymmetric	Ductile metal	Diamond turn [†]
Greater than 2.0	Any	Any	Large polishing machine

*Can generate shape or figure on a diamond-turning machine with a grinding head replacing the diamond tool. Subaperture polishing techniques, including techniques such as MRF, may be applied to advantage.

[†]Diamond-turning machines up to 2-m diameter have been built.

As indicated above, diamond turning has some unique characteristics. In some IR (and even visible) imaging systems, considerable improvements in optical performance have been obtained by combining a refractive aspheric surface and a diffractive surface in a single element. For IR applications, it is hard to produce such a component by any other fabrication process; for visible applications, such optics can be produced in volume from diamond-turned molds.

Another unique capability of diamond turning is to provide datums or alignment features machined in the same setup as the optical surface. “Snap-together” optical systems requiring no alignment adjustments after assembly are very attractive in some applications.

Over the last decade, there have been considerable advances in the ability to produce aspheric optics using computer-controlled generators and pad polishers. These technologies, combined with ion polishing, magnetorheological finishing, and computer-controlled polishing have enabled a new generation of aspheric optics. Ultimately the choice of manufacturing process requires a careful analysis of the options and the system requirements.

10.5 DIAMOND-TURNABLE MATERIALS

Selection of appropriate materials is, necessarily, a trade-off between application-specific requirements and optimization of the manufacturing process. This trade-off may drive the selection of a plated surface, for example, or the choice of fabrication steps.

Historically materials have been described as either “diamond turnable,” or not, as if this were an inherent material property. This shorthand covers two different situations. One is that, in practice, some materials cause very rapid wear of the diamond; for example, it is widely known that ferrous materials cause rapid tool wear. The other is that, particularly for certain plastics, tool-workpiece interactions produce unacceptable optical surfaces.

A number of listings of diamond-turnable materials, such as the one included in Table 2, have been published. Such listings should be treated with caution. Typically, they are incomplete and do not provide sufficient information on the materials that are listed. For example, good optical surfaces are not generally produced on all aluminum alloys: Aluminum Alloy 6061 (Aluminum Association, Inc. designation) is the most commonly used alloy, although certain 5000 series and 7000 series alloys have their proponents, and 2024 aluminum has been used but, in general, does not produce the best surfaces.

TABLE 2 Diamond-Turnable Materials

Metals	Nonmetals	Plastics
Aluminum	Calcium fluoride	Polymethmethacrylate
Brass	Magnesium fluoride	Polycarbonates
Copper	Cadmium telluride	Polyimide
Beryllium copper	Zinc selenide	
Bronze	Zinc sulphide	
Gold	Gallium arsenide	
Silver	Sodium chloride	
Lead	Calcium chloride	
Platinum	Germanium	
Tin	Strontium fluoride	
Zinc	Sodium fluoride	
Electroless nickel	KDP	
	KTP	
	Silicon	

Similarly, gold is considered diamond-turnable, but problems have been reported machining large gold-plated optics. Conventional electroplated nickels (and bulk nickel) give rapid tool wear, but electroless nickel with phosphorous contents above about 10 percent, if appropriately heat treated, can be machined effectively.⁴ Higher phosphorous contents (up to 15 percent) are obtainable in electroplated nickel^{5,6} which also machines extremely well. Both materials are metastable and will transform—with exposure to the necessary time/temperature conditions—to a mixture of crystalline nickel with hard phosphides. This transformation is accompanied by a volume change, a degradation of optical characteristics of the surface, and a dramatic increase in tool wear.

Platings may also be optimized to give low ductility and hence minimum burr formation when machining Fresnels or the molds for micro-optics arrays such as arrays of retroreflectors. Platings over a diamond-turned sacrificial mandrel allow production of otherwise unobtainable forms. Plated surfaces, however, have characteristics which can adversely affect both fabrication and application; at some level the resulting structure is a temperature-sensitive bimetal strip. Pits and inclusions can cause significant fabrication issues.⁷

Silicon, although included in the listing given here, should be considered marginal as tool wear can be high. Reasonably large areas of amorphous silicon cladding are reported to have been successfully machined.

Over the last decade or so there have been significant advances in understanding of diamond tool wear. Mechanisms associated with abrasion and chipping typically provide one limit to diamond tool life. For example, when machining bulk aluminums the interactions between hard inclusions and the diamond tool clearly lead to wear. Such mechanisms, however, do not explain the very rapid wear observed when machining soft, high-purity iron.

Paul et al.⁸ showed that machining metallic elements containing unpaired *d*-shell electrons results in catalyzed reactions between diamond and the work material. The same mechanism explains the role of phosphorous in electroless nickel and led to a recent breakthrough by Brinksmeier et al.⁹ They showed that, like phosphorous, nitrogen in a nitride surface layer on steel combines with the unpaired *d*-shell electrons from the iron. The result is a dramatic reduction in tool wear, suggesting that diamond-turned steel molds (e.g., for plastic optics) may become practical in the near future. Previous approaches—such as diamond turning at cryogenic temperatures¹⁰ or in methane or acetylene environments¹¹—provided evidence of the mechanisms at work but were not widely adopted (and in the case of cryogenic machining was not intended by the original researchers to be practical).

For a number of years, Moriwaki¹² has been developing ultrasonic-assisted machining, including cutting when the tool is vibrated with an elliptical motion. Significant reductions in tool wear have been demonstrated, although the mechanism remains controversial. The amplitude and frequency of oscillation in the cutting direction are generally selected such that separation between the rake face of the tool and the chip is expected. In this case, one might postulate poisoning of the catalytic process by, for example, hydrocarbon-based cutting fluids. Elliptical motion would also move the clearance face out of contact. Other measurements show significant reductions in cutting forces,¹³ suggesting a reduction of tool temperatures.

In case of some plastics, recent work by Gubbels et al.¹⁴ shows that the chemical explanations of Paul et al.⁸ do not apply, but that triboelectric effects dominate. In general, there are some suggestions that parameters, such as surface speed, are more important for successful diamond machining of plastics than for metal and crystalline substrates. Some plastics are diamond turned in volume production.

Other material characteristics, in addition to the material composition, are important. For example, large grain size results in a more pronounced “orange peel” as tools become dull and the variation in modulus of the different grain orientations leads to different deflections due to cutting forces. Residual stresses can relax over time and cause changes in figure. Because of these types of problems it is important to involve experienced personnel early in the design phase¹⁵ to ensure that the material specified is appropriate. In some projects, the part is so valuable and/or so difficult to produce by other techniques, it is worth consuming tools more rapidly than would normally be acceptable. However, such a decision should be taken consciously, not by default late in a project.

10.6 COMPARISON OF DIAMOND TURNING AND TRADITIONAL OPTICAL FABRICATION

In diamond turning, the final shape and surface of the optic produced depend on the machine tool accuracy, whereas, in traditional optical fabrication, the final shape and surface of the optical element depend on the process variables involved with using an abrasive-loaded lap. The differences between diamond turning and traditional optical fabrication can be summarized by describing diamond turning as a displacement-controlled process versus a force-controlled process for traditional optical fabrication.¹⁶ The goal in diamond turning is to have a machine tool that produces an extremely accurate path with the diamond tool, hence a displacement-controlled process. A traditional polishing machine used for optical fabrication depends on the force being constant over the area where the abrasive-loaded lap—or tool—touches the surface being worked. Selective removal of material can be produced by increasing the lap pressure in selected areas or by use of a zone lap. The stiffness of a diamond-turning machine is important because, to control the displacement, it is important that cutting forces and other influences do not cause unwanted displacements. Feeds, speeds, and depth of cut are typically much lower in diamond turning than conventional machining, thus giving lower forces. However, the displacements of concern are also much lower. Thus the stiffness required is as much, or more, of a concern than conventional machining even though the total force capability may be lower for diamond turning.

10.7 MACHINE TOOLS FOR DIAMOND TURNING

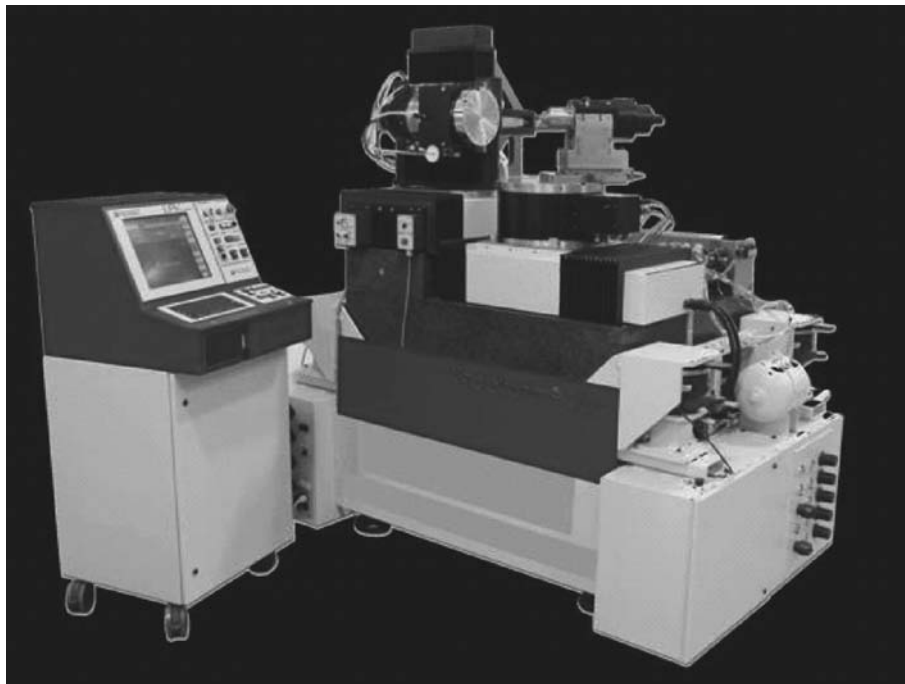
In general, the machine tools used for diamond turning are very expensive compared to the equipment needed for traditional optical fabrication. The positioning accuracy required for diamond turning is beyond the capability of conventional machine tools, thus some of the first widely adopted diamond-turning machines for fabricating optics were modified Moore measuring machines.¹⁷

Although there are some records of machine tools being used to generate optical surfaces as early as the seventeenth century, most of the effort is modern, accelerated in the 1960s and 1970s with the advent of computer-based machine tool controls and laser interferometer systems used as positional feedback devices. Evans¹⁸ has documented much of the history of diamond turning and provides an extensive reference list. Ikawa¹⁹ summarizes some of the research in metal cutting related to diamond turning and associated machine tools.

The early diamond-turning machines were two-axis lathes that could produce axisymmetric optical elements. With recent advances in computer-based control systems, and improved motion control and feedback sensors, multiaxis diamond-turning machines have become readily available. Two commercial diamond-turning machines are shown in Fig. 2. Both machines can be configured with five-axis motion control combining both linear and rotary axes. Measuring scales have replaced the laser interferometers in many diamond-turning machines and give a very reliable positioning feedback system at lower cost.

Programming of these multiaxis machines draws on the technology developed in precision machine shops for large five-axis machine tools used to make complicated parts. By adapting the multiaxis control to diamond-turning machines, a great variety of shapes can now be diamond turned which opens up the process to many new optical applications. Before judging an optical element shape to be unsuitable for diamond turning, a manufacturer of modern diamond-turning machines should be consulted.

Producing nonaxisymmetric parts—such as an off-axis parabola machined while centered on the rotating axis—has become possible with fast tool servos. These systems can rapidly move a cutting tool a short distance coordinated with the rotation of the spindle.²⁰ There are also cases where the machine's slideways or rotary motions can be used to produce nonaxisymmetric parts.



(a)



(b)

FIGURE 2 (a) Diamond-turning machine configurable for five-axis machining. (Courtesy of Precitech, Inc., Keene, New Hampshire.) (b) Diamond-turning machine configurable for five-axis machining. (Courtesy of Moore Nanotechnology Systems, LLC., Keene, New Hampshire.)

Many diamond-turning machines are used in the traditional turning lathe mode where the workpiece turns and the tool is held stationary in the tool post. Most diamond-turning machines can also be configured such that the tool rotates about the spindle axis—commonly called fly cutting—to produce components such as long flat mirror surfaces or other milled surfaces.

10.8 BASIC STEPS IN DIAMOND TURNING

Much like the traditional optical-fabrication process, the diamond-turning process can be described as a series of steps used to make an optical element. The steps used in diamond turning are

1. *Preparing the blank* with all the required features of the element with an extra thickness of material (generally 0.1-mm extra material or plating is adequate) on the surface to be diamond turned
2. *Mounting the blank* in an appropriate fixture or chuck on the diamond-turning machine
3. *Selecting the diamond tool* appropriate for the material and shape of the optical component
4. *Mounting and adjusting the diamond tool* on the machine
5. *Machining the optical surface* to final shape and surface quality
6. *Cleaning the optical surface* to remove cutting oils or solvents

Mounting the optical element blank on a diamond-turning machine is extremely important. If a blank is slightly distorted in the holding fixture, and then machined to a perfect shape on the machine, it will be a distorted mirror when released from the fixture. Therefore, fixtures and chucks to hold mirrors during diamond turning need to be carefully designed to prevent distortion. Often the best way to hold a mirror during machining is to use the same mounting method that will be used to hold the mirror in service.

It is advantageous in many applications to machine a substrate of aluminum or copper and then add a plating to be diamond turned. The design and application of platings is part science and part art. Many aspects of the platings as related to diamond turning were covered at the ASPE Spring 1991 Topical Meeting.⁷

Tool setting—the mounting and adjusting of the diamond-tipped cutting tool—is often accomplished by cutting a test surface, either on the actual mirror blank to be later machined over, or by placing a test piece on the machine just for tool setting. If the cutting tool is too high or too low, a defect at the center of a mirror is produced. It is possible, using reasonable care and patience, to set the tool height within about 0.1 μm of the exact center. Setting the tool in the feed direction after the height is correct is somewhat more difficult. For example, an error in setting will produce an ogive shape rather than a sphere which is not obvious until the figure is measured. Gerchman²¹ describes these types of defects.

The selection of the tool for diamond turning is important. Large cutting tip radii (2 mm or greater) are often used when producing flats, convex, or concave mirrors with large radius of curvature. However, small-radii diamond tools are available (in the range of 0.1 mm) for making small deep mirrors or molds. Tools with special geometries, including so-called “dead sharp” tools, can be obtained for such applications as Fresnel lenses or retroreflector arrays. In general, approximately 0° rake tools, with about 5° or 6° front clearance, are used for diamond-turning ductile metals. Negative rake tools are often good for crystalline materials and positive rakes may be beneficial when machining some plastics. The cutting edge has to be chip free to produce a good diamond-turned surface. A normal specification for edge quality is “chip free when examined at 1000 \times .” The edge sharpness is a concern for very small depths of cut—especially where the depth of cut is close to the cutting edge sharpness—because the cutting forces increase and more of a plowing than a cutting process occurs. The effect of cutting edge sharpness has been investigated by researchers, for example Lucca, et al,²² however, there is currently no convenient way to specify and inspect tools for edge sharpness.

The orientation of the diamond itself on the shank is of concern because the single-crystal diamond is anisotropic. The orientation of diamond tools has been studied, for example, by Wilks,²³ Decker,²⁴ and Hurt.²⁵ It is necessary for the tool manufacturer to mount the diamond so that it can be shaped to the required radius and produce a good cutting edge. The usual orientation for diamond tools is with the cleavage plane parallel to the rake face.

The actual diamond turning, or machining to final size and surface finish, is often the fastest part of the process. The machine-tool controller has to be programmed to move the tool along the correct path, the chip-removal system has to be positioned, and the cutting-fluid applicator needs to be adjusted to provide consistent clean cutting.

For machining of flats and spherical surfaces, the part programs that define the machine motion are straightforward. But when cutting aspherical surfaces, caution has to be exercised so that the radius of the tool is properly handled in calculating the tool path. Modern computer-aided design (CAD) systems perform the necessary calculations, but tests should be performed prior to cutting a difficult or expensive component.

In general, the cutting speeds for diamond turning are similar to those used for conventional machining: less than 1 m/min to more than 100 m/min. However, the slower cutting speeds produced by facing to the center of a workpiece do not affect the surface finish in diamond turning as is often the case with nondiamond tools. Thus, varying the spindle speed to keep the cutting speed constant is not necessary in diamond turning. The upper speed for diamond turning is often limited by the distortion of the optical element due to inertial forces, especially for larger elements. The upper spindle speed can also be limited due to any unbalance of the workpiece and fixture. The feed rate in diamond turning is usually adjusted to give a good theoretical surface finish. (See the following section.)

Cleaning of diamond-turned optics has a lot in common with cleaning conventionally polished optics. But because many of the diamond-turned elements are of soft metals, caution has to be exercised to prevent scratching. In general, a degreaser is used (soap or solvent), followed by a rinse in pure ethyl alcohol. The drag-wiping technique traditionally used on some glass optics can be used on some diamond-turned elements. Care must be taken to ensure that the lens tissue is very clean and remains wet. Some work has been done to study the best solvents to use for cleaning diamond-turned optics from an environmental-impact standpoint.²⁶

10.9 SURFACE FINISH OF DIAMOND-TURNED OPTICS

The surface structure is different for diamond-turned surfaces as compared with conventionally polished surfaces. A diamond-turned surface is produced by moving a cutting tool across the surface of the turning component, such as the facing operation illustrated in Fig. 3. Therefore, diamond-turned elements always have some periodic surface roughness, which can produce a diffraction-grating effect, whereas polished optical surfaces have a random roughness pattern. The traditional “scratch and dig” approach to describing surfaces is not meaningful for diamond-turned surfaces.

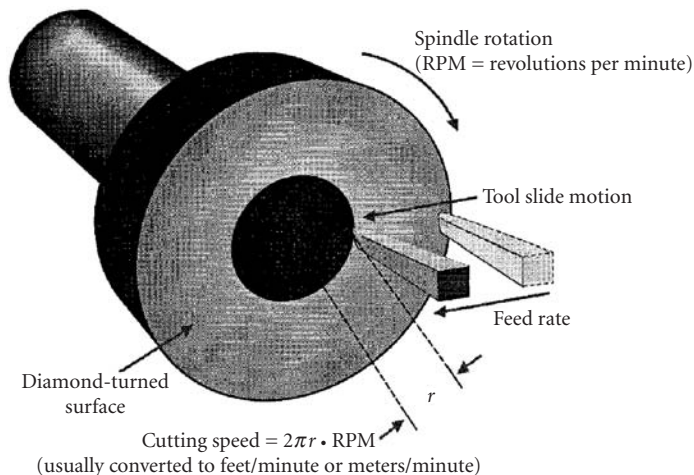


FIGURE 3 Diamond turning an optical element.

The machining process produces a periodic surface structure directly related to the tool tip radius and feed rate. The theoretical diamond-turned surface is illustrated in Fig. 4. The formula displayed in the figure for calculating the height of the cusps is

$$h = \frac{f^2}{8 \cdot R} \tag{1}$$

where h = peak-to-valley height of the periodic surface defect
 f = feed per revolution
 R = tool tip radius

For example, if a surface is diamond turned using a spindle speed of 31.4 rad/s (300 rpm), a feed of 7.5 mm/min, and a 5.0-mm tool tip radius:

$$h = \frac{(7.5/300)^2}{8 \times 5} = 1.56 \times 10^{-5} \text{ mm}$$

$$h = 15.6 \text{ nm} \tag{2}$$

In addition to the “theoretical finish” based on cusp structure, the measured surface finish on diamond-turned parts is influenced by other factors.

- Waviness within the long-wavelength cut-off for surface measurement may be correlated, for example, with slide straightness errors.
- Asynchronous error motions of the spindle can cause surface defects. If, for a given angular spindle position, there is nonrepeatability in axial, radial, or tilt directions, these errors will transfer into surface structure. Details of spindle errors are important in diamond turning. Further information can be found in the “Axis of Rotation Standard.”²⁷
- External and self-induced vibration, not at the spindle frequency or at one of its harmonics, can have the same effect on finish—measured across the lay—as asynchronous spindle motions.
- Materials effects such as the differential elastic recovery of adjacent grains can cause steps in the machined surface that have an appearance commonly referred to as “orange peel.” Impurities in the material can also degrade surface finish.
- Within each cusp, there can be a repeated surface structure related to chips in the edge of the tool.

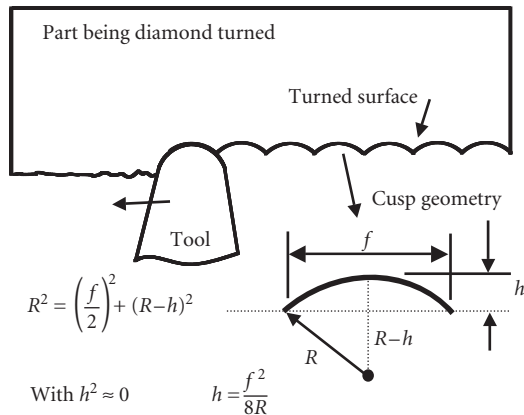
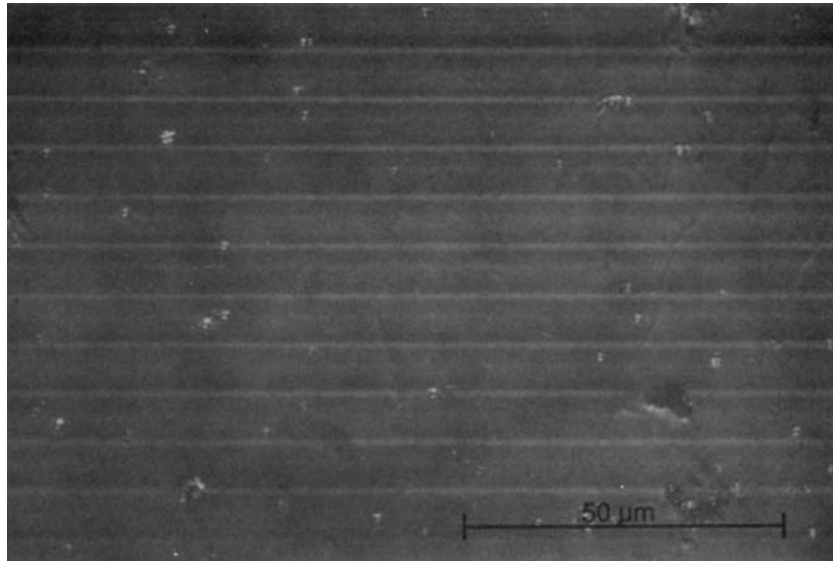
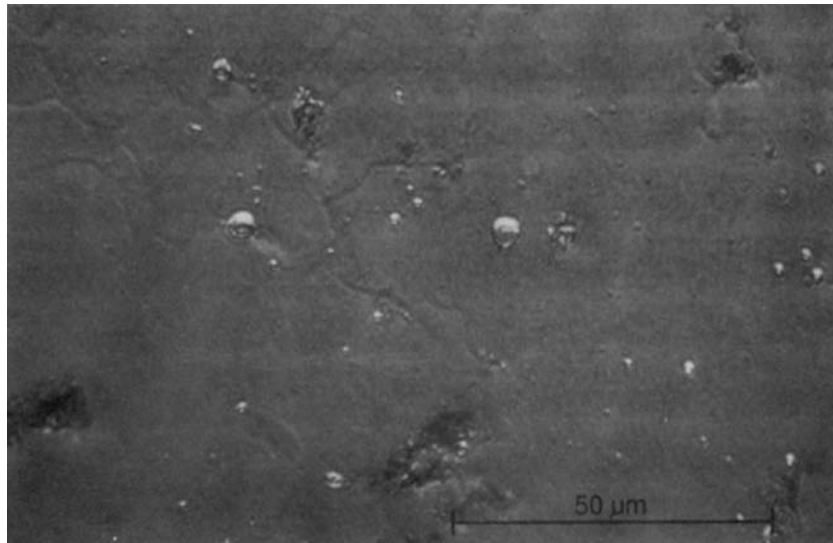


FIGURE 4 “Cusp” surface of diamond-turned optical element.

The Nomarski microscope is an excellent means of qualitatively evaluating diamond-turned surfaces. The Nomarski photo in Fig. 5a illustrates the periodic structure of a diamond-turned surface. The feed rate used in producing the surface causes the wavelength of the periodic structure to be about $8\ \mu\text{m}$. Figure 5b illustrates other defects in the diamond-turned surface when the Nomarski microscope is adjusted such that the periodic cusps are not seen.²⁸



(a)



(b)

FIGURE 5 Nomarski micrograph of a diamond-turned aluminum alloy (a) aligned so that the grooves can be seen and (b) aligned so that the grooves are canceled. (From Bennett, p. 84.²⁸)

10.10 METROLOGY OF DIAMOND-TURNED OPTICS

In general, measurement of diamond-turned optics is similar to the measurement of any other optic; figure, midspatial frequency errors, transmitted wavefront, and surface roughness may all need characterization, depending on the specification. As with other optics, the choice of figure metrology is driven by the optical surface itself. Classical null tests—especially autocollimation tests for parabolae and the related tests for other conics—are widely used. Over the last decade or so, use of in-cavity holograms in Fizeau tests has increased. What remains elusive is a general test. Over a limited range of surfaces, subaperture stitching²⁹ may be viable or, for circularly symmetric aspheres, zonal stitching.³⁰ Kuechel³¹ described a zonal technique that uses only the zone of null data and, hence, is free of retrace errors and applicable to a range of aspheres without the need for null optics. An instrument based on this technique is shown in Fig. 6.

One area in which diamond turning differs from conventional optics production is that the machine itself can be used as a measuring machine. The diamond tool can be replaced with an appropriate sensor (such as a capacitance sensor, air bearing linear variable differential transformer (LVDT), optical triangulation sensor, etc.) or the sensor can be built into an auxiliary mount. With sufficient care,³² the geometric errors of the machine can be mapped so that the limits in the metrology are the uncertainties associated with probing and with the environment. This approach is particularly advantageous when making (and measuring) radical aspheres or discontinuous, structured surfaces² such as molds for faceted automotive lighting. On multi-axis machines, it is sometimes more useful to use a different combination of axes for metrology than for machining to better decouple machine geometry errors from measurement uncertainty. For example, on a diamond-turning machine with a B axis (rotary table), near hemispheres and some aspheres can conveniently be machined using only the x and y axes, with measurement of the departure from a best-fit sphere performed using a separate probe mounted on the rotary table.



FIGURE 6 Aspheric measuring system. (Courtesy Zygo Corporation, Middlefield, CT.)



FIGURE 7 Microinterferometer. (Courtesy of Zygo Corporation).

There is little practical difference between measuring optical surfaces produced using traditional methods and by diamond turning. It is worth bearing in mind, however, that during diamond turning there is usually a monotonic progression in cutting from outside diameter to inside diameter or vice versa; hence, diamond tool wear or small edge nicks will cause a degradation in finish that depends on position on the part. The surface finish measurement sampling strategy should be adjusted accordingly. Surfaces produced using traditional 2-axis or 3-axis diamond-turning have significantly different characteristics along and transverse to the lay; scattering is isotropic, a characteristic that should be considered in both the specification and metrology of diamond-turned optics. Four-axis and 5-axis machining using methods akin to milling produce cusp structures usually at different spatial wavelengths in both directions.

Microinterferometers (Fig. 7) have become the tool of choice for characterizing optical surfaces at spatial wavelengths down to the limits posed by the instrument transfer function.³³ Microinterferometers—particularly those using scanning coherence techniques frequently referred to as scanning white light interferometry (SWLI)—can be useful, provided the surface slopes and lateral extent are compatible with the available numerical aperture of the objective and the field of view. Replication—for example, using dental replica materials, silicone-based caulks, two-part epoxies, and the like—allows sampling of large surfaces, although there is inevitably some increase in “noise” due to the replication process.

Higher spatial frequency structured surfaces, such as retroreflectors or other micro-optic arrays,³⁴ often pose metrology challenges for which there is no general solution.

10.11 CONCLUSIONS

Diamond turning has been used for many years to commercially produce infrared optics. Some visible and ultraviolet applications are now possible. Moreover, the limits of diamond turning for both figure and surface finish accuracy have not yet been reached. Taniguchi³⁵ and others have shown that precision in both conventional machining and ultraprecision machining, such as diamond turning, has steadily improved for many decades, with roughly a factor of three improvements possible every 10 years. If this trend continues, we could expect diamond-turning machines with accuracies below 10 nm and even approaching 1 nm by the year 2020. Yet, it is important to remember that it becomes increasingly difficult to push the capabilities in this regime—nor is it clear that it is cost effective to do so. Other manufacturing techniques may be more appropriate for production of the highest quality optics.

The technology developed for diamond-turning optics in some industries is now beginning to impact the precision machining of nonoptical components. In the future, the improvement of all machine tools will likely be driven by both optical and nonoptical applications, with diamond-turning machines possibly reaching the accuracy level that will allow visible and ultraviolet optics to be fabricated by machining or grinding without postpolishing.

10.12 REFERENCES

1. R. E. Parks and C. J. Evans, "Rapid Post-Polishing of Diamond-Turned Optics," *Precision Engineering* **16**: 223–227 (1994).
2. C. J. Evans and J. B. Bryan, "Structured, Engineered and Textured Surfaces," *CIRP Annals* **48/2**:541–546 (1999).
3. C. K. Syn, J. S. Taylor, and R. R. Donaldson, "Diamond Tool Wear vs. Cutting Distance on Electroless Nickel Mirrors," *Proc. SPIE* **676**:128–140 (1986).
4. J. S. Taylor, C. K. Syn, T. T. Saito, and R. R. Donaldson, "Surface Finish Measurements of Diamond Turned Electroless Nickel Plated Mirrors," *Optical Engineering* **25**(9):1013–1020 (1986).
5. A. Mayer, et al., "Electrodeposited Coatings for Diamond Turning Applications," *Proc. ASPE Spring Topical Meeting, Metal Platings for Precision Finishing Operations*, 1991.
6. C. J. Evans, R. S. Polvani, and A. Mayer, "Diamond Turned Electrodeposited Nickel Alloys," *OSA Technical Digest Series* **9**:110 (1990).
7. ASPE, "Metal Platings for Precision Finishing Operations," *Spring Topical Meeting*, Raleigh, N.C., 1991.
8. E. Paul, C. J. Evans, A. Mangamelli, M. L. McGlaufflin, and R. S. Polvani, "Chemical Aspects of Tool Wear in Single Point Diamond Turning," *Precision Engineering* **18**:4–19 (1996).
9. E. Brinksmeier, R. Glabe, and J. Osmer, "Ultra-Precision Diamond Cutting of Steel Molds," *CIRP Annals* **55**:551–554 (2006).
10. C. J. Evans, "Cryogenic Diamond Turning of Stainless Steel," *CIRP Annals* **40**:571–575 (1991).
11. J. Casstevens, "Diamond Turning of Steel in Carbon Saturated Environments," *Precision Engineering* **5**:9–15 (1983).
12. T. Moriwaki, "Ultraprecision Diamond Turning of Stainless Steel by Applying Ultrasonic Vibration," *CIRP Annals* **40**:559–562 (1991).
13. E. Shamoto and T. Moriwaki, "Ultraprecision Diamond Cutting of Hardened Steel by Applying Elliptical Vibration Cutting," *CIRP Annals* **48**:441–444 (1999).
14. G. P. H. Gubbels, G. J. F. T. van der Beek, A. L. Hoep, F. L. M. Delbressine, and H. Halewijn, "Diamond Tool Wear When Cutting Amorphous Polymers," *CIRP Annals* **53**:447–450 (2004).
15. J. S. Taylor and C. J. Evans, "Fabrication of a Metal Plated Mirror, Beginning from a Performance Specification," *Proc. ASPE Spring Topical Meeting, Metal Platings for Precision Finishing Operations*, 1991.
16. A. Gee, Cranfield Institute of Technology, in a private communication to the authors used the description of "displacement" controlled and "force" controlled.
17. W. R. Moore, *Foundations of Mechanical Accuracy*, Moore Special Tool Co., Bridgeport, CT., 1970.
18. C. J. Evans, *Precision Engineering: An Evolutionary View*, Cranfield Press, Bedford, UK, 1989, pp. 135–155.
19. N. Ikawa, et al., "Ultraprecision Metal Cutting : the Past, the Present, and the Future," *CIRP Annals* **40**: 587–594 (1991).
20. S. Patterson and E. Magrab, "Design and Testing of a Fast Tool Servo for Diamond Turning," *Precision Engineering* **7**:131–136 (1985).
21. M. C. Gerchman, "Optical Tolerancing for Diamond Turning Ogive Error," *Proc. SPIE (Reflective Optics II)*: 224–229 (1989).
22. D. A. Lucca, Y. W. Seo, R. L. Rhorer, and R. R. Donaldson, "Aspects of Surface Generation in Orthogonal Ultraprecision Machining," *CIRP Annals* **43**:43–46 (1994).
23. J. Wilks, "Performance of Diamonds as Cutting Tools for Precision Machining," *Precision Engineering* **2**: 57–71 (1980).

24. D. C. Decker, H. H. Hurt, J. H. Dancy, and C. W. Fountain, "Preselection of Diamond Single Point Tools," *Proc. SPIE* **508** (1986).
25. H. H. Hurt and D. L. Decker, "Tribological Considerations of the Diamond Single Point Tool," *Proc. SPIE* **508** (1986).
26. L. A. Theye and R. D. Day, "Evaluation of Environmentally Safe Cleaning Agents for Diamond Turned Optics," *Proc. ASPE*, ASPE, Raleigh, N.C., 1991.
27. ANSI/ASME B89.3.4M-2006 Standard: "Axes of Rotation: Methods for Specifying and Testing," ASME, New York, 2006.
28. J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, Washington, D.C., 1989.
29. P. Murphy, J. Fleig, G. Forbes, D. Miladinovic, G. DeVries, and S. O'Donohue, "Subaperture Stitching Interferometry for Testing Mild Aspheres," *Proc. SPIE* **6293** (2006).
30. M. J. Tronolone, J. F. Fleig, C. Huang, and J. H. Bruning, US Patent 5,416,586 (1995).
31. M. Kuechel, US Patent 6,781,700, (2004).
32. W. T. Estler and E. B. Magrab, "Validation Metrology of the Large Optics Diamond Turning Machine," NBSIR 85-3182(R), U.S. Department of Commerce, National Bureau of Standards (1985).
33. P. de Groot and X. Colonna de Lega, "Interpreting Interferometric Height Measurements Using the Instrument Transfer Function," *Proc. Fringe 2005: 5th International Workshop on Automatic Processing of Fringe Patterns*, 2005.
34. M. A. Davies, C. J. Evans, R. R. Vohra, B. C. Bergner, and S. R. Patterson, "Application of Precision Diamond Machining to the Manufacture of Microphotonics Components," *Proc. SPIE* **5183**:94–108 (2003).
35. N. Taniguchi, "Nanotechnology for Ultraprecision Instruments," *Precision Engineering* **16**:5–24 (1994).

This page intentionally left blank.

PART

3

TESTING

This page intentionally left blank.

ORTHONORMAL POLYNOMIALS IN WAVEFRONT ANALYSIS

Virendra N. Mahajan*

*The Aerospace Corporation
El Segundo, California*

ABSTRACT

Zernike circle polynomials are in widespread use for wavefront analysis because they are orthogonal over a unit circle and represent balanced classical aberrations for imaging systems with circular pupils. However, they are not suitable for systems with noncircular pupils. Examples of such pupils are annular as in astronomical telescopes, elliptical as in the off-axis pupil of an otherwise rotationally symmetric system with a circular on-axis pupil, hexagonal as in the hexagonal segments of a large telescope, for example, Keck, and rectangular and square as in high-power laser beams. In this chapter, we list the orthonormal circle, annular, elliptical, hexagonal, rectangular, and square polynomials. The polynomials for a noncircular pupil can be obtained by orthogonalizing the circle polynomials over the pupil using the recursive Gram-Schmidt process or a nonrecursive matrix approach. These polynomials are unique in that they are not only orthogonal across such pupils, but also represent balanced classical aberrations for such pupils, just as the Zernike circle polynomials are unique in these respects for circular pupils. The polynomials are given in terms of the circle polynomials as well as in polar and Cartesian coordinates. The orthonormal polynomials for a one-dimensional slit pupil are given as a limiting case of a rectangular pupil. The polynomials corresponding to Seidel aberrations are illustrated isometrically, interferometrically, and with the corresponding point-spread functions (PSFs).

11.1 GLOSSARY

- a half width of a unit rectangular pupil
- a_j j th expansion coefficient
- A area of pupil
- b aspect ratio of a unit elliptical pupil

*The author is also an adjunct professor at the College of Optical Sciences, University of Arizona, Tucson, Arizona and Department of Optics and Photonics, National Central University, Chung Li, Taiwan. He gratefully acknowledges helpful discussions with Drs. Guang-ming Dai and Bill Swantner.

$E_j(x, y)$	orthonormal elliptical polynomial in Cartesian coordinates (x, y)
F	focal ratio of the image-forming light cone
$F_j(x, y)$	j th orthonormal polynomial
$H_j(x, y)$	orthonormal hexagonal polynomial
j	polynomial number
N_n	number of polynomials through an order n
$P_j(x)$	orthonormal slit polynomial along the x axis
$P_n(\cdot)$	Legendre polynomial of order n
$R_j(x, y)$	orthonormal rectangular polynomial
$R_n^m(\rho)$	Zernike circle radial polynomial
$R_n^m(\rho; \epsilon)$	Zernike annular radial polynomial
$S_j(x, y)$	orthonormal square polynomial
$W(x, y)$	wave aberration at a point (x, y)
$Z_j(\rho, \theta)$	orthonormal Zernike circle polynomial in polar coordinates (ρ, θ)
$Z_j(\rho, \theta; \epsilon)$	orthonormal Zernike annular polynomial
σ	standard deviation
σ^2	variance
ϵ	obscuration ratio of an annular pupil

11.2 INTRODUCTION

Optical systems generally have a circular pupil. The imaging elements of such systems have a circular boundary. Hence they also represent circular pupils in fabrication and testing. As a result, the Zernike circle polynomials have been in widespread use since Zernike introduced them in his phase contrast method for testing circular mirrors.¹ They are used in optical design and testing to understand the aberration content of a wavefront. They have also been used for analyzing the wavefront aberration introduced by atmospheric turbulence on a wave propagating through it.² Their utility stems from the fact that they are orthogonal over a unit circle and they represent balanced classical aberrations yielding minimum variance over a circular pupil.³⁻⁶ They are unique in this respect since no other polynomials have these properties. Because of their orthogonality, when a wavefront is expanded in terms of them, the value of an expansion coefficient is independent of the number of polynomials used in the expansion. Hence, one or more polynomial terms can be added or subtracted without affecting the other coefficients. The piston coefficient represents the mean value of the aberration function and the variance of the function is given simply by the sum of the squares of the other expansion coefficients.⁷

For systems with noncircular pupils, the Zernike circle polynomials are neither orthogonal over such pupils nor do they represent balanced aberrations. Hence their special utility is lost. However, since they form a complete set, an aberration function over a noncircular wavefront can be expanded in terms of them. The expansion coefficients are no longer independent of each other and their values change as the number of polynomials used in the expansion changes. The piston coefficient does not represent the mean value of the aberration function, and the sum of the squares of the other coefficients does not yield the aberration variance.

The reflecting telescopes, such as the Hubble, have annular pupils and require polynomials that are orthogonal across an annulus to describe their aberrations.⁸⁻¹¹ The primary mirrors of large telescopes, such as the Keck, consist of hexagonal segments.¹² The wavefront analysis of such segments requires polynomials that are orthogonal over a hexagon. The pupil for off-axis imaging by a system with an axial circular pupil is vignetted, but can be approximated by an ellipse.¹³ When a flat mirror is tested by shining a circular beam on it at some angle (other than normal incidence), the illuminated spot is elliptical. Similarly, the overlap region of two circular wavefronts that are

displaced from each other, as in lateral shearing interferometry¹⁴ or in the calculation of the optical transfer function of a system,¹⁵ can also be approximated by an ellipse. In such cases we need polynomials that are orthogonal over an ellipse. In Refs. 14 and 15, the polynomials that are orthogonal over an elliptical region were obtained simply by scaling the Cartesian coordinates by its aspect ratio. However, such orthogonal polynomials cannot represent classical aberrations. For example, defocus, which varies as ρ^2 , has the same scale for both the x and y coordinates. Similarly, they cannot represent balanced classical aberrations, for example, coma balanced with tilt. High-power laser beams have rectangular or square cross sections¹⁶ and require polynomials that are orthogonal over a rectangle or a square, respectively.

The polynomials orthonormal over a unit annulus, hexagon, ellipse, rectangle, and a square inscribed inside a unit circle may be obtained from the circle polynomials by the recursive Gram-Schmidt orthogonalization process^{17,18} or a nonrecursive matrix approach.¹⁹ The orthonormal polynomials representing balanced aberrations for a slit pupil can be obtained as a limiting case of the rectangular polynomials, where one dimension of the rectangle approaches zero. They are the Legendre polynomials.²⁰ We use the circle polynomials as the basis functions for the orthogonalization process, so that the relationship of a noncircle polynomial to the circle polynomials is evident, since the former is a linear combination of the latter. We give the orthonormal form of the polynomials so that when an aberration function is expanded in terms of them, each expansion coefficient (with the exception of piston) represents the standard deviation of the corresponding expansion term. The noncircle polynomials are given not only in terms of the circle polynomials, but in polar and Cartesian coordinates as well. The circle, annular, hexagonal, and square polynomials are given up to the eighth order, and the elliptical and rectangular polynomials are given up to the fourth order. Just as the Zernike circle polynomials uniquely represent the orthogonal and balanced aberrations across circular pupils, similarly, the orthonormal polynomials for the noncircular pupils given in this chapter also uniquely represent the orthogonal and balanced aberrations across such pupils.

Orthogonal square polynomials were obtained by Bray by orthogonalizing the circle polynomials, but he chose a circle inscribed inside a square instead of the other way around.²¹ Thus his square with a full width of unity has regions that fall outside the unit circle. Defining a unit square in this manner has the disadvantage that the coefficient of a term in a certain polynomial does not represent its peak value. Products of x and y Legendre polynomials,¹⁷ which are orthogonal over a square pupil, have been suggested for analysis of square wavefronts.²² But they do not represent classical or balanced aberrations. For example, defocus is represented by a term in $x^2 + y^2$. While it can be expanded in terms of a complete set of Legendre polynomials, it cannot be represented by a single two-dimensional Legendre polynomial (i.e., as a product of x and y Legendre polynomial). The same difficulty holds for spherical aberration and coma, and the like.

Although in many imaging applications, the amplitude across the pupil is uniform, such is not always the case, for example, a system with an apodized pupil. An example of such a pupil is the Gaussian pupil, where the amplitude has the form of a Gaussian due either to an amplitude filter placed at the pupil or to the wave incident on the pupil being Gaussian, as in the case of a Gaussian laser beam. Again, the balanced aberrations for a Gaussian pupil have a form that is different from the corresponding balanced aberrations for a uniform pupil due to the amplitude weighting of the pupil.^{23–25} The amount of defocus to optimally balance spherical aberration, or the amount of wavefront tilt to optimally balance coma, for example, is different for a Gaussian pupil than its corresponding value for a uniform pupil.

11.3 ORTHONORMAL POLYNOMIALS

In Cartesian coordinates (x, y) , the aberration function $W(x, y)$ for a certain pupil may be expanded in terms of J polynomials $F_j(x, y)$ that are orthonormal over the pupil:²⁶

$$W(x, y) = \sum_{j=1}^J a_j F_j(x, y) \quad (1)$$

where a_j is an expansion or the aberration coefficient of the polynomial $F_j(x, y)$. The orthonormality of the polynomials is represented by

$$\frac{1}{A} \int_{\text{pupil}} F_j(x, y) F_{j'}(x, y) dx dy = \delta_{jj'} \quad (2)$$

where A is the area of the pupil inscribed inside a unit circle, the integration is carried out over the area of the pupil, and $\delta_{jj'}$ is a Kronecker delta. If $F_1 = 1$, then the mean value of each polynomial, except for $j = 1$, is zero, that is,

$$\frac{1}{A} \int_{\text{pupil}} F_j(x, y) dx dy = 0 \quad \text{for } j \neq 1 \quad (3)$$

as may be seen by letting $j' = 1$ in Eq. (2). The aberration coefficients are given by

$$a_j = \frac{1}{A} \int_{\text{pupil}} W(x, y) F_j(x, y) dx dy \quad (4)$$

as may be seen by substituting Eq. (1) into Eq. (4) and using the orthonormality Eq. (2).

The mean and the mean square values of the aberration function are given by

$$\langle W(x, y) \rangle = a_1 \quad (5)$$

and

$$\langle W^2(x, y) \rangle = \sum_{j=1}^J a_j^2 \quad (6)$$

Accordingly, the variance σ^2 of the aberration function is given by

$$\sigma^2 = \langle W^2(x, y) \rangle - \langle W(x, y) \rangle^2 = \sum_{j=2}^J a_j^2 \quad (7)$$

where σ is the standard deviation of the aberration function. The number of polynomials J used in the expansion is a sufficiently large that the variance obtained from Eq. (6) equals the actual value obtained from the function $W(x, y)$ within some prescribed tolerance.

11.4 ZERNIKE CIRCLE POLYNOMIALS

An aberration function $W(\rho, \theta)$, across a *unit circle* can be expanded in terms of the orthonormal *Zernike circle polynomials* $Z_j(\rho, \theta)$ in the form^{2,5}

$$W(\rho, \theta) = \sum_j a_j Z_j(\rho, \theta) \quad (8)$$

where (ρ, θ) are the polar coordinates of a point on the circle, $0 \leq \rho \leq 1$, $0 \leq \theta < 2\pi$, and a_j are the expansion coefficients. The polynomials may be written in the form

$$Z_{\text{even}j}(\rho, \theta) = \sqrt{2(n+1)} R_n^m(\rho) \cos m\theta, m \neq 0 \quad (9a)$$

$$Z_{\text{odd}j}(\rho, \theta) = \sqrt{2(n+1)} R_n^m(\rho) \sin m\theta, m \neq 0 \quad (9b)$$

$$Z_j(\rho, \theta) = \sqrt{n+1} R_n^0(\rho), m = 0 \quad (9c)$$

where n and m are positive integers (including zero) and $n - m \geq 0$ and even. It is evident from Eqs. (9) that the circle polynomials are separable in the polar coordinates ρ and θ . A radial polynomial $R_n^m(\rho)$ is given by

$$R_n^m(\rho) = \sum_{s=0}^{(n-m)/2} \frac{(-1)^s (n-s)!}{s! \left(\frac{n+m}{2} - s\right)! \left(\frac{n-m}{2} - s\right)!} \rho^{n-2s} \quad (10)$$

with a degree n in ρ containing terms in $\rho^n, \rho^{n-2}, \dots$, and ρ^m . It is even or odd in ρ depending on whether n (or m) is even or odd. Also, $R_n^n(\rho) = \rho^n$, $R_n^n(1) = 1$, and $R_n^m(0) = \delta_{m0}$ for even $n/2$ and $-\delta_{m0}$ for odd $n/2$. The polynomials $R_n^m(\rho)$ obey the orthogonality relation

$$\int_0^1 R_n^m(\rho) R_{n'}^{m'}(\rho) \rho d\rho = \frac{1}{2(n+1)} \delta_{nn'} \quad (11)$$

The orthogonality of the angular functions yields

$$\int_0^{2\pi} d\theta \begin{cases} \cos m\theta \cos m'\theta, & j \text{ and } j' \text{ are both even} \\ \cos m\theta \sin m'\theta, & j \text{ is even and } j' \text{ is odd} \\ \sin m\theta \cos m'\theta, & j \text{ is odd and } j' \text{ is even} \\ \sin m\theta \sin m'\theta, & j \text{ and } j' \text{ are both odd} \end{cases} \\ = \begin{cases} \pi(1 + \delta_{m0})\delta_{mm'}, & j \text{ and } j' \text{ are both even} \\ \pi\delta_{mm'}, & j \text{ and } j' \text{ are both odd} \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

Therefore, the Zernike polynomials are orthonormal according to

$$\int_0^1 \int_0^{2\pi} Z_j(\rho, \theta) Z_{j'}(\rho, \theta) \rho d\rho d\theta \bigg/ \int_0^1 \int_0^{2\pi} \rho d\rho d\theta = \delta_{jj'} \quad (13)$$

The expansion coefficients are given by

$$a_j = \frac{1}{\pi} \int_0^1 \int_0^{2\pi} W(\rho, \theta) Z_j(\rho, \theta) \rho d\rho d\theta \quad (14)$$

as may be seen by substituting Eq. (8) into Eq. (14) and using the orthonormality Eq. (13).

While the index n represents the radial *degree* or the *order* of a polynomial, since it represents the highest power of ρ in the polynomial, m is referred to as its *azimuthal frequency*. The index j is a *polynomial-ordering number* and is a function of both n and m . The polynomials are ordered such that an even j corresponds to a symmetric polynomial varying as $\cos m\theta$, while an odd j corresponds to an antisymmetric polynomial varying as $\sin m\theta$. A polynomial with a lower value of n is ordered first, and for a given value of n , a polynomial with a lower value of m is ordered first.

The Zernike circle polynomials are unique in that they are the only polynomials in two variables ρ and θ , which (a) are orthogonal over a circle, (b) are invariant in form with respect to rotation of the coordinate axes about the origin, and (c) include a polynomial for each permissible pair of n and m values.^{4,27}

The orthonormal Zernike circle polynomials and the names associated with some of them when identified with classical aberrations are listed in Table 1a for $n \leq 8$. The polynomials independent of θ are the spherical aberrations, those varying as $\cos\theta$ are the coma aberrations, and those varying as $\cos 2\theta$ are the astigmatism aberrations. The variation of several radial polynomials $R_n^m(\rho)$ with ρ is illustrated in Fig. 1.

TABLE 1a Orthonormal Zernike Circle Polynomials $Z_j(\rho, \theta)$ Ordered Such That an Even j Corresponds to a Symmetric Polynomial Varying as $\cos m\theta$, While an Odd j Corresponds to an Antisymmetric Polynomial Varying as $\sin m\theta$

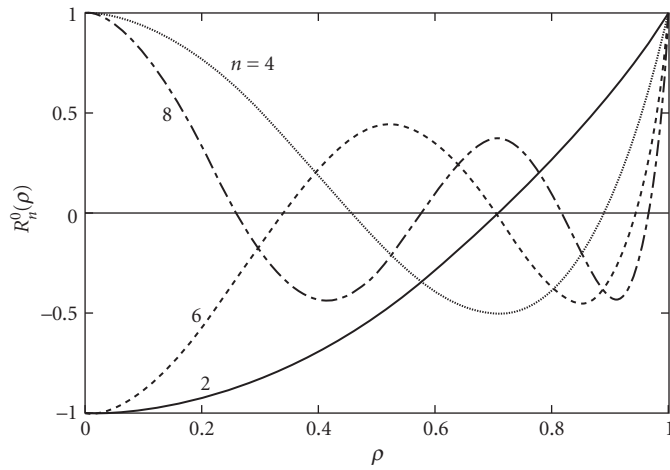
j	n	m	$Z_j(\rho, \theta)$	Aberration Name*
1	0	0	1	Piston
2	1	1	$2\rho\cos\theta$	x tilt
3	1	1	$2\rho\sin\theta$	y tilt
4	2	0	$\sqrt{3}(2\rho^2-1)$	Defocus
5	2	2	$\sqrt{6}\rho^2\sin 2\theta$	Primary astigmatism at 45°
6	2	2	$\sqrt{6}\rho^2\cos 2\theta$	Primary astigmatism at 0°
7	3	1	$\sqrt{8}(3\rho^3-2\rho)\sin\theta$	Primary y coma
8	3	1	$\sqrt{8}(3\rho^3-2\rho)\cos\theta$	Primary x coma
9	3	3	$\sqrt{8}\rho^3\sin 3\theta$	
10	3	3	$\sqrt{8}\rho^3\cos 3\theta$	
11	4	0	$\sqrt{5}(6\rho^4-6\rho^2+1)$	Primary spherical aberration
12	4	2	$\sqrt{10}(4\rho^4-3\rho^2)\cos 2\theta$	Secondary astigmatism at 0°
13	4	2	$\sqrt{10}(4\rho^4-3\rho^2)\sin 2\theta$	Secondary astigmatism at 45°
14	4	4	$\sqrt{10}\rho^4\cos 4\theta$	
15	4	4	$\sqrt{10}\rho^4\sin 4\theta$	
16	5	1	$\sqrt{12}(10\rho^5-12\rho^3+3\rho)\cos\theta$	Secondary x coma
17	5	1	$\sqrt{12}(10\rho^5-12\rho^3+3\rho)\sin\theta$	Secondary y coma
18	5	3	$\sqrt{12}(5\rho^5-4\rho^3)\cos 3\theta$	
19	5	3	$\sqrt{12}(5\rho^5-4\rho^3)\sin 3\theta$	
20	5	5	$\sqrt{12}\rho^5\cos 5\theta$	
21	5	5	$\sqrt{12}\rho^5\sin 5\theta$	
22	6	0	$\sqrt{7}(20\rho^6-30\rho^4+12\rho^2-1)$	Secondary spherical aberration
23	6	2	$\sqrt{14}(15\rho^6-20\rho^4+6\rho^2)\sin 2\theta$	Tertiary astigmatism at 45°
24	6	2	$\sqrt{14}(15\rho^6-20\rho^4+6\rho^2)\cos 2\theta$	Tertiary astigmatism at 0°
25	6	4	$\sqrt{14}(6\rho^6-5\rho^4)\sin 4\theta$	
26	6	4	$\sqrt{14}(6\rho^6-5\rho^4)\cos 4\theta$	
27	6	6	$\sqrt{14}6\rho^6\sin 6\theta$	
28	6	6	$\sqrt{14}\rho^6\cos 6\theta$	
29	7	1	$4(35\rho^7-60\rho^5+30\rho^3-4\rho)\sin\theta$	Tertiary y coma
30	7	1	$4(35\rho^7-60\rho^5+30\rho^3-4\rho)\cos\theta$	Tertiary x coma
31	7	3	$4(21\rho^7-30\rho^5+10\rho^3)\sin 3\theta$	
32	7	3	$4(21\rho^7-30\rho^5+10\rho^3)\cos 3\theta$	
33	7	5	$4(7\rho^7-6\rho^5)\sin 5\theta$	
34	7	5	$4(7\rho^7-6\rho^5)\cos 5\theta$	
35	7	7	$4\rho^7\sin 7\theta$	
36	7	7	$4\rho^7\cos 7\theta$	
37	8	0	$3(70\rho^8-140\rho^6+90\rho^4-20\rho^2+1)$	Tertiary spherical aberration
38	8	2	$\sqrt{18}(56\rho^8-105\rho^6+60\rho^4-10\rho^2)\cos 2\theta$	Quaternary astigmatism at 0°
39	8	2	$\sqrt{18}(56\rho^8-105\rho^6+60\rho^4-10\rho^2)\sin 2\theta$	Quaternary astigmatism at 45°
40	8	4	$\sqrt{18}(28\rho^8-42\rho^6+15\rho^4)\cos 4\theta$	

(Continued)

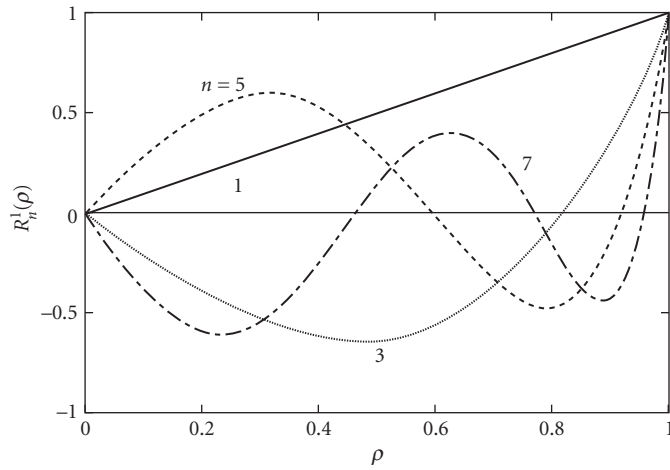
TABLE 1a Orthonormal Zernike Circle Polynomials $Z_j(\rho, \theta)$ Ordered Such That an Even j Corresponds to a Symmetric Polynomial Varying as $\cos m\theta$, While an Odd j Corresponds to an Antisymmetric Polynomial Varying as $\sin m\theta$ (Continued)

j	n	m	$Z_j(\rho, \theta)$	Aberration Name*
41	8	4	$\sqrt{18}(28\rho^8 - 42\rho^6 + 15\rho^4)\sin 4\theta$	
42	8	6	$\sqrt{18}(8\rho^8 - 7\rho^6)\cos 6\theta$	
43	8	6	$\sqrt{18}(8\rho^8 - 7\rho^6)\sin 6\theta$	
44	8	8	$\sqrt{18}\rho^8 \cos 8\theta$	
45	8	8	$\sqrt{18}\rho^8 \sin 8\theta$	

*The words *orthonormal Zernike circle* are to be associated with these names, e.g., *orthonormal Zernike circle primary astigmatism at 0°*.



(a)



(b)

FIGURE 1 Variation of a Zernike circle radial polynomial $R_n^m(\rho)$ with ρ : (a) defocus and spherical aberrations; (b) tilt and coma; and (c) astigmatism.

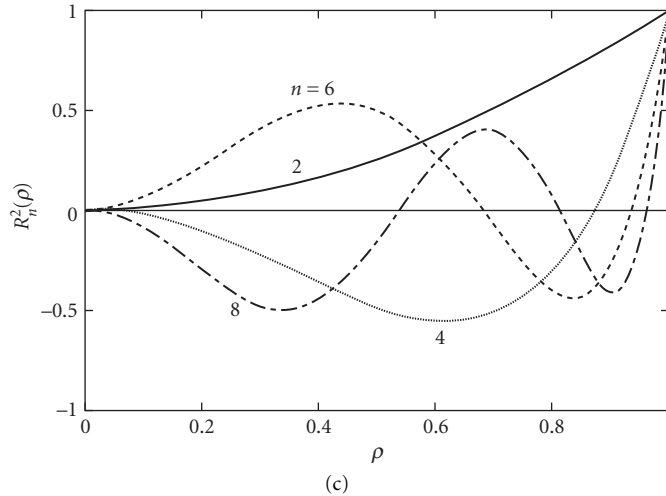


FIGURE 1 (Continued)

The number of polynomials of a given order n is $n + 1$. Their number through a certain order n is given by

$$N_n = (n+1)(n+2)/2 \quad (15)$$

For a rotationally symmetric imaging system, each of the $\sin m\theta$ terms is zero.^{4,28-32} Accordingly the number of polynomials of an even order is $(n/2) + 1$ and $(n + 1)/2$ for an odd order. Their number through an order n is given by

$$N_n = \left(\frac{n}{2} + 1\right)^2 \quad \text{for even } n \quad (16a)$$

$$= (n+1)(n+3)/4 \quad \text{for odd } n \quad (16b)$$

Relationships among the Indices n , m , and j

The number of polynomials N_n through a certain order n represents the largest value of j . Since the number of terms with the same value of n but different values of m is equal to $n + 1$, the smallest value of j for a given value of n is $N_n - n$. For a given value of n and m , there are two j values, $N_n - n + m - 1$ and $N_n - n + m$. The even value of j represents the $\cos m\theta$ term and the odd value of j represents the $\sin m\theta$ term. The value of j with $m = 0$ is $N_n - n$. For example, for $n = 5$, $N_n = 21$, and $j = 21$ represents the $\sin 5\theta$ term. The number of the corresponding $\cos 5\theta$ term is $j = 20$. The two terms with $m = 3$, for example, have j values of 18 and 19 representing the $\cos 3\theta$ and the $\sin 3\theta$ terms, respectively.

For a given value of j , n is given by

$$n = [(2j-1)^{1/2} + 0.5]_{\text{integer}} - 1 \quad (17)$$

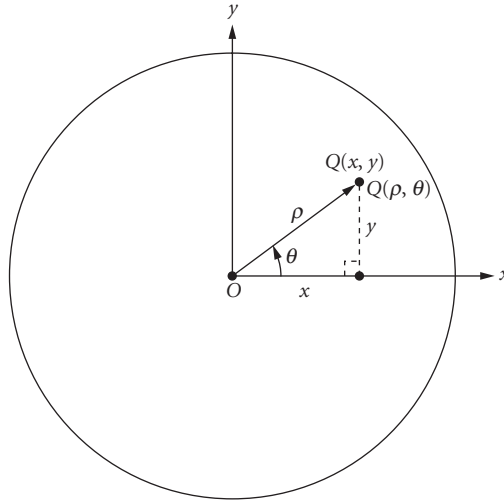


FIGURE 2 Cartesian and polar coordinates (x, y) and (ρ, θ) , respectively, of a point Q in the plane of a unit circle representing the circular exit pupil of an imaging system.

where the subscript integer implies the integer value of the number in brackets. Once n is known, the value of m is given by

$$m = \begin{cases} 2\{[2j+1-n(n+1)]/4\}_{\text{integer}} & \text{when } n \text{ is even} \\ 2\{[2(j+1)-n(n+1)]/4\}_{\text{integer}} - 1 & \text{when } n \text{ is odd} \end{cases} \quad (18a)$$

$$m = \begin{cases} 2\{[2j+1-n(n+1)]/4\}_{\text{integer}} & \text{when } n \text{ is even} \\ 2\{[2(j+1)-n(n+1)]/4\}_{\text{integer}} - 1 & \text{when } n \text{ is odd} \end{cases} \quad (18b)$$

For example, suppose we want to know the values of n and m for the term $j = 10$. From Eq. (17), $n = 3$ and from Eq. (18b), $m = 3$. Hence, it is a $\cos 3\theta$ term.

The polar coordinates (ρ, θ) and the Cartesian coordinates (x, y) of a pupil point Q , as illustrated in Fig. 2, are related to each other according to

$$(x, y) = \rho(\cos\theta, \sin\theta) \quad (19)$$

The circle polynomials in the Cartesian coordinates (x, y) of a pupil point are listed in Table 1b. It is quite common in the optics literature to consider a point object lying along the y axis when imaged by a rotationally symmetric optical system, thus making the yz plane the tangential plane.^{4,28-32} To maintain symmetry of the aberration function about this plane, the polar angle θ of a pupil point is accordingly defined as the angle made by its position vector OQ with the y axis, contrary to the standard convention as the angle with the x axis. We choose a point object along the x axis so that, for example, the coma aberration is expressed as $x(x^2 + y^2)$ and not as $y(x^2 + y^2)$. A positive value of our coma aberration yields a diffraction point spread function that is symmetric about the x axis (or symmetric in y) with its peak and centroid shifted to a positive value of x with respect to the Gaussian image point.

TABLE 1b Orthonormal Zernike Circle Polynomials $Z_j(x, y)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $0 \leq \rho = \sqrt{x^2 + y^2} \leq 1$

Polynomial	$Z_j(x, y)$
Z_1	1
Z_2	$2x$
Z_3	$2y$
Z_4	$\sqrt{3}(2\rho^2 - 1)$
Z_5	$2\sqrt{6}xy$
Z_6	$\sqrt{6}(x^2 - y^2)$
Z_7	$\sqrt{8}y(3\rho^2 - 2)$
Z_8	$\sqrt{8}x(3\rho^2 - 2)$
Z_9	$\sqrt{8}y(3x^2 - y^2)$
Z_{10}	$\sqrt{8}x(x^2 - 3y^2)$
Z_{11}	$\sqrt{5}(6\rho^4 - 6\rho^2 + 1)$
Z_{12}	$\sqrt{10}(x^2 - y^2)(4\rho^2 - 3)$
Z_{13}	$2\sqrt{10}xy(4\rho^2 - 3)$
Z_{14}	$\sqrt{10}(\rho^4 - 8x^2y^2)$
Z_{15}	$4\sqrt{10}xy(x^2 - y^2)$
Z_{16}	$\sqrt{12}x(10\rho^4 - 12\rho^2 + 3)$
Z_{17}	$\sqrt{12}y(10\rho^4 - 12\rho^2 + 3)$
Z_{18}	$\sqrt{12}x(x^2 - 3y^2)(5\rho^2 - 4)$
Z_{19}	$\sqrt{12}y(3x^2 - y^2)(5\rho^2 - 4)$
Z_{20}	$\sqrt{12}x(16x^4 - 20x^2\rho^2 + 5\rho^4)$
Z_{21}	$\sqrt{12}y(16y^4 - 20y^2\rho^2 + 5\rho^4)$
Z_{22}	$\sqrt{7}(20\rho^6 - 30\rho^4 + 12\rho^2 - 1)$
Z_{23}	$2\sqrt{14}xy(15\rho^2 - 20\rho^2 + 6)$
Z_{24}	$\sqrt{14}(x^2 - y^2)(15\rho^4 - 20\rho^2 + 6)$
Z_{25}	$4\sqrt{14}xy(x^2 - y^2)(6\rho^2 - 5)$
Z_{26}	$\sqrt{14}(8x^4 - 8x^2\rho^2 + \rho^4)(6\rho^2 - 5)$
Z_{27}	$\sqrt{14}xy(32x^4 - 32x^2\rho^2 + 6\rho^4)$
Z_{28}	$\sqrt{14}(32x^6 - 48x^4\rho^2 + 18x^2\rho^4 - \rho^6)$
Z_{29}	$4y(35\rho^6 - 60\rho^4 + 30\rho^2 - 4)$
Z_{30}	$4x(35\rho^6 - 60\rho^4 + 30\rho^2 - 4)$
Z_{31}	$4y(3x^2 - y^2)(21\rho^4 - 30\rho^2 + 10)$
Z_{32}	$4x(x^2 - 3y^2)(21\rho^4 - 30\rho^2 + 10)$
Z_{33}	$4(7\rho^2 - 6)[4x^2y(x^2 - y^2) + y(\rho^4 - 8x^2y^2)]$
Z_{34}	$4(7\rho^2 - 6)[x(\rho^4 - 8x^2y^2) - 4xy^2(x^2 - y^2)]$
Z_{35}	$8x^2y(3\rho^4 - 16x^2y^2) + 4y(x^2 - y^2)(\rho^4 - 16x^2y^2)$
Z_{36}	$4x(x^2 - y^2)(\rho^4 - 16x^2y^2) - 8xy^2(3\rho^4 - 16x^2y^2)$
Z_{37}	$3(70\rho^8 - 140\rho^6 + 90\rho^4 - 20\rho^2 + 1)$
Z_{38}	$\sqrt{18}(56\rho^6 - 105\rho^4 + 60\rho^2 - 10)(x^2 - y^2)$
Z_{39}	$2\sqrt{18}xy(56\rho^6 - 105\rho^4 + 60\rho^2 - 10)$
Z_{40}	$\sqrt{18}(28\rho^4 - 42\rho^2 + 15)(\rho^4 - 8x^2y^2)$
Z_{41}	$4\sqrt{18}xy(28\rho^4 - 42\rho^2 + 15)(x^2 - y^2)$
Z_{42}	$\sqrt{18}(x^2 - y^2)(\rho^4 - 16x^2y^2)(8\rho^2 - 7)$
Z_{43}	$2\sqrt{18}xy(3\rho^4 - 16x^2y^2)$
Z_{44}	$2\sqrt{18}(\rho^4 - 8x^2y^2)^2 - \rho^8$
Z_{45}	$8\sqrt{18}xy(x^2 - y^2)(\rho^4 - 8x^2y^2)$

11.5 ZERNIKE ANNULAR POLYNOMIALS

The aberration function $W(\rho, \theta; \epsilon)$ across a *unit annulus* with an obscuration ratio ϵ , representing the ratio of its inner and outer radii, as illustrated in Fig. 3a, can be expanded in terms of a complete set of *Zernike annular polynomials* $Z_j(\rho, \theta; \epsilon)$ that are orthonormal over the unit annulus in the form⁸⁻¹¹

$$W(\rho, \theta; \epsilon) = \sum_j a_j Z_j(\rho, \theta; \epsilon) \quad (20)$$

where a_j is an expansion coefficient of the polynomial, $\epsilon \leq \rho \leq 1$ and $0 \leq \theta < 2\pi$. The annular polynomials are written in a manner similar to the circle polynomials. Thus

$$Z_{\text{even } j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)} R_n^m(\rho; \epsilon) \cos m\theta, m \neq 0 \quad (21a)$$

$$Z_{\text{odd } j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)} R_n^m(\rho; \epsilon) \sin m\theta, m \neq 0 \quad (21b)$$

$$Z_j(\rho, \theta; \epsilon) = \sqrt{n+1} R_n^0(\rho; \epsilon), m = 0 \quad (21c)$$

where n and m are positive integers (including zero) and $n - m \geq 0$ and even. The radial annular polynomials $R_n^m(\rho; \epsilon)$ obey the orthogonality relation

$$\int_{\epsilon}^1 R_n^m(\rho; \epsilon) R_n^m(\rho; \epsilon) \rho d\rho = \frac{1-\epsilon^2}{2(n+1)} \delta_{nn'} \quad (22)$$

Accordingly, the annular polynomials obey the orthonormality condition

$$\int_{\epsilon}^1 \int_0^{2\pi} Z_j(\rho, \theta; \epsilon) Z_j(\rho, \theta; \epsilon) \rho d\rho d\theta \bigg/ \int_{\epsilon}^1 \int_0^{2\pi} \rho d\rho d\theta = \delta_{jj'} \quad (23)$$

The Zernike expansion coefficients are given by

$$a_j = \frac{1}{\pi(1-\epsilon)^2} \int_{\epsilon}^1 \int_0^{2\pi} W(\rho, \theta; \epsilon) Z_j(\rho, \theta; \epsilon) \rho d\rho d\theta \quad (24)$$

as may be seen by substituting Eq. (20) into Eq. (24) and using Eq. (23) for the orthonormality of the polynomials.

The annular polynomials are similar to the circle polynomials, except that they are orthogonal over an annular pupil. They can be obtained from the circle polynomials by the Gram-Schmidt orthogonalization process.¹⁷ The radial polynomials are accordingly given by

$$R_n^m(\rho; \epsilon) = N_n^m \left[R_n^m(\rho) - \sum_{i \geq 1}^{(n-m)/2} (n-2i+1) \langle R_n^m(\rho) R_{n-2i}^m(\rho; \epsilon) \rangle R_{n-2i}^m(\rho; \epsilon) \right] \quad (25)$$

where

$$\langle R_n^m(\rho) R_n^m(\rho; \epsilon) \rangle = \frac{2}{1-\epsilon^2} \int_{\epsilon}^1 R_n^m(\rho) R_n^m(\rho; \epsilon) \rho d\rho \quad (26)$$

and N_n^m is a normalization constant such that the radial polynomials satisfy the orthogonality Eq. (22). Thus, $R_n^m(\rho; \epsilon)$ is a radial polynomial of degree n in ρ containing terms in $\rho^n, \rho^{n-2}, \dots$, and ρ^m with coefficients that depend on ϵ . The radial polynomials are even or odd in ρ depending on whether n

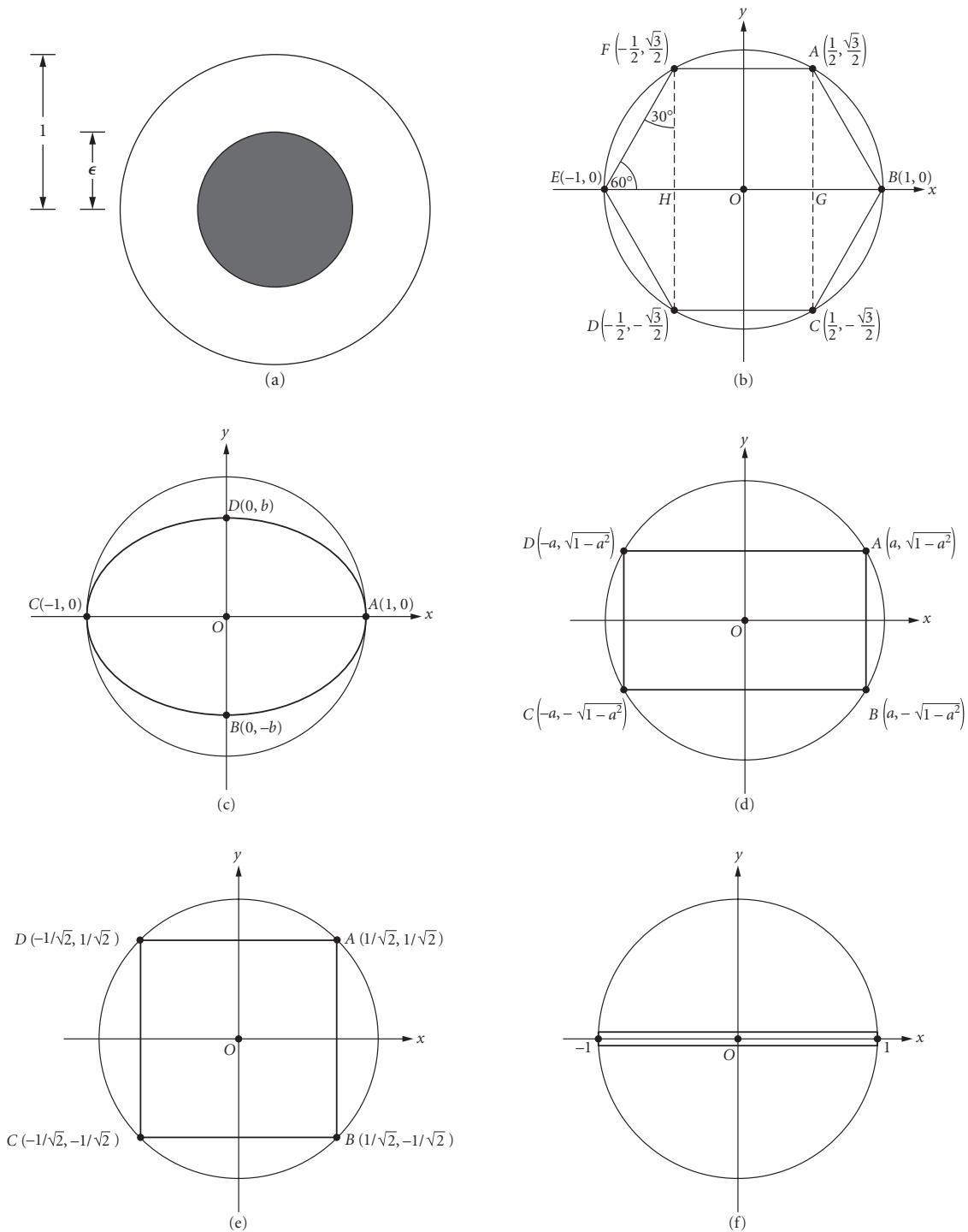


FIGURE 3 Unit pupils inscribed inside a unit circle: (a) annulus of obscuration ratio ϵ ; (b) hexagon; (c) ellipse of aspect ratio b ; (d) rectangle of half width a ; (e) square of half width $1/\sqrt{2}$; and (f) slit.

(or m) is even or odd. For $m = 0$, the radial polynomials are equal to the Legendre polynomials $P_n(\cdot)$ according to

$$R_{2n}^0(\rho; \epsilon) = P_n \left[\frac{2(\rho^2 - \epsilon^2)}{1 - \epsilon^2} - 1 \right] \quad (27)$$

Thus, they can be obtained from the circle radial polynomials $R_{2n}^0(\rho)$ by replacing ρ by $[(\rho^2 - \epsilon^2)/(1 - \epsilon^2)]^{1/2}$, that is,

$$R_{2n}^0(\rho; \epsilon) = R_{2n}^0 \left[\left(\frac{\rho^2 - \epsilon^2}{1 - \epsilon^2} \right)^{1/2} \right] \quad (28)$$

It can be seen from Eqs. (22) and (25) that

$$R_n^n(\rho; \epsilon) = \rho^n / \left(\sum_{i=0}^n \epsilon^{2i} \right)^{1/2} \quad (29)$$

$$= \rho^n \{ (1 - \epsilon^2) / [1 - \epsilon^{2(n+1)}] \}^{1/2} \quad (30)$$

Moreover,

$$R_2^{n-2}(\rho; \epsilon) = \frac{n\rho^n - (n-1)[(1 - \epsilon^{2n})/(1 - \epsilon^{2(n-1)})]\rho^{n-2}}{\{(1 - \epsilon^2)^{-1}[n^2(1 - \epsilon^{2(n+1)}) - (n^2 - 1)(1 - \epsilon^{2n})^2 / (1 - \epsilon^{2(n-1)})]\}^{1/2}} \quad (31)$$

It is evident that the radial polynomial $R_n^m(\rho; \epsilon)$ differs from the corresponding circle polynomial $R_n^m(\rho)$ only in its normalization. We also note that

$$\begin{aligned} R_n^m(1; \epsilon) &= 1, \quad m = 0 \\ &\neq 1, \quad m \neq 0 \end{aligned} \quad (32)$$

The variation of several Zernike annular radial polynomials with ρ is shown in Fig. 4 for $\epsilon = 0.5$.

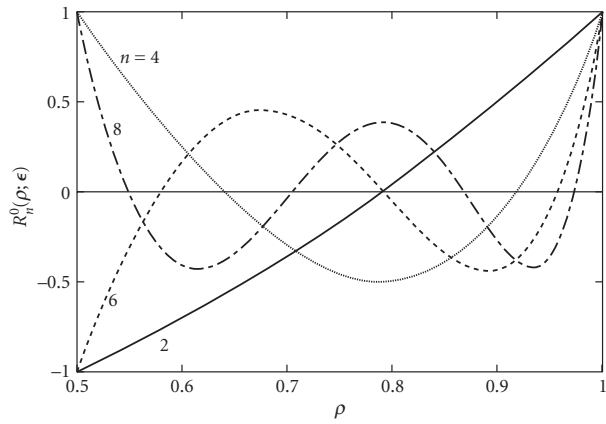
It is evident from Eqs. (21) that the annular polynomials, like the circle polynomials, are separable in the polar coordinates ρ and θ . This is a consequence of the radial symmetry of the annular pupil. As may be evident from the Gram-Schmidt orthogonalization process, each annular polynomial is a linear combination of the circle polynomials.³³ Accordingly, each radial polynomial $R_n^m(\rho; \epsilon)$ can be written as a linear combination of the polynomials $R_n^m(\rho)$, $R_{n-2}^m(\rho)$, \dots , and $R_m^m(\rho)$. For example,

$$R_3^1(\rho; \epsilon) = \frac{1}{(1 - \epsilon^2)(1 + 5\epsilon^2 + 5\epsilon^4 + \epsilon^6)^{1/2}} [(1 + \epsilon^2)R_3^1(\rho) - 2\epsilon^4 R_1^1(\rho)] \quad (33a)$$

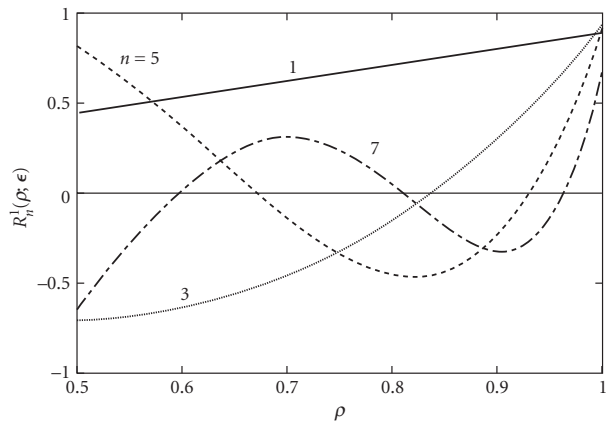
and

$$R_4^0(\rho; \epsilon) = \frac{1}{(1 - \epsilon^2)^2} [R_4^0(\rho) - 3\epsilon^2 R_2^0(\rho) + \epsilon^2(1 + \epsilon^2)R_0^0(\rho)] \quad (33b)$$

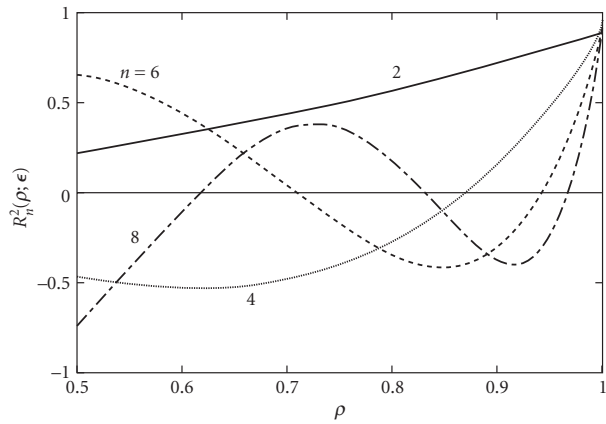
The Zernike annular radial polynomials for $n \leq 8$ are listed in Table 2a. The number polynomials of a certain order or through a certain order n is given by the same expressions as in the case of Zernike circle polynomials. Table 2b lists the full annular polynomials illustrating their ordering. In Table 2c, they are given in the Cartesian coordinates.



(a)



(b)



(c)

FIGURE 4 Variation of a Zernike *annular* radial polynomial $R_n^m(\rho; \epsilon)$ with ρ for $\epsilon = 0.5$: (a) defocus and spherical aberrations; (b) tilt and coma; and (c) astigmatism.

TABLE 2a Zernike Annular Radial Polynomials $R_n^m(\rho; \epsilon)$, Where ϵ Is the Obscuration Ratio of Annular Pupil and $\epsilon \leq \rho \leq 1$

n	m	$R_n^m(\rho; \epsilon)$
0	0	1
1	1	$\rho/(1+\epsilon^2)^{1/2}$
2	0	$(2\rho^2-1-\epsilon^2)/(1-\epsilon^2)$
2	2	$\rho^2/(1+\epsilon^2+\epsilon^2)^{1/2}$
3	1	$\frac{3(1+\epsilon^2)\rho^3-2(1+\epsilon^2+\epsilon^4)\rho}{(1-\epsilon^2)[(1+\epsilon^2)(1+4\epsilon^2+\epsilon^4)]^{1/2}}$
3	3	$\rho^3/(1+\epsilon^2+\epsilon^4+\epsilon^6)^{1/2}$
4	0	$[6\rho^4-6(1+\epsilon^2)\rho^2+1+4\epsilon^2+\epsilon^4]/(1-\epsilon^2)^2$
4	2	$\frac{4\rho^4-3[(1-\epsilon^8)/(1-\epsilon^6)]\rho^2}{\{(1-\epsilon^2)^{-1}[16(1-\epsilon^{10})-15(1-\epsilon^8)^2/(1-\epsilon^6)]\}^{1/2}}$
4	4	$\rho^4/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8)^{1/2}$
5	1	$\frac{10(1+4\epsilon^2+\epsilon^4)\rho^5-12(1+4\epsilon^2+4\epsilon^4+\epsilon^6)\rho^3+3(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)\rho}{(1-\epsilon^2)^2[(1+4\epsilon^2+\epsilon^4)(1+9\epsilon^2+9\epsilon^4+\epsilon^6)]^{1/2}}$
5	3	$\frac{5\rho^5-4[(1-\epsilon^{10})/(1-\epsilon^8)]\rho^3}{\{(1-\epsilon^2)^{-1}[25(1-\epsilon^{12})-24(1-\epsilon^{10})^2/(1-\epsilon^8)]\}^{1/2}}$
5	5	$\rho^5/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10})^{1/2}$
6	0	$[20\rho^6-30(1+\epsilon^2)\rho^4+12(1+3\epsilon^2+\epsilon^4)\rho^2-(1+9\epsilon^2+9\epsilon^4+\epsilon^6)]/(1-\epsilon^2)^3$ $15(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)\rho^6-20(1+4\epsilon^2+10\epsilon^4+10\epsilon^6+4\epsilon^8+\epsilon^{10})\rho^4$
6	2	$\frac{+6(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})\rho^2}{(1+\epsilon^2)^2[(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})]^{1/2}}$
6	4	$\frac{6\rho^6-5[(1-\epsilon^{12})/(1-\epsilon^{10})]\rho^4}{\{(1-\epsilon^2)^{-1}[36(1-\epsilon^{14})-35(1-\epsilon^{12})^2/(1-\epsilon^{10})]\}^{1/2}}$
6	6	$\rho^6/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})^{1/2}$
7	1	$a_7^1\rho^7+b_7^1\rho^5+c_7^1\rho^3+d_7^1\rho$
7	3	$a_7^3\rho^7+b_7^3\rho^5+c_7^3\rho^3$
7	5	$\frac{7\rho^7-6[(1-\epsilon^{14})/(1-\epsilon^{12})]\rho^5}{\{(1-\epsilon^2)^{-1}[49(1-\epsilon^{16})-48(1-\epsilon^{14})^2/(1-\epsilon^{12})]\}^{1/2}}$
7	7	$\rho^7/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14})^{1/2}$
8	0	$\frac{70\rho^8-140(1+\epsilon^2)\rho^6+30(3+8\epsilon^2+3\epsilon^4)\rho^4-20(1+6\epsilon^2+6\epsilon^4+\epsilon^6)\rho^2+\epsilon_8^0}{(1-\epsilon^2)^4}$
8	2	$a_8^2\rho^8+b_8^2\rho^6+c_8^2\rho^4+d_8^2\rho^2$
8	4	$a_8^4\rho^8+b_8^4\rho^6+c_8^4\rho^4$
8	6	$a_8^6\rho^8+b_8^6\rho^6$
8	8	$\rho^8/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14}+\epsilon^{16})^{1/2}$

(Continued)

TABLE 2a Zernike Annular Radial Polynomials $R_n^m(\rho; \epsilon)$, Where ϵ Is the Obscuration Ratio of Annular Pupil and $\epsilon \leq \rho \leq 1$ (Continued)

$$a_7^1 = 35(1+9\epsilon^2+9\epsilon^4+\epsilon^6)/A_7^1$$

$$b_7^1 = -60(1+9\epsilon^2+15\epsilon^4+9\epsilon^6+\epsilon^8)/A_7^1$$

$$c_7^1 = 30(1+9\epsilon^2+25\epsilon^4+25\epsilon^6+9\epsilon^8+\epsilon^{10})/A_7^1$$

$$d_7^1 = -4(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})/A_7^1$$

$$A_7^1 = (1-\epsilon^2)^3(1+9\epsilon^2+9\epsilon^4+\epsilon^6)^{1/2}(1+16\epsilon^2+36\epsilon^4+16\epsilon^6+\epsilon^8)^{1/2}$$

$$a_7^3 = 21(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})/A_7^3$$

$$b_7^3 = -30(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+20\epsilon^8+10\epsilon^{10}+4\epsilon^{12}+\epsilon^{14})/A_7^3$$

$$c_7^3 = 10(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})/A_7^3$$

$$A_7^3 = (1-\epsilon^2)^2(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})^{1/2}$$

$$\quad \times (1+9\epsilon^2+45\epsilon^4+165\epsilon^6+270\epsilon^8+270\epsilon^{10}+165\epsilon^{12}+45\epsilon^{14}+9\epsilon^{16}+\epsilon^{18})^{1/2}$$

$$e_8^0 = 1+16\epsilon^2+36\epsilon^4+16\epsilon^6+\epsilon^8$$

$$a_8^2 = 56(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})/A_8^2$$

$$b_8^2 = -105(1+9\epsilon^2+45\epsilon^4+85\epsilon^6+85\epsilon^8+45\epsilon^{10}+9\epsilon^{12}+\epsilon^{14})/A_8^2$$

$$c_8^2 = 60(1+9\epsilon^2+45\epsilon^4+115\epsilon^6+150\epsilon^8+115\epsilon^{10}+45\epsilon^{12}+9\epsilon^{14}+\epsilon^{16})/A_8^2$$

$$d_8^2 = -10(1+9\epsilon^2+45\epsilon^4+165\epsilon^6+270\epsilon^8+270\epsilon^{10}+165\epsilon^{12}+45\epsilon^{14}+9\epsilon^{16}+\epsilon^{18})/A_8^2$$

$$A_8^2 = (1-\epsilon^2)^3(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})^{1/2}$$

$$\quad \times (1+16\epsilon^2+136\epsilon^4+416\epsilon^6+626\epsilon^8+416\epsilon^{10}+136\epsilon^{12}+16\epsilon^{14}+\epsilon^{16})^{1/2}$$

$$a_8^4 = 28(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})/A_8^4$$

$$b_8^4 = -42(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+35\epsilon^{10}+20\epsilon^{12}+10\epsilon^{14}+4\epsilon^{16}+\epsilon^{18})/A_8^4$$

$$c_8^4 = 15(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+56\epsilon^{10}+35\epsilon^{12}+20\epsilon^{14}+10\epsilon^{16}+4\epsilon^{18}+\epsilon^{20})/A_8^4$$

$$A_8^4 = (1-\epsilon^2)^2(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})^{1/2}$$

$$\quad \times (1+9\epsilon^2+45\epsilon^4+165\epsilon^6+495\epsilon^8+846\epsilon^{10}+994\epsilon^{12}+846\epsilon^{14}+495\epsilon^{16}+165\epsilon^{18}+45\epsilon^{20}+9\epsilon^{22}+\epsilon^{24})^{1/2}$$

$$a_8^6 = 8(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})/A_8^6$$

$$b_8^6 = -7(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14})/A_8^6$$

$$A_8^6 = (1-\epsilon^2)(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})^{1/2}$$

$$\quad \times (1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+56\epsilon^{10}+84\epsilon^{12}+56\epsilon^{14}+35\epsilon^{16}+20\epsilon^{18}+10\epsilon^{20}+4\epsilon^{22}+\epsilon^{24})^{1/2}$$

TABLE 2b Orthonormal Zernike *Annular* Polynomials $Z_j(\rho, \theta; \epsilon)$, Ordered in the Same Manner as the Zernike Circle Polynomials in Table 1a

j	n	m	$Z_j(\rho, \theta; \epsilon)^*$	Aberration Name*
1	0	0	$R_0^0(\rho; \epsilon)=1$	Piston
2	1	1	$2R_1^1(\rho; \epsilon)\cos\theta$	x tilt
3	1	1	$2R_1^1(\rho; \epsilon)\sin\theta$	y tilt
4	2	0	$\sqrt{3}R_2^0(\rho; \epsilon)$	Defocus
5	2	2	$\sqrt{6}R_2^2(\rho; \epsilon)\sin 2\theta$	Primary astigmatism at 45°
6	2	2	$\sqrt{6}R_2^2(\rho; \epsilon)\cos 2\theta$	Primary astigmatism at 0°
7	3	1	$\sqrt{8}R_3^1(\rho; \epsilon)\sin\theta$	Primary y coma
8	3	1	$\sqrt{8}R_3^1(\rho; \epsilon)\cos\theta$	Primary x coma
9	3	3	$\sqrt{8}R_3^3(\rho; \epsilon)\sin 3\theta$	
10	3	3	$\sqrt{8}R_3^3(\rho; \epsilon)\cos 3\theta$	
11	4	0	$\sqrt{5}R_4^0(\rho; \epsilon)$	Primary spherical aberration
12	4	2	$\sqrt{10}R_4^2(\rho; \epsilon)\cos 2\theta$	Secondary astigmatism at 0°
13	4	2	$\sqrt{10}R_4^2(\rho; \epsilon)\sin 2\theta$	Secondary astigmatism at 45°
14	4	4	$\sqrt{10}R_4^4(\rho; \epsilon)\cos 4\theta$	
15	4	4	$\sqrt{10}R_4^4(\rho; \epsilon)\sin 4\theta$	
16	5	1	$\sqrt{12}R_5^1(\rho; \epsilon)\cos\theta$	Secondary x coma
17	5	1	$\sqrt{12}R_5^1(\rho; \epsilon)\sin\theta$	Secondary y coma
18	5	3	$\sqrt{12}R_5^3(\rho; \epsilon)\cos 3\theta$	
19	5	3	$\sqrt{12}R_5^3(\rho; \epsilon)\sin 3\theta$	
20	5	5	$\sqrt{12}R_5^5(\rho; \epsilon)\cos 5\theta$	
21	5	5	$\sqrt{12}R_5^5(\rho; \epsilon)\sin 5\theta$	
22	6	0	$\sqrt{7}R_6^0(\rho; \epsilon)$	Secondary spherical aberration
23	6	2	$\sqrt{14}R_6^2(\rho; \epsilon)\sin 2\theta$	Tertiary astigmatism at 45°
24	6	2	$\sqrt{14}R_6^2(\rho; \epsilon)\cos 2\theta$	Tertiary astigmatism at 0°
25	6	4	$\sqrt{14}R_6^4(\rho; \epsilon)\cos 4\theta$	
26	6	4	$\sqrt{14}R_6^4(\rho; \epsilon)\sin 4\theta$	
27	6	6	$\sqrt{14}R_6^6(\rho; \epsilon)\sin 6\theta$	
28	6	6	$\sqrt{14}R_6^6(\rho; \epsilon)\cos 6\theta$	
29	7	1	$4R_7^1(\rho; \epsilon)\sin\theta$	
30	7	1	$4R_7^1(\rho; \epsilon)\cos\theta$	
31	7	3	$4R_7^3(\rho; \epsilon)\sin 3\theta$	
32	7	3	$4R_7^3(\rho; \epsilon)\cos 3\theta$	
33	7	5	$4R_7^5(\rho; \epsilon)\sin 5\theta$	
34	7	5	$4R_7^5(\rho; \epsilon)\cos 5\theta$	
35	7	7	$4R_7^7(\rho; \epsilon)\sin 7\theta$	
36	7	7	$4R_7^7(\rho; \epsilon)\cos 7\theta$	
37	8	0	$3R_8^0(\rho; \epsilon)$	Tertiary spherical aberration
38	8	2	$\sqrt{18}R_8^2(\rho; \epsilon)\cos 2\theta$	Quaternary astigmatism at 0°
39	8	2	$\sqrt{18}R_8^2(\rho; \epsilon)\sin 2\theta$	Quaternary astigmatism at 45°
40	8	4	$\sqrt{18}R_8^4(\rho; \epsilon)\cos 4\theta$	
41	8	4	$\sqrt{18}R_8^4(\rho; \epsilon)\sin 4\theta$	
42	8	6	$\sqrt{18}R_8^6(\rho; \epsilon)\cos 6\theta$	
43	8	6	$\sqrt{18}R_8^6(\rho; \epsilon)\sin 6\theta$	
44	8	8	$\sqrt{18}R_8^8(\rho; \epsilon)\cos 8\theta$	
45	8	8	$\sqrt{18}R_8^8(\rho; \epsilon)\sin 8\theta$	

*The words “orthonormal Zernike annular” should be added to the name, e.g., orthonormal Zernike annular primary spherical aberration.

TABLE 2c Orthonormal Zernike *Annular* Polynomials $Z_j(x, y; \epsilon)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $\epsilon \leq \rho = \sqrt{x^2 + y^2} \leq 1$

Polynomial	$Z_j(x, y; \epsilon)$
Z_1	1
Z_2	$2x/(1+\epsilon^2)^{1/2}$
Z_3	$2y/(1+\epsilon^2)^{1/2}$
Z_4	$\sqrt{3}(2\rho^2 - 1 - \epsilon^2)/(1 - \epsilon^2)$
Z_5	$2\sqrt{6}xy/(1 + \epsilon^2 + \epsilon^4)^{1/2}$
Z_6	$\sqrt{6}(x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4)^{1/2}$
Z_7	$\frac{\sqrt{8}y[3(1 + \epsilon^2)\rho^2 - 2(1 + \epsilon^2 + \epsilon^4)]}{(1 - \epsilon^2)[1 + \epsilon^2 + (1 + 4\epsilon^2 + \epsilon^4)]^{1/2}}$
Z_8	$\frac{\sqrt{8}x[3(1 + \epsilon^2)\rho^2 - 2(1 + \epsilon^2 + \epsilon^4)]}{(1 - \epsilon^2)[1 + \epsilon^2 + (1 + 4\epsilon^2 + \epsilon^4)]^{1/2}}$
Z_9	$\sqrt{8}y(3x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6)^{1/2}$
Z_{10}	$\sqrt{8}x(x^2 - 3y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6)^{1/2}$
Z_{11}	$\sqrt{5}[6\rho^4 - 6(1 + \epsilon^2)\rho^2 + (1 + 4\epsilon^2 + \epsilon^4)]/(1 - \epsilon^2)^2$
Z_{12}	$\frac{\sqrt{10}(x^2 - y^2)[4\rho^2 - 3(1 - \epsilon^8)]/(1 - \epsilon^6)}{\{(1 - \epsilon^2)^{-1}[16(1 - \epsilon^{10}) - 15(1 - \epsilon^8)^2/(1 - \epsilon^{12})]\}^{1/2}}$
Z_{13}	$\frac{2\sqrt{10}xy[4\rho^2 - 3(1 - \epsilon^8)]/(1 - \epsilon^6)}{\{(1 - \epsilon^2)^{-1}[16(1 - \epsilon^{10}) - 15(1 - \epsilon^8)^2/(1 - \epsilon^6)]\}^{1/2}}$
Z_{14}	$\sqrt{10}(\rho^4 - 8x^2y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8)^{1/2}$
Z_{15}	$4\sqrt{10}xy(x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8)^{1/2}$
Z_{16}	$\frac{\sqrt{12}x[10(1 + 4\epsilon^2 + \epsilon^4)\rho^4 - 12(1 + 4\epsilon^2 + 4\epsilon^4 + \epsilon^6)\rho^2 + 3(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + \epsilon^4)(1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]^{1/2}}$
Z_{17}	$\frac{\sqrt{12}y[10(1 + 4\epsilon^2 + \epsilon^4)\rho^4 - 12(1 + 4\epsilon^2 + 4\epsilon^4 + \epsilon^6)\rho^2 + 3(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + \epsilon^4)(1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]^{1/2}}$
Z_{18}	$\frac{\sqrt{12}x(x^2 - 3y^2)[5\rho^2 - 4(1 - \epsilon^{10})/(1 - \epsilon^8)]}{\{(1 - \epsilon^2)^{-1}[25(1 - \epsilon^{12}) - 24(1 - \epsilon^{10})^2/(1 - \epsilon^8)]\}^{1/2}}$
Z_{19}	$\frac{\sqrt{12}y(3x^2 - y^2)[5\rho^2 - 4(1 - \epsilon^{10})/(1 - \epsilon^8)]}{\{(1 - \epsilon^2)^{-1}[25(1 - \epsilon^{12}) - 24(1 - \epsilon^{10})^2/(1 - \epsilon^8)]\}^{1/2}}$
Z_{20}	$\sqrt{12}x(16x^4 - 20x^2\rho^2 + 5\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10})^{1/2}$
Z_{21}	$\sqrt{12}y(16y^4 - 20y^2\rho^2 + 5\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10})^{1/2}$
Z_{22}	$\sqrt{7}[20\rho^6 - 30(1 + \epsilon^2)\rho^4 + 12(1 + 3\epsilon^2 + \epsilon^4)\rho^2 - (1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]/(1 - \epsilon^2)^3$
Z_{23}	$\frac{2\sqrt{14}xy[15(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)\rho^4 - 20(1 + 4\epsilon^2 + 10\epsilon^4 + 10\epsilon^6 + 4\epsilon^8 + \epsilon^{10})\rho^2 + 6(1 + 4\epsilon^2 + 10\epsilon^4 + 20\epsilon^6 + 10\epsilon^8 + 4\epsilon^{10} + \epsilon^{12})]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)(1 + 9\epsilon^2 + 45\epsilon^4 + 65\epsilon^6 + 45\epsilon^8 + 9\epsilon^{10} + \epsilon^{12})]^{1/2}}$
Z_{24}	$\frac{\sqrt{14}(x^2 - y^2)[15(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)\rho^4 - 20(1 + 4\epsilon^2 + 10\epsilon^4 + 10\epsilon^6 + 4\epsilon^8 + \epsilon^{10})\rho^2 + 6(1 + 4\epsilon^2 + 10\epsilon^4 + 20\epsilon^6 + 10\epsilon^8 + 4\epsilon^{10} + \epsilon^{12})]}{(1 - \epsilon^2)^2(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)(1 + 9\epsilon^2 + 45\epsilon^4 + 65\epsilon^6 + 45\epsilon^8 + 9\epsilon^{10} + \epsilon^{12})^{1/2}}$

(Continued)

TABLE 2c Orthonormal Zernike Annular Polynomials $Z_j(x, y; \epsilon)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $\epsilon \leq \rho = \sqrt{x^2 + y^2} \leq 1$ (Continued)

Polynomial	$Z_j(x, y; \epsilon)$
Z_{25}	$\frac{4\sqrt{14}xy(x^2 - y^2)[6\rho^2 - 5(1 - \epsilon^{12})/(1 - \epsilon^{10})]}{\{(1 - \epsilon^2)^{-1}[36(1 - \epsilon^{14}) - 35(1 - \epsilon^{12})^2/(1 - \epsilon^{10})]\}^{1/2}}$
Z_{26}	$\frac{\sqrt{14}(8x^4 - 8x^2\rho^2 + \rho^4)[6\rho^2 - 5(1 - \epsilon^{12})/(1 - \epsilon^{10})]}{\{(1 - \epsilon^2)^{-1}[36(1 - \epsilon^{14}) - 35(1 - \epsilon^{12})^2/(1 - \epsilon^{10})]\}^{1/2}}$
Z_{27}	$\sqrt{14}xy(32x^4 - 32x^2\rho^2 + 6\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10} + \epsilon^{12})^{1/2}$
Z_{28}	$\sqrt{14}(32x^6 - 48x^4\rho^2 + 18x^2\rho^4 - \rho^6)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10} + \epsilon^{12})^{1/2}$

11.6 HEXAGONAL POLYNOMIALS

Figure 3b shows a *unit hexagon* inscribed inside a unit circle. Each side of the hexagon has a length of unity. The area of the hexagon is $A = 3\sqrt{3}/2$. The orthonormality of the hexagonal polynomials $H_j(x, y)$ implies that²⁶

$$\frac{2}{3\sqrt{3}} \int_{\text{hexagon}} H_j(x, y) H_{j'}(x, y) dx dy = \delta_{jj'} \quad (34)$$

The orthonormal hexagonal polynomials are given in Tables 3 up to the eighth order in three different but equivalent forms. In Table 3a, each hexagonal polynomial is written in terms of the circle polynomials, thus illustrating the relationship between the two. In particular, it helps determine the potential error made when a hexagonal aberration function is expanded in terms of the circle polynomials.³⁴ The polynomials up to H_{19} are given in their analytical form, but those with $j > 19$ are written in a numerical form because of the increasing complexity of the coefficients of the circle polynomials. In Table 3b, the hexagonal polynomials are given in polar coordinates, showing one-to-one correspondence with the circle polynomials, but illustrating the difference from them. This form is convenient for analytical calculations because of the integration of trigonometric functions over symmetric limits. Finally, in Table 3c, they are given in Cartesian coordinates, as they would be used for any quantitative numerical analysis of, say, an interferogram.

From Table 3a, we note that each hexagonal polynomial consists of cosine or sine terms, but not both. Unlike the circle,³⁻⁶ annular,⁸⁻¹¹ or Gauss^{23,24} polynomials, the hexagonal polynomials are generally not separable in ρ and θ due to the lack of radial symmetry of the hexagonal pupil. The first 13 polynomials, that is, up to H_{13} , are separable, but H_{14} and H_{15} are not; H_{16} through H_{19} are separable, but H_{20} and H_{21} are not. Accordingly, the notion of two indices n and m with dependence on m in the form of $\cos m\theta$ or $\sin m\theta$, as in the case of circle polynomials, loses significance. For example, the Zernike polynomial Z_{14} for $n = 4$ and $m = 4$ varies as $\cos 4\theta$, but H_{14} has a term in $\cos 2\theta$ also. Hence, the hexagonal polynomials can be ordered by a single index only. While the polynomials H_{11} and H_{22} representing the balanced primary and secondary spherical aberrations are radially symmetric, the polynomial H_{37} representing the balanced tertiary spherical aberration is not, since it consists of an angle-dependent term in Z_{28} or $\cos 6\theta$ also. If this term is not included in the polynomial H_{37} , the standard deviation of the aberration increases from a value of unity to 1.3339.

11.7 ELLIPTICAL POLYNOMIALS

Figure 3c shows a *unit ellipse* of an aspect ratio b inscribed inside a unit circle. The semimajor and semiminor axes of the ellipse have lengths of unity and b , respectively. Of course, a unit ellipse is not unique, since b can have any value between 0 and 1. It is represented by an equation

$$x^2 + y^2/b^2 = 1 \quad (35a)$$

TABLE 3a Orthonormal Hexagonal Polynomials H_j in Terms of Zernike Circle Polynomials Z_j

$$\begin{aligned}
H_1 &= Z_1 \\
H_2 &= \sqrt{6/5}Z_2 \\
H_3 &= \sqrt{6/5}Z_3 \\
H_4 &= \sqrt{5/43}Z_1 + (2\sqrt{15/43})Z_4 \\
H_5 &= \sqrt{10/7}Z_5 \\
H_6 &= \sqrt{10/7}Z_6 \\
H_7 &= 16\sqrt{14/11055}Z_3 + 10\sqrt{35/2211}Z_7 \\
H_8 &= 16\sqrt{14/11055}Z_2 + 10\sqrt{35/2211}Z_8 \\
H_9 &= (2\sqrt{5/3})Z_9 \\
H_{10} &= 2\sqrt{35/103}Z_{10} \\
H_{11} &= (521/\sqrt{1072205})Z_1 + 88\sqrt{15/214441}Z_4 + 14\sqrt{43/4987}Z_{11} \\
H_{12} &= 225\sqrt{6/492583}Z_6 + 42\sqrt{70/70369}Z_{12} \\
H_{13} &= 225\sqrt{6/492583}Z_5 + 42\sqrt{70/70369}Z_{13} \\
H_{14} &= -2525\sqrt{14/297774543}Z_6 - (1495\sqrt{70/99258181/3})Z_{12} + (\sqrt{378910/18337/3})Z_{14} \\
H_{15} &= 2525\sqrt{14/297774543}Z_5 + (1495\sqrt{70/99258181/3})Z_{13} + (\sqrt{378910/18337/3})Z_{15} \\
H_{16} &= 30857\sqrt{2/3268147641}Z_2 + (49168/\sqrt{3268147641})Z_8 + 42\sqrt{1474/1478131}Z_{16} \\
H_{17} &= 30857\sqrt{2/3268147641}Z_3 + (49168/\sqrt{3268147641})Z_7 + 42\sqrt{1474/1478131}Z_{17} \\
H_{18} &= 386\sqrt{770/295894589}Z_{10} + 6\sqrt{118965/2872763}Z_{18} \\
H_{19} &= 6\sqrt{10/97}Z_9 + 14\sqrt{5/291}Z_{19} \\
H_{20} &= -0.71499593Z_2 - 0.72488884Z_8 - 0.46636441Z_{16} + 1.72029850Z_{20} \\
H_{21} &= 0.71499594Z_3 + 0.72488884Z_7 + 0.46636441Z_{17} + 1.72029850Z_{21} \\
H_{22} &= 0.58113135Z_1 + 0.89024136Z_4 + 0.89044507Z_{11} + 1.32320623Z_{22} \\
H_{23} &= 1.15667686Z_5 + 1.10775599Z_{13} + 0.43375081Z_{15} + 1.39889072Z_{23} \\
H_{24} &= 1.15667686Z_6 + 1.10775599Z_{12} - 0.43375081Z_{14} + 1.39889072Z_{24} \\
H_{25} &= 1.31832566Z_5 + 1.14465174Z_{13} + 1.94724032Z_{15} + 0.67629133Z_{23} + 1.75496998Z_{25} \\
H_{26} &= -1.31832566Z_6 - 1.14465174Z_{12} + 1.94724032Z_{14} - 0.67629133Z_{24} + 1.75496998Z_{26} \\
H_{27} &= 2\sqrt{77/93}Z_{27} \\
H_{28} &= -1.07362889Z_1 - 1.52546162Z_4 - 1.28216588Z_{11} - 0.70446308Z_{22} + 2.09532473Z_{28} \\
H_{29} &= 0.97998834Z_3 + 1.16162002Z_7 + 1.04573775Z_{17} + 0.40808953Z_{21} + 1.36410394Z_{29} \\
H_{30} &= 0.97998834Z_2 + 1.16162002Z_8 + 1.04573775Z_{16} - 0.40808953Z_{20} + 1.36410394Z_{30} \\
H_{31} &= 3.63513758Z_9 + 2.92084414Z_{19} + 2.11189625Z_{31} \\
H_{32} &= 0.69734874Z_{10} + 0.67589740Z_{18} + 1.22484055Z_{32} \\
H_{33} &= 1.56189763Z_3 + 1.69985309Z_7 + 1.29338869Z_{17} + 2.57680871Z_{21} + 0.67653220Z_{29} \\
&\quad + 1.95719339Z_{33} \\
H_{34} &= -1.56189763Z_2 - 1.69985309Z_8 - 1.29338869Z_{16} + 2.57680871Z_{20} - 0.67653220Z_{30} \\
&\quad + 1.95719339Z_{34} \\
H_{35} &= -1.63832594Z_3 - 1.74759886Z_7 - 1.27572528Z_{17} - 0.77446421Z_{21} - 0.60947360Z_{29} \\
&\quad - 0.36228537Z_{33} + 2.24453237Z_{35} \\
H_{36} &= -1.63832594Z_2 - 1.74759886Z_8 - 1.27572528Z_{16} + 0.77446421Z_{20} - 0.60947360Z_{30} \\
&\quad + 0.36228537Z_{34} + 2.24453237Z_{36} \\
H_{37} &= 0.82154671Z_1 + 1.27988084Z_4 + 1.32912377Z_{11} + 1.11636637Z_{22} - 0.54097038Z_{28} \\
&\quad + 1.37406534Z_{37} \\
H_{38} &= 1.54526522Z_6 + 1.57785242Z_{12} - 0.89280081Z_{14} + 1.28876176Z_{24} - 0.60514082Z_{26} \\
&\quad + 1.43097780Z_{38} \\
H_{39} &= 1.54526522Z_5 + 1.57785242Z_{13} + 0.89280081Z_{15} + 1.28876176Z_{23} + 0.60514082Z_{25} \\
&\quad + 1.43097780Z_{39} \\
H_{40} &= -2.51783502Z_6 - 2.38279377Z_{12} + 3.42458933Z_{14} - 1.69296616Z_{24} + 2.56612920Z_{26} \\
&\quad - 0.85703819Z_{38} + 1.89468756Z_{40} \\
H_{41} &= 2.51783502Z_5 + 2.38279377Z_{13} + 3.42458933Z_{15} + 1.69296616Z_{23} + 2.56612920Z_{25} \\
&\quad + 0.85703819Z_{39} + 1.89468756Z_{41}
\end{aligned}$$

(Continued)

TABLE 3a Orthonormal Hexagonal Polynomials H_j in Terms of Zernike Circle Polynomials Z_j
 (Continued)

$$\begin{aligned}
 H_{42} &= -2.72919646Z_1 - 4.02313214Z_4 - 3.69899239Z_{11} - 2.49229315Z_{22} + 4.36717121Z_{28} \\
 &\quad - 1.13485132Z_{37} + 2.52330106Z_{42} \\
 H_{43} &= 1362\sqrt{77/20334667}Z_{27} + (260/3)\sqrt{341/655957}Z_{43} \\
 H_{44} &= -2.76678413Z_6 - 2.50005278Z_{12} + 1.48041348Z_{14} - 1.62947374Z_{24} + 0.95864121Z_{26} \\
 &\quad - 0.69034812Z_{38} + 0.40743941Z_{40} + 2.56965299Z_{44} \\
 H_{45} &= -2.76678413Z_5 - 2.50005278Z_{13} - 1.48041348Z_{15} - 1.62947374Z_{23} - 0.95864121Z_{25} \\
 &\quad - 0.69034812Z_{39} - 0.40743941Z_{41} + 2.56965299Z_{45}
 \end{aligned}$$

TABLE 3b Orthonormal Hexagonal Polynomials $H_j(\rho, \theta)$ in Polar Coordinates

$$\begin{aligned}
 H_1 &= 1 \\
 H_2 &= 2\sqrt{6/5}\rho \cos \theta \\
 H_3 &= 2\sqrt{6/5}\rho \sin \theta \\
 H_4 &= \sqrt{5/43}(-5 + 12\rho^2) \\
 H_5 &= 2\sqrt{15/7}\rho^2 \sin 2\theta \\
 H_6 &= 2\sqrt{15/7}\rho^2 \cos 2\theta \\
 H_7 &= 4\sqrt{42/3685}(-14\rho + 25\rho^3) \sin \theta \\
 H_8 &= 4\sqrt{42/3685}(-14\rho + 25\rho^3) \cos \theta \\
 H_9 &= (4\sqrt{10/3})\rho^3 \sin 3\theta \\
 H_{10} &= 4\sqrt{70/103}\rho^3 \cos 3\theta \\
 H_{11} &= (3/\sqrt{1072205})(737 - 5140\rho^2 + 6020\rho^4) \\
 H_{12} &= (30/\sqrt{492583})(-249\rho^2 + 392\rho^4) \cos 2\theta \\
 H_{13} &= (30/\sqrt{492583})(-249\rho^2 + 392\rho^4) \sin 2\theta \\
 H_{14} &= (10/3)\sqrt{7/99258181}[10(297 - 598\rho^2)\rho^2 \cos 2\theta + 5413\rho^4 \cos 4\theta] \\
 H_{15} &= (10/3)\sqrt{7/99258181}[-10(297 - 598\rho^2)\rho^2 \sin 2\theta + 5413\rho^4 \sin 4\theta] \\
 H_{16} &= 2\sqrt{6/1089382547}(70369\rho - 322280\rho^3 + 309540\rho^5) \cos \theta \\
 H_{17} &= 2\sqrt{6/1089382547}(70369\rho - 322280\rho^3 + 309540\rho^5) \sin \theta \\
 H_{18} &= 4\sqrt{385/295894589}(-3322\rho^3 + 4635\rho^5) \cos 3\theta \\
 H_{19} &= 4\sqrt{5/97}(-22\rho^3 + 35\rho^5) \sin 3\theta \\
 H_{20} &= (-2.17600248\rho + 13.23551876\rho^3 - 16.15533716\rho^5) \cos \theta + 5.95928883\rho^5 \cos 5\theta \\
 H_{21} &= (2.17600248\rho - 13.23551876\rho^3 + 16.15533716\rho^5) \sin \theta + 5.95928883\rho^5 \sin 5\theta \\
 H_{22} &= -2.47059083 + 33.14780774\rho^2 - 93.07966445\rho^4 + 70.01749250\rho^6 \\
 H_{23} &= (23.72919095\rho^2 - 90.67126833\rho^4 + 78.51254738\rho^6) \sin 2\theta + 1.37164051\rho^4 \sin 4\theta \\
 H_{24} &= (23.72919095\rho^2 - 90.67126833\rho^4 + 78.51254738\rho^6) \cos 2\theta - 1.37164051\rho^4 \cos 4\theta \\
 H_{25} &= (7.55280798\rho^2 - 36.13018255\rho^4 + 37.95675688\rho^6) \sin 2\theta + (-26.67476754\rho^4 \\
 &\quad + 39.39897852\rho^6) \sin 4\theta \\
 H_{26} &= (-7.55280798\rho^2 + 36.13018255\rho^4 - 37.95675688\rho^6) \cos 2\theta + (-26.67476754\rho^4 \\
 &\quad + 39.39897852\rho^6) \cos 4\theta \\
 H_{27} &= 14\sqrt{22/93}\rho^6 \sin 6\theta \\
 H_{28} &= 0.56537219 - 10.44830313\rho^2 + 38.71296332\rho^4 - 37.27668254\rho^6 + 7.83998727\rho^6 \cos 6\theta \\
 H_{29} &= (-15.56917599 + 130.07864353\rho^2 - 288.33220017\rho^4 + 190.97455178\rho^6)\rho \sin \theta \\
 &\quad + 2.82732724\rho^5 \sin 3\theta + 1.41366362\rho^5 \sin 5\theta \\
 H_{30} &= (-15.56917599 + 130.07864353\rho^2 - 288.33220017\rho^4 + 190.97455178\rho^6)\rho \cos \theta \\
 &\quad + 2.82732724\rho^5 \cos 3\theta + 1.41366362\rho^5 \cos 5\theta \\
 H_{31} &= (54.28516840 - 202.83704634\rho^2 + 177.39928561\rho^4)\rho^3 \sin 3\theta \\
 H_{32} &= (41.60051295 - 135.27397959\rho^2 + 102.88660624\rho^4)\rho^3 \cos 3\theta \\
 H_{33} &= (-3.87525156 + 41.84243767\rho^2 - 193.65605837\rho^4 + 204.31733848\rho^6)\rho \sin \theta + (76.09262860 \\
 &\quad - 109.60283027\rho^2)\rho^5 \sin 3\theta + (38.04631430 - 54.80141514\rho^2)\rho^5 \sin 5\theta \\
 H_{34} &= (3.87525156 - 41.84243767\rho^2 + 117.56342977\rho^4 - 94.71450820\rho^6)\rho \cos \theta + (-76.09262860 \\
 &\quad + 109.60283027\rho^2)\rho^5 \cos 3\theta + (38.04631430 - 54.80141514\rho^2)\rho^5 \cos 5\theta
 \end{aligned}$$

(Continued)

TABLE 3b Orthonormal Hexagonal Polynomials $H_j(\rho, \theta)$ in Polar Coordinates (Continued)

$$\begin{aligned}
H_{35} &= (3.10311187 - 34.93479698\rho^2 + 114.10529848\rho^4 - 87.65802721\rho^6)\rho \sin \theta + (12.02405243 \\
&\quad - 2.33172188\rho^2)\rho^5 \sin 3\theta + (12.02405243 + 3.68030434\rho^2)\rho^5 \sin 5\theta + 6.01202622\rho^7 \sin 7\theta \\
H_{36} &= (3.10311187 - 34.93479698\rho^2 + 114.10529848\rho^4 - 87.65802721\rho^6)\rho \cos \theta + (12.02405243 \\
&\quad - 2.33172188\rho^2)\rho^5 \cos 3\theta + (12.02405243 + 3.68030434\rho^2)\rho^5 \sin 5\theta + 6.01202622\rho^7 \cos 7\theta \\
H_{37} &= 2.74530738 - 60.39881618\rho^2 + 300.22087475\rho^4 - 518.03488742\rho^6 + 288.55372176\rho^8 \\
&\quad - 2.02412582\rho^6 \cos 6\theta \\
H_{38} &= (-42.96232789 + 287.78381063\rho^2 - 565.13651608\rho^4 + 339.98298180\rho^4)\rho^2 \cos 2\theta \\
&\quad + (8.49786414 - 13.58537785\rho^2)\rho^4 \cos 4\theta \\
H_{39} &= (-42.96232789 + 287.78381063\rho^2 - 565.13651608\rho^4 + 339.98298180\rho^4)\rho^2 \sin 2\theta \\
&\quad + (8.49786414 - 13.58537785\rho^2)\rho^4 \sin 4\theta \\
H_{40} &= (14.79181046 - 121.61654135\rho^2 + 286.77354559\rho^4 - 203.62188574\rho^6)\rho^2 \cos 2\theta \\
&\quad + (83.39879886 - 280.00664075\rho^2 + 225.07739907\rho^4)\rho^4 \cos 4\theta \\
H_{41} &= (-14.79181046 + 121.61654135\rho^2 - 286.77354559\rho^4 + 203.62188574\rho^6)\rho^2 \sin 2\theta \\
&\quad + (83.39879886 - 280.00664075\rho^2 + 225.07739907\rho^4)\rho^4 \sin 4\theta \\
H_{42} &= -0.84269170 + 24.65387703\rho^2 - 158.21741244\rho^4 + 344.75780000\rho^6 - 238.31877895\rho^8 \\
&\quad + (-58.59775991 + 85.64367812\rho^2)\rho^6 \cos 6\theta \\
H_{43} &= 2\sqrt{22/20334667(-23443 + 32240\rho^2)}\rho^6 \sin 6\theta \\
H_{44} &= (9.64776957 - 85.41873843\rho^2 + 216.08041438\rho^4 - 164.01834750\rho^6)\rho^2 \cos 2\theta \\
&\quad + (12.67622930 - 51.08055822\rho^2 + 48.40133344\rho^4)\rho^4 \cos 4\theta + 10.90211434\rho^8 \cos 8\theta \\
H_{45} &= (9.64776957 - 85.41873843\rho^2 + 216.08041438\rho^4 - 164.01834750\rho^6)\rho^2 \sin 2\theta \\
&\quad - (12.67622930 - 51.08055822\rho^2 + 48.40133344\rho^4)\rho^4 \sin 4\theta + 10.90211434\rho^8 \sin 8\theta
\end{aligned}$$

TABLE 3c Orthonormal Hexagonal Polynomials $H_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$

$$\begin{aligned}
H_1 &= 1 \\
H_2 &= 2\sqrt{6/5}x \\
H_3 &= 2\sqrt{6/5}y \\
H_4 &= \sqrt{5/43}(-5 + 12\rho^2) \\
H_5 &= 4\sqrt{15/7}xy \\
H_6 &= 2\sqrt{15/7}(x^2 - y^2) \\
H_7 &= 4\sqrt{42/3685}(-14 + 25\rho^2)y \\
H_8 &= 4\sqrt{42/3685}(-14 + 25\rho^2)x \\
H_9 &= (4/3)\sqrt{10}(3x^2y - y^3) \\
H_{10} &= 4\sqrt{70/103}(x^3 - 3xy^2) \\
H_{11} &= (3/\sqrt{1072205})(737 - 5140\rho^2 + 6020\rho^4) \\
H_{12} &= (30/\sqrt{492583})(392\rho^2 - 249)(x^2 - y^2) \\
H_{13} &= (60/\sqrt{492583})(392\rho^2 - 249)xy \\
H_{14} &= -(10/3)\sqrt{7/99258181}[567x^4 + 32478x^2y^2 - 11393y^4 - 2970(x^2 - y^2)] \\
H_{15} &= (40/3)\sqrt{7/99258181}(-1485 + 8403x^2 - 2423y^2)xy \\
H_{16} &= 2\sqrt{2/3268147641}(211107 - 966840\rho^2 + 928620\rho^4)x \\
H_{17} &= 2\sqrt{2/3268147641}(211107 - 966840\rho^2 + 928620\rho^4)y \\
H_{18} &= 4\sqrt{385/295894589}(-3322 + 4635\rho^2)(x^3 - 3xy^2) \\
H_{19} &= 4\sqrt{5/97}(-22 + 35\rho^2)(3x^2y - y^3) \\
H_{20} &= (-2.17600247 + 13.23551876\rho^2 - 10.19604832x^4 - 91.90356268x^2y^2 + 13.64110702y^4)x \\
H_{21} &= (2.17600247 - 13.23551876\rho^2 + 45.95178134x^4 - 27.28221405x^2y^2 + 22.11462599y^4)y \\
H_{22} &= -2.47059083 + 33.14780774\rho^2 - 93.07966445\rho^4 + 70.01749250\rho^6 \\
H_{23} &= (47.45838189 - 175.85597460x^2 - 186.82909872y^2 + 157.02509476x^4 \\
&\quad + 314.05018953x^2y^2 + 157.02509476y^4)xy \\
H_{24} &= (23.72919094 - 92.04290884x^2 + 78.51254738x^4) + (-23.72919094 + 8.22984309x^2 \\
&\quad + 89.29962781y^2 + 78.51254738x^4 - 78.51254738x^2y^2 - 78.51254738y^4)y^2
\end{aligned}$$

TABLE 3c Orthonormal Hexagonal Polynomials $H_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$ (Continued)

$$\begin{aligned}
 H_{25} &= (15.10561596 - 178.95943525x^2 + 34.43870505y^2 + 233.50942786x^4 \\
 &\quad + 151.82702751x^2y^2 - 81.68240034y^4)xy \\
 H_{26} &= (-7.55280798 + 9.45541501x^2 + 1.44222164x^4)x^2 + (7.55280798 + 160.04860523x^2 \\
 &\quad - 62.80495008y^2 - 234.95164950x^4 - 159.03813574x^2y^2 + 77.35573540y^4)y^2 \\
 H_{27} &= (40.85537039x^4 - 136.18456799x^2y^2 + 40.85537039y^4)xy \\
 H_{28} &= 0.56537219 - 10.44830312\rho^2 + 38.71296332x^4 + 77.42592664x^2y^2 + 38.71296332y^4 \\
 &\quad - 29.43669525x^6 - 229.42985678x^4y^2 + 5.76976155x^2y^4 - 45.11666981y^6 \\
 H_{29} &= (-15.56917599 + 7.06831810x^4 - 14.13663621x^2y^2 + 1.41366362y^4 + 130.07864353\rho^2 \\
 &\quad - 291.15952741\rho^4 + 190.97455178\rho^6)y \\
 H_{30} &= (-15.56917599 - 1.41366362x^4 + 14.13663621x^2y^2 - 7.06831810y^4 + 130.07864353\rho^2 \\
 &\quad - 291.15952741\rho^4 + 190.97455178\rho^6)x \\
 H_{31} &= 162.85550520x^2 - 54.28516840y^2 - 608.51113904x^2\rho^2 + 202.83704634y^2\rho^2 \\
 &\quad + 532.19785685x^2\rho^4 - 177.39928561y^2\rho^4)y \\
 H_{32} &= [(41.60051295 - 135.27397959x^2 + 102.88660624x^4)x^2 + (-124.80153887 + 270.54795919x^2 \\
 &\quad + 405.82193879y^2 - 102.88660624x^4 - 514.43303123x^2y^2 - 308.65981874y^4)y^2]x \\
 H_{33} &= [-3.87525156 + (41.84243767 - 307.79500129x^2 + 368.72158389x^4)x^2 + (41.84243767 \\
 &\quad + 145.33628349x^2 - 155.60974407y^2 + 10.13644892x^4 - 209.06921162x^2y^2 + 149.51592334y^4)y^2]y \\
 H_{34} &= [3.87525156 + (-41.84243767 + 79.51711547x^2 - 39.91309306x^4)x^2 + (-41.84243767 \\
 &\quad + 615.59000259x^2 - 72.66814174y^2 - 777.35626084x^4 - 558.15060029x^2y^2 + 179.29256748y^4)y^2]x \\
 H_{35} &= [3.10311187 + (-34.93479698 + 132.14137712x^2 - 73.19935100x^4)x^2 + (-34.93479698 \\
 &\quad + 144.04222993x^2 + 108.09327226y^2 - 519.49349681x^4 + 23.85771799x^2y^2 - 104.44842531y^4)y^2]y \\
 H_{36} &= [3.10311187 + (-34.93479698 + 96.06921983x^2 - 66.20418535x^4)x^2 + (-34.93479698 \\
 &\quad + 264.28275425x^2 + 72.02111496y^2 - 535.81555000x^4 + 7.53566481x^2y^2 - 97.45325965y^4)y^2]x \\
 H_{37} &= 2.74530738 - 60.39881618\rho^2 + 300.22087475\rho^4 + 288.55372176\rho^8 - 520.05901324x^6 \\
 &\quad - 1523.74277487x^4y^2 - 1584.46654966x^2y^4 - 516.01076159y^6 \\
 H_{38} &= (-42.96232789 + 296.28167478x^2 - 578.72189394x^4 + 339.98298180x^6)x^2 + (42.96232789 \\
 &\quad - 50.98718488x^2 - 279.28594648y^2 - 497.20962679x^4 + 633.06340537x^2y^2 + 551.55113822y^4 \\
 &\quad + 679.96596360x^6 - 679.96596360x^2y^4 - 339.98298180y^6)y^2 \\
 H_{39} &= [-85.92465579 + (541.57616468 - 1075.93152073x^2 + 679.96596360x^4)x^2 + (609.55907786 \\
 &\quad - 2260.54606433x^2 - 1184.61454360y^2 + 2039.89789081x^4 + 2039.89789081x^2y^2 + 679.96596360y^4)y^2]xy \\
 H_{40} &= (14.79181046 - 38.21774249x^2 + 6.76690483x^4 + 21.45551332x^6)x^2 + (-14.79181046 \\
 &\quad - 500.39279319x^2 + 205.01534022y^2 + 1686.80674937x^4 + 1113.25965819x^2y^2 - 566.78018634y^4 \\
 &\quad - 1307.55336779x^6 - 2250.77399075x^4y^2 - 493.06582480x^2y^4 + 428.69928482y^6)y^2 \\
 H_{41} &= [-29.58362093 + (576.82827818 - 1693.57365421x^2 + 1307.55336779x^4)x^2 + (-90.36211274 \\
 &\quad - 1147.09418236x^2 + 546.47947184y^2 + 2122.04091078x^4 + 321.42171817x^2y^2 - 493.06582480y^4)y^2]xy \\
 H_{42} &= -0.84269170 + (24.65387703 - 158.21741244x^2 + 286.16004008x^4 - 152.67510082x^6)x^2 \\
 &\quad + (24.65387703 - 316.43482489x^2 - 158.21741244y^2 + 1913.23979875x^4 + 155.30700127x^2y^2 \\
 &\quad + 403.3555992y^4 - 2152.28660953x^6 - 1429.91267370x^4y^2 + 245.73637792x^2y^4 - 323.96245707y^6)y^2 \\
 H_{43} &= 2\sqrt{22/20334667}(6x^5y - 20x^3y^3 + 6xy^5)(-23443 + 32240\rho^2) \\
 H_{44} &= (9.64776957 - 72.74250912x^2 + 164.99985615x^4 - 104.71489971x^6)x^2 + (-9.64776957 \\
 &\quad - 76.05737585x^2 + 98.09496774y^2 + 471.48320551x^4 + 39.32237674x^2y^2 - 267.16097261y^4 \\
 &\quad - 826.90123032x^6 + 279.13466933x^4y^2 - 170.82784030x^2y^4 + 223.32179529y^6)y^2 \\
 H_{45} &= [19.29553915 + (-221.54239411 + 636.48306167x^2 - 434.42511407x^4)x^2 + (-120.13255963 \\
 &\quad + 864.32165754x^2 + 227.83859586y^2 - 1788.23382186x^4 - 179.98634818x^2y^2 - 221.64827593y^4)y^2]xy
 \end{aligned}$$

or

$$y = \pm b\sqrt{1-x^2} \tag{35b}$$

Its area is equal to πb . The orthonormality of the elliptical polynomials $E_j(x, y)$ is represented by²⁶

$$\frac{1}{\pi b} \int_{-1}^1 dx \int_{-b\sqrt{1-x^2}}^{b\sqrt{1-x^2}} E_j(x, y) E_{j'}(x, y) dy = \delta_{jj'} \tag{36}$$

The orthonormal elliptical polynomials up to the fourth order are given in Tables 4 in three different but equivalent forms, as in the case of hexagonal polynomials. As in the case of a hexagonal

TABLE 4a Orthonormal *Elliptical* Polynomials E_j in terms of Zernike Circle Polynomials Z_j , Which Reduce to the Corresponding Circle Polynomials as the Aspect Ratio $b \rightarrow 1$

$$\begin{aligned}
E_1 &= Z_1 \\
E_2 &= Z_2 \\
E_3 &= Z_3/b \\
E_4 &= (1/\sqrt{3-2b^2+3b^4})[\sqrt{3}(1-b^2)Z_1 + 2Z_4] \\
E_5 &= Z_5/b \\
E_6 &= [1/(2\sqrt{2b^3\sqrt{3-2b^2+3b^4}})] [-\sqrt{3}(3-4b^2+b^4)Z_1 - 3(1-b^4)Z_4 + \sqrt{2}(3-2b^2+3b^4)Z_6] \\
E_7 &= [1/(b\sqrt{5-6b^2+9b^4})][6(1-b^2)Z_3 + 2\sqrt{2}Z_7] \\
E_8 &= (2/\sqrt{9-6b^2+5b^4})[(1-b^2)Z_2 + \sqrt{2}Z_8] \\
E_9 &= [1/(2\sqrt{2b^3\sqrt{5-6b^2+9b^4}})] [-2\sqrt{2}(5-8b^2+3b^4)Z_3 - (5-2b^2-3b^4)Z_7 + (5-6b^2+9b^4)Z_9] \\
E_{10} &= [1/(2\sqrt{2b^3\sqrt{9-6b^2+5b^4}})] [-2\sqrt{2}(3-4b^2+b^4)Z_2 - (3+2b^2-5b^4)Z_8 + (9-6b^2+5b^4)Z_{10}] \\
E_{11} &= (1/\alpha)[\sqrt{5}(7-10b^2+3b^4)Z_1 + 4\sqrt{15}(1-b^2)Z_4 - 2\sqrt{30}(1-b^2)Z_6 + 8Z_{11}] \\
E_{12} &= -\sqrt{5/8}b^{-2}(195-475b^2+558b^4-422b^6+159b^8-15b^{10})\beta^{-1}Z_1 \\
&\quad -\sqrt{15/8}b^{-2}(105-205b^2+194b^4-114b^6+5b^8+15b^{10})\beta^{-1}Z_4 \\
&\quad +\sqrt{15/4}(75-155b^2+174b^4-134b^6+55b^8-15b^{10})\beta^{-1}Z_6 \\
&\quad -10\sqrt{2}b^{-2}(3-2b^2+2b^6-3b^8)\beta^{-1}Z_{11} + b^{-2}\alpha\gamma^{-1}Z_{12} \\
E_{13} &= [1/(b\sqrt{5-6b^2+5b^4})][\sqrt{15}(1-b^2)Z_5 + 2Z_{13}] \\
E_{14} &= (\sqrt{5/2}/4)(1-b^2)^2b^{-4}(35-10b^2-b^4)\gamma^{-1}Z_1 + (5\sqrt{15/2}/8)(1-b^2)^2b^{-4}(7+2b^2-b^4)\gamma^{-1}Z_4 \\
&\quad -(\sqrt{15}/8)(35-70b^2+56b^4-26b^6+5b^8)\gamma^{-1}Z_6 + [5/(8\sqrt{2})](1-b^2)^2b^{-4}(7+10b^2+7b^4)\gamma^{-1}Z_{11} \\
&\quad - (5/8)b^{-4}(7-6b^2+6b^6-7b^8)\gamma^{-1}Z_{12} + [\gamma/(8b^4)]Z_{14} \\
E_{15} &= -(\sqrt{15}/4)b^{-3}(5-8b^2+3b^4)\delta^{-1}Z_5 - (5/4)(1-b^4)b^{-3}\beta^{-1}Z_{13} + [\delta/(2b^3)]Z_{15} \\
\alpha &= (45-60b^2+94b^4-60b^6+45b^8)^{1/2} \\
\beta &= (1575-4800b^2+12020b^4-17280b^6+21066b^8-17280b^{10}+12020b^{12}-4800b^{14}+1575b^{16})^{1/2} \\
\gamma &= (35-60b^2+114b^4-60b^6+35b^8)^{1/2} \\
\delta &= (5-6b^2+5b^4)^{1/2}
\end{aligned}$$

TABLE 4b Orthonormal *Elliptical* Polynomials $E_j(\rho, \theta)$ in Polar Coordinates

$$\begin{aligned}
E_1 &= 1 \\
E_2 &= 2\rho \cos \theta \\
E_3 &= (2\rho \sin \theta)/b \\
E_4 &= \sqrt{3/(3-2b^2+3b^4)}(-1-b^2+4\rho^2) \\
E_5 &= (\sqrt{6}/b)\rho^2 \sin 2\theta \\
E_6 &= [1/(2b^2)]\sqrt{6/(3-2b^2+3b^4)}[2b^2(1-b^2)-3(1-b^4)\rho^2 + (3-2b^2+3b^4)\rho^2 \cos 2\theta] \\
E_7 &= [4/(b\sqrt{5-6b^2+9b^4})][-(1+3b^2)\rho + 6\rho^3] \sin \theta \\
E_8 &= (4/\sqrt{9-6b^2+5b^4})[-(3+b^2)\rho + 6\rho^3] \cos \theta \\
E_9 &= [1/(b^3\sqrt{5-6b^2+9b^4})]\{3[4b^2(1-b^2)\rho - (5-2b^2-3b^4)\rho^3] \sin \theta + (5-6b^2+9b^4)\rho^3 \sin 3\theta\} \\
E_{10} &= [1/(b^2\sqrt{9-6b^2+5b^4})]\{3[4b^2(1-b^2)\rho - (3+2b^2-5b^4)\rho^3] \cos \theta + (9-6b^2+5b^4)\rho^3 \cos 3\theta\} \\
E_{11} &= \sqrt{5}[3+2b^2+3b^4-24(1+b^2)\rho^2+48\rho^4-12(1-b^2)\rho^2 \cos 2\theta]\alpha \\
E_{12} &= [\sqrt{10}\alpha/(\gamma b^2)](-3\rho^2+4\rho^4) \cos 2\theta + [\sqrt{5/2}/(2b^2\beta)][-12b^2(5-2b^2+2b^6-5b^8) \\
&\quad +6(15+125b^2-194b^4+194b^6-125b^8-15b^{10})\rho^2+240(-3+2b^2-2b^6+3b^8)\rho^4 \\
&\quad +6(75-155b^2+174b^4-134b^6+55b^8-15b^{10})\rho^2 \cos 2\theta] \\
E_{13} &= (\sqrt{10}/b)\delta^{-1}[-3(1+b^2)\rho^2+8\rho^4] \sin 2\theta \\
E_{14} &= [\sqrt{10}/(8b^4\gamma)]\{3(1-b^2)^2[8b^4-40b^2(1+b^2)\rho^2+5(7+10b^2+7b^4)\rho^4] \\
&\quad +4[6b^2(5-7b^2+7b^4-5b^6)-5(7-6b^2+6b^6-7b^8)\rho^2]\rho^2 \cos 2\theta + (35-60b^2+114b^4 \\
&\quad -60b^6+35b^8)\rho^4 \cos 4\theta\} \\
E_{15} &= (\sqrt{10}/b^3)\delta^{-1}\{[6b^2(1-b^2)-5(1-b^4)\rho^2]\rho^2 \sin 2\theta + [(5-6b^2+5b^4)/2]\rho^4 \sin 4\theta\}
\end{aligned}$$

TABLE 4c Orthonormal *Elliptical* Polynomials $E_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-1 \leq x \leq 1$, and $-\sqrt{1-b^2x^2} \leq y \leq \sqrt{1-b^2x^2}$

$$\begin{aligned}
 E_1 &= 1 \\
 E_2 &= 2x \\
 E_3 &= 2y/b \\
 E_4 &= (\sqrt{3}/\sqrt{3-2b^2+3b^4})(-1-b^2+4\rho^2) \\
 E_5 &= (2\sqrt{6}/b)xy \\
 E_6 &= [\sqrt{3}/(b^2\sqrt{6-4b^2+6b^4})][b^2(1-b^2)+b^2(3b^2-1)x^2-(3-b^2)y^2] \\
 E_7 &= [4/(b\sqrt{5-6b^2+9b^4})][-(1+3b^2)+6\rho^2]y \\
 E_8 &= (4/\sqrt{9-6b^2+5b^4})[-(3+b^2)+6\rho^2]x \\
 E_9 &= [4/(b^3\sqrt{5-6b^2+9b^4})][3b^2(3b^2-1)x^2-(5-3b^2)y^2+3b^2(1-b^2)]y \\
 E_{10} &= [4/(b^2\sqrt{9-6b^2+5b^4})][b^2(5b^2-3)x^2-3(3-b^2)y^2+3b^2(1-b^2)]x \\
 E_{11} &= (\sqrt{5}/\alpha)[48\rho^4-12(3+b^2)x^2-12(1+3b^2)y^2+3+2b^2+3b^4] \\
 E_{12} &= [\sqrt{10}\alpha/(b^2\gamma)][(x^2-y^2)(4\rho^2-3)+[\sqrt{5}/(2\sqrt{2}b^2\beta)][240(-3+2b^2-2b^6+3b^8)\rho^4 \\
 &\quad -60(-9+3b^2+2b^4-6b^6+7b^8+3b^{10})x^2-24(15-70b^2+92b^4-82b^6+45b^8)y^2 \\
 &\quad +12b^2(-5+2b^2-2b^6+5b^8)]] \\
 E_{13} &= [2\sqrt{10}/(b\delta)](8\rho^2-3-3b^2)xy \\
 E_{14} &= [\sqrt{10}/(b^4\gamma)][b^4(3-30b^2+35b^4)x^4+6b^2(5-18b^2+5b^4)x^2y^2+(35-30b^2+3b^4)y^4 \\
 &\quad -6b^4(1-6b^2+5b^4)x^2-6b^2(5-6b^2+b^4)y^2+3b^4(1-b^2)^2] \\
 E_{15} &= [4\sqrt{10}/(b^3\delta)][b^2(5b^2-3)x^2-(5-3b^2)y^2+3b^2(1-b^2)]xy
 \end{aligned}$$

pupil, each elliptical polynomial consists of either cosine or sine terms, but not both. For example, E_6 is a linear combination of Z_6 , Z_4 , and Z_1 . It also shows that the balancing defocus for (zero-degree) Seidel astigmatism is different for an elliptical pupil compared to that for a circular,³⁻⁶ annular,⁸⁻¹¹ or a Gaussian pupil.²³⁻²⁵ Moreover, E_{11} is a linear combination of Z_{11} , Z_6 , Z_4 , and Z_1 . Thus, spherical aberration ρ^4 is balanced with not only defocus ρ^2 but astigmatism $\rho^2 \cos^2\theta$ as well. The elliptical polynomials are generally more complex in that they are made up of a larger number of circle polynomials. These results are a consequence of the fact that the x and y dimensions of the elliptical pupil are not equal. As expected, the elliptical polynomials reduce to the circle polynomials as $b \rightarrow 1$, that is, as the unit ellipse approaches a unit circle.

11.8 RECTANGULAR POLYNOMIALS

Figure 3d shows a *unit rectangle* inscribed inside a unit circle. While the distance of a corner point of the rectangle, such as A , from its center O is unity, the half widths of the rectangle along the x and y axes are a and $\sqrt{1-a^2}$, respectively. Accordingly, the aspect ratio of the rectangle is $\sqrt{1-a^2}/a$, and its area is $4a\sqrt{1-a^2}$. As in the case of a unit ellipse, a unit rectangle is also not unique, since a can have any value between 0 and 1. The orthonormality of the rectangular polynomials $R_j(x, y)$ is represented by²⁶

$$\frac{1}{4a\sqrt{1-a^2}} \int_{-\sqrt{1-a^2}}^{\sqrt{1-a^2}} dy \int_{-a}^a R_j(x, y) R_j(x, y) dx = \delta_{jj} \quad (37)$$

The rectangular polynomials thus obtained up to the fourth order are given in Tables 5 in the same manner as the hexagonal and elliptical polynomials. As in the case of hexagonal and elliptical polynomials, each rectangular polynomial also consists of either cosine or sine terms, but not both. Like the elliptical polynomials, the rectangular polynomials also consist of a larger number of circle polynomials. The rectangular polynomial R_{11} , like the elliptical polynomials E_{11} , representing a balanced primary spherical aberration is not radially symmetric, since it consists of a term in astigmatism Z_6 or $\cos 2\theta$. As discussed below, the rectangular polynomials reduce to the square polynomials as $a \rightarrow 1/\sqrt{2}$, and the slit polynomials for a slit pupil parallel to the x axis as $a \rightarrow 1$.

TABLE 5a Orthonormal Rectangular Polynomials R_j in Terms of Zernike Circle Polynomials Z_j Which Reduce to the Corresponding Square Polynomials as $a \rightarrow 1/\sqrt{2}$

$$\begin{aligned}
 R_1 &= Z_1 \\
 R_2 &= [\sqrt{3}/(2a)]Z_2 \\
 R_3 &= [\sqrt{3}/(2\sqrt{1-a^2})]Z_3 \\
 R_4 &= [\sqrt{5}/(4\sqrt{1-2a^2+2a^4})](Z_1 + \sqrt{3}Z_4) \\
 R_5 &= [\sqrt{3}/2/(2a\sqrt{1-a^2})]Z_5 \\
 R_6 &= \{\sqrt{5}/[8a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[(3-10a^2+12a^4-8a^6)Z_1 + \sqrt{3}(1-2a^2)Z_4 \\
 &\quad + \sqrt{6}(1-2a^2+2a^4)Z_6] \\
 R_7 &= [\sqrt{21}/(4\sqrt{2}\sqrt{27-81a^2+116a^4-62a^6})][\sqrt{2}(1+4a^2)Z_3 + 5Z_7] \\
 R_8 &= [\sqrt{21}/(4\sqrt{2a}\sqrt{35-70a^2+62a^4})][\sqrt{2}(5-4a^2)Z_2 + 5Z_8] \\
 R_9 &= \{\sqrt{5}/2\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[16a^2(27-81a^2+116a^4-62a^6)]\}[2\sqrt{2}(9-36a^2 \\
 &\quad + 52a^4-60a^6)Z_3 + (9-18a^2-26a^4)Z_7 + (27-54a^2+62a^4)Z_9] \\
 R_{10} &= \{\sqrt{5}/2/[16a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}[2\sqrt{2}(35-112a^2+128a^4-60a^6)Z_2 \\
 &\quad + (35-70a^2+26a^4)Z_8 + (35-70a^2+62a^4)Z_{10}] \\
 R_{11} &= [1/(16\mu)][8(3+4a^2-4a^4)Z_1 + 25\sqrt{3}Z_4 + 10\sqrt{6}(1-2a^2)Z_6 + 21\sqrt{5}Z_{11}] \\
 R_{12} &= \{3\mu/[16a^2\nu\eta]\}[(105-550a^2+1559a^4-2836a^6+2695a^8-1078a^{10})Z_1 \\
 &\quad + 5\sqrt{3}(14-74a^2+205a^4-360a^6+335a^8-134a^{10})Z_4 + (5\sqrt{3}/2)(35-156a^2 \\
 &\quad + 421a^4-530a^6+265a^8)Z_6 + 21\sqrt{5}(1-4a^2+6a^4-4a^6)Z_{11} + [(7/2)\sqrt{5}/2\eta/(1-a^2)]Z_{12}] \\
 R_{13} &= [\sqrt{21}/(16\sqrt{2a}\sqrt{1-3a^2+4a^4-2a^6})](\sqrt{3}Z_5 + \sqrt{5}Z_{13}) \\
 R_{14} &= \tau[6(245-1400a^2+3378a^4-4452a^6+3466a^8-1488a^{10}+496a^{12})Z_1 \\
 &\quad + 15\sqrt{3}(49-252a^2+522a^4-540a^6+270a^8)Z_4 + 15\sqrt{6}(49-252a^2+534a^4-596a^6 \\
 &\quad + 360a^8-144a^{10})Z_6 + 3\sqrt{5}(49-196a^2+282a^4-172a^6+86a^8)Z_{11} \\
 &\quad + 147\sqrt{10}(1-4a^2+6a^4-4a^6)Z_{12} + 3\sqrt{10}\nu^2Z_{14}] \\
 R_{15} &= \{1/[32a^3(1-a^2)(1-3a^2+4a^4-2a^6)^{1/2}]\}[3\sqrt{7}/2(5-18a^2+24a^4-16a^6)Z_5 \\
 &\quad + \sqrt{105}/2(1-2a^2)Z_{13} + \sqrt{210}(1-2a^2+2a^4)Z_{15}] \\
 \mu &= (9-36a^2+103a^4-134a^6+67a^8)^{1/2} \\
 \nu &= (49-196a^2+330a^4-268a^6+134a^8)^{1/2} \\
 \tau &= 1/[128\nu a^4(1-a^2)^2] \\
 \eta &= 9-45a^2+139a^4-237a^6+210a^8-67a^{10}
 \end{aligned}$$

TABLE 5b Orthonormal *Rectangular* Polynomials $R_j(\rho, \theta)$ in Polar Coordinates

$$\begin{aligned}
 R_1 &= 1 \\
 R_2 &= (\sqrt{3}/a)\rho \cos \theta \\
 R_3 &= \sqrt{3/(1-a^2)}\rho \sin \theta \\
 R_4 &= [\sqrt{5}/(2\sqrt{1-2a^2+2a^4})](3\rho^2-1) \\
 R_5 &= [3/(2a\sqrt{1-a^2})]\rho^2 \sin 2\theta \\
 R_6 &= \{\sqrt{5}/[4a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[3(1-2a^2+2a^4)\rho^2 \cos 2\theta + 3(1-2a^2)\rho^2 \\
 &\quad - 2a^2(1-a^2)(1-2a^2)] \\
 R_7 &= [\sqrt{21}/(2\sqrt{27-81a^2+116a^4-62a^6})](15\rho^2-9+4a^2)\rho \sin \theta \\
 R_8 &= [\sqrt{21}/(2a\sqrt{35-70a^2+62a^4})](15\rho^2-5-4a^2)\rho \cos \theta \\
 R_9 &= \{\sqrt{5}\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[8a^2(27-81a^2+116a^4-62a^6)]\}\{(27-54a^2+62a^4) \\
 &\quad \times \rho^3 \sin 3\theta - 3[4a^2(3-13a^2+10a^4) - (9-18a^2-26a^4)\rho^2]\rho \sin \theta\} \\
 R_{10} &= \{\sqrt{5}/[8a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}\{(35-70a^2+62a^4)\rho^3 \cos 3\theta \\
 &\quad - 3[4a^2(7-17a^2+10a^4) - (35-70a^2+26a^4)\rho^2]\rho \cos \theta\} \\
 R_{11} &= [1/(8\mu)][315\rho^4 + 30(1-2a^2)\rho^2 \cos 2\theta - 240\rho^2 + 27 + 16a^2 - 16a^4] \\
 R_{12} &= [3\mu/(8a^2\nu\eta)][315(1-2a^2)(1-2a^2+2a^4)\rho^4 + 5(7\mu^2\rho^2 - 21 + 72a^2 - 225a^4 + 306a^6 \\
 &\quad - 153a^8)\rho^2 \cos 2\theta - 15(1-2a^2)(7+4a^2-71a^4+134a^6-67a^8)\rho^2 \\
 &\quad + a^2(1-a^2)(1-2a^2)(70-233a^2+233a^4)] \\
 R_{13} &= [\sqrt{21}/(4a\sqrt{1-3a^2+4a^4-2a^6})](5\rho^2-3)\rho^2 \sin 2\theta \\
 R_{14} &= 6\tau\{5\nu^2\rho^4 \cos 4\theta - 20(1-2a^2)[6a^2(7-16a^2+18a^4-9a^6) - 49(1-2a^2+2a^4)\rho^2]\rho^2 \cos 2\theta \\
 &\quad + 8a^4(1-a^2)^2(21-62a^2+62a^4) - 120a^2(7-30a^2+46a^4-23a^6)\rho^2 \\
 &\quad + 15(49-196a^2+282a^4-172a^6+86a^8)\rho^4\} \\
 R_{15} &= \{\sqrt{21}/[8a^3(1-a^2)^{3/2}(1-2a^2+2a^4)^{1/2}]\}[-(1-2a^2)(6a^2-6a^4-5\rho^2)\rho^2 \sin 2\theta \\
 &\quad + (5/2)(1-2a^2+2a^4)\rho^4 \sin 4\theta]
 \end{aligned}$$

TABLE 5c Orthonormal *Rectangular* Polynomials $R_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-a \leq x \leq a$, and $-\sqrt{1-a^2} \leq y \leq \sqrt{1-a^2}$

$$\begin{aligned}
 R_1 &= 1 \\
 R_2 &= (\sqrt{3}/a)x \\
 R_3 &= \sqrt{3/(1-a^2)}y \\
 R_4 &= [\sqrt{5}/(2\sqrt{1-2a^2+2a^4})](3\rho^2-1) \\
 R_5 &= [3/(a\sqrt{1-a^2})]xy \\
 R_6 &= \{\sqrt{5}/[2a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[3(1-a^2)^2x^2 - 3a^4y^2 - a^2(1-3a^2+2a^4)] \\
 R_7 &= [\sqrt{21}/(2\sqrt{27-81a^2+116a^4-62a^6})](15\rho^2-9+4a^2)y \\
 R_8 &= [\sqrt{21}/(2a\sqrt{35-70a^2+62a^4})](15\rho^2-5-4a^2)x \\
 R_9 &= \{\sqrt{5}\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[2a^2(27-81a^2+116a^4-62a^6)]\}[27(1-a^2)^2x^2 \\
 &\quad - 35a^4y^2 - a^2(9-39a^2+30a^4)]y \\
 R_{10} &= \{\sqrt{5}/[2a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}[35(1-a^2)^2x^2 - 27a^4y^2 - a^2(21-51a^2+30a^4)]x \\
 R_{11} &= [1/(8\mu)][315\rho^4 - 30(7+2a^2)x^2 - 30(9-2a^2)y^2 + 27 + 16a^2 - 16a^4] \\
 R_{12} &= [3\mu/(8a^2\nu\eta)][35(1-a^2)^2(18-36a^2+67a^4)x^4 + 630(1-2a^2)(1-2a^2+2a^4)x^2y^2 \\
 &\quad - 35a^4(49-98a^2+67a^4)y^4 - 30(1-a^2)(7-10a^2-12a^4+75a^6-67a^8)x^2 \\
 &\quad - 30a^2(7-77a^2+189a^4-193a^6+67a^8)y^2 + a^2(1-a^2)(1-2a^2)(70-233a^2+233a^4)] \\
 R_{13} &= [\sqrt{21}/(2a\sqrt{1-3a^2+4a^4-2a^6})](5\rho^2-3)xy \\
 R_{14} &= 16\tau[735(1-a^2)^4x^4 - 540a^4(1-a^2)^2x^2y^2 + 735a^8y^4 - 90a^2(1-a^2)^3(7-9a^2)x^2 \\
 &\quad + 90a^6(1-a^2)(2-9a^2)y^2 + 3a^4(1-a^2)^2(21-62a^2+62a^4)] \\
 R_{15} &= \{\sqrt{21}/[2a^3(1-a^2)\sqrt{1-3a^2+4a^4-2a^6}]\}[5(1-a^2)^2x^2 - 5a^4y^2 - a^2(3-9a^2+6a^4)]xy
 \end{aligned}$$

11.9 SQUARE POLYNOMIALS

Figure 3e shows a *unit square* inscribed inside a unit circle, as in the case of a rectangle. The distance of a corner point of the square, such as A , from its center O is unity, but each of its sides has a length of $\sqrt{2}$ and its area is 2. The orthonormality of the square polynomials $S_j(x, y)$ is represented by²⁶

$$\frac{1}{2} \int_{-1/\sqrt{2}}^{1/\sqrt{2}} dy \int_{-1/\sqrt{2}}^{1/\sqrt{2}} S_j(x, y) S_{j'}(x, y) dx = \delta_{jj'} \quad (38)$$

The square polynomials through the eighth order are given in terms of the Zernike polynomials in Table 6a. The first 15 polynomials are given in their analytical form, but those with $j > 15$ are written in a numerical form because of the increasing complexity of the coefficients of the circle polynomials. The corresponding polynomials in polar and Cartesian coordinates are given in Tables 6b and 6c, respectively. Of course, up to the fourth order, they can be obtained simply from the rectangular polynomials $R_j(x, y)$ given in Tables 5 by letting $a=1/\sqrt{2}$. The square polynomial S_{11} representing the balanced primary spherical aberration is radially symmetric, but the polynomial S_{22} representing the balanced secondary spherical aberration is not, since it consists of a term in Z_{14} or $\cos 4\theta$ also. Similarly, the polynomial S_{37} representing the balanced tertiary spherical aberration is also not radially symmetric, since it consists of terms in Z_{14} and Z_{26} both varying as $\cos 4\theta$.

11.10 SLIT POLYNOMIALS

By letting $a \rightarrow 1$ in the rectangular pupil, we obtain a *unit slit* pupil that is parallel to the x axis, as illustrated in Figure 3f. The corresponding orthonormal polynomials representing balanced aberrations for such pupils can be obtained from the rectangular polynomials $R_j(x, y)$ given in Table 5c by letting $y \rightarrow 0$ and $a \rightarrow 1$. Half of the rectangular polynomials thus reduce to zero. Some of the other polynomials are redundant. For example, the one-dimensional defocus and astigmatism can not be distinguished from each other. The slit polynomials are orthonormal according to²⁶

$$\frac{1}{2} \int_{-1}^1 P_j(x) P_{j'}(x) dx = \delta_{jj'} \quad (39)$$

The relevant orthonormal slit polynomials are listed in Table 7. They are the Legendre polynomials,¹⁷ which represent balanced aberrations uniquely.^{20,26} Since the pupil is one dimensional along the x axis, the aberrations vary with x only.

11.11 ABERRATION BALANCING AND TOLERANCING, AND DIFFRACTION FOCUS

For small aberrations, the Strehl ratio of the image of a point object is approximately given by $1 - \sigma^2$ or $\exp(-\sigma^2)$ when the standard deviation σ of the aberration is in units of radians.^{4,5,35} The Zernike circle and annular polynomials are separable in ρ and θ . The balanced spherical aberrations for these radially symmetric pupils are radially symmetric, and the balanced primary astigmatism for them has the same form. This is also true of a Gaussian circular or annular pupil, again because of the radial symmetry of the pupil and the amplitude across it.²³⁻²⁵ From the orthonormal form H_4 of defocus for a hexagonal pupil, the sigma of the defocus aberration ρ^2 is given by $(1/12)\sqrt{43/5}$. The hexagonal polynomials H_5 and H_6 show that the balanced astigmatism has the same form as the circle polynomials Z_5 and Z_6 , respectively. Thus the relative amount of defocus ρ^2 that balances classical or Seidel astigmatism $\rho^2 \cos^2 \theta$ is the same for a hexagonal pupil as for a circular pupil. Hence, for

TABLE 6a Orthonormal Square Polynomials S_j in Terms of Zernike Circle Polynomials Z_j

$$S_1 = Z_1$$

$$S_2 = \sqrt{3/2}Z_2$$

$$S_3 = \sqrt{3/2}Z_3$$

$$S_4 = (\sqrt{5/2/2})Z_1 + (\sqrt{15/2/2})Z_4$$

$$S_5 = \sqrt{3/2}Z_5$$

$$S_6 = (\sqrt{15/2})Z_6$$

$$S_7 = (3\sqrt{21/31/2})Z_3 + (5\sqrt{21/62/2})Z_7$$

$$S_8 = (3\sqrt{21/31/2})Z_2 + (5\sqrt{21/62/2})Z_8$$

$$S_9 = -(7\sqrt{5/31/2})Z_3 - (13\sqrt{5/62/4})Z_7 + (\sqrt{155/2/4})Z_9$$

$$S_{10} = (7\sqrt{5/31/2})Z_2 + (13\sqrt{5/62/4})Z_8 + (\sqrt{155/2/4})Z_{10}$$

$$S_{11} = (8/\sqrt{67})Z_1 + (25\sqrt{3/67/4})Z_4 + (21\sqrt{5/67/4})Z_{11}$$

$$S_{12} = (45\sqrt{3/16})Z_6 + (21\sqrt{5/16})Z_{12}$$

$$S_{13} = (3\sqrt{7/8})Z_5 + (\sqrt{105/8})Z_{13}$$

$$S_{14} = 261/(8\sqrt{134})Z_1 + (345\sqrt{3/134/16})Z_4 + (129\sqrt{5/134/16})Z_{11} + (3\sqrt{335/16})Z_{14}$$

$$S_{15} = (\sqrt{105/4})Z_{15}$$

$$S_{16} = 1.71440511Z_2 + 1.71491497Z_8 + 0.65048499Z_{10} + 1.52093102Z_{16}$$

$$S_{17} = 1.71440511Z_3 + 1.71491497Z_7 - 0.65048449Z_9 + 1.52093102Z_{17}$$

$$S_{18} = 4.10471345Z_2 + 3.45884077Z_8 + 5.34411808Z_{10} + 1.51830574Z_{16} + 2.80808005Z_{18}$$

$$S_{19} = -4.10471345Z_3 - 3.45884078Z_7 + 5.34411808Z_9 - 1.51830575Z_{17} + 2.80808005Z_{19}$$

$$S_{20} = 5.57146696Z_2 + 4.44429264Z_8 + 3.00807599Z_{10} + 1.70525179Z_{16} + 1.16777987Z_{18}$$

$$+ 4.19716701Z_{20}$$

$$S_{21} = 5.57146696Z_3 + 4.44429264Z_7 - 3.00807599Z_9 + 1.70525179Z_{17} - 1.16777988Z_{19}$$

$$+ 4.19716701Z_{21}$$

$$S_{22} = 1.33159935Z_1 + 1.94695912Z_4 + 1.74012467Z_{11} + 0.65624211Z_{14} + 1.50989174Z_{22}$$

$$S_{23} = 0.95479991Z_5 + 1.01511643Z_{13} + 1.28689496Z_{23}$$

$$S_{24} = 9.87992565Z_6 + 7.28853095Z_{12} + 3.38796312Z_{24}$$

$$S_{25} = 5.61978925Z_{15} + 2.84975327Z_{25}$$

$$S_{26} = 11.00650275Z_1 + 14.00366597Z_4 + 9.22698484Z_{11} + 13.55765720Z_{14} + 3.18799971Z_{22}$$

$$+ 5.11045000Z_{26}$$

$$S_{27} = 4.24396143Z_5 + 2.70990074Z_{13} + 0.84615108Z_{23} + 5.17855026Z_{27}$$

$$S_{28} = 17.58672314Z_6 + 11.15913268Z_{12} + 3.57668869Z_{24} + 6.44185987Z_{28}$$

$$S_{29} = 2.42764289Z_3 + 2.69721906Z_7 - 1.56598064Z_9 + 2.12208902Z_{17} - 0.93135653Z_{19}$$

$$+ 0.25252773Z_{21} + 1.59017528Z_{29}$$

$$S_{30} = 2.42764289Z_2 + 2.69721906Z_8 + 1.56598064Z_{10} + 2.12208902Z_{16} + 0.93135653Z_{18}$$

$$+ 0.25252773Z_{20} + 1.59017528Z_{30}$$

$$S_{31} = -9.10300982Z_3 - 8.79978208Z_7 + 10.69381427Z_9 - 5.37383385Z_{17} + 7.01044701Z_{19}$$

$$- 1.26347272Z_{21} - 1.90131756Z_{29} + 3.07960207Z_{31}$$

$$S_{32} = 9.10300982Z_2 + 8.79978208Z_8 + 10.69381427Z_{10} + 5.37383385Z_{16} + 7.01044701Z_{18}$$

$$+ 1.26347272Z_{20} + 1.90131756Z_{30} + 3.07960207Z_{32}$$

$$S_{33} = 21.39630883Z_3 + 19.76696884Z_7 - 12.70550260Z_9 + 11.05819453Z_{17} - 7.02178756Z_{19}$$

$$+ 15.80286172Z_{21} + 3.29259996Z_{29} - 2.07602718Z_{31} + 5.40902889Z_{33}$$

$$S_{34} = 21.39630883Z_2 + 19.76696884Z_8 + 12.70550260Z_{10} + 11.05819453Z_{16} + 7.02178756Z_{18}$$

$$+ 15.80286172Z_{20} + 3.29259996Z_{30} + 2.07602718Z_{32} + 5.40902889Z_{34}$$

$$S_{35} = -16.54454462Z_3 - 14.89205549Z_7 + 22.18054997Z_9 - 7.94524849Z_{17} + 11.85458952Z_{19}$$

$$- 6.18963457Z_{21} - 2.19431441Z_{29} + 3.24324400Z_{31} - 1.72001172Z_{33} + 8.16384008Z_{35}$$

$$S_{36} = 16.54454462Z_2 + 14.89205549Z_8 + 22.18054997Z_{10} + 7.94524849Z_{16} + 11.85458952Z_{18}$$

$$+ 6.18963457Z_{20} + 2.19431441Z_{30} + 3.24324400Z_{32} + 1.72001172Z_{34} + 8.16384008Z_{36}$$

(Continued)

TABLE 6a Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Terms of Zernike Circle Polynomials (Continued)

$$S_{37} = 1.75238960Z_1 + 2.72870567Z_4 + 2.76530671Z_{11} + 1.43647360Z_{14} + 2.12459170Z_{22} \\ + 0.92450043Z_{26} + 1.58545010Z_{37}$$

$$S_{38} = 19.24848143Z_6 + 16.41468913Z_{12} + 9.76776798Z_{24} + 1.47438007Z_{28} + 3.83118509Z_{38}$$

$$S_{39} = 0.46604820Z_5 + 0.84124290Z_{13} + 1.00986774Z_{23} - 0.42520747Z_{27} + 1.30579570Z_{39}$$

$$S_{40} = 28.18104531Z_1 + 38.52219208Z_4 + 30.18363661Z_{11} + 36.44278147Z_{14} + 15.52577202Z_{22} \\ + 19.21524879Z_{26} + 4.44731721Z_{37} + 6.00189814Z_{40}$$

$$S_{41} = (369/4)\sqrt{35/3574}Z_{15} + [11781/(32\sqrt{3574})]Z_{25} + (2145/32)\sqrt{7/3574}Z_{41}$$

$$S_{42} = 85.33469748Z_6 + 64.01249391Z_{12} + 30.59874671Z_{24} + 34.09158819Z_{28} + 7.75796322Z_{38} \\ + 9.37150432Z_{42}$$

$$S_{43} = 14.30642479Z_5 + 11.17404702Z_{13} + 5.68231935Z_{23} + 18.15306055Z_{27} + 1.54919583Z_{39} \\ + 5.90178984Z_{43}$$

$$S_{44} = 36.12567424Z_1 + 47.95305224Z_4 + 35.30691679Z_{11} + 56.72014548Z_{14} + 16.36470429Z_{22} \\ + 26.32636277Z_{26} + 3.95466397Z_{37} + 6.33853092Z_{40} + 12.38056785Z_{44}$$

$$S_{45} = 21.45429746Z_{15} + 9.94633083Z_{25} + 2.34632890Z_{41} + 10.39130049Z_{45}$$

TABLE 6b Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Polar Coordinates

$$S_1 = 1$$

$$S_2 = \sqrt{6}\rho \cos \theta$$

$$S_3 = \sqrt{6}\rho \sin \theta$$

$$S_4 = \sqrt{5/2}(3\rho^2 - 1)$$

$$S_5 = 3\rho^2 \sin 2\theta$$

$$S_6 = 3\sqrt{5/2}\rho^2 \cos 2\theta$$

$$S_7 = \sqrt{21/31}(15\rho^2 - 7)\rho \sin \theta$$

$$S_8 = \sqrt{21/31}(15\rho^2 - 7)\rho \cos \theta$$

$$S_9 = (\sqrt{5/31/2})[31\rho^3 \sin 3\theta - 3(13\rho^2 - 4)\rho \sin \theta]$$

$$S_{10} = (\sqrt{5/31/2})[31\rho^3 \cos 3\theta + 3(13\rho^2 - 4)\rho \cos \theta]$$

$$S_{11} = [1/(2\sqrt{67})](315\rho^4 - 240\rho^2 + 31)$$

$$S_{12} = [15/(2\sqrt{2})](7\rho^2 - 3)\rho^2 \cos 2\theta$$

$$S_{13} = \sqrt{21/2}(5\rho^2 - 3)\rho^2 \sin 2\theta$$

$$S_{14} = [3/(8\sqrt{134})](335\rho^4 \cos 4\theta + 645\rho^4 - 300\rho^2 + 22)$$

$$S_{15} = (5/2)\sqrt{21/2}\rho^4 \sin 4\theta$$

$$S_{16} = \sqrt{55/1966}[11\rho^3 \cos 3\theta + 3(19 - 97\rho^2 + 105\rho^4)\rho \cos \theta]$$

$$S_{17} = \sqrt{55/1966}[-11\rho^3 \sin 3\theta + 3(19 - 97\rho^2 + 105\rho^4)\rho \sin \theta]$$

$$S_{18} = (1/4)\sqrt{3/844397}[5(-10099 + 20643\rho^2)\rho^3 \cos 3\theta + 3(3128 - 23885\rho^2 + 37205\rho^4)\rho \cos \theta]$$

$$S_{19} = (1/4)\sqrt{3/844397}[5(-10099 + 20643\rho^2)\rho^3 \sin 3\theta - 3(3128 - 23885\rho^2 + 37205\rho^4)\rho \sin \theta]$$

$$S_{20} = (1/16)\sqrt{7/859}[2577\rho^5 \cos 5\theta - 5(272 - 717\rho^2)\rho^3 \cos 3\theta + 30(22 - 196\rho^2 + 349\rho^4)\rho \cos \theta]$$

$$S_{21} = (1/16)\sqrt{7/859}[2577\rho^5 \sin 5\theta + 5(272 - 717\rho^2)\rho^3 \sin 3\theta + 30(22 - 196\rho^2 + 349\rho^4)\rho \sin \theta]$$

$$S_{22} = (1/4)\sqrt{65/849}(1155\rho^6 + 30\rho^4 \cos 4\theta - 1395\rho^4 + 453\rho^2 - 31)$$

$$S_{23} = (1/2)\sqrt{33/3923}(471 - 1820\rho^2 + 1575\rho^4)\rho^2 \sin 2\theta$$

$$S_{24} = (21/4)\sqrt{65/1349}(27 - 140\rho^2 + 165\rho^4)\rho^2 \cos 2\theta$$

$$S_{25} = (7/4)\sqrt{33/2}(9\rho^2 - 5)\rho^4 \sin 4\theta$$

$$S_{26} = [1/(16\sqrt{849})][5(-98 + 2418\rho^2 - 12051\rho^4 + 15729\rho^6) + 3(-8195 + 17829\rho^2)\rho^4 \cos 4\theta]$$

$$S_{27} = [1/(16\sqrt{7846})][27461\rho^6 \sin 6\theta + 15(348 - 2744\rho^2 + 4487\rho^4)\rho^2 \sin 2\theta]$$

$$S_{28} = [21/(32\sqrt{1349})][1349\rho^6 \cos 6\theta + 5(196 - 1416\rho^2 + 2247\rho^4)\rho^2 \cos 2\theta]$$

(Continued)

TABLE 6b Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Polar Coordinates (Continued)

$$\begin{aligned}
 S_{29} &= (-13.79189793\rho + 125.49411319\rho^3 - 308.13074909\rho^5 + 222.62454035\rho^7) \sin \theta \\
 &\quad + (8.47599260\rho^3 - 16.13156842\rho^5) \sin 3\theta + 0.87478174\rho^5 \sin 5\theta \\
 S_{30} &= (-13.79189793\rho + 125.49411319\rho^3 - 308.13074909\rho^5 + 222.62454035\rho^7) \cos \theta \\
 &\quad + (-8.47599260\rho^3 + 16.13156842\rho^5) \cos 3\theta + 0.87478174\rho^5 \cos 5\theta \\
 S_{31} &= (6.14762642\rho - 79.44065626\rho^3 + 270.16115026\rho^5 - 266.18445920\rho^7) \sin \theta \\
 &\quad + (56.29115383\rho^3 - 248.12774426\rho^5 + 258.68657393\rho^7) \sin 3\theta - 4.37679791\rho^5 \sin 5\theta \\
 S_{32} &= (-6.14762642\rho + 79.44065626\rho^3 - 270.16115026\rho^5 + 266.18445920\rho^7) \cos \theta \\
 &\quad + (56.29115383\rho^3 - 248.12774426\rho^5 + 258.68657393\rho^7) \cos 3\theta + 4.37679791\rho^5 \cos 5\theta \\
 S_{33} &= (-6.78771487\rho + 103.15977419\rho^3 - 407.15689696\rho^5 + 460.96399558\rho^7) \sin \theta \\
 &\quad + (-21.68093294\rho^3 + 127.50233381\rho^5 - 174.38628345\rho^7) \sin 3\theta \\
 &\quad + (-75.07397471\rho^5 + 151.45280913\rho^7) \sin 5\theta \\
 S_{34} &= (-6.78771487\rho + 103.15977419\rho^3 - 407.15689696\rho^5 + 460.96399558\rho^7) \cos \theta \\
 &\quad + (21.68093294\rho^3 - 127.50233381\rho^5 + 174.38628345\rho^7) \cos 3\theta \\
 &\quad + \rho^5(-75.07397471 + 151.45280913\rho^2) \cos 5\theta \\
 S_{35} &= (3.69268433\rho - 59.40323317\rho^3 + 251.40397826\rho^5 - 307.20401818\rho^7) \sin \theta \\
 &\quad + (28.20381860\rho^3 - 183.86176738\rho^5 + 272.43249673\rho^7) \sin 3\theta \\
 &\quad + (19.83875817\rho^5 - 48.16032819\rho^7) \sin 5\theta + 32.65536033\rho^7 \sin 7\theta \\
 S_{36} &= (-3.69268433\rho + 59.40323317\rho^3 - 251.40397826\rho^5 + 307.20401818\rho^7) \cos \theta \\
 &\quad + (28.20381860\rho^3 - 183.86176738\rho^5 + 272.43249673\rho^7) \cos 3\theta \\
 &\quad + (-19.83875817\rho^5 + 48.16032819\rho^7) \cos 5\theta + 32.65536033\rho^7 \cos 7\theta \\
 S_{37} &= 2.34475558 - 55.32128002\rho^2 + 296.53777290\rho^4 - 553.46621887\rho^6 + 332.94452229\rho^8 \\
 &\quad + (-12.75329096\rho^4 + 20.75498320\rho^6) \cos 4\theta \\
 S_{38} &= (-51.83202694\rho^2 + 451.93890159\rho^4 - 1158.49126888\rho^6 + 910.24313983\rho^8) \cos 2\theta \\
 &\quad + 5.51662508\rho^6 \cos 6\theta \\
 S_{39} &= (-39.56789598\rho^2 + 267.47071204\rho^4 - 525.02362247\rho^6 + 310.24123146\rho^8) \sin 2\theta \\
 &\quad - 1.59098067\rho^6 \sin 6\theta \\
 S_{40} &= 1.21593465 - 45.42224477\rho^2 + 373.41167834\rho^4 - 1046.32659847\rho^6 + 933.93661610\rho^8 \\
 &\quad + (137.71626496\rho^4 - 638.10242034\rho^6 + 712.98912399\rho^8) \cos 4\theta \\
 S_{41} &= (9/8)\sqrt{7/1787}(1455 - 5544\rho^2 + 5005\rho^4)\rho^4 \sin 4\theta \\
 S_{42} &= (-40.45171657\rho^2 + 494.75561036\rho^4 - 1738.64589491\rho^6 + 1843.19802390\rho^8) \cos 2\theta \\
 &\quad + (-150.76043598\rho^6 + 318.07940431\rho^8) \cos 6\theta \\
 S_{43} &= (-9.12193686\rho^2 + 110.47679089\rho^4 - 371.21215287\rho^6 + 368.07015240\rho^8) \sin 2\theta \\
 &\quad + (-107.35168289\rho^6 + 200.31338972\rho^8) \sin 6\theta \\
 S_{44} &= 0.58427150 - 25.29433513\rho^2 + 242.54313549\rho^4 - 795.02011474\rho^6 + 830.47943579\rho^8 \\
 &\quad + (90.22533813\rho^4 - 538.44320774\rho^6 + 752.97905752\rho^8) \cos 4\theta + 52.52630092\rho^8 \cos 8\theta \\
 S_{45} &= (31.08509142\rho^4 - 194.79990628\rho^6 + 278.72965314\rho^8) \sin 4\theta + 44.08655427\rho^8 \sin 8\theta
 \end{aligned}$$

TABLE 6c Orthonormal Square Polynomials $S_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-1/\sqrt{2} \leq x \leq 1/\sqrt{2}$, and, $-1/\sqrt{2} \leq y \leq 1/\sqrt{2}$

$$\begin{aligned}
 S_1 &= 1 \\
 S_2 &= \sqrt{6}x \\
 S_3 &= \sqrt{6}y \\
 S_4 &= \sqrt{5/2}(3\rho^2 - 1) \\
 S_5 &= 6xy \\
 S_6 &= 3\sqrt{5/2}(x^2 - y^2) \\
 S_7 &= \sqrt{21/31}(15\rho^2 - 7)y \\
 S_8 &= \sqrt{21/31}(15\rho^2 - 7)x \\
 S_9 &= \sqrt{5/31}(27x^2 - 35y^2 + 6)y \\
 S_{10} &= \sqrt{5/31}(35x^2 - 27y^2 - 6)x \\
 S_{11} &= [1/(2\sqrt{67})](315\rho^4 - 240\rho^2 + 31) \\
 S_{12} &= [15/(2\sqrt{2})](x^2 - y^2)(7\rho^2 - 3) \\
 S_{13} &= \sqrt{42}(5\rho^2 - 3)xy
 \end{aligned}$$

(Continued)

TABLE 6c Orthonormal Square Polynomials $S_j(x, y)$, in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-1/\sqrt{2} \leq x \leq 1/\sqrt{2}$, and, $-1/\sqrt{2} \leq y \leq 1/\sqrt{2}$ (Continued)

$$S_{14} = [3/(4\sqrt{134})][10(49x^4 - 36x^2y^2 + 49y^4) - 150\rho^2 + 11]$$

$$S_{15} = 5\sqrt{42}(x^2 - y^2)xy$$

$$S_{16} = \sqrt{55/1966}(315\rho^4 - 280x^2 - 324y^2 + 57)x$$

$$S_{17} = \sqrt{55/1966}(315\rho^4 - 324x^2 - 280y^2 + 57)y$$

$$S_{18} = (1/2)\sqrt{3/844397}[105(1023x^4 + 80x^2y^2 - 943y^4) - 61075x^2 + 39915y^2 + 4692]x$$

$$S_{19} = (1/2)\sqrt{3/844397}[105(943x^4 - 80x^2y^2 - 1023y^4) - 39915x^2 + 61075y^2 - 4692]y$$

$$S_{20} = (1/4)\sqrt{7/859}[6(693x^4 - 500x^2y^2 + 525y^4) - 1810x^2 - 450y^2 + 165]x$$

$$S_{21} = (1/4)\sqrt{7/859}[6(525x^4 - 500x^2y^2 + 693y^4) - 450x^2 - 1810y^2 + 165]y$$

$$S_{22} = (1/4)\sqrt{65/849}[1155\rho^6 - 15(91x^4 + 198x^2y^2 + 91y^4) + 453\rho^2 - 31]$$

$$S_{23} = \sqrt{33/3923}(1575\rho^4 - 1820\rho^2 + 471)xy$$

$$S_{24} = (21/4)\sqrt{65/1349}(165\rho^4 - 140\rho^2 + 27)(x^2 - y^2)$$

$$S_{25} = 7\sqrt{33/2}(9\rho^2 - 5)xy(x^2 - y^2)$$

$$S_{26} = [1/(8\sqrt{849})][42(1573x^6 - 375x^4y^2 - 375x^2y^4 + 1573y^6) - 60(707x^4 - 225x^2y^2 + 707y^4) + 6045\rho^2 - 245]$$

$$S_{27} = [1/(2\sqrt{7846})][14(2673x^4 - 2500x^2y^2 + 2673y^4) - 10290\rho^2 + 1305]xy$$

$$S_{28} = [21/(8\sqrt{1349})][3146x^6 - 2250x^4y^2 + 2250x^2y^4 - 3146y^6 - 1770(x^4 - y^4) + 245(x^2 - y^2)]$$

$$S_{29} = (-13.79189793 + 150.92209099x^2 + 117.01812058y^2 - 352.15154565x^4 - 657.27245247x^2y^2 - 291.12439892y^4 + 222.62454035x^6 + 667.87362106x^4y^2 + 667.87362106x^2y^4 + 222.62454035y^6)y$$

$$S_{30} = (-13.79189793 + 117.01812058x^2 + 150.92209099y^2 - 291.12439892x^4 - 657.27245247x^2y^2 - 352.15154565y^4 + 222.62454035x^6 + 667.87362106x^4y^2 + 667.87362106x^2y^4 + 222.62454035y^6)x$$

$$S_{31} = (6.14762642 + 89.43280522x^2 - 135.73181009y^2 - 496.10607212x^4 + 87.83479115x^2y^2 + 513.91209661y^4 + 509.87526260x^6 + 494.87949207x^4y^2 - 539.86680367x^2y^4 - 524.87103314y^6)y$$

$$S_{32} = (-6.14762642 + 135.73181009x^2 - 89.43280522y^2 - 513.91209661x^4 - 87.83479115x^2y^2 + 496.10607212y^4 + 524.87103314x^6 + 539.86680367x^4y^2 - 494.87949207x^2y^4 - 509.87526260y^6)x$$

$$S_{33} = (-6.78771487 + 38.11697536x^2 + 124.84070714y^2 - 400.01976911x^4 + 191.43062089x^2y^2 - 609.73320550y^4 + 695.06919087x^6 - 246.30347616x^4y^2 - 154.56957886x^2y^4 + 786.80308817y^6)y$$

$$S_{34} = (-6.78771487 + 124.84070714x^2 + 38.11697536y^2 - 609.73320550x^4 + 191.43062089x^2y^2 - 400.01976911y^4 + 786.80308817x^6 - 154.56957886x^4y^2 - 246.30347616x^2y^4 + 695.06919087y^6)x$$

$$S_{35} = (3.69268433 + 25.20822264x^2 - 87.60705178y^2 - 200.98753298x^4 - 63.30315999x^2y^2 + 455.10450382y^4 + 497.87935336x^6 - 461.58554163x^4y^2 + 470.02596297x^2y^4 - 660.45220344y^6)y$$

$$S_{36} = (-3.69268433 + 87.60705178x^2 - 25.20822264y^2 - 455.10450382x^4 + 63.30315999x^2y^2 + 200.98753298y^4 + 660.45220344x^6 - 470.02596297x^4y^2 + 461.58554163x^2y^4 - 497.87935336y^6)x$$

$$S_{37} = 9.37902233 - 221.28512011\rho^2 + 1186.15109160\rho^4 - 2213.86487550\rho^6 + 1331.77808917\rho^8 + 0.0190064(x^4 - 6x^2y^2 + y^4)(-671 + 1092\rho^2)$$

$$S_{38} = (-51.83202694 + 451.93890159x^2 - 1152.97464379x^4 + 910.24313983x^6)x^2 + (51.83202694 - 451.93890159y^2 - 1241.24064523x^4 + 1241.24064523x^2y^2 + 1152.97464379y^4 + 1820.48627967x^6 - 1820.48627967x^2y^2 - 910.24313983y^6)y^2$$

$$S_{39} = (-79.13579197 + 534.94142408x^2 + 534.94142408y^2 - 1059.59312899x^4 - 2068.27487642x^2y^2 - 1059.59312899y^4 + 620.48246292x^6 + 1861.44738877x^4y^2 + 1861.44738877x^2y^4 + 620.48246292y^6)xy$$

$$S_{40} = 1.21593465 + (-45.42224477 + 511.12794331x^2 - 1684.42901882x^4 + 1646.92574009x^6)x^2 + (-45.42224477 - 79.47423312x^2 + 511.12794331y^2 + 51.53230630x^4 + 51.53230630x^2y^2 - 1684.42901882y^4 + 883.78996844x^6 - 1526.27154329x^4y^2 + 883.78996844x^2y^4 + 1646.92574009y^6)y^2$$

$$S_{41} = (409.79084415x^2 - 409.79084415y^2 - 1561.42985567x^4 + 1561.42985567y^4 + 1409.62417525x^6 + 1409.62417525x^4y^2 - 1409.62417525x^2y^4 - 1409.62417525y^6)xy$$

$$S_{42} = (-40.45171657 + 494.75561036x^2 - 1889.40633090x^4 + 2161.27742821x^6)x^2 + (40.45171657 - 494.75561036y^2 + 522.76064491x^4 - 522.76064491x^2y^2 + 1889.40633090y^4 - 766.71561254x^6 + 766.71561254x^4y^2 - 2161.27742821y^6)y^2$$

$$S_{43} = (-18.24387372 + 220.95358178x^2 + 220.95358178y^2 - 1386.53440310x^4 + 662.18504631x^2y^2 - 1386.53440310y^4 + 1938.02064313x^6 - 595.96654168x^4y^2 - 595.96654168x^2y^4 + 1938.02064313y^6)xy$$

$$S_{44} = 0.58427150 + (-25.29433513 + 332.76847363x^2 - 1333.46332249x^4 + 1635.98479424x^6)x^2 + (-25.29433513 - 56.26575785x^2 + 332.76847363y^2 + 307.15569451x^4 + 307.15569451x^2y^2 - 1333.46332249y^4 - 1160.73491284x^6 + 1129.92710444x^4y^2 - 1160.73491284x^2y^4 + 1635.98479424y^6)y^2$$

$$S_{45} = (124.34036571x^2 - 124.34036571y^2 - 779.19962514x^4 + 779.19962514y^4 + 1467.61104674x^6 - 1353.92842666x^4y^2 + 1353.92842666x^2y^4 - 1467.61104674y^6)xy$$

TABLE 7 Orthonormal Polynomials for a Unit *Slit* Pupil

j	Aberration	Orthonormal Polynomials
1	Piston	1
2	Tilt	$\sqrt{3}x$
3	Defocus	$(\sqrt{5}/2)(3x^2 - 1)$
4	Coma	$(\sqrt{7}/2)(5x^3 - 3x)$
5	Spherical aberration	$(3/8)(35x^4 - 30x^2 + 3)$
6	Secondary coma	$(\sqrt{11}/8)(63x^5 - 70x^3 + 15x)$
7	Secondary spherical aberration	$(\sqrt{13}/16)(231x^6 - 315x^4 + 105x^2 - 5)$

a small amount of astigmatism, the diffraction focus for an inscribed hexagonal pupil is the same as for a circular pupil.^{4,5} For an image with a focal ratio of F , it lies along the z axis at a distance of $-8F^2$ times the amount of the balancing defocus from the Gaussian image point. However, the hexagonal polynomials H_7 and H_8 show that the relative amount of tilt $\rho \cos \theta$ that optimally balances classical or Seidel coma $\rho^3 \cos \theta$ is $-14/25 \approx -0.56$ compared to $-2/3$ for a circular pupil. The diffraction focus in this case lies along the x axis at a distance of $-2F$ times the amount of tilt from the Gaussian image point. Similarly, the hexagonal polynomial H_{11} shows that the relative amount of defocus that optimally balances classical primary or Seidel spherical aberration ρ^4 is $-257/301 \approx -0.85$ compared to a value of -1 for a circular pupil. It has the consequence that the diffraction focus lies closer to the Gaussian image point in the case of coma, and closer to the Gaussian image plane in the case of spherical aberration, compared to their corresponding locations for a circular pupil. While the balanced primary and secondary spherical aberrations H_{11} and H_{22} are radially symmetric, the balanced tertiary spherical aberration H_{37} is not. The tertiary spherical aberration ρ^8 is balanced not only by defocus ρ^2 and primary and secondary spherical aberrations ρ^4 and ρ^6 , but by a term in Z_{28} or $\rho^6 \cos 6\theta$ as well.

In the case of an elliptical pupil, the sigma of Seidel astigmatism $\rho^2 \cos \theta$ is given by $\sigma_a = 1/4$, independent of its aspect ratio b , and thus equal to that for a circular pupil. Since Seidel astigmatism x^2 varies only along the x axis for which the unit ellipse has the same length as a unit circle, the sigma is independent of b . The amount of balancing defocus ρ^2 for astigmatism is different in the case of an elliptical or a rectangular pupil from the value of $-1/2$ for a circular pupil. Moreover, for these pupils, spherical aberration ρ^4 is balanced not only by defocus ρ^2 but astigmatism $\rho^2 \cos^2 \theta$ as well. This is a consequence of the fact that the x and y dimensions of these pupils are not equal.

A square pupil is a special case of a rectangular pupil for which $a = 1/\sqrt{2}$. It is evident from the square polynomials S_5 and S_6 that they have the same form as the corresponding circle polynomials. Thus there is no additional defocus for balancing astigmatism, as may be seen by the absence of a Z_4 term in the expression for S_6 . Hence, the diffraction focus of a system does not change when its circular pupil is replaced by an inscribed square pupil. Unlike an elliptical or a rectangular pupil, the primary spherical aberration in a square pupil is balanced by defocus only, as is evident from the radially symmetric expression for S_{11} . However, the balanced secondary and tertiary spherical aberrations are not radially symmetric, since they contain angle-dependent terms varying as $\cos 4\theta$. From the polynomials S_7 , S_8 , and S_{11} , the diffraction foci in the case of coma and spherical aberration are closer to the Gaussian image point and the Gaussian image plane, respectively, compared to their corresponding locations for a circular pupil.

The sigma of Seidel aberrations with and without balancing are listed in Table 8 for elliptical and rectangular pupils. The corresponding values for a circular, hexagonal, square, and a slit pupil are listed in Table 9.²⁶ As expected, the results for an elliptical pupil reduce to those for a circular pupil as $b \rightarrow 1$, and the results for a rectangular pupil reduce to those for a square pupil as $a \rightarrow 1/\sqrt{2}$. As the area of a unit pupil decreases in going from a circular to a hexagonal to a square pupil (from π to $3\sqrt{3}/2 \approx 2.6$ to 2), the sigma of an aberration decreases and its tolerance for a certain Strehl ratio

TABLE 8 Standard Deviation or Sigma of a Primary and a Balanced Primary Aberration for Elliptical and Rectangular Pupils

Sigma	Elliptical	Rectangular
σ_d	$(1/4)[(3 - 2b^2 + 3b^4)/3]^{1/2}$	$(2/3)[(1 - 2a^2 + 2a^4)/5]^{1/2}$
σ_a	1/4	$2a^2/(3\sqrt{5})$
σ_{ba}	$b^2/[6(3 - 2b^2 + 3b^4)]^{1/2}$	$2a^2(1 - a^2)/\{3[5(1 - 2a^2 + 2a^4)]^{1/2}\}$
σ_c	$(5 + 2b^2 + b^4)^{1/2}/8$	$a[(7 + 8a^4)/105]^{1/2}$
σ_{bc}	$(9 - 6b^2 + 5b^4)^{1/2}/24$	$2a(35 - 70a^2 + 62a^4)^{1/2}/(15\sqrt{21})$
σ_s	$(225 + 60b^2 - 58b^4 + 60b^6 + 225b^8)^{1/2}/(24\sqrt{10})$	$4(63 - 162a^2 + 206a^4 - 88a^6 + 44a^8)^{1/2}/(45\sqrt{7})$
σ_{bc}	$(45 - 60b^2 + 94b^4 - 60b^6 + 45b^8)^{1/2}/(48\sqrt{5})$	$(8/315)(9 - 36a^2 + 103a^4 - 134a^6 + 67a^8)^{1/2}$

TABLE 9 Standard Deviation or Sigma of a Primary and a Balanced Primary Aberration for Circular, Hexagonal, Square, and Slit Pupils

Sigma	Circle	Hexagon	Square	Slit
σ_d	$1/(2\sqrt{3})$ = 1/3.464	$(1/12)\sqrt{43/5}$ = 1/4.092	$(1/3)\sqrt{2/5}$ = 1/4.743	$2/(3\sqrt{5})$ = 1/3.354
σ_a	1/4	$(1/24)\sqrt{127/5}$ = 1/4.762	$1/(3\sqrt{5})$ = 1/6.708	-
σ_{ba}	$1/(2\sqrt{6})$ = 1/4.899	$(1/4)\sqrt{7/15}$ = 1/5.855	$1/(3\sqrt{10})$ = 1/9.487	-
σ_c	$1/(2\sqrt{2})$ = 1/2.828	$(1/4)\sqrt{83/70}$ = 1/3.673	$\sqrt{3/70}$ = 1/4.831	$1\sqrt{7}$ = 1/2.646
σ_{bc}	$1/(6\sqrt{2})$ = 1/8.485	$(1/20)\sqrt{737/210}$ = 1/10.676	$(1/15)\sqrt{31/21}$ = 1/12.346	$2/(5\sqrt{7})$ = 1/6.614
σ_s	$2/(3\sqrt{5})$ = 1/3.354	$(1/6)\sqrt{59/35}$ = 1/4.621	$(2/45)\sqrt{101/7}$ = 1/5.923	4/15 = 1/3.750
σ_{bs}	$1/(6\sqrt{5})$ = 1/13.416	$(1/84)\sqrt{4987/215}$ = 1/17.441	$(2/315)\sqrt{67}$ = 1/19.242	8/105 = 1/13.125

increases. The slit pupil is more sensitive compared to a circular pupil, except for spherical aberration for which it is slightly less sensitive. To obtain the Seidel coefficients from the orthonormal coefficients of a noncircular wavefront, all significant coefficients that contain a Seidel term must be taken into account, just as in the case of Zernike coefficients.³⁶

11.12 ISOMETRIC, INTERFEROMETRIC, AND PSF PLOTS FOR ORTHONORMAL ABERRATIONS

The aberration-free point-spread functions (PSFs) for unit pupils considered in this chapter are shown in Fig. 5, illustrating their symmetry, for example, 6-fold symmetry for a hexagonal pupil. Their linear scale is such that the first zero of the PSF, for example, for a square pupil occurs at unity in units of λF (corresponding to 1.22 for a circular pupil). Here λ is the wavelength of the object radiation and F is the focal ratio of the image-forming light cone. These PSFs are the ultimate goal of fabrication and testing. The obscuration ratio of the annular pupil in Fig. 5b is $\epsilon = 0.5$; the aspect ratio of the elliptical pupil in Fig. 5d is $b = 0.85$; and the half width of the rectangular pupil in Fig. 5e

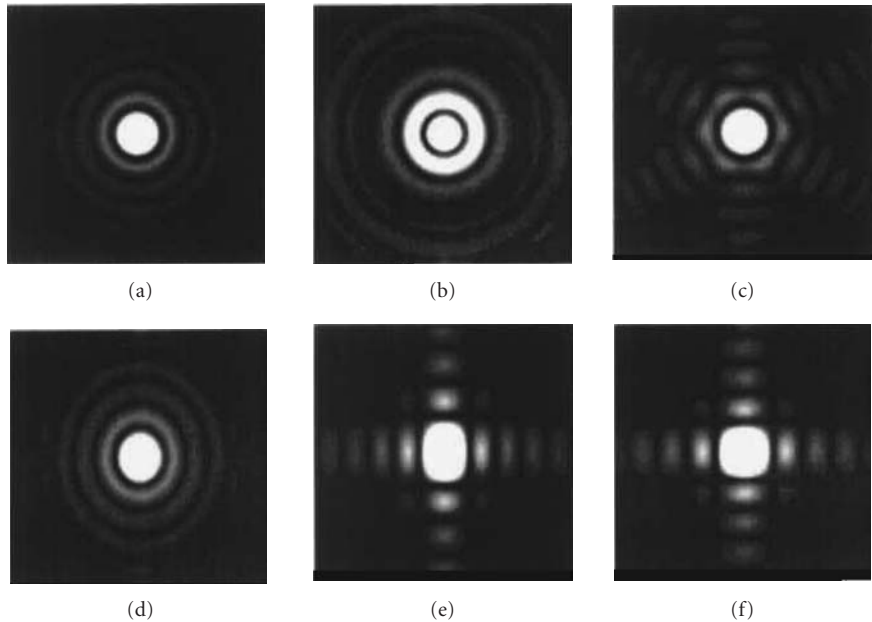


FIGURE 5 Aberration-free PSFs for different unit pupils: (a) Circular; (b) annular with obscuration ratio $\epsilon = 0.5$; (c) hexagonal; (d) elliptical with aspect ratio $b = 0.85$; (e) rectangular with half width $a = 0.8$; and (f) square.

is $a = 0.8$. The orthonormal polynomials corresponding to a Seidel aberration for a hexagonal, elliptical, rectangular, and square pupils are illustrated in three different but equivalent ways in Fig. 6.⁷ In Fig. 6d, as in Fig. 5a the aspect ratio of the elliptical pupil is $b = 0.85$. In Fig. 6e, as in Fig. 5e, the half width of the rectangular pupil is $a = 0.8$. For each polynomial, the isometric plot at the top illustrates its shape as produced, for example, in a deformable mirror. The standard deviation of each polynomial aberration in the figure is one wave. An interferogram, as in optical testing, is shown on the left. The number of fringes, which is equal to the number of times the aberration changes by one wave as we move from the center to the edge of a pupil, is different for the different polynomials. Each fringe represents a contour of constant phase or aberration. The fringe is dark when the phase is an odd multiple of π or the aberration is an odd multiple of $\lambda/2$. On the right for each polynomial are shown the PSFs, which represent the images of a point object in the presence of a polynomial aberration.

11.13 USE OF CIRCLE POLYNOMIALS FOR NONCIRCULAR PUPILS

Since the Zernike circle polynomials form a complete set, any wavefront, regardless of the shape of the pupil (which defines the perimeter of the wavefront) can be expanded in terms of them.³⁴ However, unless the pupil is circular, advantages of orthogonality and aberration balancing are lost. For example, the mean value of a Zernike circle polynomial across a noncircular pupil is not zero, the Zernike piston coefficient does not represent the mean value of the aberration, the other Zernike coefficients do not represent the standard deviation of the corresponding aberration terms, and the variance of the aberration is not equal to the sum of the squares of these other coefficients. Moreover, the value of a Zernike coefficient changes as the number of polynomials used in the

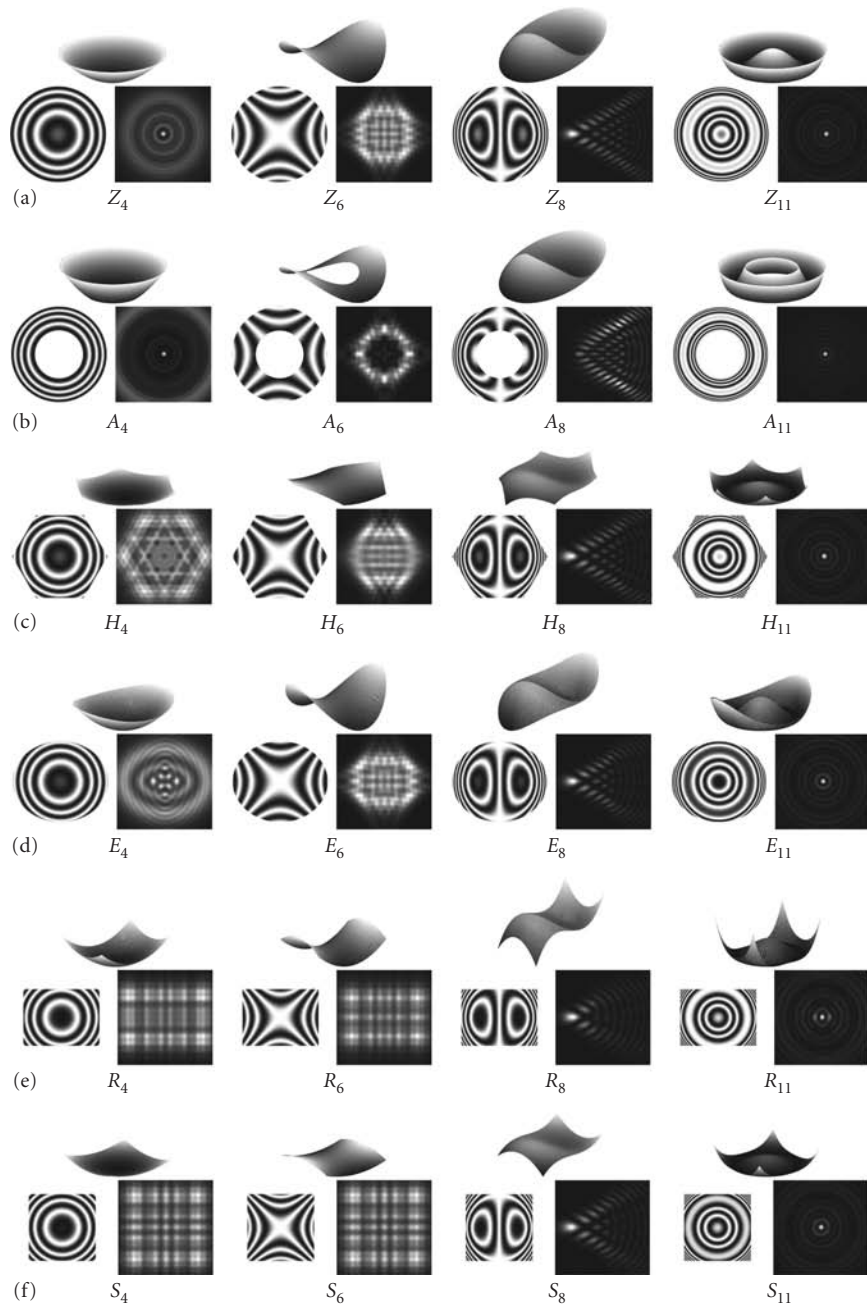


FIGURE 6 Isometric plots, interferograms, and PSFs for defocus ($j = 4$), astigmatism ($j = 6$), coma ($j = 8$), and spherical aberration ($j = 11$) in unit pupils. (a) Circular; (b) annular with $\epsilon = 0.5$; (c) hexagonal; (d) elliptical with aspect ratio $b = 0.85$; (e) rectangular with half width $a = 0.8$; and (f) square.

expansion of an aberration function changes. Hence, the circle polynomials are not appropriate for analysis of noncircular wavefronts. The polynomials given in this chapter for various pupils uniquely represent balanced classical aberrations that are also orthogonal across those pupils, just like the Zernike circle polynomials are for a circular pupil. Since each orthonormal polynomial is a linear combination of the Zernike circle polynomials, the wavefront fitting is as complete with the latter as it is with the former. However, since the circle polynomials do not represent the balanced classical aberrations for a noncircular pupil, the Zernike coefficients do not have the physical significance of their orthonormal counterparts. But the tip/tilt and defocus values in an interferometrically obtained aberration function, representing the lateral and longitudinal errors of an interferometer setting, obtained from the corresponding Zernike circle coefficients when the function is approximated with only the first four circle polynomials in a least square sense are identically the same as those obtained from the corresponding orthonormal coefficients. Accordingly, the aberration function obtained by subtracting the tip/tilt and defocus values from the measured aberration function is independent of the nature of the polynomials used in the expansion, regardless of the domain of the function or the shape of the pupil, so long as the nonorthogonal expansion is in terms of only the first four circle polynomials. The difference function is what is provided to the optician to zero out from the surface under fabrication by polishing.

11.14 DISCUSSION AND CONCLUSIONS

The Zernike circle polynomials are in widespread use for wavefront analysis in optical design and testing, because they are orthogonal over a unit circle and represent balanced aberrations of systems with circular pupils. When an aberration function of a circular wavefront is expanded in terms of them, the value of an expansion coefficient is independent of the number of polynomials used in the expansion. Accordingly, one or more terms can be added or subtracted without affecting the other coefficients. The piston coefficient represents the mean value of the aberration function and the other coefficients represent the standard deviation of the corresponding terms. The variance of the aberration is given simply by the sum of the squares of those other aberration coefficients.

We have also listed the orthonormal polynomials for analyzing the wavefronts across noncircular pupils, such as annular, hexagonal, elliptical, rectangular, and square. These polynomials are for unit pupils inscribed inside a unit circle. Such a choice keeps the maximum value of the distance of a point on the pupil from its center to be unity, thus easily identifying the peak of value of a classical aberration across it. Each orthonormal polynomial for the pupils considered consists of either the cosine or the sine terms, but not both due to the biaxial symmetry of the pupils. Whereas the circle and annular polynomials are separable in their dependence on the polar coordinates ρ and θ of a pupil point due to the radial symmetry of the pupils, only some of the polynomials for other pupils are separable. Hence polynomial numbering with two indices n and m , as for circular and annular polynomials, loses significance, and must be numbered with a single index j . The hexagonal polynomials H_{11} and H_{22} representing the balanced primary and secondary spherical aberrations are radially symmetric, but the polynomial H_{37} representing the balanced tertiary spherical aberration is not, since it contains an angle-dependent term in Z_{28} or $\cos 6\theta$ also. A hexagonal pupil has two distinct configurations where the hexagon in one is rotated by 30° with respect to that in the other. Only some of the polynomials are the same for the two configurations.²⁶ While the balancing defocus to optimally balance Seidel astigmatism for a hexagonal or a square pupil is the same as that for circular and annular pupils, it is different for the elliptical and rectangular pupils. For the elliptical and rectangular pupils, the Seidel or primary aberration ρ^4 is balanced not only by defocus but astigmatism as well. The square polynomial S_{11} representing the balanced primary spherical aberration is radially symmetric, but the square polynomial S_{22} representing the balanced secondary spherical aberration is not, since it contains a term in Z_{14} or $\cos 4\theta$ also. Similarly, the polynomial S_{37} representing the balanced tertiary spherical aberration is also not radially symmetric, since it consists of terms in Z_{14} and Z_{26} both varying as $\cos 4\theta$. We have illustrated orthonormal polynomials in three different but equivalent ways: isometrically, interferometrically, and by the corresponding aberrated PSFs.

The sigma of a Seidel aberration with and without balancing decreases as the area of a unit pupil decreases in going from a circular to a hexagonal to a square pupil. The sigma for Seidel astigmatism $\rho^2 \cos \theta$ for an elliptical pupil is independent of its aspect ratio and, therefore, is the same as for a circular pupil. This is due to the fact that the aberration is one dimensional along the dimension for which the unit ellipse has the same length as the unit circle. Since a slit pupil is one dimensional, there is no distinction between defocus and astigmatism. It is more sensitive to a Seidel aberration with or without balancing compared to a circular pupil, except for spherical aberration for which it is slightly less sensitive.

When the aberration function is known only at a discrete set of points, as in a digitized interferogram, the integral for determining the aberration coefficients reduces to a sum and the orthonormal coefficients thus obtained may be in error, since the polynomials are not orthonormal over the discrete points of the aberration data set. The magnitude of the error decreases as the number of points increases. This is not a serious problem when the wavefront errors are determined by, say, phase-shifting interferometry,³⁷ since the number of points can be very large. However, when the number of data points is small, or the pupil is irregular in shape due to vignetting, then ray tracing or testing of the system yields wavefront error data at an array of points across a region for which closed-form orthonormal polynomials are not available. In such cases, we can determine the coefficients of an expansion in terms of numerical polynomials that are orthogonal over the data set, obtained by the Gram-Schmidt orthogonalization process.^{7,38} However, if we just want to determine the values of tip/tilt and defocus terms, yielding the errors in interferometer settings, they can be obtained by least squares fitting the aberration function data with only these terms.

11.15 REFERENCES

1. F. Zernike, "Diffraction Theory of Knife-Edge Test and Its Improved Form, the Phase Contrast Method," *Mon. Not. R. Astron. Soc.* **94**: 377–384 (1934).
2. R. J. Noll, "Zernike Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **66**: 207–211 (1976).
3. B. R. A. Nijboer, "The Diffraction Theory of Optical Aberrations. Part II: Diffraction Pattern in the Presence of Small Aberrations," *Physica*. **13**: 605–620 (1947).
4. M. Born and E. Wolf, *Principles of Optics*, 7th ed., Oxford, New York (1999).
5. V. N. Mahajan, *Optical Imaging and Aberrations*, Part II: Wave Diffraction Optics, SPIE Press, Bellingham, Washington (Second Printing 2004).
6. V. N. Mahajan, "Zernike Polynomials and Aberration Balancing," *SPIE Proc.* **5173**: 1–17 (2003).
7. V. N. Mahajan, "Zernike Polynomials and Wavefront Fitting," in *Optical Shop Testing*, 3rd ed., D. Malacara, ed., Wiley, New York, pp. 498–546 (2007).
8. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 75–85 (1981).
9. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 1408 (1981).
10. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am. A* **1**: 685 (1984).
11. V. N. Mahajan, "Zernike Annular Polynomials and Optical Aberrations of Systems with Annular Pupils," *Appl. Opt.* **33**: 8125–8127 (1994).
12. <http://scikits.com/KFacts.html>
13. W. B. King, "The Approximation of Vignetted Pupil Shape by an Ellipse," *Appl. Opt.* **7**: 197–201 (1968).
14. G. Harbers, P. J. Kunst, and G. W. R. Leibbrandt, "Analysis of Lateral Shearing Interferograms by Use of Zernike Polynomials," *Appl. Opt.* **35**: 6162–6172 (1996).
15. H. Sumita, "Orthogonal Expansion of the Aberration Difference Function and Its Application to Image Evaluation," *Jpn. J. Appl. Phys.* **8**: 1027–1036 (1969).

16. K. N. LaFortune, R. L. Hurd, S. N. Fochs, M. D. Rotter, P. H. Pax, R. L. Combs, S. S. Olivier, J. M. Brase, and R. M. Yamamoto, "Technical Challenges for the Future of High Energy Lasers," *SPIE Proc.* **6454**: 1–11 (2007).
17. G. A. Korn and T. M. Korn, *Mathematical Handbook for Scientists and Engineers*, McGraw–Hill, New York, (1968).
18. V. N. Mahajan and G.-m. Dai, "Orthonormal Polynomials for Hexagonal Pupils," *Opt. Lett.* **31**: 2462–2465 (2006).
19. G.-m. Dai and V. N. Mahajan, "Nonrecursive Orthonormal Polynomials with Matrix Formulation," *Opt. Lett.* **32**: 74–76 (2007).
20. R. Barakat and L. Riseberg, "Diffraction Theory of the Aberrations of a Slit Aperture," *J. Opt. Soc. Am.* **55**: 878–881 (1965). There is an error in their polynomial S_2 , which should read as $x^2 - 1/3$.
21. M. Bray, "Orthogonal Polynomials: A Set for Square Areas," *SPIE Proc.* **5252**: 314–320 (2004).
22. J. L. Rayces, "Least-Squares Fitting of Orthogonal Polynomials to the Wave-Aberration Function," *Appl. Opt.* **31**: 2223–2228 (1992).
23. V. N. Mahajan, "Uniform Versus Gaussian Beams: a Comparison of the Effects of Diffraction, Obscuration, and Aberrations," *J. Opt. Soc. Am.* **A3**: 470–485 (1986).
24. V. N. Mahajan, "Zernike-Gauss Polynomials and Optical Aberrations of Systems with Gaussian Pupils," *Appl. Opt.* **34**: 8057–8059 (1995).
25. S. Szapiel, "Aberration Balancing Techniques for Radially Symmetric Amplitude Distributions; a Generalization of the Maréchal Approach," *J. Opt. Soc. Am.* **72**: 947–956 (1982).
26. V. N. Mahajan and G.-m. Dai, "Orthonormal Polynomials in Wavefront Analysis: Analytical Solution," *J. Opt. Soc. Am.* **A24**: 2994–3016 (2007).
27. A. B. Bhatia and E. Wolf, "On the Circle Polynomials of Zernike and Related Orthogonal Sets," in *Proc. Camb. Phil. Soc.* **50**: 40–48 (1954).
28. V. N. Mahajan, *Optical Imaging and Aberrations, Part I: Ray Geometrical Optics*, SPIE Press, Bellingham, Washington (Second Printing 2001).
29. W. T. Welford, *Aberrations of the Symmetrical Optical System*, Academic Press, New York (1974).
30. R. R. Shannon, *The Art and Science of Optical Design*, Cambridge University Press, New York (1997).
31. P. Mouroulis and J. Macdonald, *Geometrical Optics and Optical Design*, Oxford, New York (1997).
32. D. Malacara and Z. Malacara, *Handbook of Lens Design*, Dekkar, New York (1994).
33. G.-m. Dai and V. N. Mahajan, "Zernike Annular Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **A24**: 139–155 (2007).
34. G.-m. Dai and V. N. Mahajan, "Orthonormal Polynomials in Wavefront Analysis: Error Analysis," *Appl. Opt.* **47**: 3433–3445 (2008).
35. V. N. Mahajan, "Strehl Ratio for Primary Aberrations in Terms of Their Aberration Variance," *J. Opt. Soc. Am.* **73**: 860–861 (1983).
36. V. N. Mahajan and W. H. Swantner, "Seidel Coefficients in Optical Testing," *Asian J. Phys.* **15**: 203–209 (2006).
37. K. Creath, "Phase-Measurement Interferometry Techniques," *Progress in Optics*, E. Wolf, ed., Elsevier, New York, **26**: 349–393 (1988).
38. D. Malacara, J. M. Carpio-Valdéz, and J. Javier Sánchez-Mondragón, "Wavefront Fitting with Discrete Orthogonal Polynomials in a Unit Radius Circle," *Opt. Eng.* **29**: 672–675 (1990).

This page intentionally left blank.

Zacarías Malacara and Daniel Malacara-Hernández

*Centro de Investigaciones en Óptica, A. C.
León, Gto., México*

12.1 GLOSSARY

B	baseline length
f	focal length
f	signal frequency
f_b	back focal length
I	irradiance
N	average group refractive index
R	radius of curvature of an optical surface
r	radius of curvature of a spherometer ball
R	range
α	attenuation coefficient
λ	wavelength of light
Λ	synthetic wavelength
τ	delay time

In the optical shop, the measuring process has the purpose to obtain a comparison of physical variables using optical means. In this chapter we describe the most common procedures for the measurements of length, angle, curvature, and focal length of lenses and mirrors. The reader may obtain some more details about these procedures in the book *Optical Shop Testing* by D. Malacara,¹ or in Chap. 32, “Interferometers,” in Vol. I.

*Note: Figures 27, 30 and 31 are from *Optical Shop Testing*, 3d ed., edited by D. Malacara. (Reprinted with permission John Wiley and Sons, Inc. New York, 2007.)

12.2 INTRODUCTION AND DEFINITIONS

Lens parameter measuring in the optical shop has been a permanent problem in respect to the unit's agreement. Usually, rather than using SI units for length specifications, wavelength units is still a common reference for precision length measurements. Although the meter is obtained practically from the wavelength of a krypton source, helium-neon lasers are a common reference wavelength. The generally accepted measurement system of units is the International System or *Système International* (SI).

Time is the fundamental standard and is defined as follows: "The second is the duration of 9,192,631,770 periods of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium 133 atom." Formerly, the meter was a fundamental standard defined as 1,650,763.73 wavelengths in vacuum of the orange-red spectral line from the $2p_{10}$ and $5d_5$ levels of the krypton 86 atom. Shortly after the invention of the laser, it was proposed to use a laser line as a length standard.² In 1986, the speed of light was defined as 299,792,458 m/s, thus the meter is now a derived unit defined as "the distance traveled by light during $1/299,792,458$ of a second." The advantages of this definition, compared with the former meter definition, lie in the fact that it uses a relativistic constant, not subjected to physical influences, and is accessible and invariant. To avoid a previous definition of a time standard, the meter could be defined as "The length equal to $9,192,631,770/299,792,458$ wavelengths in vacuum of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium 133 atom."^{3,4} Modern description of the time and distance standards are described by Cundiff et al.⁵

The measurement process is prone to errors. They may be systematic or stochastic. A systematic error occurs in a poorly calibrated instrument. An instrument low in systematic error is said to be accurate. Accuracy is a measure of the amount of systematic errors. The accuracy is improved by adequate tracing to a primary standard. Stochastic errors appear due to random noise and other time-dependent fluctuations that affect the instrument. Stochastic errors may be reduced by taking several measurements and averaging them. A measurement from an instrument is said to be reproducible when the magnitude of stochastic errors is low. Reproducibility is a term used to define the repeatability for the measurements of an instrument. In measurements, a method for data acquisition and analysis has to be developed. Techniques for experimentation, planning, and data reduction can be found in the references.^{6,7}

12.3 LENGTH AND STRAIGHTNESS MEASUREMENTS

Length measurements may be performed by optical methods, since the definition of the meter is in terms of a light wavelength. Most of the length measurements are, in fact, comparisons to a secondary standard. Optical length measurements are made by comparisons to an external or internal scale a light time of flight, or by interferometric fringe counting.

Stadia and Range Finders

A *stadia* is an optical device used to determine distances. The principle of the measurement is a bar of known length W set at one end of the distance to be measured (Fig. 1). At the other end, a prism superimposes two images, one coming directly and the other after a reflection from a mirror, which are then observed through a telescope. At this point, the image of one end of the bar is brought in coincidence with that of the other end by rotating the mirror an angle θ . The mirror rotator is calibrated in such a way that for a given bar length W , a range R can be read directly in a dial, according to the equation:⁸

$$R = \frac{W}{\theta} \quad (1)$$

where θ is small and expressed in radians.

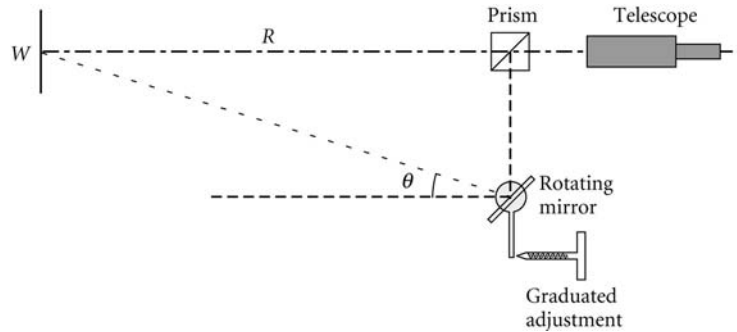


FIGURE 1 A stadia range meter. (From Patrick.⁸)

Another stadia method uses a graduated bar and a calibrated reticle in the telescope. For a known bar length W imaged on a telescope with focal length f , the bar on the focal plane will have a size i , and the range can be calculated approximately by:⁹

$$R = \left(\frac{f}{i} \right) W \quad (2)$$

Most surveying instruments have stadia markings usually in a 1:100 ratio, to measure distances from the so-called anallatic point, typically the instrument's center. A theodolite may be used for range measurements using the subtense bar method. In this method, a bar of known length is placed at the distance to be measured from the theodolite. By measuring the angle from one end to another end of the bar, the range can be easily measured.

A range finder is different from the stadia, in that the reference line is self-contained within the instrument. A general layout for a range finder is shown in Fig. 2. Two pentaprisms are separated a known baseline B ; two telescopes form an image, each through a coincidence prism. The images from the same reference point are superimposed in a split field formed by the coincidence prism. In one branch, a range compensator is adjusted to permit an effective coincidence from the reference images.

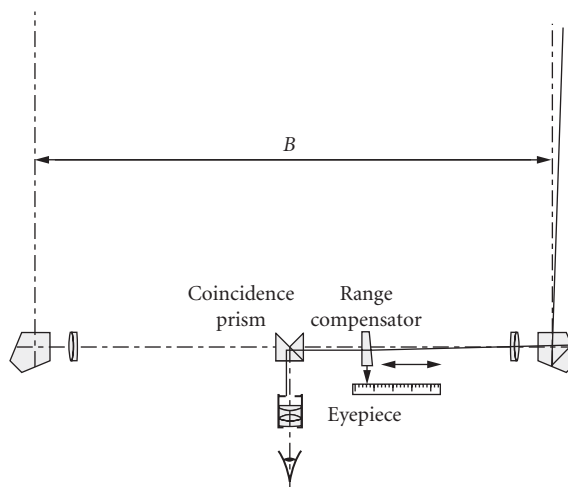


FIGURE 2 A range finder.

Assuming a baseline B and a range R , for small angles (large distances compared with the baseline), the range is

$$R = \frac{B}{\theta} \quad (3)$$

For an error $\Delta\theta$ in the angle measurement, the corresponding error ΔR in the range determination would be

$$\Delta R = -B\theta^{-2} \Delta\theta \quad (4)$$

and by substituting in Eq. (3),

$$\Delta R = -\frac{R^2}{B} \Delta\theta \quad (5)$$

From this last equation, it can be seen that the range error increases with the square of the range. Also, it is inversely proportional to the baseline. The angle error $\Delta\theta$ is a function of the eye's angular acuity, about 10 arcsec or 0.00005 rad.¹⁰

Pentaprisms permit a precise 90° deflection, independent of the alignment, and are like mirror systems with an angle of 45° between them. The focal length for the two telescopes in a range finder must be closely matched to avoid any difference in magnification. There are several versions of the range compensator. In one design, a sliding prism makes a variable deviation angle (Fig. 3a); in another system (Fig. 3b), a sliding prism displaces the image, without deviating it. A deviation can be also made with a rotating glass block (Fig. 3c). The Risley prisms (Fig. 3d) are a pair of counter rotating prisms located on the entrance pupil for one of the arms. Since the light beam is collimated, there is no astigmatism contribution from the prisms.¹⁰ A unit magnification Galilean telescope (Fig. 3e) is made with two weak lenses. A sliding lens is used to form a variable wedge to deviate the image path.⁸

Time-Based and Optical Radar

Distance measurements can also be done by the time-of-flight method. An application of the laser, obvious at the time of its advent, is the measurement of range. Distance determination by precise timing is known as optical radar. Optical radar has been used to measure the distance to the moon.

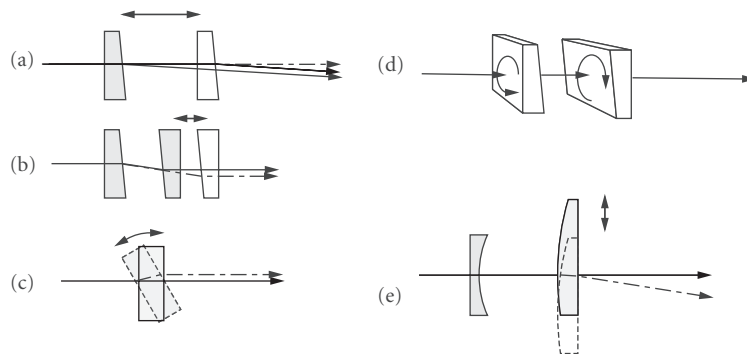


FIGURE 3 Range compensators for range finders (a) and (b) sliding prisms; (c) rotating glass block; (d) counterrotating prisms; and (e) sliding lens.

Since a small timing is involved, optical radar is applicable for distances from about 10 km. For distances from about a meter up to 50 km, modulated beams are used.

Laser radars measure the time of flight for a pulsed laser. Since accuracy depends on the temporal response of the electronic and detection system, optical radars are limited to distances larger than 1 km. Whenever possible, a cat's eye retro reflector is set at the range distance, making possible the use of a low power-pulsed laser. High-power lasers can be used over small distances without the need of a reflector. Accuracies of the order of 10^{-6} can be obtained.¹¹

The beam modulation method requires a high-frequency signal, about 10 to 30 MHz (modulating wavelength between 30 and 10 m) to modulate a laser beam carrier. The amplitude modulation is usually applied, but polarization modulation may also be used. With beam modulation distance measurements, the phase for the returning beam is compared with that of the output beam. The following description is for an amplitude modulation distance meter, although the same applies for polarizing modulation. Assuming a sinusoidally modulated beam, the returning beam can be described by

$$I_R = \alpha I_0 [1 + A \sin \omega(t - \tau)] \quad (6)$$

where α is the attenuation coefficient for the propagation media, I_0 is the output intensity for the exit beam; ω is 2π times the modulated beam frequency, and τ is the delay time. By measuring the relative phase between the outgoing and the returning beam, the quantity ω is measured. In most electronic systems, the delay time τ is measured, so, the length in multiples of the modulating wavelength is

$$L = \frac{c}{2N_g} \tau \quad (7)$$

where N_g is the average group refractive index for the modulating frequency.¹¹

Since the measured length is a multiple of the modulating wavelength, one is limited in range to one-half of the modulating wavelength. To measure with acceptable precision, and at the same time measure over large distances, several modulating frequencies are used, sometimes at multiples of ten to one another. The purpose is to obtain a synthetic wavelength Λ obtained from the mixing of two wavelengths λ_1, λ_2 as follows:^{12,13}

$$\Lambda = \frac{\lambda_1 \lambda_2}{\lambda_1 - \lambda_2} \quad (8)$$

To increment the range, several wavelengths are used to have several synthetic wavelengths. A frequency comb can increase the accuracy up to 8 nm in a range of 800 mm.¹⁴ The use of a femto-second mode locked laser produces the desired frequency comb.¹⁵

The traveling time is measured by comparing the phase of the modulating signal for the exiting and the returning beams. This phase comparison is sometimes made using a null method. In this method, a delay is introduced in the exiting signal, until it has the same phase as the returning beam, as shown in Fig. 4. Since the introduced delay is known, the light traveling time is thus determined.

Another method uses the down conversion of the frequency of both signals. These signals, with frequency f , are mixed with an electrical signal with frequency f_o , in order to obtain a signal with lower frequency f_L , in the range of a few kHz. The phase difference between the two low-frequency signals is the same as that between the two original signals. The lowering of the frequencies permits us to use conventional electronic methods to determine the phase difference between the two signals. A broad study on range finders has been done by Stitch.¹⁶

Interferometric Measurement of Small Distances

Interferometric methods may be used to measure small distances with a high degree of accuracy. This is usually done with a Michelson interferometer by comparing the thickness of the lens or glass

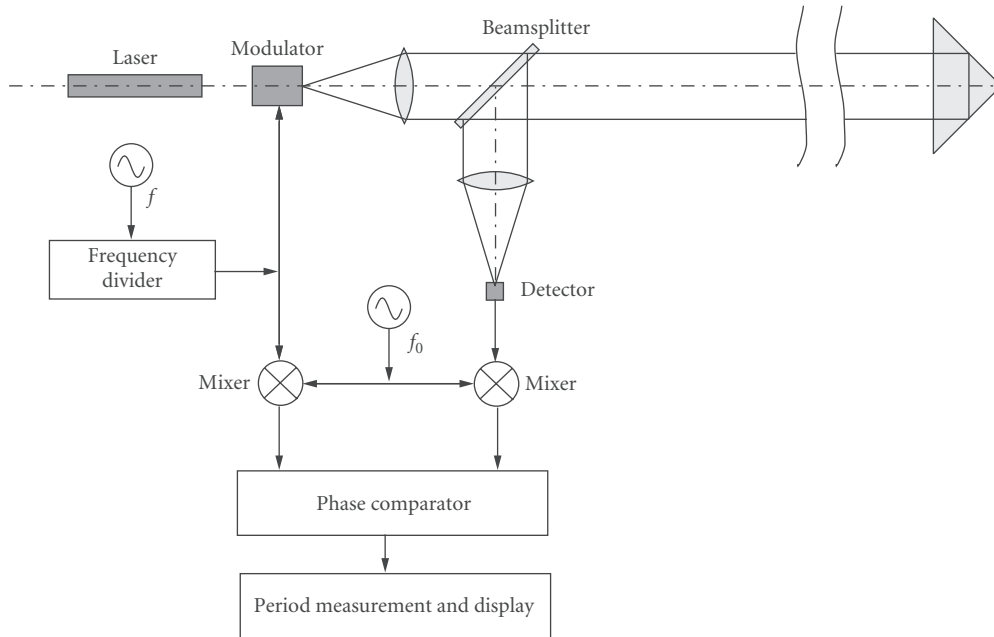


FIGURE 4 A wave modulation distance meter.

plate with that of a calibrated reference glass plate. Both plates must have approximately the same thickness and should be made with the same material.¹⁷ The two mirrors of a dispersion-compensated Michelson interferometer are replaced by the glass plate to be measured and by a reference plane-parallel plate of the same material as the lens. (The next step is to adjust the interferometer, to produce white-light Newton rings with the front surfaces of the lens and the plate.) Then, the plate is translated along the arm of the interferometer until the rear surface produces white-light rings. The displacement is the optical thickness Nt of the measured plate.

Interferometric Measurement of Medium Distances

Long and medium distances may also be measured by interferometric methods.^{18,19,20} Basically, the method counts the fringes in an interferometer while increasing or decreasing the optical path difference. The low temporal coherence or monochromaticity of most light sources limits this procedure to short distances. Lasers, however, have a much longer coherence length, due to their higher monochromaticity. Using lasers, it has been possible to make interferometric distance measurements over several meters.

In these interferometers, three things should be considered during their design. The first is that the laser light illuminating the interferometer should not be reflected back to the laser because that would cause instabilities in the laser, resulting in large irradiance fluctuations. As the optical path difference is changed by moving one of the mirrors, an irradiance detector in the pattern will detect a sinusoidally varying signal, but the direction of this change cannot be determined. Therefore, the second thing to be considered is that there should be a way to determine if the fringe count is going up or down; that is, if the distance is increasing or decreasing. There are two basic approaches to satisfy this last requirement. One is by producing two sinusoidal signals in phase quadrature (phase difference of 90° between them). The direction of motion of the moving prism may be sensed by determining the phase of which signal leads or lags the phase of the other signal. This information

is used to make the fringe counter increase or decrease. The alternative method uses a laser with two different frequencies. Finally, the third thing to consider in the interferometer design is that the number of fringes across its aperture should remain low and constant while moving the reflector in one of the two interferometer arms. This last condition is easily satisfied by using retroreflectors instead of mirrors. Then the two interfering wavefronts will always be almost flat and parallel. A typical retroreflector is a cube corner prism of reasonable quality to keep the number of fringes over the aperture low.

One method of producing two signals in quadrature is to have only a small number of fringes over the interferogram aperture. One possible method is by deliberately using an imperfect retroreflector. Then, the two desired signals are two small slits parallel to the fringes and separated by one-fourth of the distance between the fringes.¹¹ To avoid illuminating back the laser, the light from this laser should be linearly polarized. Then, a $\lambda/4$ phase plate is inserted in front of the beam, with its slow axis set at 45° to the interferometer plane to transform it into a circularly polarized beam.

Another method used to produce the two signals in quadrature phase is to take the signals from the two interference patterns that are produced in the interferometer. If the beam splitters are dielectric (no energy losses), the interference patterns will be complements of each other and, thus, the signals will be 90° apart. By introducing appropriate phase shifts in the beam splitter using metal coatings, the phase difference between the two patterns may be made 90° , as desired.²¹ This method was used by Rowley,²² as illustrated in Fig. 5. In order to separate the two patterns from the incident light beam, a large beam splitter and large retroreflectors are used. This configuration has the advantage that the laser beam is not reflected back. This is called a nonreacting interferometer configuration.

One more method, illustrated in Fig. 6, is the nonreacting interferometer designed at the Perkin-Elmer Corporation by Minkowitz and Vanir.²³ A circularly polarized light beam, produced by a linearly polarized laser and a $\lambda/4$ phase plate, illuminates the interferometer. This beam is divided by a beam splitter into two beams going to both arms of the interferometer. Upon reflection by the retro reflector, one of the beams changes its state of polarization from right to left circularly polarized. The two beams with opposite circular polarization are recombined at the beam splitter, thus producing linearly polarized light. The angle of the plane of polarization is determined by the phase difference between the two beams. The plane of polarization rotates 360° if the optical path difference

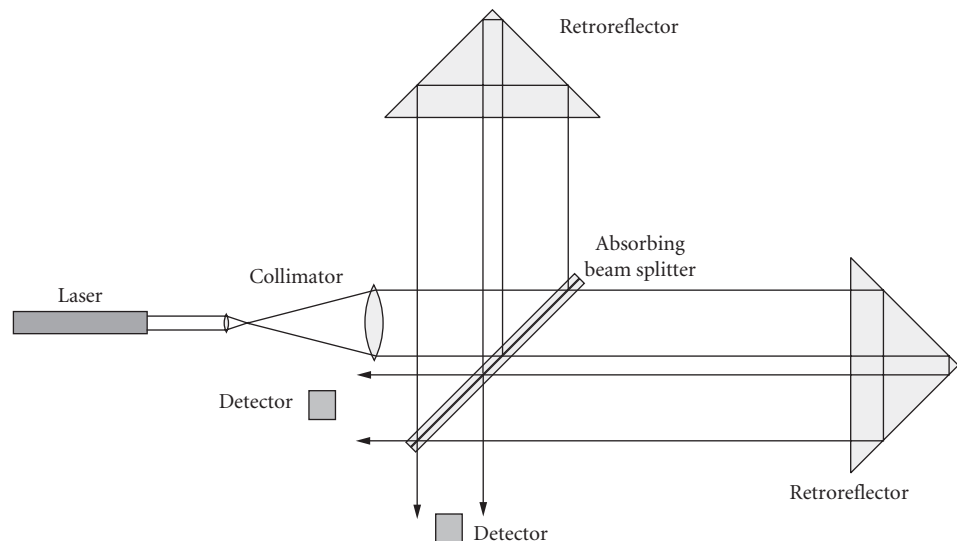


FIGURE 5 Two-interference pattern distance-measuring interferometer.

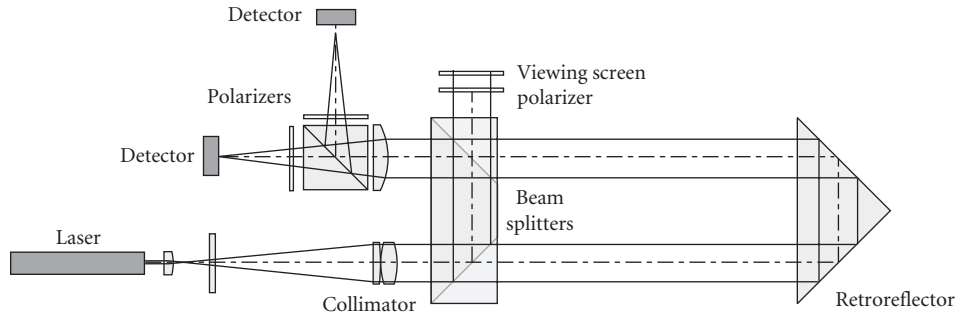


FIGURE 6 Minkowitz distance-measuring interferometer.

is changed by $\lambda/2$, the direction of rotation being given by the direction of the displacement. This linearly polarized beam is divided into two beams by a beam splitter. On each of the exiting beams, a linear polarizer is placed, one at an angle of $+45^\circ$ and the other at an angle of -45° . Then the two beams are in quadrature to each other.

In still another method, shown in Fig. 7, a beam of light, linearly polarized at 45° (or circularly polarized), is divided at a beam splitter, the p and s components. Then, both beams are converted to circular polarization with a $\lambda/4$ phase plate in front of each of them, with their axis at 45° . Upon reflection on the retro-reflectors, the handedness of the polarization is reversed. Thus, the linearly polarized beams exiting from the phase plates on the return to the beam splitter will have a plane of polarization orthogonal to that of the incoming beams. It is easily seen that no light returns to the laser. Here, the nonreacting configuration is not necessary but it may be used for additional protection. After recombination on the beam splitter, two orthogonal polarizations are present. Each plane of polarization contains information about the phase from one arm only so that no interference between the two beams has occurred. The two desired signals in quadrature are then generated by producing two identical beams with a nonpolarizing beam splitter with a polarizer on each exiting beam with their axes at $+45^\circ$ and -45° with respect to the vertical plane. The desired phase difference between the two beams is obtained by introducing a

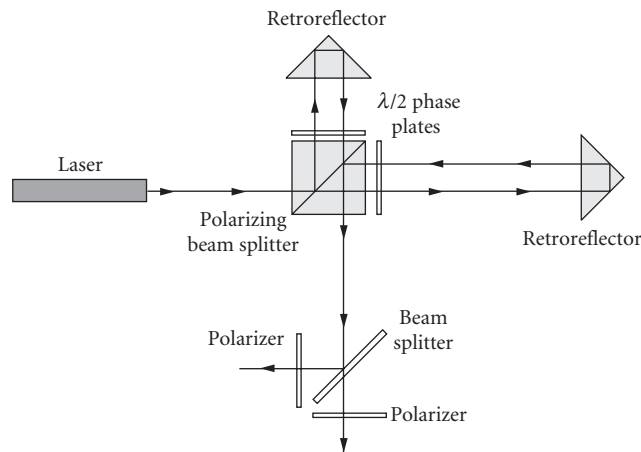


FIGURE 7 Brunning distance-measuring interferometer.

$\lambda/4$ phase plate after the beam splitter but before one of the polarizers with its slow axis vertical or horizontal. These two polarizers may be slightly rotated to make the two irradiances equal, at the same time preserving the 90° angle between their axes. If the prism is shifted a distance x , the fringe count is

$$\Delta_{\text{count}} = \pm \left[\frac{2x}{\lambda} \right] \quad (9)$$

where $[]$ denotes the integer part of the argument.

These interferometers are problematic in that any change in the irradiance may be easily interpreted as a fringe monitoring the light source. A more serious problem is the requirement that the static interference pattern be free of, or with very few, fringes. Fringes may appear because of multiple reflections or turbulence.

A completely different method uses two frequencies. It was developed by the Hewlett Packard Co.^{24,25} and is illustrated in Fig. 8. The light source is a frequency-stabilized He-Ne laser whose light beam is Zeeman split into two frequencies f_1 and f_2 by application of an axial magnetic field. The frequency difference is several megahertz and both beams have circular polarization, but with opposite sense. A $\lambda/4$ phase plate transforms the signals f_1 and f_2 into two orthogonal linearly polarized beams, one in the horizontal and the other in the vertical plane. A sample of this mixed signal is deviated by a beam splitter and detected at photo detector A, by using a polarizer at 45° . The amplitude modulation of this signal, with frequency $f_1 - f_2$, is detected and passed to a counter. Then, the two orthogonally polarized beams with frequencies f_1 and f_2 are separated at a polarizing beam splitter. Each is transformed into a circularly polarized beam by means of $\lambda/4$ phase plates. After reflection by the prisms, the handedness of these polarizations is changed. Then they go through the same phase plates where they are converted again to orthogonal linearly polarized beams. There is no light reflecting back to the laser. After recombination at the beam splitter, a polaroid at 45° will take the components of both beams in this plane. This signal is detected at the photo detector B. As with the other signal, the modulation with frequency $f_1 - f_2 + \Delta f$ is extracted from the carrier and sent to

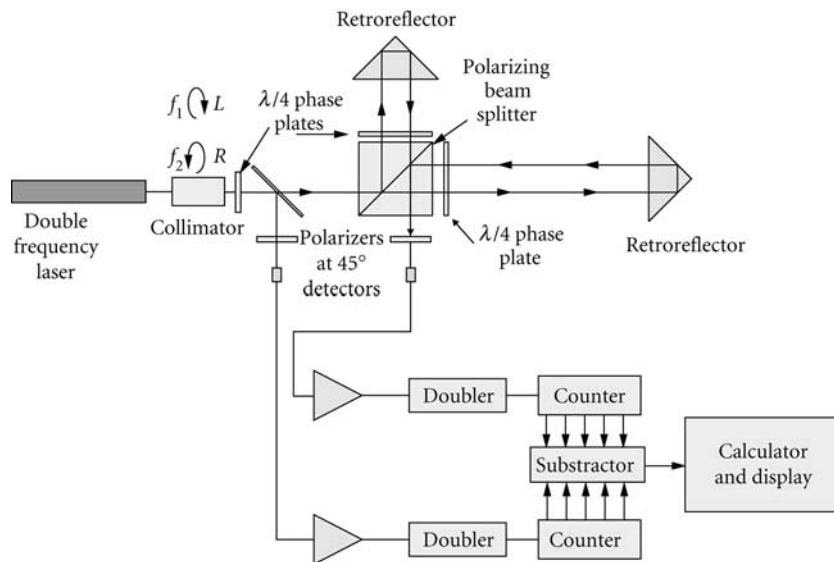


FIGURE 8 Hewlett Packard double-frequency distance-measuring interferometer.

another counter. The shift Δf in the frequency of this signal comes from a Doppler shift due to the movement of one of the retroreflectors,

$$\Delta f = \frac{2}{\lambda_2} \frac{dx}{dt} \quad (10)$$

where dx/dt is the cube corner prism velocity and λ_1 is the wavelength corresponding to the frequency f_1 . The difference between the results of the two counters is produced by the displacement of the retroreflector. If the prism moves a distance x , the number of pulses detected is given by

$$\Delta_{\text{count}} = \pm \left[\frac{2x}{\lambda_2} \right] \quad (11)$$

The advantage of this method compared with the first is that fringe counting is not subject to drift. These signals may be processed to obtain a better signal-to-noise ratio and higher resolution.²⁵

Straightness Measurements

Light propagation is assumed to be rectilinear in a homogeneous medium. This permits the use of a propagating light beam as a straightness reference. Besides the homogeneous medium, it is necessary to get a truly narrow pencil of light to improve accuracy. Laser light is an obvious application because of the high degree of spatial coherence. Beam divergence is usually less than 1 mrad for a He-Ne laser. One method uses a position-sensing detector to measure the centroid of the light spot despite its shape. In front of the laser, McLeod²⁶ used an axicon as an aligning device. When a collimated light beam is incident on an axicon, it produces a bright spot on a circular field. An axicon can give as much as 0.01 arcsec.

Another method to measure the deviation from an ideal reference line uses an autocollimator. A light beam leaving an autocollimator is reflected by a mirror. The surface slope is measured at the mirror. By knowing the distance to the mirror, one can determine the surface's profile by integration, as in a curvature measurement (see "Optical Methods for Measuring Curvature," later in this chapter). A method can be designed for measuring flatness for tables on lathe beds.²⁷

12.4 ANGLE MEASUREMENTS

Angle measurements, as well as distance measurements, require different levels of accuracy. For cutting glass, the required accuracy can be as high as several degrees, while for test plates, an error of less than a second of arc may be required. For each case, different measurement methods are developed.

Mechanical Methods

The easiest way to measure angles with medium accuracy is by means of mechanical nonoptical methods. These are

Sine Plate Essentially it is a table with one end higher than the other by a fixed amount, as shown in Fig. 9. Accuracy close to 30 arcmin may be obtained.

Goniometer This is a precision spectrometer. It has a fixed collimator and a moving telescope pointing to the center of a divided circle. Accuracies close to 20 arcsec may be obtained.

Bevel Gauge Another nonoptical method is by the use of a bevel gauge. This is made of two straight bars hinged at their edges by a pivot, as shown in Fig. 10. This device may be used to measure angle

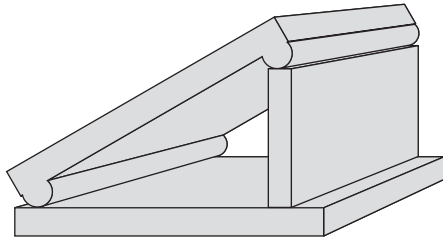


FIGURE 9 Sine plate.

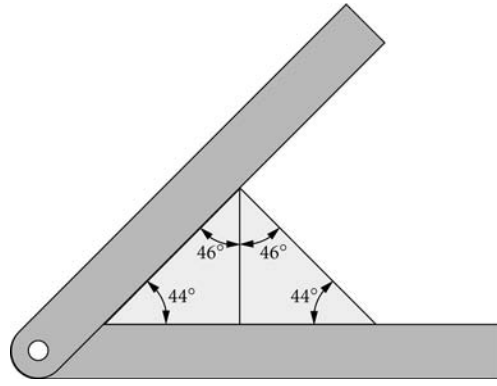


FIGURE 10 Bevel gauge.

prisms whose angle accuracies are from 45 to about 20 arcsec.²⁸ For example, if the measured prism has a 50-mm hypotenuse, a space of 5 μm at one end represents an angle of 0.0001 rad or 20 arcsec.

Numerically Controlled Machines (CNC) The advent of digitally controlled machines has brought to the optical shop machines to work prisms and angles with a high level of accuracy. Most machines can be adjusted within 0.5 arcmin of error.

Autocollimators

As shown in Fig. 11, an autocollimator is essentially a telescope focused at infinity with an illuminated reticle located at the focal plane. A complete description of autocollimators is found in Hume.²⁹ A flat reflecting surface, perpendicular to the exiting light beam, forms an image of the reticle on the same plane as the original reticle. Then, both the reticle and its image are observed through the eyepiece. When the reflecting surface is not exactly perpendicular to the exiting light beam, the reticle image is laterally displaced in the focal plane with respect to the object reticle. The magnitude of this displacement d is

$$d = 2\alpha f \quad (12)$$

where α is the tilt angle for the mirror in radians and f is the focal length of the telescope.

Autocollimator objective lenses are usually corrected doublets. Sometimes a negative lens is included to form a telephoto lens to increase the effective focal length while maintaining compactness. The collimating lens adjustment is critical for the final accuracy. Talbot interferometry can be used for a precise focus adjustment.³⁰

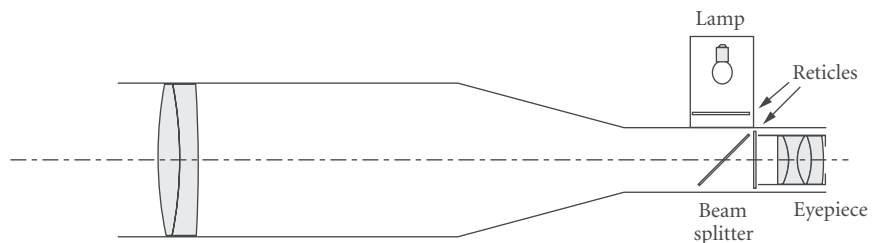


FIGURE 11 An autocollimator.

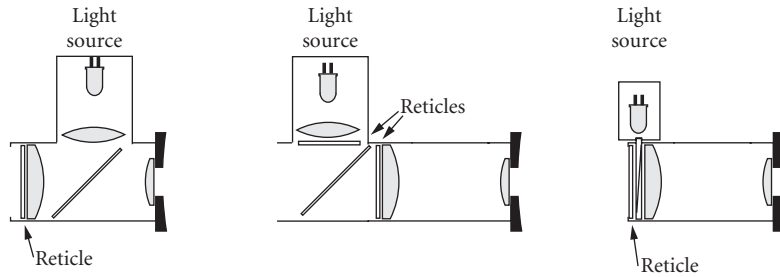


FIGURE 12 Illuminated eyepieces for autocollimators and microscopes. (a) Gauss; (b) bright line; and (c) Abbe.

The focal plane is observed through an eyepiece. Several types of illuminated reticles and eyepieces have been developed. Figure 12³¹ illustrates some illuminated eyepieces, in all of which the reticle is calibrated to measure the displacement. Gauss and Abbe illuminators show a dark reticle on a bright field. A bright field³² may be more appropriate for low reflectance surfaces. Rank³³ modified a Gauss eyepiece to produce a dark field. In other systems, a drum micrometer displaces a reticle to position it at the image plane of the first reticle. To increase sensitivity some systems, called microoptic autocollimators, use a microscope to observe the image.

Direct-reading autocollimators have a field of view of about 1° . Precision in an autocollimator is limited by the method for measuring the centroid of the image. In a diffraction-limited visual system, the diffraction image size sets the limit of precision. In a precision electronic measuring system, the accuracy of the centroid measurement is limited by the electronic detector, independent of the diffraction image itself, and can exceed the diffraction limit. In some photoelectric systems, the precision is improved by more than an order of magnitude.

Autocollimators are used for angle measurements in prisms and glass polygons. But they also have other applications; for example, to evaluate the parallelism between faces in optical flats or to manufacture divided circles.³⁴ By integrating measured slope values with an autocollimator, flatness deviations for a machine tool for an optical bed can also be evaluated.²⁷

The reflecting surface in autocollimation measurements must be kept close to the objective in order to make the alignment easier and to be sure that all of the reflected beam enters the system. The reflecting surface must be of high quality. A curved surface is equivalent to introducing another lens in the system with a change in the effective focal length.²⁷

Several accessories for autocollimators have been designed. For single-axis angle measurement, a pentaprism is used. An optical square permits angle measurements for surfaces at right angles. Perpendicularity is measured with a pentaprism and a mirror, as shown in Fig. 13. A handy horizontal reference can be produced with an oil pool, but great care must be taken with the surface stability.

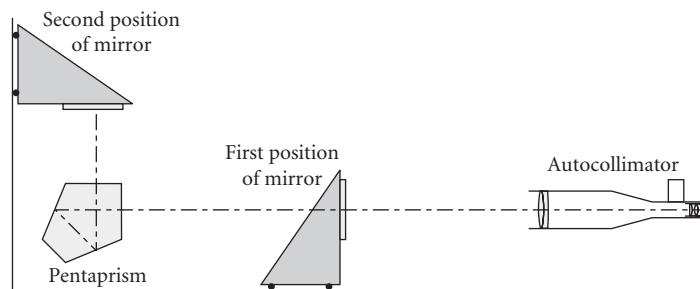


FIGURE 13 Perpendicularity measurement with an autocollimator.

Theodolites

Theodolites are surveying instruments to measure vertical and horizontal angles. A telescope with reticle is the reference to direct the instrument axis. In some theodolites, the telescope has an inverting optical system to produce an erect image. The reticle is composed of a crosswire and has a couple of parallel horizontal lines, called stadia lines, used for range measurements (see “Stadia and Range Finders” earlier in this chapter). The telescope has two focus adjustments: one to sharply focus the reticle according to the personal setting for each observer, and the other to focus the objective on the reticle. This later focus adjustment is performed by moving the whole eyepiece.

The theodolite telescope has a tree-screw base altitude-azimuth mounting on a base made horizontal with a spirit level. Divided circles attached to both telescope axes measure vertical and horizontal angles. In older instruments, an accuracy of 20 arcmin was standard. Modern instruments are accurate to within 20 arcsec, the most expensive of which can reach an accuracy of 1 arcsec. To remove errors derived from eccentric scales as well as orthogonality errors, both axes are rotated 180° and the measurement repeated. Older instruments had a provision for reading at opposite points of the scale. Scales for theodolites can be graduated in sexagesimal degrees or may use a centesimal system that divides a circle into 400 grades.

Some of the accessories available for theodolites include

1. An illuminated reticle that can be used as an autocollimator when directed to a remote retroreflector. The observer adjusts the angles until both the reflected and the instrument's reticle are superposed. This increases the pointing accuracy.
2. A solar filter, which can be attached to the eyepiece or objective side of the telescope. This is used mainly for geographic determination in surveying.
3. An electronic range meter, which is superposed to the theodolite to measure the distance. Additionally, some instruments have electronic position encoders that allow a computer to be used as an immediate data-gathering and reducing device.
4. A transverse micrometer for measuring angular separation in the image plane.

Accuracy in a theodolite depends on several factors in its construction. Several of these errors can be removed by a careful measuring routine. Some of the systematic or accuracy limiting errors are

1. Perpendicularity—deviation between vertical and horizontal scales. This error can be nulled by plunging and rotating the telescope, then averaging.
2. Concentricity deviation of scales. When scales are not concentric, they are read at opposite ends to reduce this error. Further accuracy can be obtained by rotating the instrument 90° , 180° , and 270° and averaging measurements.

Level

Levels are surveying instruments for measuring the deviation from a horizontal reference line. A level is a telescope with an attached spirit level. The angle between the telescope axis and the one defined by the spirit level must be minimized. It can be adjusted by a series of measurements to a pair of surveying staves (Fig. 14). Once the bubble is centered in the tube, two measurements are taken on

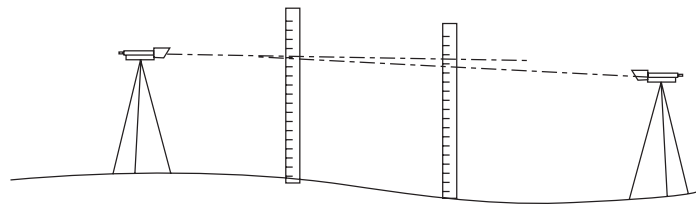


FIGURE 14 Level adjustment. (After Kingslake.³⁵)

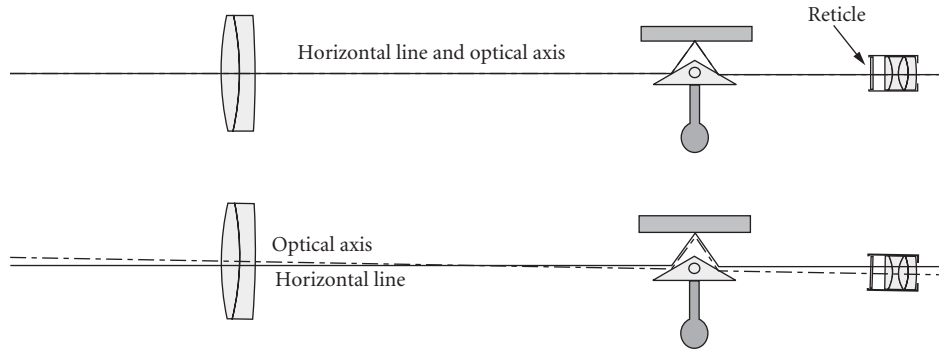


FIGURE 15 The autoset level.

each side of the staves.³⁵ The level differences between the two staves must be equal if the telescope axis is parallel with the level horizontal axis.

The autoset level (Fig. 15) uses a suspended prism and a fixed mirror on the telescope tube. The moving prism maintains the line aimed at the horizon and passing through the center of the reticle, despite the tube orientation, as long as it is within about 15 arcmin. Typical precision for this automatic level can go up to 1 arcsec.²⁷

Interferometric Measurements

Interferometric methods find their main applications in measuring very small wedge angles in glass slabs^{36,37} and in parallelism evaluation³⁸ by means of the Fizeau or Haidinger interferometers.³⁹

Interferometric measurements of large angles may also be performed. In one method, a collimated laser beam is reflected from the surfaces by a rotating glass slab. The resulting fringes can be considered as coming from a Murty lateral shear interferometer.⁴⁰ This device can be used as a secondary standard to produce angles from 0° to 360° with accuracy within a second of arc. Further analysis of this method has been done by Tentori and Celaya.⁴¹ In another system, a Michelson interferometer is used with an electronic counter to measure over a range of $\pm 5^\circ$ with a resolution of 10° .^{42,43} An interferometric optical sine bar for angles in the milliseconds of arc was built by Chapman.⁴⁴

Angle Measurements in Prisms

A problem frequently encountered in the manufacture of prisms is the precise measurement of angle. In most cases, prism angles are 90° , 45° , and 30° . These angles are easily measured by comparison with a standard but it is not always necessary.

An important aspect of measuring angles in a prism is to determine if the prism is free of pyramidal error. Consider a prism with angles A, B, and C (Fig. 16a). Let OA be perpendicular to plane ABC. If line AP is perpendicular to segment BC, then the angle AOP is a measurement of the pyramidal error. In a prism with pyramidal error, the angles between the faces, as measured in planes perpendicular to the edges between these faces, add up to over 180° . To simply detect pyramidal error in a prism, Johnson⁴⁵ and Martin⁴⁶ suggest that both the refracted and the reflected images from a straight line be examined (Fig. 16b). When pyramidal error is present, the line appears to be broken. A remote target could be graduated to measure directly in minutes of arc. A sensitivity of up to 3 arcmin may be obtained.

During the milling process in the production of a prism, a glass blank is mounted in a jig collinear with a master prism (Fig. 17). An autocollimator aimed at the master prism accurately sets the

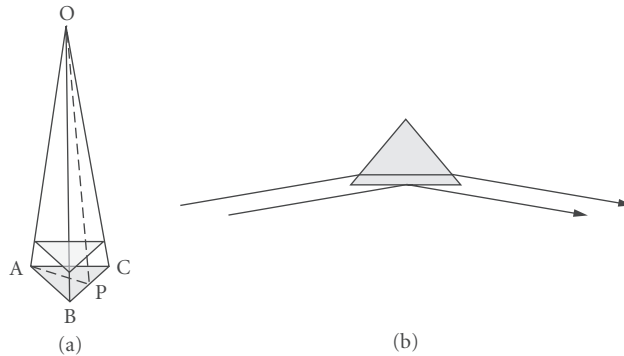


FIGURE 16 Pyramidal error in a prism (a) nature of the error and (b) test of the error.

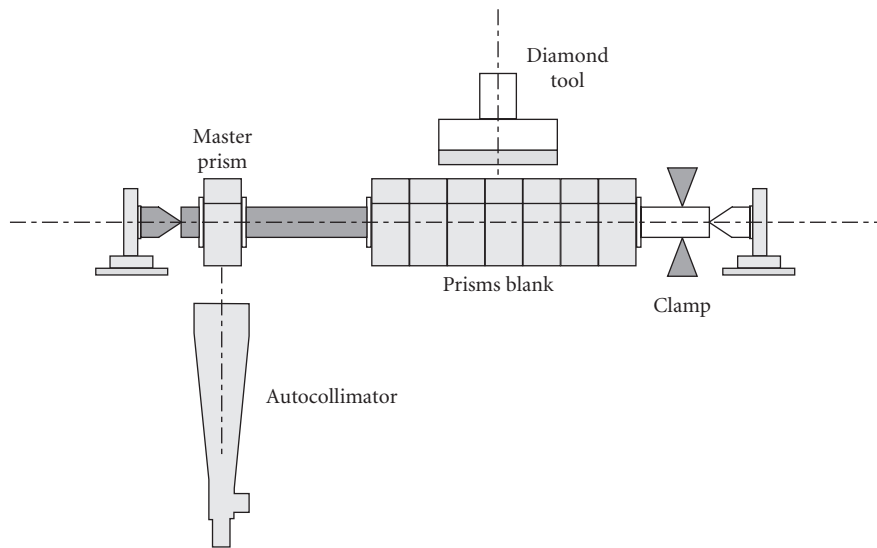


FIGURE 17 Milling prisms for replication.

position for each prism face.^{47,48} With a carefully set diamond lap, pyramidal error is minimized. In a short run, angles can be checked with a bevel gauge. Visual tests for a prism in a bevel gauge can measure an error smaller than a minute of arc.³¹

A 90° angle in a prism can be measured by internal reflection, as shown in Fig. 18a. At the autocollimator image plane, two images are seen with an angular separation of $2N\alpha$, where α is the magnitude of the prism angle error, and its sign is unknown. Since the hypotenuse face has to be polished and the glass must be homogeneous, the measurement of the external angle with respect to a reference flat is preferred (Fig. 18b). In this case, the sign of the angle error is determined by a change in the angle by tilting the prism. If the external angle is decreased and the images separate further, then the external angle is less than 90° . Conversely, if the images separate by tilting in such way that the external angle increases, then the external angle is larger than 90° .

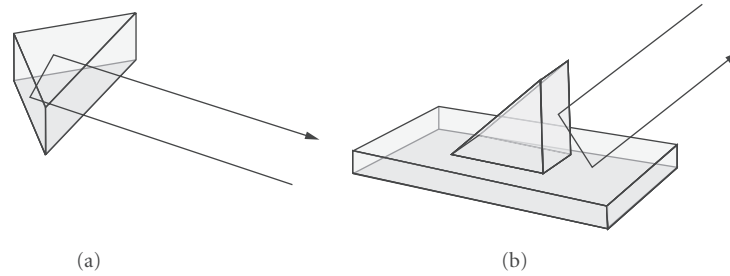


FIGURE 18 Right angle measurement in prisms: (a) internal measurement and (b) external measurement.

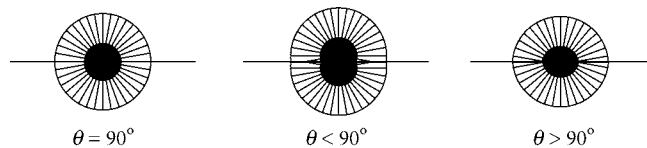


FIGURE 19 Retroreflected images of the observer's pupil in a 90° prism.

To determine the sign of the error, several other methods have been proposed. DeVany⁴⁹ suggested that when looking at the double image from the autocollimator, the image should be defocused inward. If the images tend to separate, then the angle in the prism is greater than 90° . Conversely, an outward defocusing will move the images closer to each other for an angle greater than 90° . Another way to eliminate the sign of the error in the angle is by introducing, between the autocollimator and the prism, a glass plate with a small wedge whose orientation is known. The wedge should cover only one-half of the prism aperture. Ratajczyk and Bodner⁵⁰ suggested a different method using polarized light.

Right-angle prisms can be measured using an autocollimator with acceptable precision.⁵¹ With some practice, perfect cubes with angles more accurate than 2 arcsec can be obtained.⁴⁹ An extremely simple test for the 90° angle in prisms⁴⁵ is performed by looking to the retroreflected image of the observer's pupil without any instrument. The shape of the image of the pupil determines the error, as shown in Fig. 19. The sensitivity of this test is not very great and may be used only as a coarse qualitative test. As shown by Malacara and Flores,⁵² a small improvement in the sensitivity of this test may be obtained if a screen with a small hole is placed in front of the eye, as in Fig. 20a. A cross centered on the small hole is painted on the front face of the screen. The observed images are as shown in the same Fig. 20b. As opposed to the collimator test, there is no uncertainty in the sign of the error in the tests just described, since the observed plane is located where the two prism surfaces intersect. An improvement described by Malacara and Flores,⁵² combining these simple tests with an autocollimator, is obtained with the instrument in Fig. 21. In this system, the line defining the intersection between the two surfaces is out of focus and barely visible while the reticle is in perfect focus at the eyepiece.

Corner cube prisms are a real challenge to manufacture, since besides the large precision required in the angles, all surfaces should be exempt of any curvature. The dihedral angle in pentaprisms is tested usually with an interferometer. An error in the prism alignment results in an error in angle determination.⁵³ A simple geometric method for angle measurement in corner cube reflector has been described by Rao.⁵⁴ Also, the calculations for the electric field in a corner cube are performed by Schöll.⁵⁵ These calculations include the effect by nonhomogeneities, angle of incidence, and errors in the surface finish.

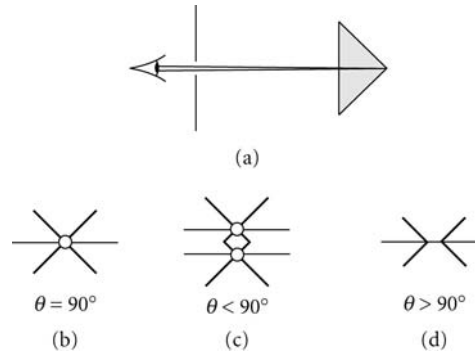


FIGURE 20 Testing a right-angle prism: (a) screen in front of the eye and (b) to (d) its observed images.

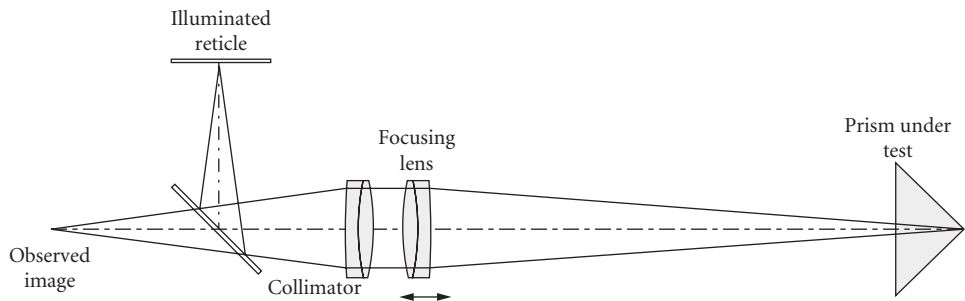


FIGURE 21 Modified autocollimator for testing the right angle in prisms without sign uncertainty in measured error.

12.5 CURVATURE AND FOCAL LENGTH MEASUREMENTS

The curvature of a spherical optical surface or the local curvature of an aspherical surface may be measured by means of mechanical or optical methods. Some methods measure the sagitta, some the surface slope, and some others directly locate the position of the center of curvature.

Mechanical Methods for Measuring Curvature

Templates The simplest and most common way to measure the radius of curvature is by comparing it with metal templates with known radii of curvature until a good fit is obtained. The template is held in contact with the optical surface with a bright light source behind the template and the optical surface. If the surface is polished, gaps between the template and the surface may be detected to an accuracy of one wavelength. If the opening is very narrow, the light passing through the gap becomes blue due to diffraction.

Test Plates This method uses a glass test plate with a curvature opposite to that of the glass surface to be measured. The accuracy is much higher than in the template method, but the surface must be polished.

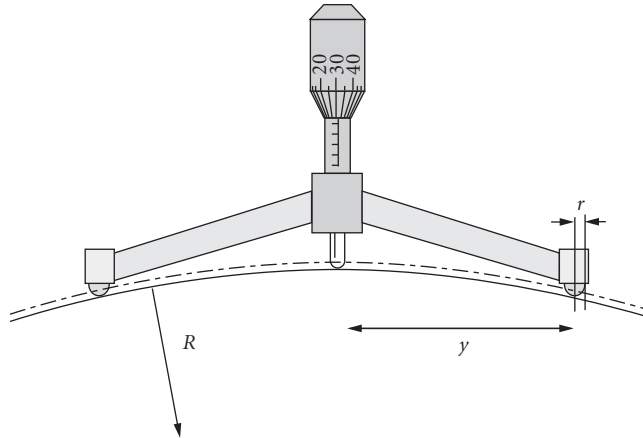


FIGURE 22 Three-leg spherometer.

Spherometers This is probably the most popular mechanical device used to measure radii of curvature. It consists of three equally spaced feet with a central moving plunger. The value of the radius of curvature is calculated after measuring the sagitta, as shown in Fig. 22. The spherometer must first be calibrated by placing it on top of a flat surface. Then it is placed on the surface to be measured. The difference in the position of the central plunger for these two measurements is the sagitta of the spherical surface being measured. Frequently, a steel ball is placed at the end of the legs as well as at the end of the plunger to avoid the possibility of scratching the surface with sharp points. In this case, if the measured sagitta is z , the radius of curvature R of the surface is given by

$$R = \frac{z}{2} + \frac{y^2}{2z} \pm r \quad (13)$$

where r is the radius of curvature of the balls. The plus sign is used for concave surfaces and the minus sign for convex surfaces. The precision of this instrument in the measurement of the radius of curvature for a given uncertainty in the measured sagitta may be obtained by differentiating Eq. (13)

$$\frac{dR}{dz} = \frac{1}{2} - \frac{y^2}{2z^2} \quad (14)$$

obtaining

$$\Delta R = \frac{\Delta z}{2} \left(1 - \frac{y^2}{z^2} \right) \quad (15)$$

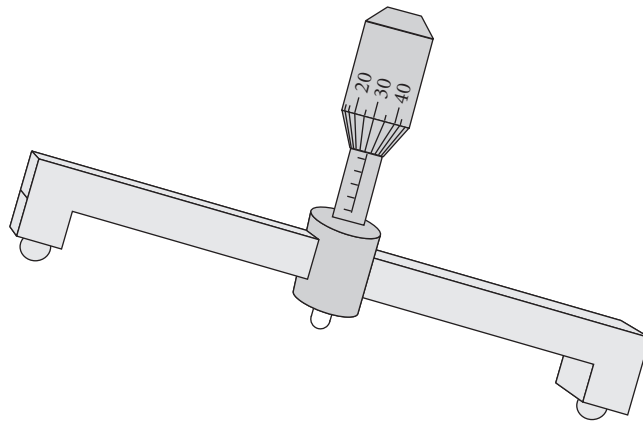
This accuracy assumes that the spherometer is perfectly built and that its dimensional parameters y and r are well known. The uncertainty comes only from human or instrumental errors in the measurement of the sagitta. Noble³¹ has evaluated the repeatability for a spherometer with $y = 50$ mm and an uncertainty in the sagitta reading equal to $5 \mu\text{m}$, and has reported the results in Table 1 where it can be seen that the precision is better than 2 percent. An extensive analysis of the precision and accuracy of several types of spherometers is given by Jurek.⁵⁶

A ring may be used instead of the three legs in a mechanical spherometer. A concave surface contacts the external edge of the cup, and a convex surface is contacted by the internal edge of the ring. Thus, Eq. (5) may be used if a different value of y is used for concave and convex surfaces, and r is taken as zero. Frequently in spherometers of this type, the cups are interchangeable in order to have different diameters for different surface diameters and radii of curvature. The main advantage of the

TABLE 1 Spherometer precision*

Radius of Sphere R (mm)	Sagitta Z (mm)	Fractional Precision ΔR (mm)	Precision $\Delta R/R$
10,000	0.125	-400	-0.040
5,000	0.250	-100	-0.020
2,000	0.625	-16	-0.008
1,000	1.251	-4	-0.004
500	2.506	-1	-0.002
200	6.351	-0.15	-0.0008

* $y = 50$ mm; $\Delta z = 5$ μm .
Source: From Noble.³¹

**FIGURE 23** Bar spherometer.

use of a ring instead of three legs is that an astigmatic deformation of the surface is easily detected, although it cannot be measured.

A spherometer that permits the evaluation of astigmatism is the bar spherometer, shown in Fig. 23. It can measure the curvature along any diameter. A commercial version of a small bar spherometer for the specific application in optometric work is the Geneva gauge, where the scale is directly calibrated in diopters assuming that the refractive index of the glass is 1.53.

Automatic spherometers use a differential transformer as a transducer to measure the plunger displacement. This transformer is coupled to an electronic circuit and produces a voltage that is linear with respect to the plunger displacement. This voltage is fed to a microprocessor which calculates the radius of curvature or power in any desired units and displays it.

Optical Methods for Measuring Curvature

Foucault Test Probably the oldest and easiest method to measure the radius of curvature of a concave surface is the knife-edge test. In this method, the center of curvature is first located by means of the knife edge. Then, the distance from the center of curvature to the optical surface is measured.

Autocollimator The radius of curvature may also be determined through measurements of the slopes of the optical surface with an autocollimator as described by Horne.⁵⁷ A pentaprism producing a 90° deflection of a light beam independent of small errors in its orientation is used in this technique,

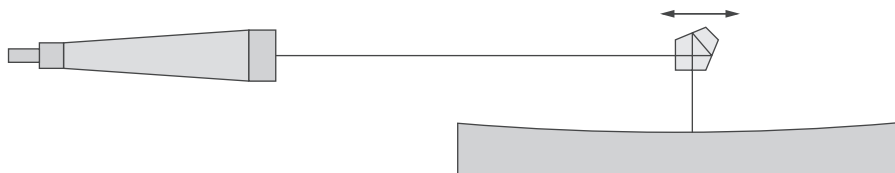


FIGURE 24 Autocollimator and pentaprism used to determine radius of curvature by measuring surface slopes.

as illustrated in Fig. 24, where the pentaprism travels over the optical surface to be measured along one diameter. First the light on the reticle of the autocollimator is centered on the vertex of the surface being examined. Then the pentaprism is moved toward the edge of the surface in order to measure any slope variations. From these slope measurements, the radius of curvature may be calculated. This method is useful only for large radii of curvature for either concave or convex surfaces.

Confocal Cavity Technique Gerchman and Hunter^{58,59} have described the so-called optical cavity technique that permits the interferometric measurement of very long radii of curvature with an accuracy of 0.1 percent. The cavity of a Fizeau interferometer is formed, as illustrated in Fig. 25. This is a confocal cavity of n th order, where n is the number of times the path is folded. The radius of curvature is equal to approximately $2n$ times the cavity length Z .

Traveling Microscope This instrument is used to measure the radius of curvature of small concave optical surfaces with short radius of curvature. As illustrated in Fig. 26, a point light source is produced at the front focus of a microscope objective. This light source illuminates the concave optical surface to be measured near its center of curvature. Then this concave surface forms an image which is also close to its center of curvature. This image is observed with the same microscope used to illuminate the surface. During this procedure, the microscope is focused both at the center of curvature and at the surface to be measured. A sharp image of the light source is observed at both places. The radius of curvature is the distance between these two positions for the microscope.

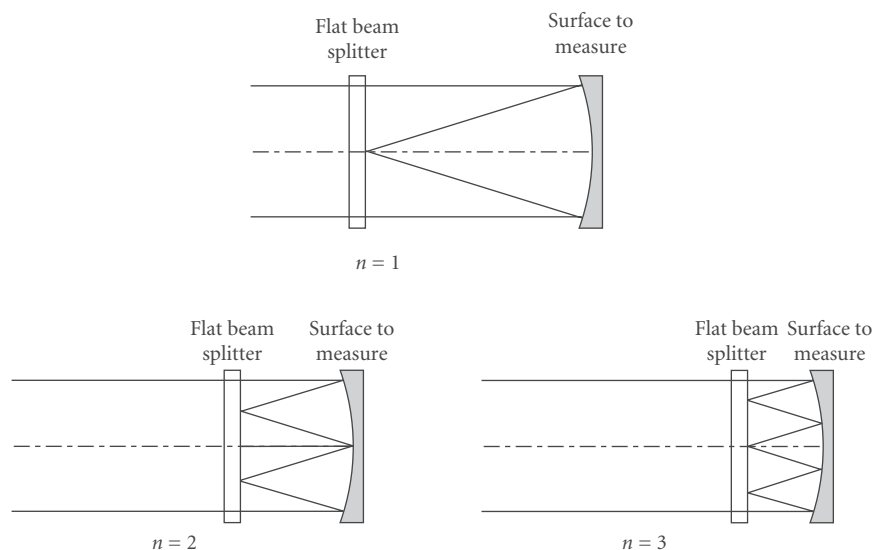


FIGURE 25 Confocal cavity arrangements used to measure radius of curvature.

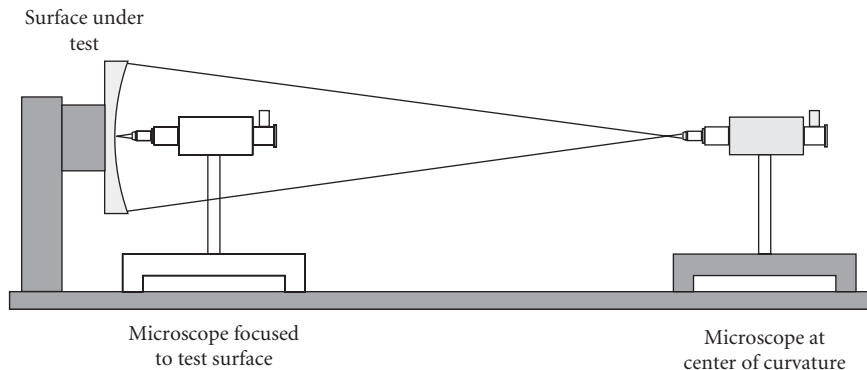


FIGURE 26 Traveling microscope to measure radii of curvature.

This distance traveled by the microscope may be measured on a vernier scale, obtaining a precision of about 0.1 mm. If a bar micrometer is used, the precision may be increased by an order of magnitude. In this case, two small convex buttons are required: one fixed to the microscope carriage and the other to the stationary part of the bench. They must face each other when the microscope carriage is close to the optical bench fixed component.

Carnell and Welford³² describe a method that requires only one measurement. The microscope is focused only at the center of curvature. Then the radius of curvature is measured by inserting a bar micrometer with one end touching the vertex of the optical surface. The other end is adjusted until it is observed to be in focus on the microscope. Accuracies of a few microns are obtained with this method.

In order to focus the microscope properly, the image of an illuminated reticle must fall, after reflection, back on itself, as in the Gauss eyepiece shown in Fig. 12. The reticle and its image appear as dark lines in a bright field. The focusing accuracy may be increased with a dark field. Carnell and Welford obtained a dark field with two reticles, as in Fig. 12, one illuminated with bright lines and the other with dark lines.

A convex surface may also be measured with this method if a well-corrected lens with a conjugate longer than the radius of curvature of the surface under test is used. Another alternative for measuring convex surfaces is by inserting an optical device with prisms in front of the microscope, as described by Jurek.⁵⁶

Some practical aspects of the traveling microscope are examined by Rank,³³ who obtained a dark field at focus with an Abbe eyepiece which introduces the illumination with a small prism. This method has been implemented using a laser light source by O'Shea and Tilstra.⁶⁰

Additional optical methods to measure the radius of curvature of a spherical surface have been described. Evans^{61,62,63} determines the radius by measuring the lateral displacements on a screen of a laser beam reflected on the optical surface when this optical surface is laterally displaced. Cornejo-Rodriguez and Cordero-Dávila,⁶⁴ Klingsporn,⁶⁵ and Diaz-Urbe et al.⁶⁶ rotate the surface about its center of curvature on a nodal bench.

Focal Length Measurements

There are two focal lengths in an optical system: the back focal length and the effective focal length. The back focal length is the distance from the last surface of the system to the focus. The effective focal length is the distance from the principal plane to the focus. The back focal length is easily measured, following the same procedure used for measuring the radius of curvature, using a microscope and the lens bench. On the other hand, the effective focal length requires the previous location of the principal plane.

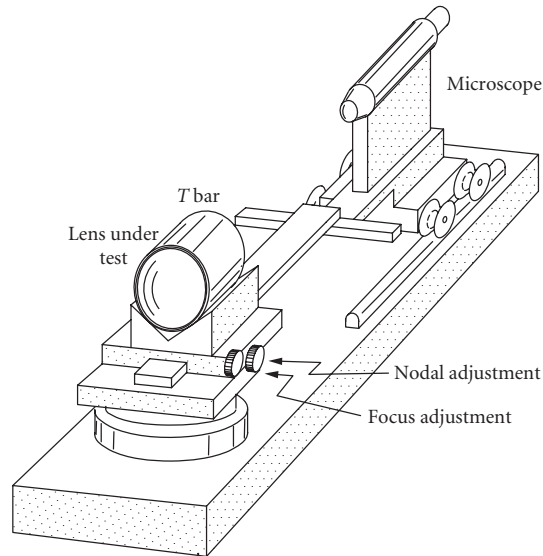


FIGURE 27 Nodal slide bench. (From Malacara.¹)

Nodal Slide Bench In an optical system in air, the principal points (intersection of the principal plane and the optical axis) coincide with the nodal points. Thus, to locate this point we may use the well-known property that small rotations of the lens about an axis perpendicular to the optical axis and passing through the nodal point do not produce any lateral shift of the image. The instrument used to perform this procedure, shown in Fig. 27, is called an optical nodal slide bench.⁶⁷ This bench has a provision for slowly moving the lens under test longitudinally in order to find the nodal point.

The bench is illuminated with a collimated light source and the image produced by the lens under test is examined with a microscope. The lens is then displaced slightly about a vertical axis as it is being displaced longitudinally. This procedure is stopped until a point is found in which the image does not move laterally while rotating the lens. This axis of rotation is the nodal point. Then, the distance from the nodal point to the image is the effective focal length.

Focimeters A focimeter is an instrument designed to measure the focal length of lenses in a simple manner. The optical scheme for the classical version of this instrument is shown in Fig. 28. A light source illuminates a reticle and a convergent lens, with focal length f , displaced at a distance x from the reticle. The lens to be measured is placed at a distance d from the convergent lens. The magnitude

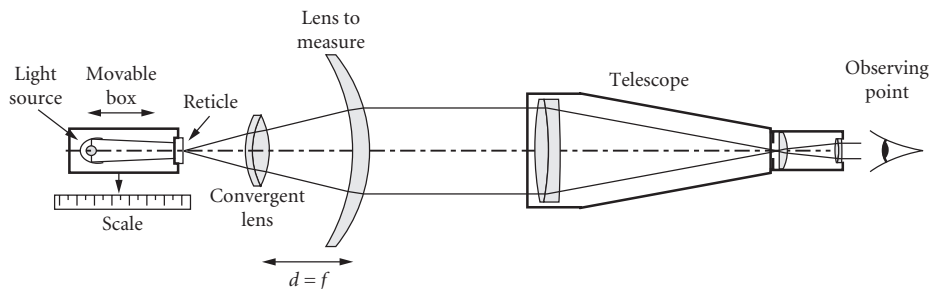


FIGURE 28 Focimeter schematics.

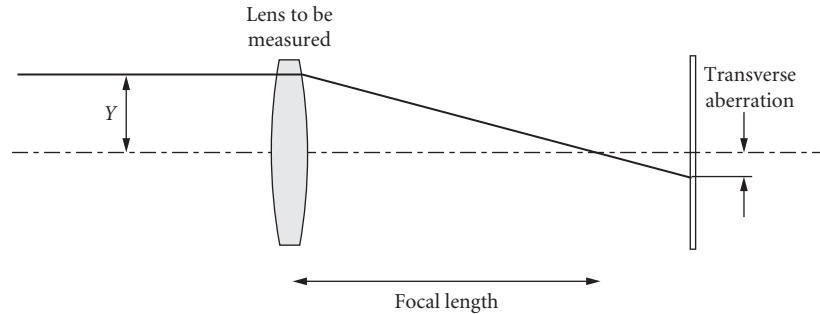


FIGURE 29 Focal length determination by transverse aberration measurements.

of x is variable and is adjusted until the light beam going out from the lens under test becomes collimated. This collimation is verified by means of a small telescope in front of this lens focused at infinity. The values of d and the focal length f are set to be equal. Then, the back focal length f_b of the lens under test is given by

$$\frac{1}{f_b} = \frac{1}{d} - \frac{x}{d^2} = P_v \quad (16)$$

where its inverse P_v is the vertex power. As can be seen, the power of the lens being measured is linear with respect to the distance x . There are many variations of this instrument. Some modern focimeters measure the lateral deviation of a light ray from the optical axis (transverse aberration), as in Fig. 29, when a defocus is introduced.^{61,62,63} This method is mainly used in some modern automatic focimeters for optometric applications. To measure the transverse aberration, a position-sensing detector is frequently used. The power error of a focimeter can be obtained by derivation of Eq. (16)

$$\delta P_v = -\frac{\delta x}{f_c^2} \quad (17)$$

Thus, the power error is a linear function of the target position error and decreases with the square of the collimating lens focal distance.⁷⁰

Other Focal Length Measurements

Moiré Deflectometry Moiré deflectometry method for focal length determination sends a collimated beam over a pair of Ronchi rulings (Fig. 30) depending on the convergence or divergence of the beam, the resulting moiré pattern rotates according to the convergence (Fig. 31). The rotation has an angle α that is related to the focal distance by^{71,72}

$$f \approx \frac{d}{\theta \tan \alpha} \quad (18)$$

d being the ruling's pitch, θ is the angle between the ruling's lines, and α is the rotation angle of the moiré pattern.

Talbot Autoimages Talbot autoimages method for focal length determination is performed by sending a coherent beam of light into a Ronchi ruling. An image of the grating will be produced periodically and evenly spaced along the light beam. Every Talbot autoimage is an object for the lens

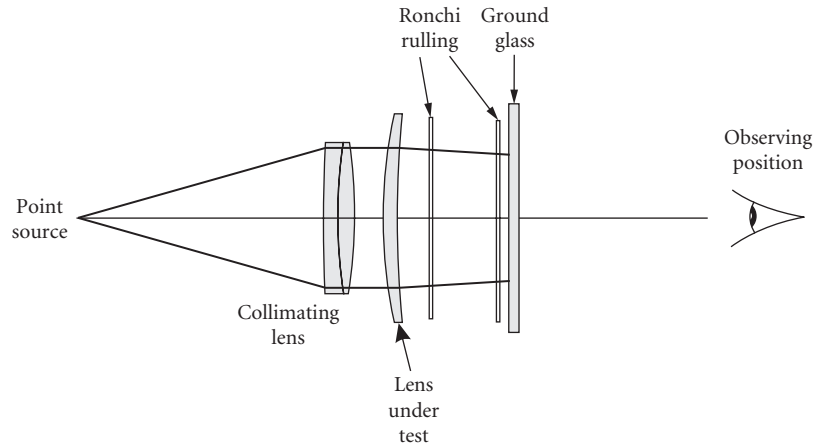


FIGURE 30 Moiré deflectometry lens power measurement. (From Malacara.¹)

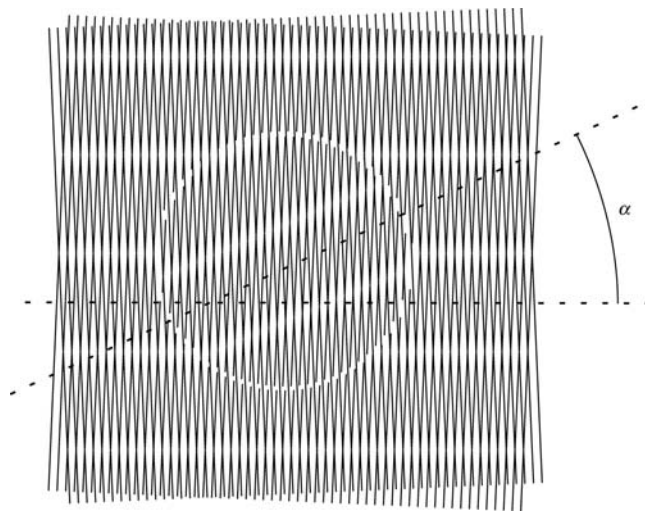


FIGURE 31 Moiré pattern as produced in a moiré deflectometer. (From Malacara.¹)

under test and produces a set of autoimages at the image side of the lens. To determine the lens focal length, another ruling coincident to the autoimages at the image plane is used.^{73,74}

Fourier Transforms The Fourier-transforming property of a lens can be used in the focal-length determination. Horner⁷⁵ measured the diffraction pattern at the focal plane produced by a slit. For this method, the light beam does not have to be perfectly collimated.

Microlenses Micro-lenses applications require new methods for small focal length determination. The propagated Gaussian beam of a laser can be analyzed.^{76,77} In this method, a lens is placed at the laser beam waist, then the propagating beam is measured to determine the focal length.

Fiber Optics A clever method used to automatically find the position of the focus has been described by Howland and Proll.⁷⁸ They used optical fibers to illuminate the lens in an autocollimating configuration, and the location of the image was also determined using optical fibers.

12.6 REFERENCES

1. D. Malacara, (ed.) *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
2. A. L. Bloom, "Gas Lasers and Their Application to Precise Length Measurements," in E. Wolf (ed.), *Progress in Optics*, vol. IX, North Holland, Amsterdam, 1971.
3. D. T. Goldman, "Proposed New Definition of the Meter," *J. Opt. Soc. Am.* **70**: 1640–1641 (1980).
4. P. Giacomo, "Metrology and Fundamental Constants." *Proc. Int. School of Phys. "Enrico Fermi," course 68*, North-Holland, Amsterdam, 1980.
5. S. Cundiff, J. Ye, and J. Hall, "Rulers of Light," *Sci. Am.* **298**(4): 52–59 (2008).
6. D. C. Baird, *Experimentation*, Prentice-Hall, New Jersey, 1962.
7. J. C. Gibbins, *The Systematic Experiment*, Cambridge Univ. Press, Cambridge, 1986.
8. F. B. Patrick, "Military Optical Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, chap. 7.
9. M. S. Dickson and D. Harkness, "Surveying and Tracking Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, chap. 8.
10. W. J. Smith, *Modern Optical Engineering*, 2d ed., McGraw-Hill, New York, 1990.
11. A. Sona, "Lasers in Metrology," in F. T. Arecchi and E. O. Schulz-Dubois (eds.), *Laser Handbook*, vol. 2, North Holland, Amsterdam, 1972.
12. R. Dändliker, R. Thalman, and D. Prongué, "Two-Wavelength Laser Interferometry Using Superheterodyne Detection," *Opt. Lett.* **13**: 339–342 (1988).
13. R. Dändliker, Y. Salvadé, and E. Zimmermann, "Distance Measurement by Multiple-Wavelength Interferometry," *J. Opt.* **29**: 105–114 (1998).
14. Y. Salvadé, N. Schuler, S. Lévêque, and S. Le Floch, "High-Accuracy Absolute Distance Measurement Using Frequency Comb Referenced Multiwavelength Source," *Appl. Opt.* **47**(14): 2715–2720 (2008).
15. K. Minoshima and H. Matsumoto, "High-Accuracy Measurement of 240-m Distance in an Optical Tunnel by Use of a Compact Femtosecond Laser," *Appl. Opt.* **39**: 5512–5517 (2000).
16. M. L. Stitch, "Laser Rangefinding," in F. T. Arecchi and E. O. Schulz-Dubois (eds.), *Laser Handbook*, vol. 2, North Holland, Amsterdam, 1972.
17. T. Tsuruta and Y. Ichihara, "Accurate Measurement of Lens Thickness by Using White-Light Fringes," *Jpn. J. Appl. Phys. Suppl.* **14-1**: 369–372 (1975).
18. J. Bruning, "Fringe Scanning," in D. Malacara (ed.), *Optical Shop Testing*, 1st ed., John Wiley and Sons, New York, 1978.
19. C. Steinmetz, R. Burgoon, and J. Herris, "Accuracy Analysis and Improvements for the Hewlett-Packard Laser Interferometer System," *Proc. SPIE* **816**: 79–94 (1987).
20. N. A. Massie, and J. Caulfield, "Absolute Distance Interferometry," *Proc. SPIE* **816**: 149–157 (1987).
21. E. R. Peck and S. W. Obetz, "Wavelength or Length Measurement by Reversible Fringe Counting," *J. Opt. Soc. Am.* **43**: 505–507 (1953).
22. W. R. C. Rowley, "Some Aspects of Fringe Counting in Loser Interferometers," *IEEE Trans. on Instr. and Measur.* **15**(4): 146–149 (1966).
23. S. Minkowitz and W. Reid Smith Vanir, "Laser Interferometer," *Proc. 1st Congress on Laser Applications (Paris)*, *J. Quantum Electronics* **3**: 237 (1967).
24. G. M. Burgwald and W. P. Kruger, "An Instant-On Laser for Length Measurements," *HPJ* **21**: 2 (1970).
25. J. N. Dukes and G. B. Gordon, "A Two-Hundred-Foot Yardstick with Graduations Every Microinch," *HPJ* **21**: 2 (1970).
26. J. H. McLeod, "The Axicon: A New Type of Optical Element," *J. Opt. Soc. Am.* **44**: 592–597 (1954).

27. A. W. Young, "Optical Workshop Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. 4, Academic Press, New York, 1967, chap. 7.
28. C. Deve, *Optical Workshop Principles*, T. L. Tippell (transl.), Hilger and Watts, London, 1945.
29. K. J. Hume, *Metrology with Autocollimators*, Hilger and Watts, London, 1965.
30. M. P. Kothiyal and R. S. Sirohi, "Improved Collimation Testing Using Talbot Interferometry," *Appl. Opt.* **26**: 4056–4057 (1987).
31. R. E. Noble, "Some Parameter Measurements," in D. Malacara (ed.), *Optical Shop Testing*, 1st ed., John Wiley and Sons, New York, 1978.
32. K. H. Carnell and W. T. Welford, "A Method for Precision Spherometry of Concave Surfaces," *J. Phys. E.* **4**: 1060–1062 (1971).
33. D. H. Rank, "Measurement of the Radius of Curvature of Concave Spheres," *J. Opt. Soc. Am.* **36**: 108–110 (1946).
34. D. F. Horne, *Dividing, Ruling and Mask Making*, Adam Hilger, London, 1974, chap. VII.
35. R. Kingslake, *Optical System Design*, Academic Press, New York, 1983, chap. 13.
36. V. Met, "Determination of Small Wedge Angles Using a Gas Laser," *Appl. Opt.* **5**: 1242–1244 (1966).
37. G. W. Leppelmeier and D. J. Mullenhoff, "A Technique to Measure the Wedge Angle of Optical Flats," *Appl. Opt.* **9**: 509–510 (1970).
38. J. H. Wasilik, T. V. Blomquist, and C. S. Willett, "Measurement of Parallelism of the Surfaces of a Transparent Sample Using Two-Beam Non-Localized Fringes Produced by a Laser," *Appl. Opt.* **10**: 2107–2112 (1971).
39. D. Malacara, (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1992.
40. D. Malacara and O. Harris, "Interferometric Measurement of Angles," *Appl. Opt.* **9**: 1630–1633 (1970).
41. D. Tentori and M. Celaya, "Continuous Angle Measurement with a Jamin Interferometer," *Appl. Opt.* **25**: 215–220 (1986).
42. E. Stijns, "Measuring Small Rotation Rates with a Modified Michelson Interferometer," *Proc. SPIE* **661**: 264–266 (1986).
43. P. Shi and E. Stijns, "New Optical Method for Measuring Small Angle Rotations," *Appl. Opt.* **27**: 4342–4344 (1988).
44. G. D. Chapman, "Interferometric Angular Measurement," *Appl. Opt.* **13**: 1646–1651 (1974).
45. B. K. Johnson, *Optics and Optical Instruments*, Dover, New York, 1947, chaps. II and VIII.
46. L. C. Martin, *Optical Measuring Instruments*, Blackie and Sons Ltd., London, 1924.
47. F. Twyman, *Prisms and Lens Making*, 2d ed., Hilger and Watts, London, 1957.
48. A. S. DeVany, "Reduplication of a Penta-Prism Angle Using Master Angle Prisms and Plano Interferometer," *Appl. Opt.* **10**: 1371–1375 (1971).
49. A. S. DeVany, "Testing Glass Reflecting-Angles of Prisms," *Appl. Opt.* **17**: 1661–1662 (1978).
50. F. Ratajczyk and Z. Bodner, "An Autocollimation Measurement of the Right Angle Error with the Help of Polarized Light," *Appl. Opt.* **5**: 755–758 (1966).
51. A. M. Tareev, "Testing the Angles of High-Precision Prisms by Means of an Autocollimator and a Mirror Unit," *Sov. J. Opt. Technol.* **52**: 50–52 (1985).
52. D. Malacara and R. Flores, "A Simple Test for the 90 Degrees Angle in Prisms," *Proc. SPIE* **1332**: 678 (1990).
53. C. Ai and K. L. Smith, "Accurate Measurement of the Dihedral Angle of a Corner Cube," *Appl. Opt.* **31**: 519–527 (1992).
54. S. M. Rao, "Method for the Measurement of the Angles of a Tetragonal or Corner Cube Prism," *Opt. Eng.* **41**: 1612–1614 (2002).
55. M. S. Scholl, "Ray Trace through a Corner-Cube Retroreflector with Complex Reflection Coefficients," *J. Opt. Soc. Am. A.* **12**(7): 1589–1592 (1995).
56. B. Jurek, *Optical Surfaces*, Elsevier Scient. Pub. Co., New York, 1977.
57. D. F. Horne, *Optical Production Technology*, Adam Hilger, London, and Crane Russak, New York, 1972, chap. XI.
58. M. C. Gerchman and G. C. Hunter, "Differential Technique for Accurately Measuring the Radius of Curvature of Long Radius Concave Optical Surfaces," *Proc. SPIE* **192**: 75–84 (1979).

59. M. C. Gerchman and G. C. Hunter, "Differential Technique for Accurately Measuring the Radius of Curvature of Long Radius Concave Optical Surfaces," *Opt. Eng.* **19**: 843–848 (1980).
60. D. C. O'Shea and S. A. Tilstra, "Non-Contact Measurements of Refractive Index and Surface Curvature," *Proc. SPIE* **966**: 172–176 (1988).
61. J. D. Evans, "Method for Approximating the Radius of Curvature of Small Concave Spherical Mirrors Using a He-Ne Laser," *Appl. Opt.* **10**: 995–996 (1971).
62. J. D. Evans, "Equations for Determining the Focal Length of On-Axis Parabolic Mirrors by He-Ne Laser Reflection," *Appl. Opt.* **11**: 712–714 (1972).
63. J. D. Evans, "Error Analysis to: Method for Approximating the Radius of Curvature of Small Concave Spherical Mirrors Using a He-Ne Laser," *Appl. Opt.* **11**: 945–946 (1972).
64. A. Cornejo-Rodriguez and A. Cordero-Dávila, "Measurement of Radii of Curvature of Convex and Concave Surfaces Using a Nodal Bench and a He-Ne Laser," *Appl. Opt.* **19**: 1743–1745 (1980).
65. P. E. Klingsporn, "Use of a Laser Interferometric Displacement-Measuring System for Noncontact Positioning of a Sphere on a Rotation Axis through Its Center and for Measuring the Spherical Contour," *Appl. Opt.* **18**: 2881–2890 (1979).
66. R. Díaz-Urribe, J. Pedraza-Contreras, O. Cardona-Nuñez, A. Cordero-Dávila, and A. Cornejo Rodriguez, "Cylindrical Lenses: Testing and Radius of Curvature Measurement," *Appl. Opt.* **25**: 1707–1709 (1986).
67. R. Kingslake, "A New Bench for Testing Photographic Lenses," *J. Opt. Soc. Am.* **22**: 207–222 (1932).
68. P. Bouchaud, and J. A. Cogno, "Automatic Method for Measuring Simple Lens Power," *Appl. Opt.* **21**: 3068 (1982).
69. D. Malacara and Z. Malacara, "Testing and Centering of Lenses by Means of Hartmann Test with Four Holes," *Opt. Eng.* **31**: 1551–1555 (1996).
70. M. Martínez-Corral, W. D. Furlan, A. Pons, and G. Saavedra, *Instrumentos Ópticos y Optométricos. Teoría y Prácticas*, Universitat de Valencia, Valencia, (1998).
71. O. Kafri and I. Glatt, *The Physics of Moiré Metrology*, Wiley Interscience, New York, 1990.
72. I. Glatt and O. Kafri, "Determination of the Focal Length of Non-Paraxial Lenses by Moiré Deflectometry," *Appl. Opt.* **26**: 2507–2508, (1987).
73. D. Malacara-Doblado and D. Malacara-Hernández, "Measuring Convergence or Divergence Power with Moiré fringes," *Proc. SPIE* **2860**: 390–393 (1996).
74. Y. Nakano, and K. Murata, "Talbot Interferometry for Measuring the Focal Length of a Lens," *Appl. Opt.* **24**:19 (1985).
75. J. L. Horner, Collimation Invariant Technique for Measuring the Focal Length of a Lens," *Appl. Opt.* **28**: 1047–1047 (1989).
76. A. A. Camacho, C. Solano, M. Cywiak, G. Martínez-Ponce, and R. Baltazar, "Method for the Determination of the Focal Length of a Microlens," *Opt. Eng.* **39**: 2149–2152 (2000).
77. A. A. Camacho, C. Solano, G. Martínez-Ponce, and R. Baltazar, "Simple Method to Measure the Focal Length of a Lens," *Opt. Eng.* **41**: 2899–2902 (2002).
78. B. Howland and A. F. Proll, "Apparatus for the Accurate Determination of Flange Focal Distance," *Appl. Opt.* **11**: 1247–1251 (1970).

This page intentionally left blank.

Daniel Malacara-Hernández

*Centro de Investigaciones en Optica, A. C.
León, Gto., Mexico*

13.1 GLOSSARY

E	electric field strength
k	radian wave number
r	position
t	time
λ	wavelength
φ	phase
ω	radian frequency

13.2 INTRODUCTION

The requirements for high-quality optical surfaces are more demanding every day. They should be tested in an easier, faster, and more accurate manner. Optical surfaces usually have a flat or a spherical *shape*, but they also may be toroidal or generally aspheric. Frequently, an aspherical surface is a conic of revolution, but an aspherical surface can only be made as good as it can be tested. Here, the field of optical testing will be reviewed. There are some references that the reader may consult for further details.¹

13.3 CLASSICAL NONINTERFEROMETRIC TESTS

Some classical tests will never be obsolete, because they are cheap, simple, and provide qualitative results about the shape of the optical surface or wavefront almost instantly. These are the Foucault or knife-edge test, the Ronchi test, and the Hartmann test. They will be described next.

Foucault Test

The Foucault or knife-edge test was invented by Leon Foucault² in France, to evaluate the quality of spherical surfaces. This test described by Ojeda-Castañeda³ detects the presence of transverse aberrations by intercepting the reflected rays deviated from their ideal trajectory, as Fig. 1 shows. The observer is behind the knife, looking at the illuminated optical surface, with the reflected rays entering the eye. The regions corresponding to the intercepted rays will appear dark, as in Fig. 2.

This test is extremely sensitive. If the wavefront is nearly spherical, irregularities as small as a fraction of the wavelength of the light may be easily detected. This is the simplest and most powerful qualitative test for observing small irregularities and evaluating the general smoothness of the spherical

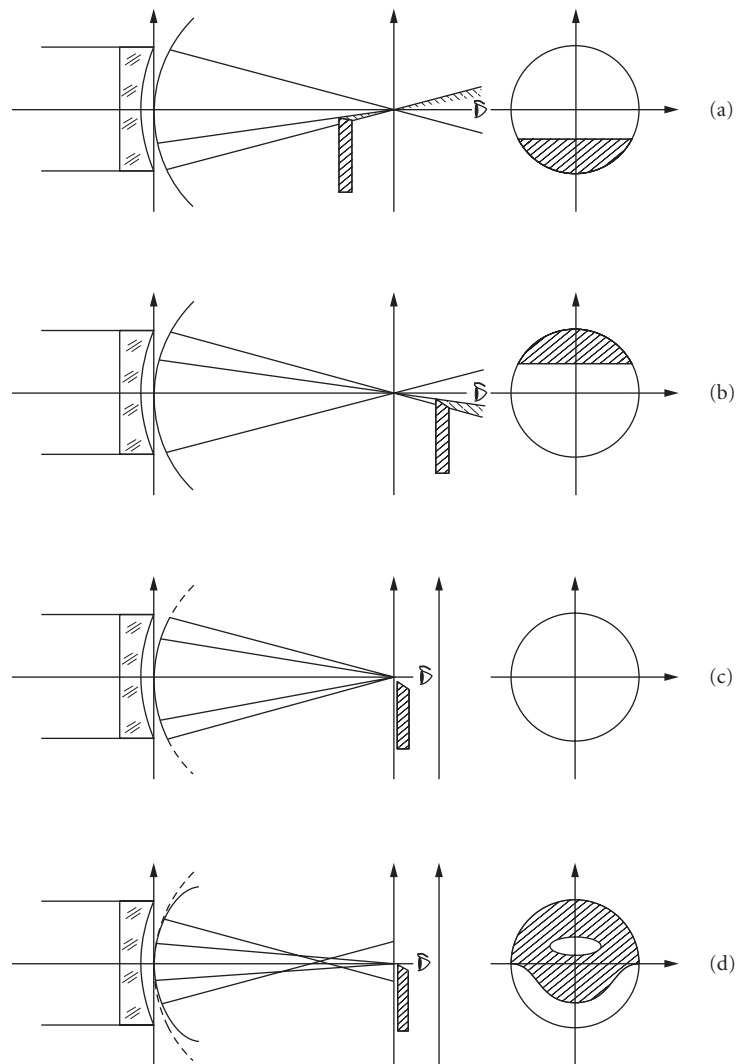


FIGURE 1 Optical schematics for the Foucault test of a spherical mirror at several positions of the knife edge.



FIGURE 2 An optical surface being examined by the Foucault test. (From Ojeda-Castañeda.³)

surface under test. Any other surface or lens may be tested, as long as it produces an almost spherical wavefront, otherwise, an aberration compensator must be used, as will be described later. Very often a razor blade makes a good, straight, sharp edge that is large enough to cover the focal region.

Ronchi Test

Vasco Ronchi⁴ invented his famous test in Italy in 1923. A coarse ruling (50–100 lines per inch) is placed in the convergent light beam reflected from the surface under test, near its focus. The observer is behind the ruling, as Fig. 3 shows, with the light entering the eye. The dark bands in the ruling intercept light, forming shadows on the illuminated optical surface. These shadows will be straight and parallel only if the reflected wavefront is perfectly spherical. Otherwise, the fringes will be curves whose shape and separation depends on the wavefront deformations. The Ronchi test

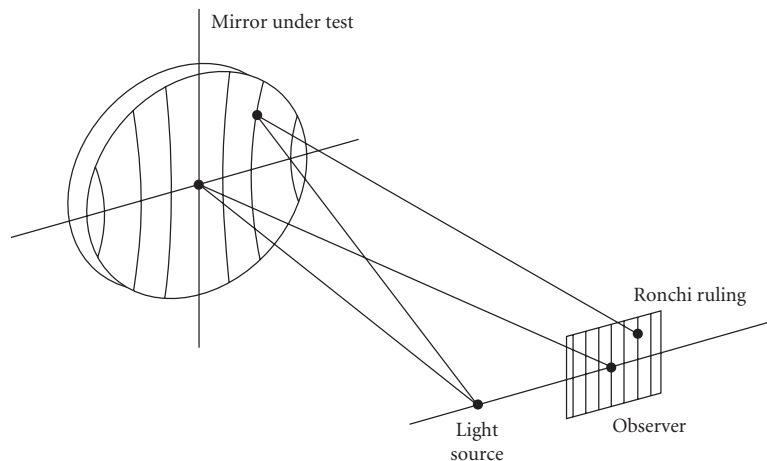


FIGURE 3 Testing a concave surface by means of the Ronchi test.

measures the transverse aberrations in the direction perpendicular to the slits on the grating. The wavefront deformations $W(x, y)$ are related to the transverse aberrations $TA_x(x, y)$ and $TA_y(x, y)$ by the following well-known relations:

$$TA_x(x, y) = -r \frac{\partial W(x, y)}{\partial x} \quad (1)$$

and

$$TA_y(x, y) = -r \frac{\partial W(x, y)}{\partial y} \quad (2)$$

where r is the radius of curvature of the wavefront $W(x, y)$. Thus, if we assume a ruling with period d , the expression describing the m th fringe on the optical surface is given by

$$\frac{\partial W(x, y)}{\partial x} = -\frac{md}{r} \quad (3)$$

Each type of aberration wavefront has a characteristic Ronchi pattern, as shown in Fig. 4; thus, the aberrations in the optical system may be easily identified, and their magnitude estimated. We may interpret the Ronchi fringes not only as geometrical shadows, but also as interferometric fringes, identical with those produced by a lateral shear interferometer.

Hartmann Test

J. Hartmann⁵ invented his test in Germany. It is one of the most powerful methods to determine the figure of a concave spherical or aspherical mirror. Figure 5 shows the optical configuration used

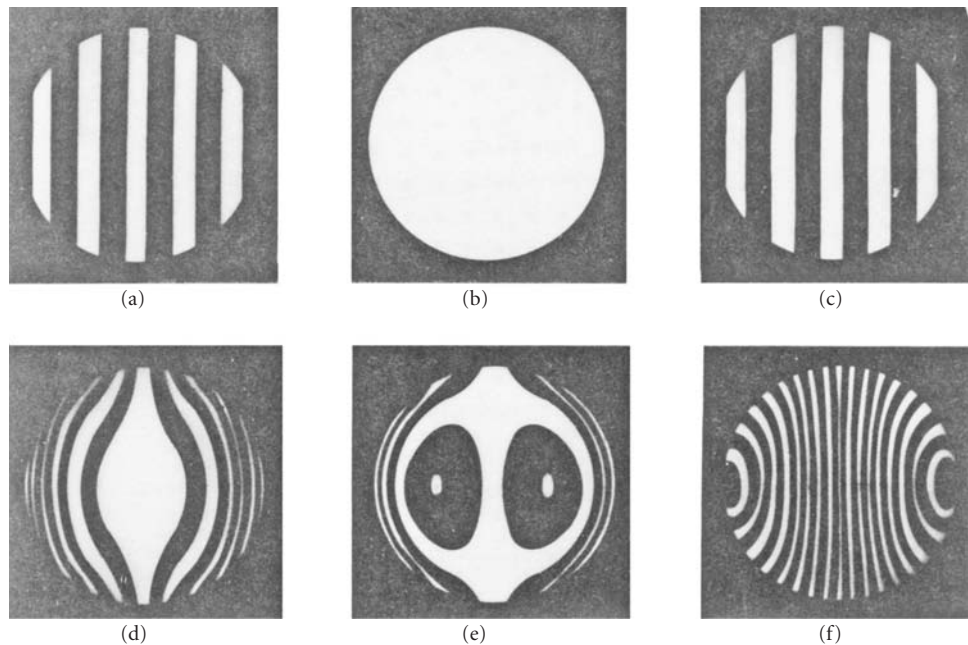


FIGURE 4 Typical Ronchi patterns for a spherical and a parabolic mirror for different positions of the Ronchi ruling.

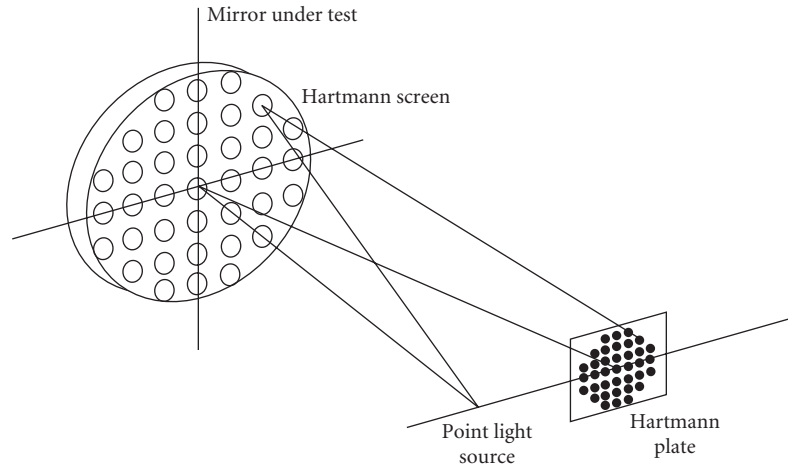


FIGURE 5 Optical arrangement to perform the Hartmann test.

in this test, where a point light source illuminates the optical surface, with its Hartmann screen in front of it. The light beams reflected through each hole on the screen are intercepted on a photographic plate near the focus. Then, the position of the recorded spots is measured to find the value of the transverse aberration on each point. If the screen has a rectangular array of holes, the typical Hartmann plate image for a parabolic mirror looks like that in Fig. 6. The wavefront $W(x, y)$ may be obtained from integration of Eqs. (1) and (2) as follows:

$$W(x, y) = -\frac{1}{r} \int_0^x TA_x(x, y) dx \quad (4)$$

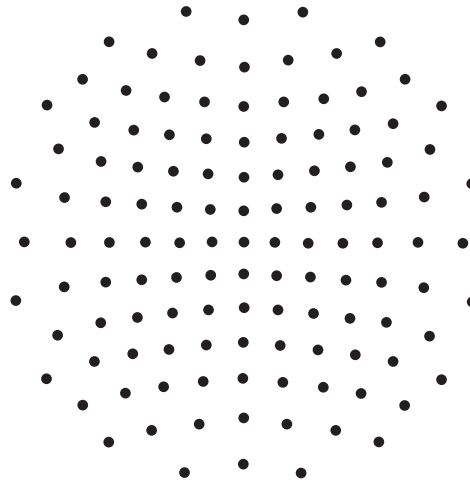


FIGURE 6 Array of spots in a Hartmann plate of a parabolic mirror.

and

$$W(x, y) = -\frac{1}{r} \int_0^y TA_y(x, y) dy \quad (5)$$

After numerical integration of the values of the transverse aberrations, this test provides the concave surface shape with very high accuracy. If the surface is not spherical, the transverse aberrations to be integrated are the difference between the measured values and the ideal values for a perfect surface. Extended, localized errors, as well as asymmetric errors like astigmatism, are detected with this test. The two main problems of this test are that small, localized defects are not detected if they are not covered by the holes on the screen. Not only is this information lost, but the integration results will be false if the localized errors are large. The second important problem of the Hartmann test is that it is very time consuming, due to the time used in measuring all the data points on the Hartmann plate. These problems are avoided by complementing this test with the Foucault test, using an Offner compensator, in order to be sure about the smoothness of the surface (discussed under “Measuring Aspherical Wavefronts”). Various stratagems are available to speed the process. These include modulating the light at different frequencies at each of the holes. Variations also include measuring in front of, behind, or at the focus to get slope information. This technique can be considered an experimental ray trace.

Hartmann-Shack Test

Platt and Shack,⁶ proposed using a lenticular screen, instead of a screen with an array of holes, as illustrated in Fig. 7. This is a simple but important modification from the classic Hartmann. Some differences are

- (a) In the Hartmann test the pattern is obtained in focused convergent light beam, near the focus. On the other hand, in the Hartmann-Shack the test is made in a nearly collimated beam of light.
- (b) A practical advantage of the Hartmann-Shack method is that any positive or negative power can be easily detected and measured.

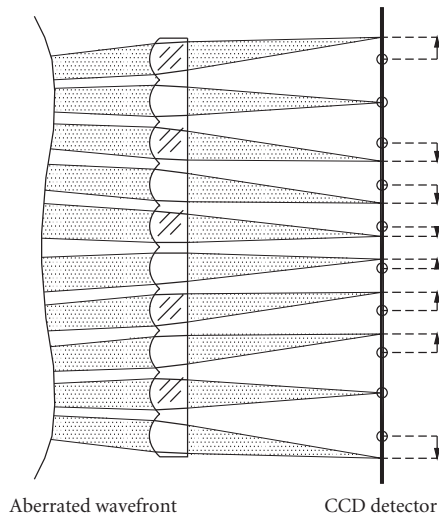


FIGURE 7 Hartmann-Shack Test.

- (c) A second advantage is that each of the spots is individually focused on the detector, making the light energy density of the spot higher than in the Hartmann test.
- (d) The Hartmann-Shack lenticular screen can be made with two identical layers of cylindrical lenses perpendicular to each other, or with a lenslet arrays in molded plastic, glass, or fused silica.

If the wavefront is flat, the light beam passing through each lens is focused close to the optical axis of each lenslet. Since the lens array is not perfect, the lenticular array must be previously calibrated with a reference well-known flat wavefront.

The spot displacement on the detector is equal to the wavefront slope multiplied by the focal length of the lenslet, thus, a shorter focal length will give a greater dynamic range but a reduced angular sensitivity. The optimum focal length depends on the application.

13.4 INTERFEROMETRIC TESTS

Classical geometrical tests are very simple, but they do not provide the accuracy of the interferometric tests. Quite generally, an interferometric test produces an interferogram by producing the interference between two wavefronts. One of these two wavefronts is the wavefront under test. The other wavefront is either a perfectly spherical or flat wavefront, or a copy of the wavefront under test.

When the second wavefront is perfectly spherical or flat, this wavefront acts as a reference. The separation between the two wavefronts, or optical path difference $OPD(x, y)$, is a direct indication of the deformations $W(x, y)$ of the wavefront under test. Then, we may simply write $W(x, y) = OPD(x, y)$. There are many types of interferometers producing interferograms of these type of interferograms, for example, the Twyman-Green and the Fizeau interferometers. Some other interferometers can be considered as modifications of these two basic interferometers, such as the Point Diffraction and the Burch interferometers, and many others that will not be described.

Twyman-Green Interferometer

The Twyman-Green interferometer is illustrated in Fig. 8. The light from a monochromatic point light source is collimated to produce a flat wavefront. Then, the two interfering wavefronts are generated by means of a partially reflective and partially transmitting glass plate, called beam splitter.

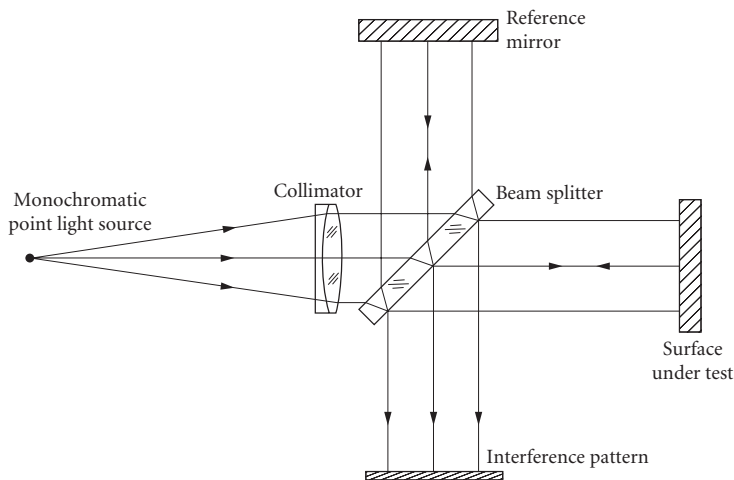


FIGURE 8 Twyman-Green interferometer.

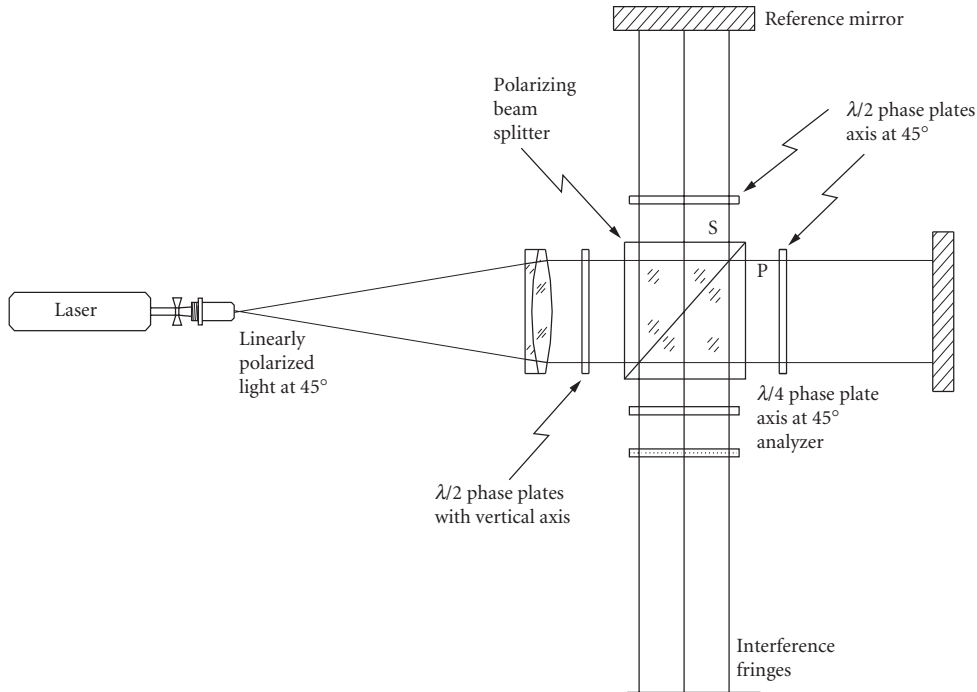


FIGURE 9 Twyman-Green interferometer with polarized beam splitter.

After the beam splitter two flat wavefronts travel in orthogonal directions, one of them to the flat reference mirror and the other to the surface or optical element under test. After returning to the beam splitter, the two wavefronts are recombined to produce an interference pattern. The beam splitter can be oriented at 45° as in Fig. 8, but sometimes a different angle is chosen, for example, a Brewster angle to avoid the reflected beam from the second face on the beam splitter.

Instead of a plane parallel beam splitter, sometimes a polarizing cube beam splitter is used, as in Fig. 9. In this system the plane wavefront entering the beam splitter is linearly polarized at an angle of 45° with respect to the plane of the interferometer. Then, the two wavefronts exiting the cube in orthogonal directions will also be linearly polarized but one of them in the vertical plane and the other in the horizontal plane. A $\lambda/4$ phase plate is located at each of the two exiting faces on the cube. These phase plates produce circularly polarized beams, one going to the reference mirror and the other to the surface under test, but one is right handed and the other is left handed. When arriving back to the cube, after passing twice through the phase plates, the two beams will be again linearly polarized, one in the vertical plane and the other in the horizontal plane. However, these planes of polarization will be orthogonal to the planes of polarization when the wavefront exited the cube. Thus, the two wavefronts are sent to the observation plane and not back to the light source.

Two important facts should be noticed: (a) that the two wavefronts are orthogonally polarized and hence they can not interfere and (b) that there are not any light beams going back to the light source to produce a complementary pattern as in the classic Twyman-Green interferometer. In order to produce observable fringes a linear polarizer should be placed just before the observing plane.

This interferometer with a polarizing beam splitter has many practical advantages.

Fizeau Interferometer

A Fizeau interferometer, illustrated in Fig. 10, is also a two-beam system like the Twyman-Green interferometer. The main difference is that the plate beam splitter is not at 45° with the illuminated

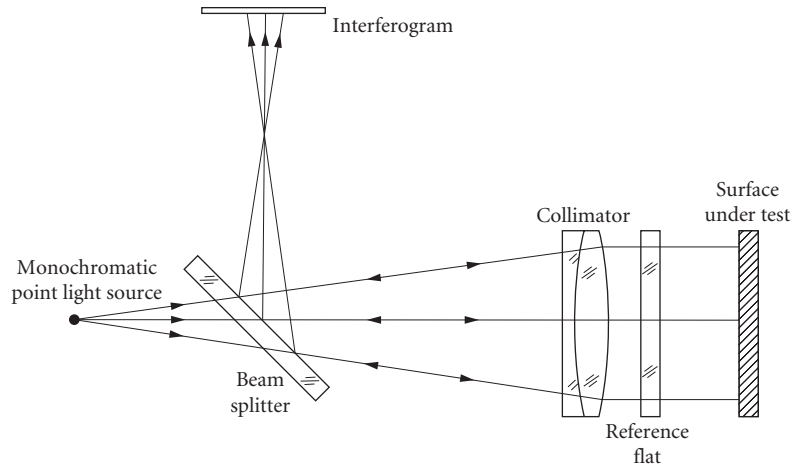


FIGURE 10 Fizeau interferometer.

collimated beam, but perpendicular to its incidence path. Some important practical advantages are that (a) the beam splitter can be smaller for the same aperture, (b) it is more compact, and (c) it is easier to align.

Figure 11 shows some typical interferograms for the Seidel primary aberrations obtained with a Twyman-Green or a Fizeau interferometer.¹ The mathematical analysis of interferograms for wavefronts with arbitrary deformations is a research topic of great interest that has been described in a large number of publications.⁷

Common Path Interferometer

A common path interferometer is one for which the two interfering beams travel the same paths. The optical path difference at the center of the optical axis is zero and cannot be modified. Thus interference with a white light source can be easily achieved. An example of this kind of interferometer is the point-diffraction interferometer first described by Linnik in 1933 and later rediscovered by Smart and Strong.⁸ The lens of optical element under test focuses the light at the center of a small diffracting plate as illustrated in Fig. 12. This diffracting plate is coated with a thin partially transmitting film with a small uncoated disk at its center. The diameter of the central clear disk is about the size of the diffraction air disk produced by a perfect optical system. If the wavefront from the system is not spherical but distorted, the focused spot would be larger than the central disk. Then, two wavefronts are produced. One is a reference spherical wavefront arising from the light diffracted at the central disk. The undiffracted light passing outside of the disk is the wavefront under test.

The fringe patterns in the point diffraction interferometer are identical to those produced by the Twyman-Green interferometer.

Lateral Shearing Interferometers

When the second wavefront is not perfectly flat or spherical, but a copy of the wavefront under test, its relative dimensions or orientation must be changed (sheared) in some way with respect to the wavefront under test. Otherwise, no information about the wavefront deformations is obtained, because the fringes will always be straight and parallel independent of any aberrations. There are several kinds of shearing interferometers, depending on the kind of transformation applied to the reference wavefront.

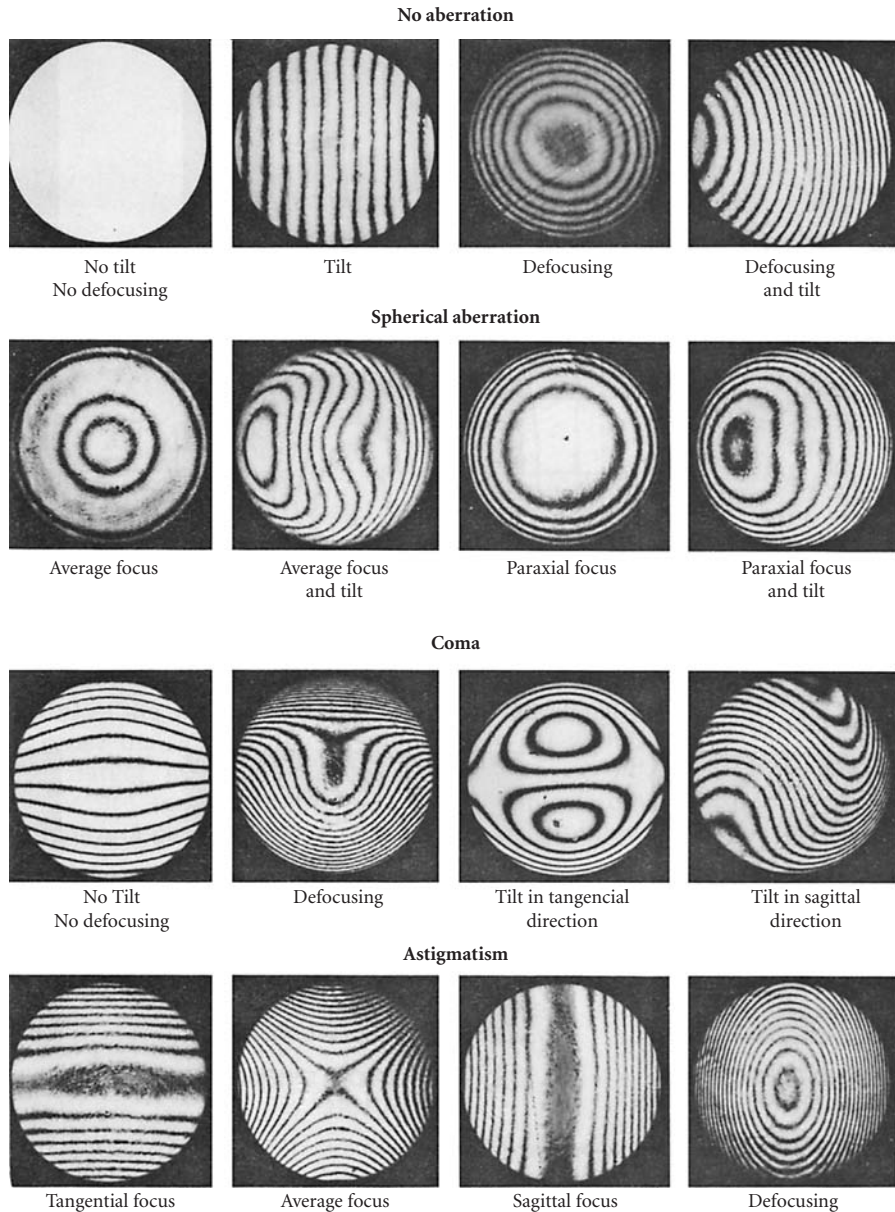


FIGURE 11 Twyman-Green interferograms. (From Malacara.¹)

The most popular of these instruments is the lateral shearing interferometer, with the reference wavefront laterally displaced with respect to the other, as in the interferograms in Fig. 13 shows. The optical path difference $OPD(x, y)$ and the wavefront deformations $W(x, y)$ are related by

$$OPD(x, y) = W(x, y) - W(x - S, y) \quad (6)$$

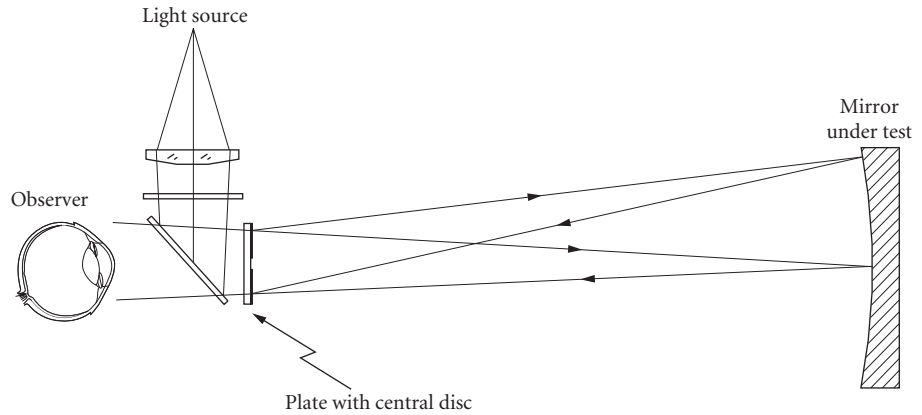


FIGURE 12 Point diffraction interferometer.

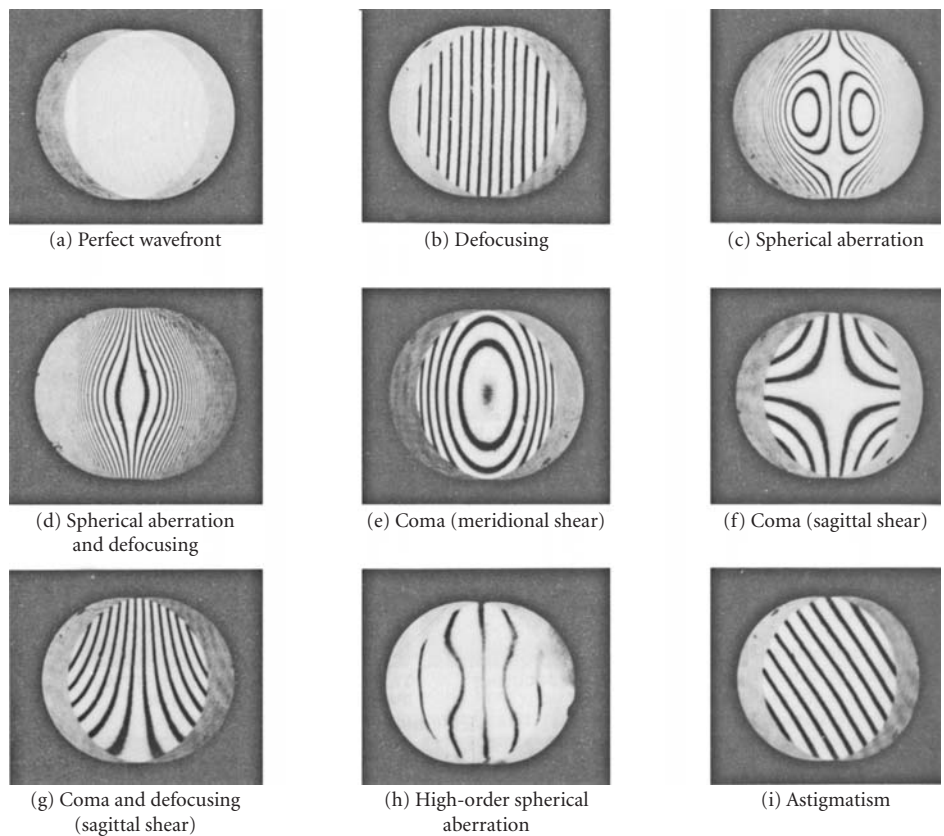


FIGURE 13 Laterally sheared interferograms. (From Malacara.¹)

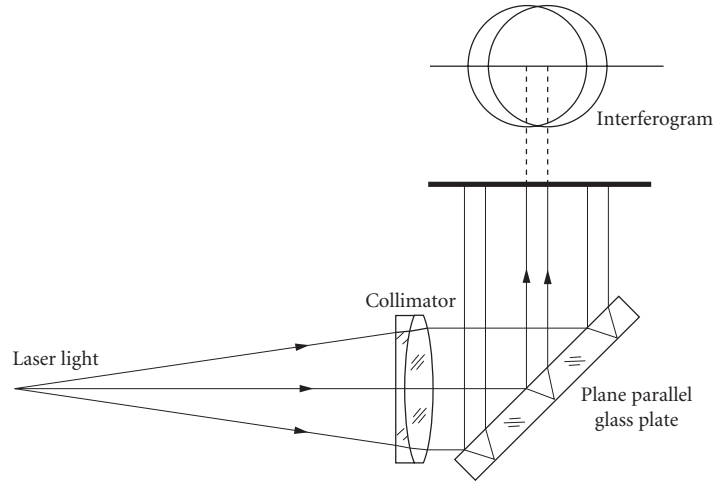


FIGURE 14 Murty's lateral shear interferometer.

where S is the lateral shear of one wavefront with respect to the other. If the shear is small with respect to the diameter of the wavefront, this expression may be approximated by

$$\text{OPD}(x, y) = -S \frac{\partial W(x, y)}{\partial x} = -\frac{S}{r} T A_x(x, y) \quad (7)$$

This relation suggests that the parameter being directly measured is the slope in the x direction of the wavefront (x component $T A_x$ of the transverse aberration). An example of a lateral shear interferometer is the Murty interferometer, illustrated in Fig. 14.

Radial, Rotational, and Reversal Shearing Interferometers

There are also radial, rotational, and reversal shearing interferometers, where the interfering wavefronts are as illustrated in Fig. 15. A radial shear interferometer with a large shear approaches an interferometer with a perfect reference wavefront. These interferometers procedure fringe patterns by the interference of two wavefronts with the same aberrations and deformations. The difference between the two interfering wavefronts is their size or orientation.

The optical path in the radial shear interferometer can be represented by

$$\text{OPD}(x, y) = W(x, y) - W(\rho x, \rho y) \quad (8)$$

A typical radial shear interferogram is in Fig. 16.

In the rotational shear interferometer. The two interfering wavefronts have the same size, but one of those is rotated with respect to the other. In the particular case on a 180° rotation the sensitivity of the interferometer is zero for symmetrical (even) aberrations, like spherical aberration. However, the sensitivity is doubled for antisymmetrical (odd) aberrations, like coma.

In a reversing shear interferometer one of the two wavefronts is reversed with respect to the other about any diameter on the wavefront's pupil. As in the rotational shear interferometer, the sensitivity to aberrations is symmetric with respect to the axis of reversion. Also, the sensitivity to aberration antisymmetric about the axis of reversion is doubled.

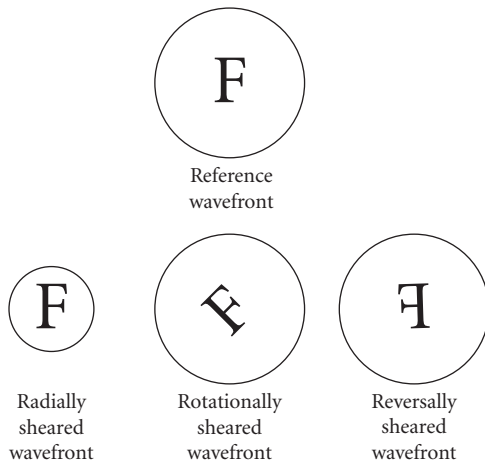


FIGURE 15 Wavefronts in radial, rotational, and reversal shear interferometers.

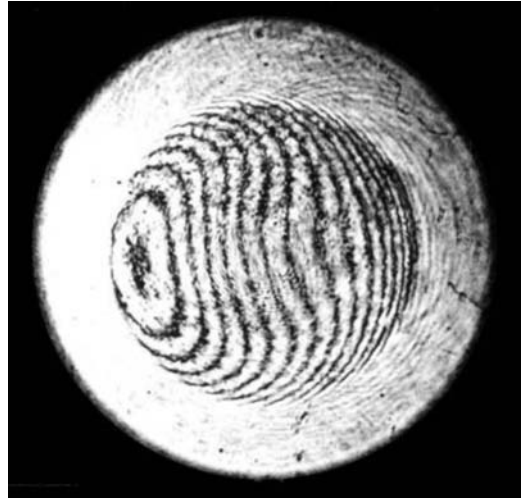


FIGURE 16 A radial shear interferogram.

13.5 INCREASING THE SENSITIVITY OF INTERFEROMETERS

The sensitivity of interferometers is a small fraction of the wavelength being used (about $\lambda/20$). There are several methods to increase this sensitivity, but the most common methods will now be described.

Multiple-Reflection Interferometers

A method to increase the sensitivity of interferometric methods is to use multiple reflections, as in the Fabry-Perot interferometer. The Newton as well as the Fizeau interferometers can be made multiple-reflection interferometers by coating the reference surface and the surface under test with a high-reflection film. Then, the fringes are greatly narrowed and their deviations from straightness are more accurately measured.⁹

Multiple-Pass Interferometers

Another method to increase the sensitivity of interferometers is by double, or even multiple, pass. An additional advantage of double-pass interferometry is that the symmetrical and antisymmetrical parts of the wavefront aberration may be separated. This makes their identification easier, as Hariharan and Sen¹⁰ have proved. Several arrangements have been devised to use multiple pass.¹¹

Zernike Tests

The Zernike phase-contrast method is another way to improve the sensitivity of an interferometer to small aberrations. It was suggested by Zernike as a way to improve the knife-edge test.¹² There are several versions of this test. The basic principle in all of them is the introduction of a

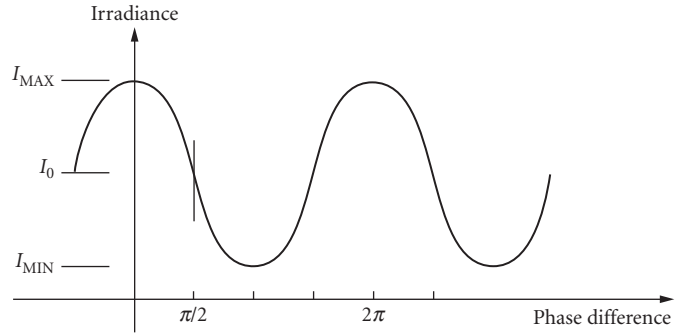


FIGURE 17 Irradiance in an interference pattern, as a function of the phase difference between the two interfering waves.

phase difference equal to $\lambda/2$ between the wavefront under test and the reference wavefront. To understand why this phase difference is convenient, let us consider two interfering beams and irradiances $I_1(x, y)$ and $I_2(x, y)$ and a phase $\phi(x, y)$ between them. The final irradiance $I(x, y)$ in the interferogram is given by

$$I(x, y) = I_1(x, y) + I_2(x, y) + 2\sqrt{I_1(x, y)I_2(x, y)} \cos\phi(x, y) \quad (9)$$

Thus, the irradiance $I(x, y)$ of the combination would be a sinusoidal function of the phase, as illustrated in Fig. 17. If the phase difference is zero for a perfect wavefront, deformations of the wavefront smaller than the wavelength of the light will not be easy to detect, because the slope of the function is zero for a phase near zero. The slope of this function is larger and linear for a phase value of 90° . Thus, the small wavefront deformations are more easily detected if the interferometer is adjusted, so that the wavefronts have a phase difference equal to 90° when the wavefront under test is perfect.

13.6 INTERFEROGRAM EVALUATION

An interferogram may be analyzed in several manners. One way begins by measuring several points on the interferogram, on top of the fringes. Then, the wavefront values between the fringes are interpolated. Another way uses a Fourier analysis of the interferogram. A third method interprets the fringe deformations as a phase modulation.

Fixed Interferogram Evaluation

Once the interferogram has been formed, a quantitative evaluation of it is a convenient method to find the wavefront deformations. The fixed interferogram evaluation by fringe measurements is done by measuring the position of several data points located on top of the fringes. These measurements are made in many ways, for example, with a measuring microscope, with a digitizing tablet, or with a video camera connected to a computer.

The fringe centers can be located either manually, using a digitizing tablet, or automatically, with the computer directly examining a single fringe image that has been captured using a digital frame grabber. After locating the fringe centers, fringe order numbers must be assigned to each point. The wavefront can then be characterized by direct analysis of the fringe centers. If desired, instead of global interpolation, a local interpolation procedure may be used.

To analyze the fringes by a computer, they must first be digitized by locating the fringe centers, and assigning fringe order numbers to them. The optical path difference (OPD) at the center of any fringe is a multiple m of the wavelength ($OPD = m\lambda$), where m is the fringe order. To obtain the

wavefront deformation, only the relative values of the fringe order are important. So any value of the fringe order may be assigned to the first fringe being measured. However, for the second fringe, it may be increased or decreased by one. This choice affects only the sign of the OPD. An important disadvantage of the fixed interferogram analysis is that the sign of the OPD cannot be obtained from the interferogram alone. This information can be retrieved if the sign of any term in the wavefront deformation expression, like defocusing or tilt, is previously determined when taking the interferogram.

Fringes have been digitized using scanners,¹³ television cameras,¹⁴ photoelectric scanners, and digitizing tablets. Review articles by Reid^{15,16} give useful references for fringe digitization using television cameras.

Global and Local Interpolation of Interferograms

After the measurements are made, the wavefront is computed with the measured points. The data density depends on the density of fringes in the interferogram. Given a wavefront deformation, the ratio of the fringe deviations from straightness to the separation between the fringes remains a constant, independently of the number of fringes introduced by tilting of the reference wavefront. If the number of fringes is large due to a large tilt, the fringes look more straight than if the number of fringes is small. Thus, the fringe deviations may more accurately be measured if there are few fringes in the interferogram. Thus, information about many large zones is lost. A way to overcome this problem is to interpolate intermediate values by any of several existing methods. One method is to fit the wavefront data to a two-dimensional polynomial with a least-squares fitting, as described by Malacara et al.¹⁷ or by using splines as described by Hayslett and Swantner¹⁸ and Becker et al.¹⁹ Unfortunately, this procedure has many problems if the wavefront is very irregular. The values obtained with the polynomial may be wrong, especially near the edge, or between fringes if the wavefront is too irregular.

The main disadvantage of global fits is that they smooth the measured surface more than desired. Depending on the degree of the polynomial, there will only be a few degrees of freedom to fit many data points. It is even possible that the fitted surface will pass through none of the measured points. If the surface contains irregular features that are not well described by the chosen polynomial, such as steps or small bumps, the polynomial fit will smooth these features. Then, they will not be visible in the fitted surface.

Global interpolation is done by least-squares fitting the measured data to a two-dimensional polynomial in polar coordinates. The procedure to make the least-squares fitting begins by defining the variance of the discrete wavefront fitting as follows:

$$\sigma = \frac{1}{N} \sum_{i=1}^N [W_i' - W(\rho_i, \theta_i)]^2 \quad (10)$$

where N is the number of data points, W_i is the measured wavefront deviation for data point i , and $W(\rho_i, \theta_i)$ is the functional wavefront deviation after the polynomial fitting. The only requirement is that this variance or fit error is minimized. It is well known that the normal least-squares procedure leads to the inversion of an almost singular matrix. Then, the round-off errors will be so large that the results will be useless. To avoid this problem, the normal approach is to fit the measured points to a linear combination of polynomials that are orthogonal over the discrete set of data points. Thus, the wavefront is represented by

$$W(\rho_i, \theta_i) = \sum_{n=1}^L B_n V_n(\rho_i, \theta_i) \quad (11)$$

$V(\rho, \theta)$ are polynomials of degree r and not the monomials x . These polynomials satisfy the orthogonality condition

$$\sum_{i=1}^N V_n(\rho_i, \theta_i) V_m(\rho_i, \theta_i) = F_n \rho_{nm} \quad (12)$$

where $F_n = \sum V_n^2$.

The advantage of using these orthogonal polynomials is that the matrix of the system becomes diagonal and there is no need to invert it.

The only problem that remains is to obtain the orthogonal polynomials by means of the Gram-Schmidt orthogonalization procedure. It is important to notice that the set of orthogonal polynomials is different for every set of data points. If only one data point is removed or added, the orthogonal polynomials are modified. If the number of data points tends to infinity and they are uniformly distributed over a circular pupil with unit radius, these polynomials $V_r(\rho, \theta)$ approach the Zernike polynomials.²⁰

Several properties of orthogonal polynomials make them ideal for representing wavefronts, but the most important of them is that we may add or subtract one or more polynomial terms without affecting the fit coefficients of the other terms. Thus, we can subtract one or more fitted terms—defocus, for example—without having to recalculate the least-squares fit. In an interferometric optical testing procedure the main objective is to determine the shape of the wavefront measured with respect to a best-fit sphere. Nearly always it will be necessary to add or subtract some terms.

The only problem with these orthogonal polynomials over the discrete set of data points is that they are different for every set of data points. A better choice for the wavefront representation is the set of Zernike polynomials, which are orthogonal on the circle with unit radius, as follows:

$$\int_0^1 \int_0^{2\pi} U_n(\rho, \theta) U_m(\rho, \theta) \rho d\rho d\theta = F_{nm} \delta_{nm} \quad (13)$$

These polynomials are not exactly orthogonal on the set of data points, but they are close to satisfying this condition. Therefore, it is common to transform the wavefront representation in terms of the polynomials V_n to another similar representation in terms of Zernike polynomials $U_n(\rho, \theta)$, as

$$W(\rho, \theta) = \sum_{n=1}^L A_n U_n(\rho, \theta) \quad (14)$$

Fourier Analysis of Interferograms

A completely different way to analyze an interferogram without having to make any interpolation between the fringes is by a Fourier analysis of the interferogram. An interpolation procedure is not needed because the irradiance at a fine two-dimensional array of points is measured and not only at the top of the fringes. The irradiance should be measured directly on the interferogram with a two-dimensional detector or television camera, and not on a photographic picture. Womack,²¹ Macy,²² Takeda et al.,²³ and Roddier and Roddier²⁴ have studied in detail the Fourier analysis of interferograms to obtain the wavefront deformations.

Consider an interferogram produced by the interference of the wavefront under test and a flat reference wavefront, with a large tilt between them, as in the interferogram in Fig. 18. The tilt is about the y axis, increasing the distance between the wavefronts in the x direction. The picture of this interferogram may be thought of as a hologram reconstructing the wavefront. Thus, three wavefronts (images) are generated when this hologram is illuminated with a flat wavefront. In order to have complete separation between these images, the tilt between the wavefronts must be large enough, so that the angle between them is not zero at any point over the interferogram. This is equivalent to saying that the fringes must be open, and never cross any line parallel to the x axis more than once. One image is the wavefront under test and another is the conjugate of this wavefront.

If the tilt between the wavefront is θ and the wavefront shape is $W(x, y)$, the irradiance is given by

$$I(x, y) = I_1(x, y) + I_2(x, y) + 2\sqrt{I_1(x, y)I_2(x, y)} \cos(\phi_0 + kx \sin\theta + kW(x, y)) \quad (15)$$

where $k = 2\pi/\lambda$. This expression may be rewritten as

$$I = [I_1 + I_2] + \sqrt{I_1 I_2} e^{i(kx \sin\theta + kW)} + \sqrt{I_1 I_2} e^{-i(kx \sin\theta + kW)} \quad (16)$$

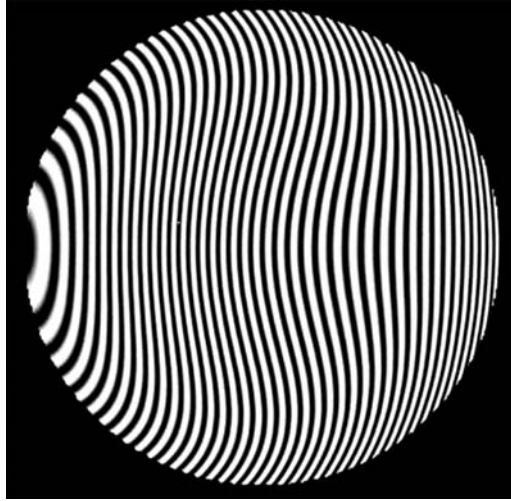


FIGURE 18 Interferogram with a large tilt (linear carrier) to avoid closed fringes.

The first term represents the zero order, the second is the real image, and the third is the virtual image. We also may say that the Fourier transform of the interferogram is formed by a Dirac impulse $\delta(f)$ at the origin and two terms shifted from the origin, at frequencies $+f_0$ and $-f_0$. The quantity f is the spatial frequency, defined by the tilt between the reference wavefront and the wavefront under test ($f = \sin \theta / \lambda$). These terms may be found by taking the Fourier transform of the interferogram. The term at f_0 is due to the wavefront under test. This wavefront may be obtained by taking the Fourier transform of this term, mathematically isolated from the others. This method is performed in a computer by using the fast Fourier transform. The undesired terms are simply eliminated before taking the second fast Fourier transform in order to obtain the wavefront.

Direct Interferometry

This is another method to obtain the wavefront from an interferogram without the need of any interpolation. As in the Fourier method, the image of the interferogram is directly measured with a two-dimensional detector or television camera. The interferogram must have many fringes, produced with a large tilt between the wavefronts. The requirements for the magnitude of this tilt are the same as in the Fourier method.

Consider the irradiance in the interferogram in Fig. 18 along a line parallel to the x axis. This irradiance plotted versus the coordinate x is a perfectly sinusoidal function only if the wavefront is perfect, that is, if the fringes are straight, parallel, and equidistant. Otherwise, this function appears as a wave with a phase modulation. The phase-modulating function is the wavefront shape $W(x, y)$. If the tilt between the wavefronts is θ , the irradiance function is described by Eq. (15). If φ_0 is a multiple of 2π , this expression may be rewritten as

$$I(x, y) = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(kx \sin \theta + kW) \quad (17)$$

Multiplying this phase-modulated function by a sinusoidal signal with the same frequency as the carrier $\sin(kx \sin \theta)$ a new signal S is obtained. Similarly, multiplying by a cosinusoidal signal

$\cos(kx \sin \theta)$ a new signal C is obtained. If all terms in the signals S and C with frequencies equal to or greater than the carrier frequency are removed with a low-pass filter, they become

$$S(x, y) = -\sqrt{I_1 I_2} \sin kW(x, y) \quad (18)$$

$$C(x, y) = \sqrt{I_1 I_2} \cos kW(x, y) \quad (19)$$

then, the wavefront $W(x, y)$ is given by

$$W(x, y) = -\frac{1}{k} \tan^{-1} \left[\frac{S(x, y)}{C(x, y)} \right] \quad (20)$$

which is our desired result.

13.7 PHASE-SHIFTING INTERFEROMETRY

All the methods just described are based on the analysis of a single static interferogram. Static fringe analysis is generally less precise than phase-shifting interferometry, by more than one order of magnitude. However, fringe analysis has the advantage that a single image of the fringes is needed. On the other hand, phase-shifting interferometry requires several images, acquired over a long time span during which the fringes must be stable. This is the main reason why phase-shifting interferometry has seldom been used for the testing of astronomical optics.

Phase-shifting interferometry^{25,26} is possible, thanks to modern tools like array detectors and microprocessors. Figure 19 shows a Twyman-Green interferometer adapted to perform phase-shifting interferometry. Most conventional interferometers, like the Fizeau and the Twyman-Green, have been used to do phase shifting. A good review about these techniques may be found in the review article by Creath.²⁷

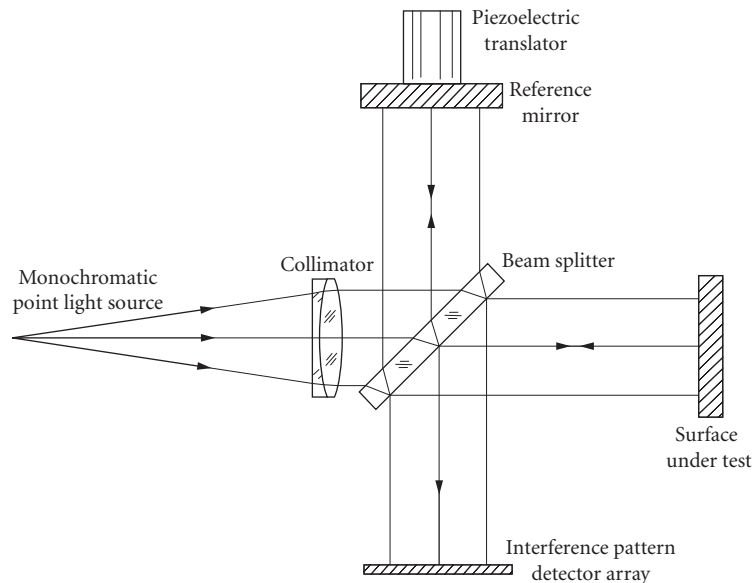


FIGURE 19 Twyman-Green interferogram adapted to do phase shifting.

In phase-shifting interferometers, the reference wavefront is moved along the direction of propagation, with respect to the wavefront under test, changing in this manner their phase difference. This phase shifting is made in steps or in a continuous manner. Of course, this relative displacement of one wavefront with respect to the other may only be achieved through a momentary or continuous change in the frequency of one of the beams, for example, by Doppler shift, moving one of the mirrors in the interferometer. In other words, this change in the phase is accomplished when the frequency of one of the beams is modified in order to form beats.

By measuring the irradiance changes for different values of the phase shifts, it is possible to determine the initial difference in phase between the wavefront under test and the reference wavefront, for that measured point over the wavefront. By obtaining this initial phase difference for many points over the wavefront, the complete wavefront shape is thus determined. If we consider any fixed point in the interferogram, the initial phase difference between the two wavefronts has to be changed in order to make several measurements.

One method that can be used to shift this phase is by moving the mirror for the reference beam along the light trajectory, as in Fig. 19. This can be done in many ways, for example, with a piezoelectric crystal or with a coil in a magnetic field. If the mirror moves with a speed V , the frequency of the reflected light is shifted by an amount equal to $\Delta\nu = 2V/\lambda$.

Another method to shift the phase is by inserting a plane parallel glass plate in the light beam (see Fig. 20). Then the plate is rotated about an axis perpendicular to the optical axis. The phase may also be shifted by means of the device shown in Fig. 21. The first quarter-wave retarding plate is stationary, with its slow axis at 45° with respect to the plane of polarization of the incident linearly polarized light. This plate also transforms the returning circularly polarized light back to linearly

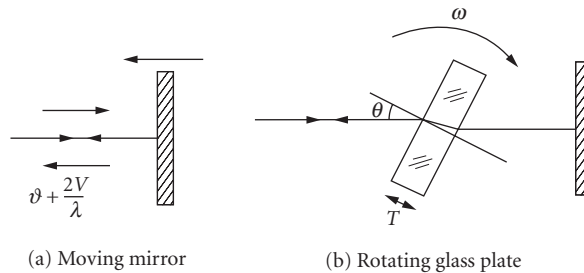


FIGURE 20 Obtaining the phase shift by means of a moving mirror or a rotating glass plate.

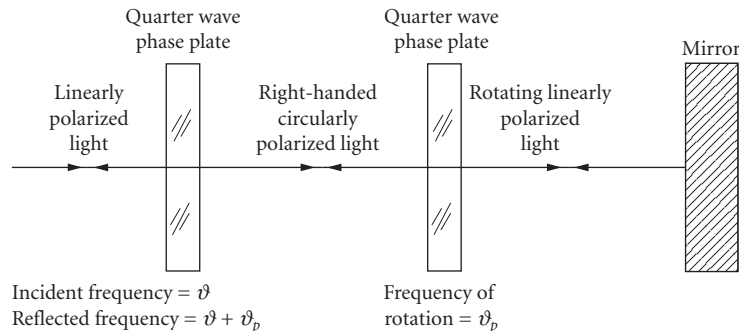


FIGURE 21 Obtaining the phase shift by means of phase plates and polarized light, with a double pass of the light beam.

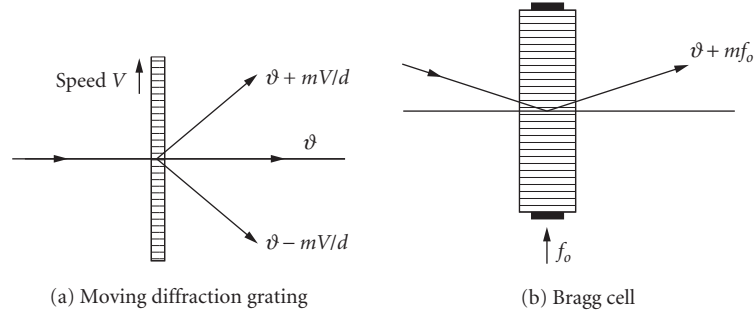


FIGURE 22 Obtaining the phase shift by means of diffraction: (a) with a diffraction grating and (b) with an acousto-optic Bragg cell.

polarized. The second phase retarder is also a quarter-wave plate, it is rotating, and the light goes through it twice; therefore, it is acting as a half-wave plate.

Still another manner to obtain the shift of the phase is by a diffraction grating moving perpendicularly to the light beam, as shown in Fig. 22a, or with an acousto-optic Bragg cell, as shown in Fig. 22b. The change in the frequency is equal to the frequency f of the ultrasonic wave times the order of diffraction m . Thus: $\Delta\nu = mf$.

The nonshifted relative phase of the two interfering wavefronts is found by measuring the irradiance with several predefined and known phase shifts. Let us assume that the irradiance of each of the two interfering light beams at the point x, y in the interference patterns are $I_1(x, y)$ and $I_2(x, y)$ and that their phase difference is $\phi(x, y)$. It was shown before, in Eq. (9), that the resultant irradiance $I(x, y)$ is a sinusoidal function describing the phase difference between the two waves. The basic problem is to determine the nonshifted phase difference between the two waves, with the highest possible precision. This may be done by any of several different procedures.

Phase-Stepping and Four Steps Algorithms

This method²⁷ consists of measuring the irradiance values for several known increments of the phase. There are several versions of this method, which will be described later. The measurement of the irradiance for any given phase takes some time, since there is a time response for the detector. Therefore, the phase has to be stationary during a short time in order to take the measurement. Between two consecutive measurements, the phase is changed by an increment α_i . For those values of the phase, the irradiance becomes

$$I(x, y) = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(\phi + \alpha_i) \quad (21)$$

There are many different algorithms, with many different phase steps, as shown in Fig. 23. The minimum number of steps needed to reconstruct this sinusoidal function is three. As an example with four steps,

$$I_A = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos\phi \quad (22)$$

$$I_B = I_1 + I_2 - 2\sqrt{I_1 I_2} \sin\phi \quad (23)$$

$$I_C = I_1 + I_2 - 2\sqrt{I_1 I_2} \cos\phi \quad (24)$$

$$I_D = I_1 + I_2 + 2\sqrt{I_1 I_2} \sin\phi \quad (25)$$

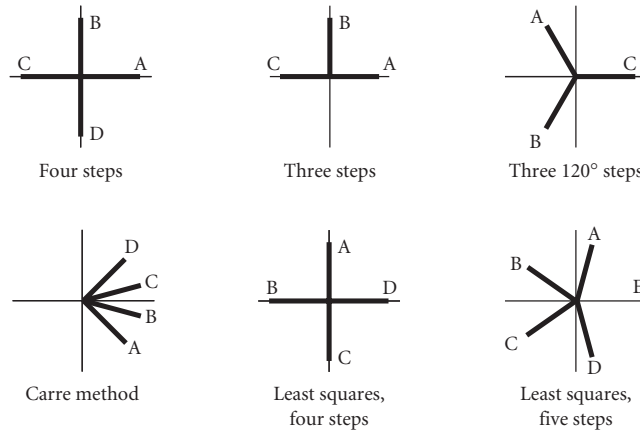


FIGURE 23 Six different ways to shift the phase using phase steps.

Integrating Bucket

In the integrating phase-shifting method the detector continuously measures the irradiance during a fixed time interval, without stopping the phase. Since the phase changes continuously, the average value of the irradiance during the measuring time interval is measured. Thus, the integrating phase-stepping method may be mathematically considered a particular case of the phase-stepping method if the detector has an infinitely short time response. Then, the measurement time interval is reduced to zero. If the measurement is taken, as in Fig. 24, from $\alpha_i + \Delta/2$ to $\alpha_i - \Delta/2$ with center at α_i , then

$$I = \frac{1}{\Delta} \int_{\alpha_i - \Delta/2}^{\alpha_i + \Delta/2} [I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(\phi + \alpha_i)] d\alpha \quad (26)$$

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} \text{sinc}(\Delta/2) \cos(\phi + \alpha_i) \quad (27)$$

In general, in the phase-stepping as well as in the integrating phase-shifting methods, the irradiance is measured at several different values of the phase α_p , and then the phase is calculated.

Two Steps Plus One Method

As pointed out before, phase-shifting interferometry is not useful for testing systems with vibrations or turbulence because the three or four interferograms are taken at different times. An attempt

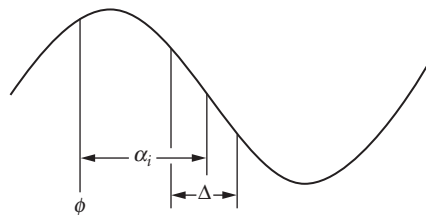


FIGURE 24 Averaged signal measurements with the integrating phase-shifting method.

to reduce this time is the so-called two steps plus one method, in which only two measurements separated by 90° are taken.²⁸ A third reading is taken any time later, of the sum of the irradiance of the beams, independently of their relative phase. This last reading may be taken using an integrating interval $\Delta = 2\pi$. Thus,

$$I_A = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos\phi \quad (28)$$

$$I_B = I_1 + I_2 + 2\sqrt{I_1 I_2} \sin\phi \quad (29)$$

$$I_C = I_1 + I_2 \quad (30)$$

Therefore,

$$\phi = \tan^{-1} \left\{ \frac{I_B - I_C}{I_A - I_C} \right\} \quad (31)$$

Other Phase-Stepping Algorithms

When calculating the phase with several different three or more phase steps, several sources of error affect the accuracy of the result, for example:

- (a) A line miscalibration of the phase shifter. Then the introduced phase steps will have a linear error, which can be interpreted as a change of the reference frequency f_r , making it different from the signed frequency as it should be.
- (b) A nonlinearity in the phase shifter or on the detector that introduces light order harmonics in the detected signed.

These phase errors can be reduced or minimized by properly selecting the number phase steps and their phase increments. Using Fourier theory as shown by Freischlad and Koliopolus,²⁹ a large number of different algorithms with different contributors and number of phase steps have been described in the literature. Depending on the kind of source error and the maximum number of phase steps desired, the proper algorithm can be selected.

Simultaneous Measurement

It has been said several times that the great disadvantage of phase-shifting interferometry is its great sensitivity to vibrations and atmospheric turbulence. To eliminate this problem, it has been proposed that the different interferograms corresponding to different phases be taken simultaneously.^{30,31} To obtain the phase-shifted interferogram, they have used polarization-based interferometers. The great disadvantage of these interferometers is their complexity. To measure the images these interferometers have to use several television cameras.

Heterodyne Interferometer

When the phase shift is made in a continuous manner rather than in steps, the frequency of the shifting beam is permanently modified, and a beating between the two interferometer beams is formed.³²

The phase of the modulated or beating wave may be determined in many ways. One way is by electronic analog techniques, for example, using leading-edge detectors. Another way is by detecting when the irradiance passes through zero, that is, through the axis of symmetry of the irradiance function.

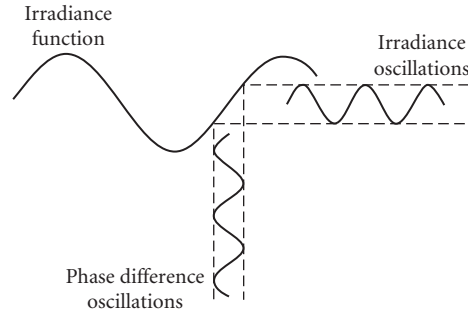


FIGURE 25 Phase-lock method to find the phase with a small sinusoidal modulation of the phase.

Phase Lock

The phase-lock method^{31–35} can be explained with the help of Fig. 25. Assume that an additional phase difference is added to the initial phase $\phi(x, y)$. The additional phase being added has two components: one of them with a fixed value and the other with a sinusoidal time shape. Both components can have any predetermined desired value. Thus, the resultant phase ϕ_r is given by

$$\phi_r = \phi(x, y) + \delta(x, y) + a \sin \omega t \quad (32)$$

then, the irradiance $I(x, y)$ would be given by

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos[\phi + \delta + a \sin \omega t] \quad (33)$$

The amplitude of the phase oscillations $a \sin t$ is much smaller than π . We may now adjust the fixed phase to a value such that $\phi + \delta = \pi/2 + n\pi$. Then the value of $\cos(\phi + \delta)$ is zero. The curve is antisymmetric at this point; hence, only odd harmonics remain on the irradiance signal. This is done in practice by slowly changing the value of the phase δ , while maintaining the oscillation $a \sin \omega t$, until the maximum amplitude of the first harmonic, or fundamental frequency, is obtained. At this point, then, we have $\delta + \phi = \pi/2 + n\pi$, and since the value of δ is known, the value ϕ has been determined.

13.8 MEASURING ASPHERICAL WAVEFRONTS

The most common type of interferometer, with the exception of lateral or rotational shearing interferometers, produces interference patterns in which the fringes are straight, equidistant, and parallel, when the wavefront under test is perfect and spherical with the same radius of curvature as the reference wavefront.

If the surface under test does not have a perfect shape, the fringes will not be straight and their separations will be variable. The deformations of the wavefront may be determined by a mathematical examination of the shape of the fringes. By introducing a small spherical curvature on the reference wavefront (focus shift) or by changing its angle with respect to the wavefront under test (tilt), the number of fringes in the interferogram may be changed. This is done to reduce the number of fringes as much as possible, since the greater the number of fringes, the smaller the sensitivity of the test. However, for aspherical surfaces this number of fringes cannot be smaller than a certain minimum. The larger the asphericity is, the greater is this minimum number of fringes. Since the fringe

separations are not constant, in some places the fringes will be widely spaced, but in some others the fringes will be too close together.

The sensitivity of the test depends on the separation between the fringes, because an error of one wavelength in the wavefront distorts the fringe shape by an amount equal to the separation between the fringes. Thus, the sensitivity is directly proportional to the fringe separation. When the fringes are widely separated, the sampled points will be quite separated from each other, leaving many zones without any information. On the other hand, where the fringes are very close to each other, there is a high density of sampled data points, but the sensitivity is low.

Then, it is desirable that the spherical aberration of the wavefront under test is compensated in some way, so that the fringes appear straight, parallel, and equidistant, for a perfect wavefront. This is called a null test and may be accomplished by means of some special configurations. These special configurations may be used to conduct a null test of a conic surface. These are described in several books.¹ Almost all of these surfaces have rotational symmetry.

If no testing configuration can be found to get rid of the spherical aberration, additional optical components, called null compensators, have to be used. Many different types of compensators have been invented. The compensators may be refractive (lenses), reflective (mirrors), or diffractive (real or computer-generated holograms).

Refractive or Reflective Compensators

The simplest way to compensate the spherical aberration of a paraboloid or a hyperboloid tested at the center of curvature is a single convergent lens placed near the point of convergence of the rays, as Fig. 26 shows. This lens is called a Dall compensator. Unfortunately, the correction due to a single lens is not complete, so a system of two lenses must be used to obtain a better compensation. This system is called an Offner compensator and is shown in Fig. 27. The field lens L is used to image the surface under test on the plane of the compensating lens L . Mirrors may also be used to design a null compensator.

As the sad experience of the Hubble space telescope proves, the construction parameters in a lens compensator have to be very carefully measured and adjusted, otherwise an imperfect correction is

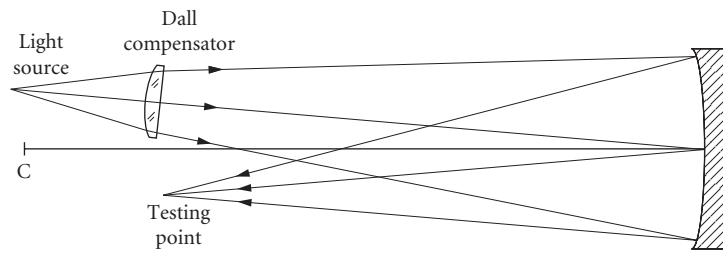


FIGURE 26 The Dall compensator.

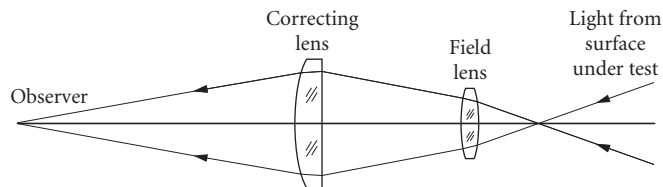


FIGURE 27 The Offner compensator. Only the reflected beam is shown.

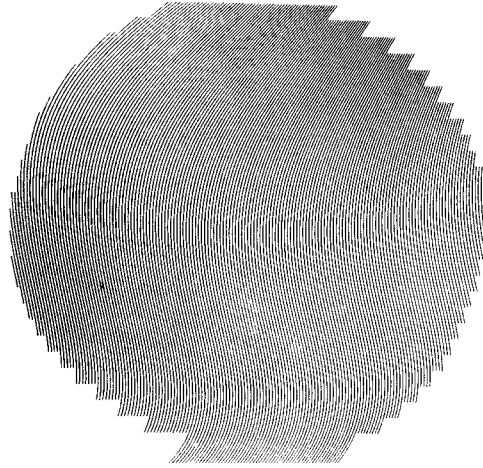


FIGURE 28 Computer-generated hologram for testing an aspherical wavefront. (From Wyant.³⁷)

obtained either by undercorrection or overcorrection. The distance from the compensator to the surface under test is one of those parameters to be carefully measured. A way around this problem would be to assume that the compensator detects smoothness imperfections but not the exact degree of asphericity. This degree of asphericity may then be measured with some independent measurement like the Hartmann test.

Holographic Compensators

Diffractive holographic elements also may be used to compensate the spherical aberration of the system and to obtain a null test. The hologram may be real, produced by photographing an interferometric pattern. This pattern has to be formed by superimposing on the screen a wavefront like the one we have to test and a perfectly flat or spherical wavefront. The only problem with this procedure is that a perfect wavefront with the same shape as the wavefront to be tested has first to be produced. This is not always easy.

A better approach is to simulate the holographic interference pattern in a computer,³⁶ as in Fig. 28. Then this image is transferred to a small photographic plate, with the desired dimensions. There are many experimental arrangements to compensate the aspherical wavefront aberration with a hologram. One of these is illustrated in Fig. 29.

Infrared Interferometry

Another simple approach to reduce the number of fringes in the interferogram is to use a long infrared wavelength. Light from a CO₂ laser has been used with this purpose. It can also be used when the surface is still quite rough.

Two-Wavelength Interferometry

In phase-shifting interferometry, each detector must have a phase difference smaller than π from the closest neighboring detector, in order to avoid 2π phase ambiguities and ensure phase continuity. In

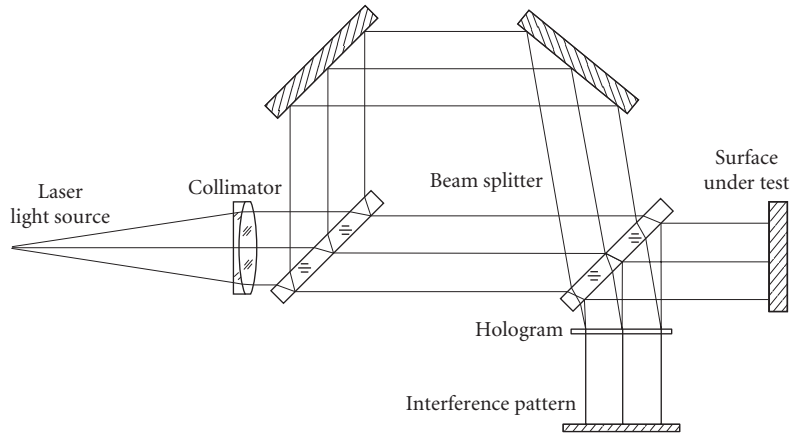


FIGURE 29 An optical arrangement for testing an aspherical wavefront with a computer-generated hologram.

other words, there should be at least two detector elements for each fringe. If the slope of the wavefront is very large, the fringes will be too close together and the number of detector elements would be extremely large.³⁷

A solution to this problem is to use two different wavelengths λ_1 and λ_2 simultaneously. The group wavelength or equivalent wavelength λ_{eq} is longer than any of the two components and is given by

$$\lambda_{\text{eq}} = \frac{\lambda_1 \lambda_2}{|\lambda_1 - \lambda_2|} \quad (34)$$

Under these conditions, the requirement in order to avoid phase uncertainties is that there should be at least two detectors for each fringe produced if the wavelength is λ_{eq} . The great advantage of this method is that we may test wavefronts with large asphericities, limited in asphericity by the group wavelength, and accuracy limited by the shortest wavelength of the two components.

Moiré Tests

An interferogram in which a large amount of tilt has been introduced is an ideal periodic structure to form moiré patterns. A moiré pattern represents the difference between two periodic structures. Thus, a moiré formed by two interferograms represents the difference between the two interferograms. There are several possibilities for the use in optical testing of this technique, as shown by Patorski.³⁸

Let us assume that the two interferograms are taken from the same optical system producing an aspherical wavefront, but with two different wavelengths λ_1 and λ_2 . The moiré obtained represents the interferogram that would be obtained with an equivalent wavelength λ_{eq} given by Eq. (31). If the tilt is of different magnitude in the two interferograms, the difference appears as a tilt in the moiré between them. Strong aspheric wavefronts may be tested with this method.

A second possibility is to produce the moiré between the ideal interferogram for an aspheric wavefront and the actual wavefront. Any differences between both would be easily detected.

Another possibility of application is for eliminating the wavefront imperfections in a low-quality interferometer. One interferogram is taken with the interferometer alone, without any optical piece under test. The second interferogram is taken with the optical component being tested. The moiré

represents the wavefront deformations due to the piece being tested, without the interferometer imperfections.

Sub-Nyquist Interferometry

It was pointed out before that in phase-shifting interferometry each detector must have a phase difference smaller than π from the closest neighboring detector, in order to avoid 2π phase ambiguities and to ensure phase continuity. In other words, there should be at least two detector elements for each fringe. This condition is known as the Nyquist condition.

Since there is a minimum practical distance between detectors, the maximum asphericity in a surface to be tested by phase-shifting interferometry is only a few wavelengths. This condition may be relaxed³⁹ if the wavefront and its slope are assumed to be continuous on the whole aperture. Then, optical surfaces with larger asphericities may be tested.

Wavefront Stitching

When an aspheric and a flat reference wavefronts interfere, the fringe spacing is minimum where the angle between the two wavefronts is larger. When the aspheric wavefront has rotational symmetry and there is no angle between them (no tilt) near the optical axis, the minimum fringe spacing occurs at the edge of the pupil. The maximum fringe spacing is at the center, where the two wavefronts are parallel to each other.

Most times the fringe pattern is imaged on a CCD detector with a rectangular array of small square pixels. According to the Nyquist sampling condition, the fringes can be detected only if the detector has more than two pixels per fringe spacing. With strongly aspheric wavefronts this condition can not be satisfied near the pupil edge.

If the flat reference wavefront is tilted or made slightly spherical, the zone of the interferogram with maximum fringe spacing can be located at any desired. Then the interferogram can be measured at a zone around this point with maximum fringe spacing where the Nyquist condition is satisfied. By moving this point the whole interferogram can thus be measured in small pieces. Then, all pieces should be joined together in a process called wavefront stitching,⁴⁰ as illustrated in Fig. 30.

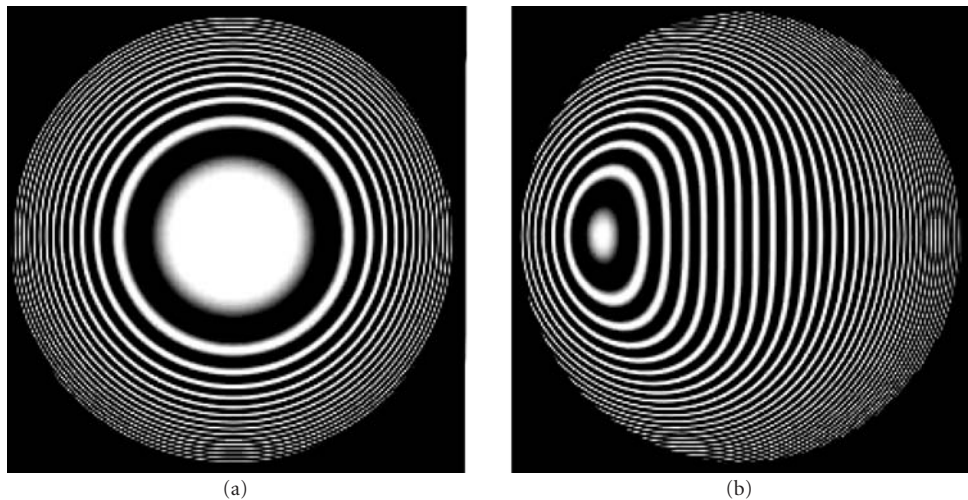


FIGURE 30 Two different interferograms of the same wavefront with different tilts to do wavefront stitching.

13.9 REFERENCES

1. D. Malacara, *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
2. L. M. Foucault, "Description des Procédes Employes pour Reconnaître la Configuration des Surfaces Optiques," *C.R. Acad. Sci. Paris* **47**:958 (1852); reprinted in Armand Colin, *Classiques de al Science*, vol. II.
3. J. Ojeda-Castañeda, "Foucault, Wire and Phase Modulation Tests," in D. Malacara (ed.), *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
4. V. Ronchi, "Le Franque di Combinazione Nello Studio Delle Superficie e Dei Sistemi Ottici," *Ri. Ottica mecc. Precis.* **2**:9 (1923).
5. J. Hartmann, "Bemerkungen über den Bann und die Justirung von Spektrographen," *Zt. Instrumentenk.* **20**:47 (1900).
6. B. C. Platt and R. V. Shack, "Lenticular Hartmann Screen," *Opt. Sci. Newsl.* **5**: 15–16 (1971).
7. D. Malacara, *Interferogram Analysis for Optical Testing*, 2d ed., CRC Press, Taylor and Francis Group, Boca Raton, FL, 2005.
8. R. N. Smart and J. Strong, "Point Diffraction Interferometer," (Abstract only) *J. Opt. Soc. Am.* **62**:737 (1972).
9. C. Roychoudhuri, "Multiple-Beam Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
10. P. Hariharan and D. Sen, "The Separation of Symmetrical and Asymmetrical Wave-Front Aberrations in the Twyman Interferometer," *Proc. Phys. Soc.* **77**:328 (1961).
11. P. Hariharan, "Multiple-Pass Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1991.
12. F. Zernike, "Diffraction Theory of Knife Edge Test and Its Improved Form: The Phase Contrast," *Mon. Not. R. Astron. Soc.* **94**:371 (1934).
13. D. Rozenzweig and B. Alte, "A Facility for the Analysis of Interferograms," in A. H. Guenther, D. H. Liedbergh (eds), *Optical Interferograms—Reduction and Interpretation*, ASTM Symposium, *Am. Soc. for Test and Mat. Tech.* Publ. 666, West Conshohocken, PA, 1978.
14. K. H. Womack, J. A. Jonas, C. L. Koliopoulos, K. L. Underwood, J. C. Wyant, J. S. Loomis, and C. R. Hayslett, "Microprocessor-Based Instrument for Analysis of Video Interferograms," *Proc. SPIE* **192**:134 (1979).
15. G. T. Reid, "Automatic Fringe Pattern Analysis: A Review," *Opt. and Lasers in Eng.* **7**:37 (1986).
16. G. T. Reid, "Image Processing Techniques for Fringe Pattern Analysis," *Proc. SPIE* **954**:468 (1988).
17. D. Malacara, J. M. Carpio-Valadéz, and J. J. Sánchez-Mondragón, "Wavefront Fitting with Discrete Orthogonal Polynomials in a Unit Radius Circle," *Opt. Eng.* **29**:672 (1990).
18. C. R. Hayslett and W. Swantner, "Wave Front Derivation from Interferograms by Three Computer Programs," *Appl. Opt.* **19**:3401 (1980).
19. F. Becker, G. E. A. Maier, and H. Wegner, "Automatic Evaluation of Interferograms," *Proc. SPIE* **359**:386 (1982).
20. F. Zernike, "Begünstheorie des Schneidener-Fahrens und Seiner Verbasserten Form, der Phasenkontrastmethode," *Physica* **1**:689 (1934).
21. K. H. Womack, "Frequency Domain Description of Interferogram Analysis," *Opt. Eng.* **23**:396 (1984).
22. W. W. Macy, Jr., "Two Dimensional Fringe Pattern Analysis," *Appl. Opt.* **22**:3898 (1983).
23. M. Takeda, H. Ina, and S. Kobayashi, "Fourier Transform Method of Fringe-Pattern Analysis for Computer-Based Topography and Interferometry," *J. Opt. Soc. Am.* **72**:156 (1982).
24. C. Roddier and F. Roddier, "Interferogram Analysis Using Fourier Transform Techniques," *Appl. Opt.* **26**:1668 (1987).
25. J. H. Bruning, D. J. Herriott, J. E. Gallagher, D. P. Rosenfeld, A. D. White, and D. J. Brangaccio, "Digital Wavefront Measurement Interferometer," *Appl. Opt.* **13**:2693 (1974).
26. J. Greivenkamp and J. H. Bruning, "Phase Shifting Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1991.
27. K. Creath, "Phase-Measurement Interferometry Techniques," in E. Wolf (ed.), *Progress in Optics*, vol. XXVI, Elsevier Science Publishers, Amsterdam, 1988.

28. P. L. Wizinowich, "Systems for Phase Shifting Interferometry in the Presence of Vibration: A New Algorithm and System," *Appl. Opt.* **29**:3271–3279 (1990).
29. K. Freischlad and C. L. Koliopoulos, "Fourier Description of Digital Phase Measuring Interferometry," *J. Opt. Soc. Am. A.* **7**:542–551 (1990).
30. N. Bareket, "Three-Channel Phase Detector for Pulsed Wavefront Sensing," *Proc. SPIE* **551**:12 (1985).
31. C. L. Koliopoulos, "Simultaneous Phase Shift Interferometer," *Proc. SPIE* **1531**:119–127 (1991).
32. N. A. Massie, "Digital Heterodyne Interferometry," *Proc. SPIE* **816**:40 (1987).
33. G. W. Johnson, D. C. Leiner, and D. T. Moore, "Phase Locked Interferometry," *Proc. SPIE* **126**:152 (1977).
34. G. W. Johnson, D. C. Leiner, and D. T. Moore, "Phase Locked Interferometry," *Opt. Eng.* **18**:46 (1979).
35. D. T. Moore, "Phase-Locked Moire Fringe Analysis for Automated Contouring of Diffuse Surfaces," *Appl. Opt.* **18**:91 (1979).
36. J. C. Wyant, "Holographic and Moire Techniques," in D. Malacara (ed.), *Optical Shop Testing*, John Wiley and Sons, New York, 1978.
37. J. C. Wyant, B. F. Oreb, and P. Hariharan, "Testing Aspherics Using Two-Wavelength Holography: Use of Digital Electronic Techniques," *Appl. Opt.* **23**:4020 (1984).
38. K. Paturski, "Moire' Methods in Interferometry," *Opt. and Lasers in Eng.* **8**:147 (1988).
39. J. E. Greivenkamp, "Sub-Nyquist Interferometry," *Appl. Opt.* **26**:5245 (1987).
40. J. Liesener and H. Tiziani, "Interferometer with Dynamic Reference," *Proc. SPIE* **5252**:264–271 (2004).

This page intentionally left blank.

USE OF COMPUTER-GENERATED HOLOGRAMS IN OPTICAL TESTING

Katherine Creath

*Optineering
Tucson, Arizona, and
College of Optical Sciences
University of Arizona
Tucson, Arizona*

James C. Wyant

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

14.1 GLOSSARY

CGH	computer-generated hologram
M	linear, lateral magnification
N	diffracted order number
n	integers
P	number of distortion-free resolution points
r	radius
S	maximum wavefront slope (waves/radius)
$x, \Delta x$	distance
$\Delta\theta$	rotational angle error
$\Delta\phi$	wavefront phase error
θ	rotational angle
λ	wavelength
$\phi(\)$	wavefront phase described by hologram

14.2 INTRODUCTION

Holography is extremely useful for the testing of optical components and systems. If a master optical component or optical system is available, a hologram can be made of the wavefront produced by the component or system and this stored wavefront can be used to perform null tests of similar

optical systems. If a master optical system is not available for making a hologram, a synthetic or a computer-generated hologram (CGH) can be made to provide the reference wavefront.¹⁻⁸ When an aspheric optical element with a large departure from a sphere is tested, a CGH can be combined with null optics to perform a null test.

There are several ways of thinking about CGHs. For the testing of aspheric surfaces, it is easiest to think of a CGH as a binary representation of the ideal interferogram that would be produced by interfering the reference wavefront with the wavefront produced by a perfect sphere. In the making of the CGH the entire interferometer should be ray traced to determine the so-called perfect aspheric wavefront at the hologram plane. This ray trace is essential because the aspheric wavefront will change as it propagates, and the interferometer components may change the shape of the perfect aspheric wavefront.

Figure 1 shows an example of a CGH. Since the amplitude of the aspheric wavefront is constant across the wavefront, best results are obtained if the lines making up the hologram have approximately one-half the spacing of the lines (i.e., fringe spacing) at the location of the lines. Thus, the line width will vary across the hologram. The major difference between the binary synthetic hologram

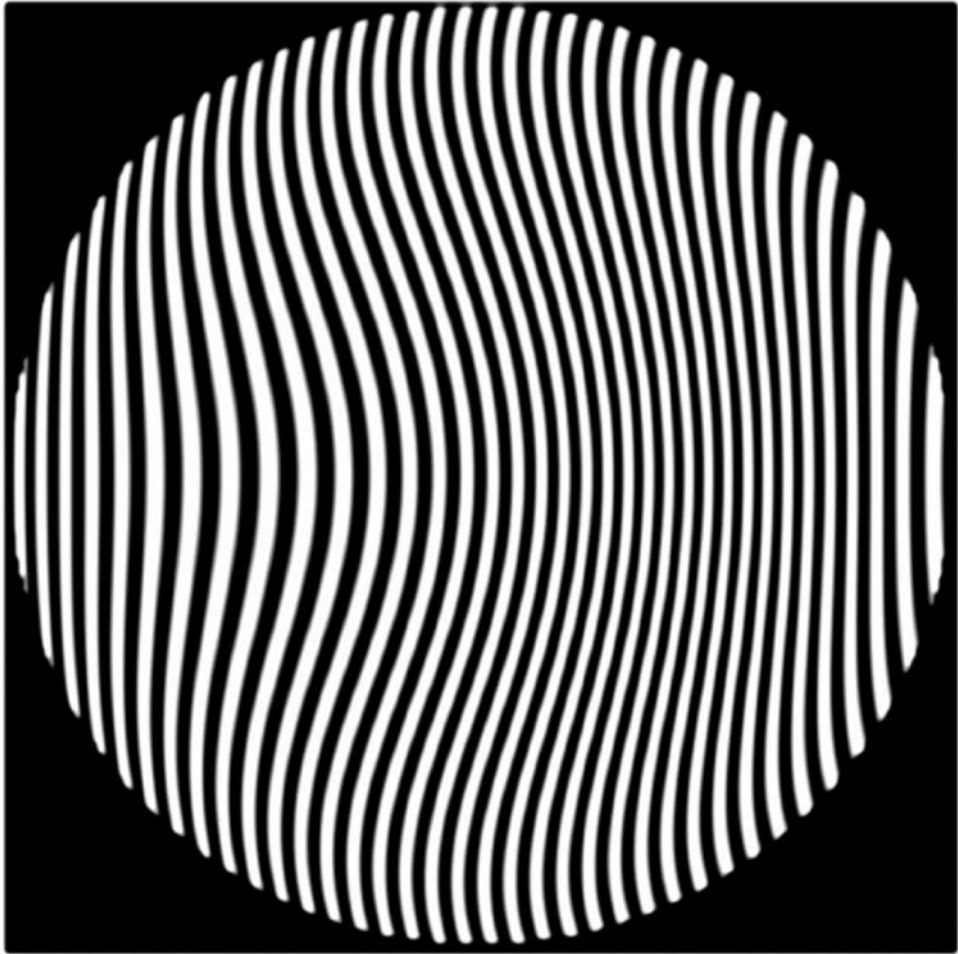


FIGURE 1 Sample computer-generated hologram (CGH).

and the real grayscale hologram that would be produced by interfering a reference wavefront and the aspheric wavefront is that additional diffraction orders are produced. These additional diffraction orders can be eliminated by spatial filtering.

14.3 PLOTTING CGHs

The largest problem in making CGHs is the plotting. The accuracy of the plot determines the accuracy of the wavefront. It is easier to see the plotting accuracy by comparing a binary synthetic hologram with an interferogram. In an interferogram, a wavefront error of $1/n$ waves causes a fringe to deviate from the ideal position by $1/n$ the fringe spacing. The same is true for CGHs. A plotting error of $1/n$ the fringe spacing will cause an error in the produced aspheric wavefront of $1/n$ wave. As an example, assume the error in drawing a line is $0.1\ \mu\text{m}$ and the fringe spacing is $20\ \mu\text{m}$, then the wavefront produced by the CGH will have an error in units of wave of $0.1/20$, or $1/200$ wave.

To minimize wavefront error due to the plotter, the fringe spacing in the CGH should be as large as possible. The minimum fringe spacing is set by the slope difference between the aspheric wavefront and the reference wavefront used in the making of the synthetic hologram. While it is not mandatory, the interferogram is cleaner if the slope difference is large enough to separate the diffraction orders so spatial filtering can be used to select out only the first order. Figure 2 shows a photograph of the diffracted orders. As shown in Fig. 2, to ensure no overlapping of the first and second orders in the Fourier plane, the tilt angle of the reference beam needs to be greater than three times the maximum slope of the aberrated wave.⁹ This means that, in general, the maximum slope difference between the reference and test beams is four times the maximum slope of the test beam. Thus, the error produced by plotter distortion is proportional to the slope of the aspheric wavefront being produced.

Many plotters have been used to plot holograms, but the best holograms are now made using either laser-beam recorders or more commonly electron-beam (e-beam) recorders of the type used for producing masks in the semiconductor industry.¹⁰ The e-beam recorders write onto photoresist

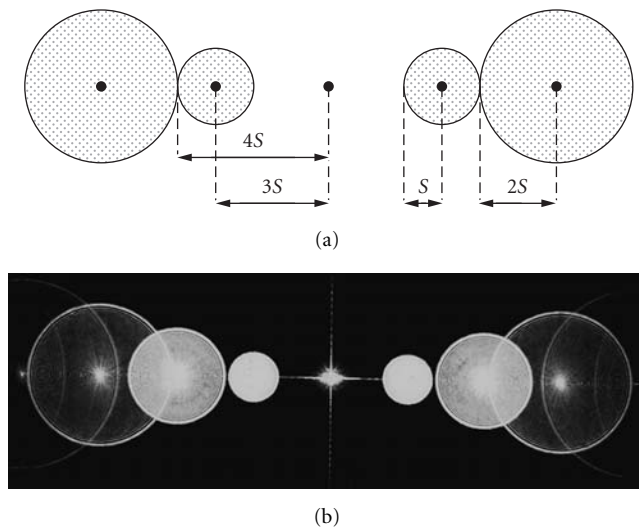


FIGURE 2 Diffracted orders in Fourier plane of CGH: (a) drawing and (b) photograph.

deposited on an optical-quality glass plate and currently produce the highest-quality CGHs. Typical e-beam recorders will write areas larger than $100\text{ mm} \times 100\text{ mm}$ with positional accuracies of less than 100 nm .¹¹

If needed, plotter distortion can be measured and calibrated out in the making of the hologram.^{12,13} The easiest way of determining plotter distortion is to draw straight lines and then treat this plot as a diffraction grating. If the computer-generated grating is illuminated with two plane waves, and the $-N$ order of beam 1 is interfaced with the $+N$ order of beam 2, the resulting interferogram gives us the plotter distortion. If the lines drawn by the plotter are spaced a distance Δx , a fringe error in the interferogram corresponds to a distortion error of $\Delta x/2N$ in the plot.

14.4 INTERFEROMETERS USING COMPUTER-GENERATED HOLOGRAMS

Many different experimental setups can be used for the holographic testing of optical elements. Figure 3 shows one common setup. The setup must be ray traced so the aberration in the hologram plane is known. While in theory there are many locations where the hologram can be placed, it is convenient to place the hologram in a plane conjugate to the asphere under test so the intensity across the image of the asphere is uniform. The longitudinal positional sensitivity for the hologram is reduced if the hologram is made in a region where the beams are collimated. Another advantage of this setup is that both the test and the reference beams pass through the hologram so errors resulting from hologram substrate thickness variations are eliminated without requiring the hologram be made on a good optical flat.

Another common setup for using a CGH to test aspheres is shown in Fig. 4.¹⁴ The largest advantage of this setup is that it works well with commercial Fizeau interferometers. The only addition to the commercial interferometer is a mount to hold the CGH between the transmission sphere and the optics under test. Since the light is diffracted by the CGH twice, the CGH must be a phase

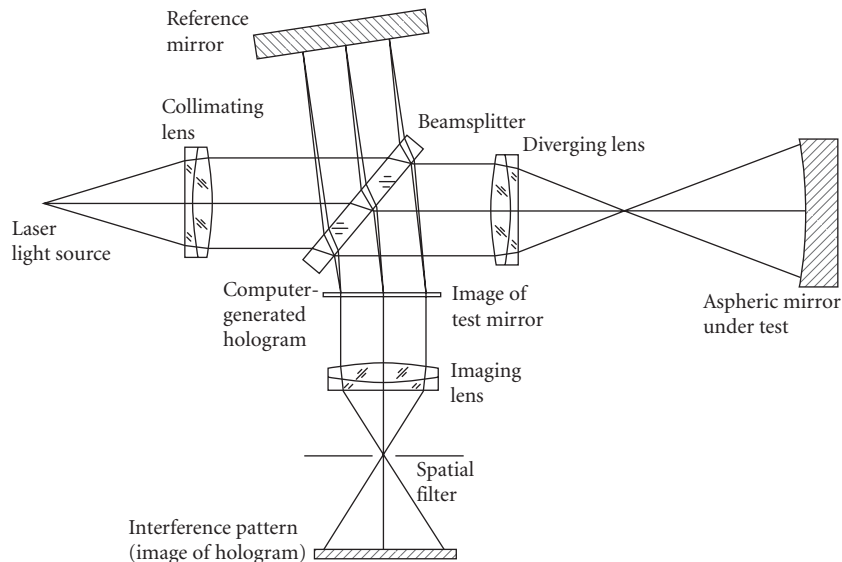


FIGURE 3 Interferometer setup using CGHs to test aspheric wavefronts.

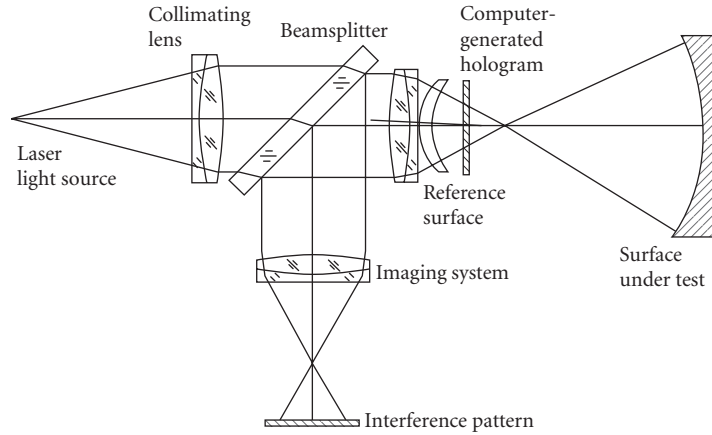


FIGURE 4 Use of CGH with Fizeau interferometer.

hologram so the diffraction efficiency is good, and since only the test beam is transmitted through the CGH, the substrate must either be high quality or thickness variations in the substrate must be measured and subtracted from the test results.

Figure 5 shows a setup for testing convex surfaces. In this case an on-axis CGH is used and the CGH is made on the concave reference surface.¹⁵ The light waves are perpendicular to the concave reference surface and then after diffraction they become perpendicular to the surface under test. The CGH pattern may be drawn exposing photoresist, ablating a metallic coating, or by creating a thin oxidation layer by heating a metal coating with a focused laser beam.¹⁶

CGHs can also be combined with partial null optics to test much more complicated aspherics than can be practically tested with either a CGH or null optics. This combination gives the real power of computer-generated holograms.¹⁷

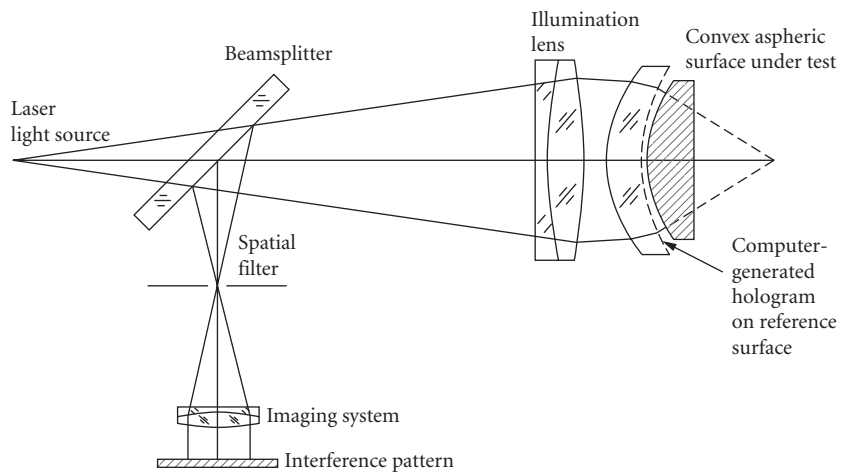


FIGURE 5 Using CGH to test convex surface.

14.5 ACCURACY LIMITATIONS

The largest source of error is the error due to plotter distortion as discussed previously. The other large sources of error are improper positioning of the hologram in the interferometer, and incorrect hologram size.

Any translation or rotation of the hologram produces error.² If the hologram is made conjugate to the exit pupil of the master optical system, the exit pupil of the system under test must coincide with the hologram. If the test wavefront in the hologram plane is described by the function $\phi(x, y)$, a displacement of the hologram a distance Δx in the x direction produces an error

$$\Delta\phi(x, y) \approx \frac{\partial\phi(x, y)}{\partial x} \Delta x \quad (1)$$

where $\partial\phi/\partial x$ is the slope of the wavefront in the x direction. Similarly, for a wavefront described by $\phi(r, \theta)$, the rotational error $\Delta\theta$ is given by

$$\Delta\phi(r, \theta) \approx \frac{\partial\phi(r, \theta)}{\partial\theta} \Delta\theta \quad (2)$$

Another source of error is incorrect hologram size. If the aberrated test wavefront in the plane of the hologram is given by $\phi(r, \theta)$, a hologram of incorrect size will be given by $\phi(r/M, \theta)$, where M is a magnification factor. The error due to incorrect hologram size will be given by the difference $\phi(r/M, \theta) - \phi(r, \theta)$, and can be written in terms of a Taylor expansion as

$$\begin{aligned} \phi\left(\frac{r}{M}, \theta\right) - \phi(r, \theta) &= \phi\left[r + \left(\frac{1}{M} - 1\right)r, \theta\right] - \phi(r, \theta) \\ &= \left[\frac{\partial\phi(r, \theta)}{\partial r}\right] \left(\frac{1}{M} - 1\right)r + \dots \end{aligned} \quad (3)$$

where terms higher than first order can be neglected if M is sufficiently close to 1, and a small region is examined. Note that this error is similar to a radial shear. When the CGH is plotted, alignment aids, which can help in obtaining the proper hologram size, can be drawn on the hologram plot. Figure 6 shows a CGH where the hologram is made in the center of the substrate and alignment aids are placed on the outer portion of the CGH.¹¹ Not only can the alignment aids help in putting the CGH in the proper position, but they can be used to help position the optics being tested. Figure 7

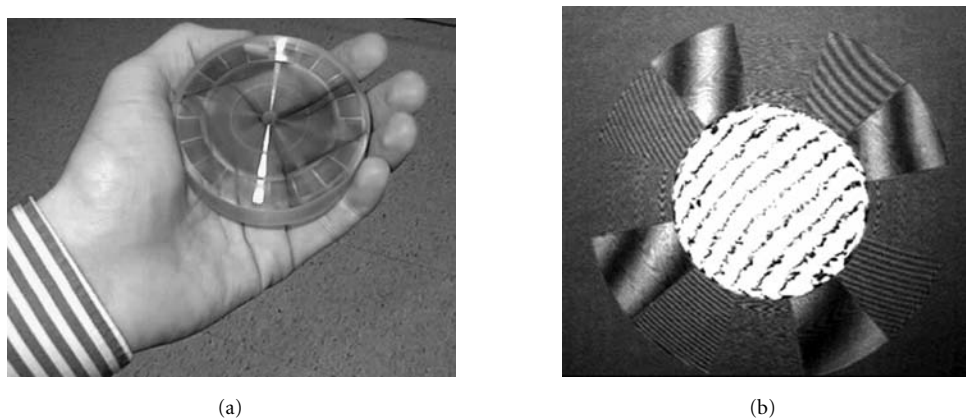


FIGURE 6 Use of CGH for alignment: (a) note structure in CGH and (b) interferogram produced with this CGH.

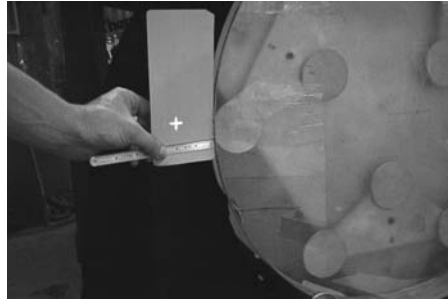


FIGURE 7 Use of crosshair produced by CGH to aid in the alignment of an off-axis parabola mirror.

shows a crosshair produced by a CGH that aids in the alignment of an off-axis parabolic mirror. The same CGH used to produce the crosshair produces the aspheric wavefront required to provide a null test of the off-axis parabola.¹¹

14.6 EXPERIMENTAL RESULTS

Figure 8 shows the results of using the setup shown in Fig. 3 to measure a 10-cm-diameter F/2 parabola using a CGH generated with an e-beam recorder. The fringes obtained in a Twyman-Green interferometer using a helium-neon source without the CGH present are shown in Fig. 8*a*. After the CGH is placed in the interferometer, a much less complicated interferogram is obtained as shown in Fig. 8*b*. The CGH corrects for about 80 fringes of spherical aberration, and makes the test much easier to perform.

To illustrate the potential of a combined CGH/null-lens test, results for a CGH/null-lens test of the primary mirror of an eccentric Cassegrain system with a departure of approximately 455 waves (at 514.5 nm) and a maximum slope of approximately 1500 waves per radius are shown.¹⁷ The mirror was a 69-cm diameter off-axis segment whose center lies 81 cm from the axis of symmetry of the parent aspheric surface. The null optics was a Maksutov sphere (as illustrated in Fig. 9), which reduces the departure and slope of the aspheric wavefront from 910 to 45 waves, and 300 to 70 waves per radius, respectively. A hologram was then used to remove the remaining asphericity.

Figure 10*a* shows interferograms of the mirror under test obtained using the CGH Maksutov test. Figure 10*b* shows the results when the same test was performed using a rather expensive refractive

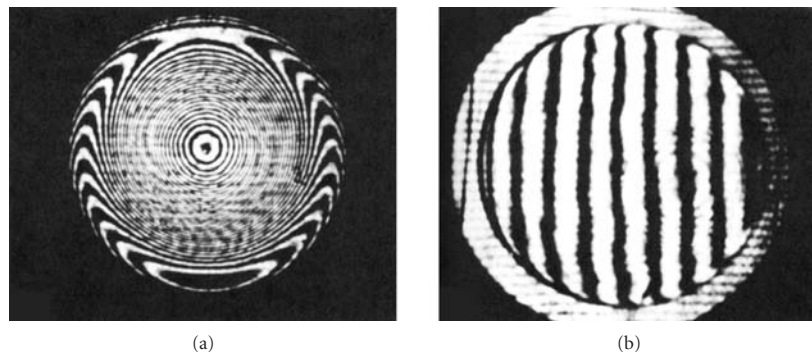


FIGURE 8 Results obtained testing a 10-cm-diameter F/2 parabola: (a) without using CGH and (b) using CGH made using an e-beam recorder.

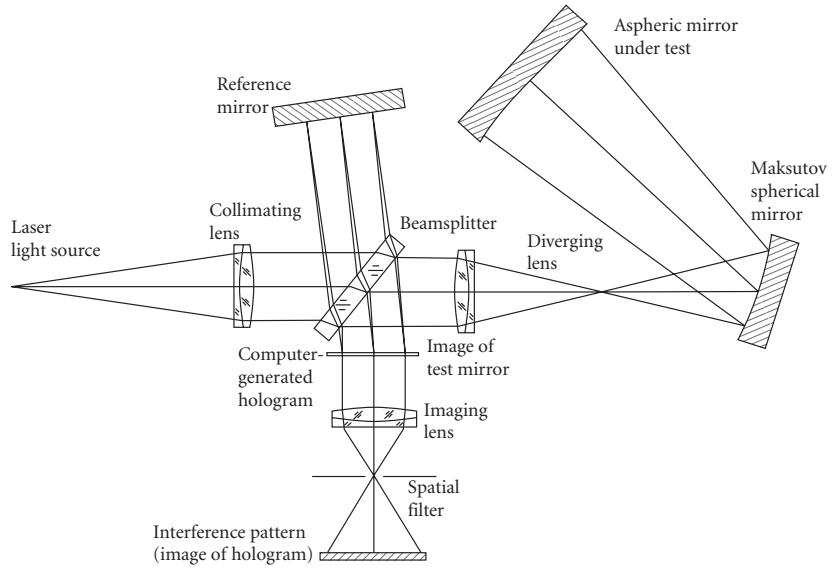


FIGURE 9 Setup to test the primary mirror of a Cassegrain telescope using a Maksutov sphere as a partial null and a CGH.

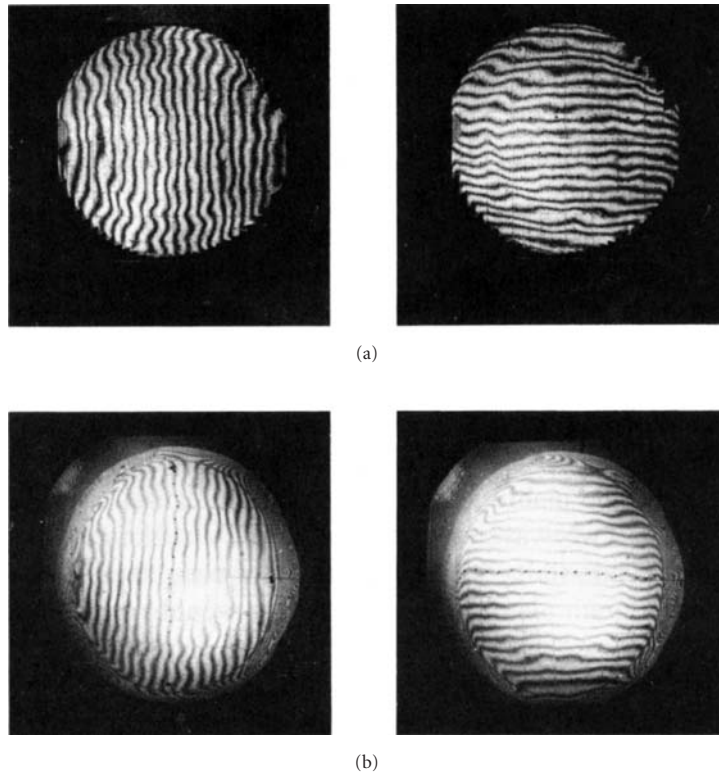


FIGURE 10 Results obtained using Fig. 9: (a) CGH-Maksutov test ($\lambda = 514.5 \text{ nm}$) and (b) using null lens ($\lambda = 632.8 \text{ nm}$).

null lens. When allowance is made for the fact that the interferogram obtained with the null lens has much more distortion than the CGH Maksutov interferogram, and for the difference in sensitivity ($\lambda = 632.8$ nm for the null-lens test and 514.5 nm for the CGH-Maksutov test), the results for the two tests are seen to be very similar. The “hills” and “valleys” on the mirror surface appear the same for both tests, as expected. The peak-to-valley surface error measured using the null lens was 0.46 waves (632.8 nm), while for the CGH-Maksutov test it was 0.39 waves (514.5 nm). The rms surface error measured was 0.06 waves (632.8 nm) for the null lens, while the CGH Maksutov test gave 0.07 wave (514.5 nm). These results certainly demonstrate that expensive null optics can be replaced by a combination of relatively inexpensive null optics and a CGH.

14.7 DISCUSSION

The difficult problem of testing aspheric surfaces, which are becoming increasingly popular in optical design, is made easier by the use of CGHs. The technology has reached the point that commercial interferometers using computer-generated holograms are now available. The main problem with testing aspheric optical elements is reducing the aberration sufficiently to ensure that light gets back through the interferometer. Combinations of simple null optics with a CGH to perform a test enable the measurement of a wide variety of optical surfaces. The making and use of a CGH are analogous to using an interferometer setup that yields a large number of interference fringes, and measuring the interferogram at a large number of data points. Difficulties involved in recording and analyzing a high-density interferogram and making a CGH are very similar. In both cases, a large number of data points are necessary, and the interferometer must be ray traced so that the aberrations due to the interferometer are well known. The advantage of the CGH technique is that once the CGH is made, it can be used for testing a single piece of optics many times or for testing several identical optical components. Additional alignment aids can be placed on the CGH to aid in the alignment of the CGH and the optics under test.

14.8 REFERENCES

1. A. J. MacGovern and J. C. Wyant, “Computer Generated Holograms for Testing Optical Elements,” *Appl. Opt.* **10**(3):619–624 (1971).
2. J. C. Wyant and V. P. Bennett, “Using Computer Generated Holograms to Test Aspheric Wavefronts,” *Appl. Opt.* **11**(12):2833–2839 (1972).
3. A. F. Fercher and M. Kriese, “Binäre Synthetische Hologramme zur Prüfung Asphärischer Optischer Elemente,” *Optik* **35**(2):168–179 (1972).
4. Y. Ichioka and A. W. Lohmann, “Interferometric Testing of Large Optical Components with Circular Computer Holograms,” *Appl. Opt.* **11**(11):2597–2602 (1972).
5. J. Schwider and R. Burrow, “The Testing of Aspherics by Means of Rotational-Symmetric Synthetic Holograms,” *Optica Applicata* **6**:83 (1976).
6. T. Yatagai and H. Saito, “Interferometric Testing with Computer-Generated Holograms: Aberration Balancing Method and Error Analysis,” *Appl. Opt.* **17**(4):558–565 (1978).
7. J. Schwider, R. Burrow, and J. Grzanna, “CGH—Testing of Rotational Symmetric Aspheric in Compensated Interferometers,” *Optica Applicata* **9**:39 (1979).
8. C.S. Pruss, S. Reichelt, H.J. Tiziani, and W. Osten, “Computer-Generated Holograms in Interferometric Testing,” *Opt. Eng.* **43**:2534–2540 (2004).
9. J. W. Goodman, *Introduction to Fourier Optics*, 3d ed. Roberts & Company: Greenwood Village, Colorado, 2004.
10. Y. C. Chang and J. H. Burge, “Error Analysis for CGH Optical Testing,” *Proc. SPIE* **3872**:358–366 (1999).
11. J. H. Burge, R. Zehnder, and Chunyu Zhao, “Optical Alignment with Computer Generated Holograms,” *Proc. SPIE* **6676**:66760C (2007).

12. J. C. Wyant, P. K. O'Neill, and A. J. MacGovern, "Interferometric Method of Measuring Plotter Distortion," *Appl. Opt.* **13**(7):1549–1551 (1974).
13. A. F. Fercher, "Computer Generated Holograms for Testing Optical Elements: Error Analysis and Error Compensation," *Optica Applicata* **23**(5):347–365 (1976).
14. H. J. Tiziani, J. S. Reichlet, C. Pruss, M. Rocktachel, and U. Hofbauer, "Testing of Aspheric Surfaces," *Proc. SPIE*, **4440**:109–119 (2001).
15. J. H. Burge and D. S. Anderson, "Full-Aperture Interferometric Test of Convex Secondary Mirrors Using Holographic Test Plates," *Proc. SPIE* **2199**:181–192 (1994).
16. J. H. Burge, M. J. Fehniger, and G. C. Cole, "Demonstration of Accuracy and Flexibility of Using CGH Test Plates for Measuring Aspheric Surfaces," *Proc. SPIE* **3134**:379–389 (1997).
17. J. C. Wyant and P. K. O'Neill, "Computer Generated Hologram; Null Lens Test of Aspheric Wavefronts," *Appl. Opt.* **13**(12):2762–2765 (1974).

PART

4

SOURCES

This page intentionally left blank.

ARTIFICIAL SOURCES

Anthony LaRocca*

*General Dynamics
Advanced Information Systems
Ypsilanti, Michigan*

15.1 GLOSSARY

λ	wavelength
$d\lambda$	differential wavelength
M_λ	spectral radiant exitance
T	absolute temperature
c_1	first radiation constant
c_2	second radiation constant
h	Planck's constant
c	velocity of light
k	Boltzmann constant
λ_{\max}	wavelength at peak of radiant exitance
K	factor
R	radius of the interior surface of the cavity
r	radius of the aperture
S	interior surface area of the cavity
s	aperture area
\mathcal{E}	emissivity
\mathcal{E}_0	uncorrected emissivity

15.2 INTRODUCTION

Whereas most of the sources described in this chapter can be used for any purpose for which one can justify their use, the emphasis is on the production of the appropriate radiation for the calibration of measurement instrumentation. This implies that the basis for their use is supported by their

*Retired

traceabilities to calibrated standards of radiation from an internationally known and respected standards laboratory such as the National Institute of Standards and Technology (NIST) in the United States or, say, the National Physical Laboratory (NPL) in the United Kingdom. Because calibration implies a high degree of accuracy, the chapter initially contains a short exposition on the so-called Planckian, or blackbody radiation standard, and the equation which describes the Planck radiation.

This chapter deals with artificial sources of radiation as subdivided into two classes: laboratory and field sources. Much of the information on commercial sources is taken from a chapter previously written by the author.¹ Where it was feasible, similar information here has been updated. When vendors failed to comply to requests for information, the older data were retained to maintain completeness, but the reader should be aware that some sources cited here may no longer exist, or perhaps may not exist in the specification presented. Normally, laboratory sources are used in some standard capacity and field sources are used as targets. Both varieties appear to be limitless. Only laboratory sources are covered here.

The sources in this chapter were chosen arbitrarily, often depending on manufacturer response to requests for information. The purpose of this chapter is to consolidate much of this information to assist the optical-systems designer in making reasonable choices. To attempt to include the hundreds of types of lasers, however, and the thousands of varieties, would be useless for several reasons, but particularly because they change often.

Complementing the material in the following chapter on Lasers, a fairly comprehensive source of information on lasers can be found in the *CRC Handbook of Laser Science and Technology*, Supplement 1: published by the CRC Press, Inc., 2000 Corporate Blvd., N.W., Boca Raton, Florida, 33431. Of course, the literature is laden with material on lasers, including the chapter on Lasers in this *Handbook*, and the reader would be wise to consult the Internet from which compilations such as the one cited above or a host of others can be obtained from companies like Amazon.com, or, better yet, by accessing a literature-rich source such as Google.

Regarding the selection of a source, Worthing² suggests that one ask the following questions:

1. Does it supply energy at such a rate or in such an amount as to make measurements possible?
2. Does it yield an irradiation that is generally constant or that may be varied with time as desired?
3. Is it reproducible?
4. Does it yield irradiations of the desired magnitudes over areas of the desired extent?
5. Has it the desired spectral distribution?
6. Has it the necessary operating life?
7. Has it sufficient ruggedness for the proposed problem?
8. Is it sufficiently easy to obtain and replace, or is its purchase price or its construction cost reasonable?

15.3 RADIATION LAW

All of the radiation sources described in this chapter span the region of the electromagnetic spectrum mainly from the visible region (starting from about 400) through the infrared (around 400 μm and beyond). Given that they have a demonstrable temperature, they relate in their own peculiar ways, depending on material properties, to the radiation called “blackbody” radiation, which is described by the Planck radiation law. Many attempts were made in the latter part of the nineteenth century and the early twentieth century to describe blackbody radiation mathematically, all doomed to failure before the recognition of quantum concepts, in particular, by Max Planck. Any attempt to describe the mechanisms surrounding the Planck theory would be superfluous here. Suffice it to say that the basis for the theory can be explained from an examination of the experimental curve shown in Fig. 1 (borrowed from Richtmeyer and Kennard³) determined from the examination of the radiation from a “blackbody” at several different temperatures. By plotting the points one concludes, as shown on the graph, that the spectral radiant exitance, M_λ , is equal to the product of the negative fifth power of λ times some function of the product, λT , where λ is the wavelength (in micrometers) of the radiation and T is the absolute temperature of the radiator (in Kelvins). Confirmation of this fact is shown in Table 1. The radiant exitance values in the table were calculated

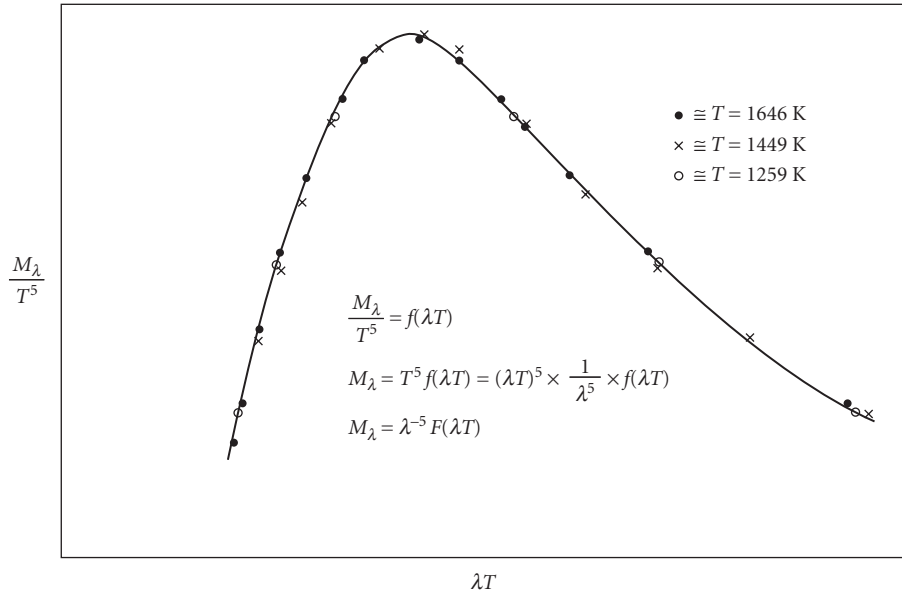


FIGURE 1 Experimental verification of the blackbody displacement law.

TABLE 1 Numerical Values to Support Fig. 1

T	λ	λT	M_λ	M_λ/T^5	
1,646	1.761	2,899	4.95	4.09×10^{-16}	
	1	1,646	1.91	1.56×10^{-16}	
	0.8	1,317	0.653	5.40×10^{-17}	
	2	3,291	4.77	3.90×10^{-16}	
	3	4,938	2.81	2.33×10^{-16}	
	5	8,230	0.803	6.65×10^{-17}	
	7	11,522	0.285	2.36×10^{-17}	
	10	16,460	0.085	7.04×10^{-18}	
	1,449	2	2,898	2.62	4.10×10^{-16}
		1.14	1,646	1.02	1.60×10^{-16}
0.909		1,317	0.346	5.42×10^{-17}	
2.27		3,291	2.52	3.95×10^{-16}	
3.41		4,938	1.49	2.32×10^{-16}	
5.68		8,230	0.425	6.65×10^{-17}	
7.95		11,522	0.151	2.36×10^{-17}	
11.36		16,460	0.045	7.06×10^{-18}	
1,259	2.3	2,896	1.3	4.11×10^{-16}	
	1.31	1,646	0.502	1.59×10^{-16}	
	1.046	1,317	0.171	5.41×10^{-17}	
	2.61	3,291	1.25	3.95×10^{-16}	
	3.91	4,938	0.741	2.34×10^{-16}	
	6.54	8,230	0.21	6.64×10^{-17}	
	9.15	11,522	0.0749	2.37×10^{-17}	
	13.07	16,460	0.0223	7.05×10^{-18}	

T = deg Kelvin; λ = μm ; M_λ = radiant exitance; $\text{W}\cdot\text{cm}^{-2}\cdot\text{ster}^{-1}\cdot\mu\text{m}^{-1}$

using the Infrared Radiance Calculator created by the author and found by choosing the term “Calculators” from the Military Sensing Information Analysis Center (SENSIAC) in a search of the Internet under www.sensiac.gatech.edu.

Postulating the quantum nature of the radiation, Planck, in a clever demonstration of the entropies resulting from small and large values of λT , was able to establish an expression for blackbody radiation as

$$M_{\lambda}(\lambda)d\lambda = c_1\lambda^{-5}(e^{c_2/\lambda T} - 1)^{-1}d\lambda$$

where

$$c_1 = 2\pi hc^2 = 3.7413 \times 10^4 \text{ W-cm}^{-2}\text{-}\mu\text{m}^4 \text{ (first radiation constant)}$$

and

$$c_2 = hc/k = 14388 \mu\text{m-K (second radiation constant)}$$

h = Planck's constant = 6.6252×10^{-34} W-s²

k = Boltzmann constant = 1.38042×10^{-23} W-s-K⁻¹

c = Velocity of light = 2.99793×10^{10} -s⁻¹

He later established the same equation from first principles.

When the Planck function is plotted on log-log paper the graph of Fig. 2 results. The special feature of this type of plot is that, regardless of the temperature of the blackbody, the shape of the curve is constant. It merely moves up and to the left (i.e., toward shorter wavelengths) as the temperature increases. The straight line, with a slope of -5 , drawn through the set of curves of Fig. 2, depicts the

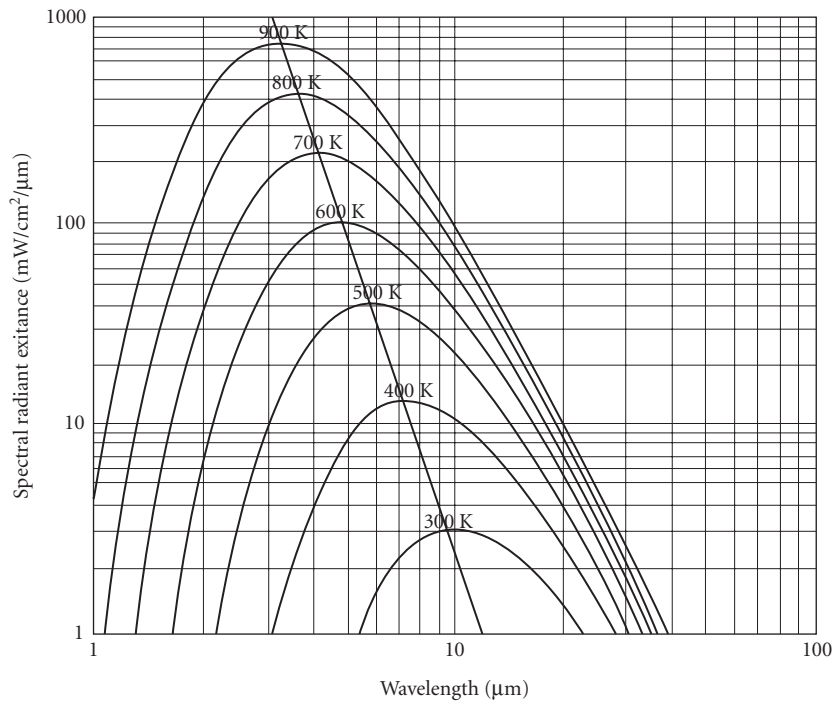


FIGURE 2 Spectral radiant exitance versus wavelength.

wavelength at which each one peaks, resulting in what is known as the Wein displacement law given for a specific temperature by

$$\lambda_{\max} T = 2897.9 \mu\text{m} - \text{K}^{-1}$$

Thus the peak of any Planck curve can be determined, given the temperature of the blackbody.

15.4 LABORATORY SOURCES

Standard Sources

The reader may be interested in the exposition by Quinn⁴ on the calculation of the emissivity of cylindrical cavities in which the method of DeVos⁵ is used. In a more recent paper Irani⁶ refers to the method of Gouffé⁷ for the construction of blackbody calibration sources. Quinn states that for certain constructions there are errors in the method of Gouffé. However, for a well-constructed source the shape of the construction is least at fault, since any heat-resistant material with a reasonably high surface emissivity will produce a resultant emissivity of better than 0.99. However, the accuracy of the value of the radiation for a given temperature depends not only on the emissivity but on generally high numerical powers of the temperature especially for high-temperature blackbodies. Therefore, very small variations of temperature over the inner surface of the source can cause relatively large errors in the radiation accuracy. Thus, great caution is used in creating a uniform temperature, resulting in the use of the fixed-point temperatures of various metals for the most basic and accurate calibration standards.

Blackbody Cavity Theory Radiation levels can be standardized by the use of a source that will emit a quantity of radiation that is both reproducible and predictable. Cavity configurations can be produced to yield radiation theoretically sufficiently close to Planckian that it is necessary only to determine what the imprecision is. Several theories have been expounded over the years to calculate the quality of a blackbody simulator.*

*The Method of Gouffé.*⁷ For the total emissivity of the cavity forming a blackbody (disregarding temperature variations) Gouffé gives

$$\epsilon_0 = \epsilon'_0(1 + K) \quad (1)$$

where

$$\epsilon'_0 = \frac{\epsilon}{\epsilon \left(1 - \frac{s}{S}\right) + \frac{s}{S}} \quad (2)$$

and $K = (1 - \epsilon) \left[\frac{s}{S} - \frac{s}{S_0} \right]$, and is always nearly zero—it can be either positive or negative.

ϵ = emissivity of materials forming the blackbody surface

s = area of aperture

S = area of interior surface

S_0 = surface of a sphere of the same depth as the cavity in the direction normal to the aperture

Figure 3 is a graph for determining the emissivities of cavities with simple geometric shapes. In the lower section, the value of the ratio s/S is given as a function of the ratio $1/r$. (Note the scale

*Generically used to describe those sources designed to produce radiation that is nearly Planckian.

change at the value for $1/r = 5$.) The values of ϵ'_0 is found by reading up from this value of the intrinsic emissivity of the cavity material. The emissivity of the cavity is found by multiplying ϵ'_0 by the factor $(1 + K)$.

When the aperture diameter is smaller than the interior diameter of the cylindrical cavity, or the base diameter of a conical cavity, it is necessary to multiply the value of s/S determined from the graph by $(r/R)^2$, which is the ratio of the squares of the aperture and cavity radii (Fig. 3).

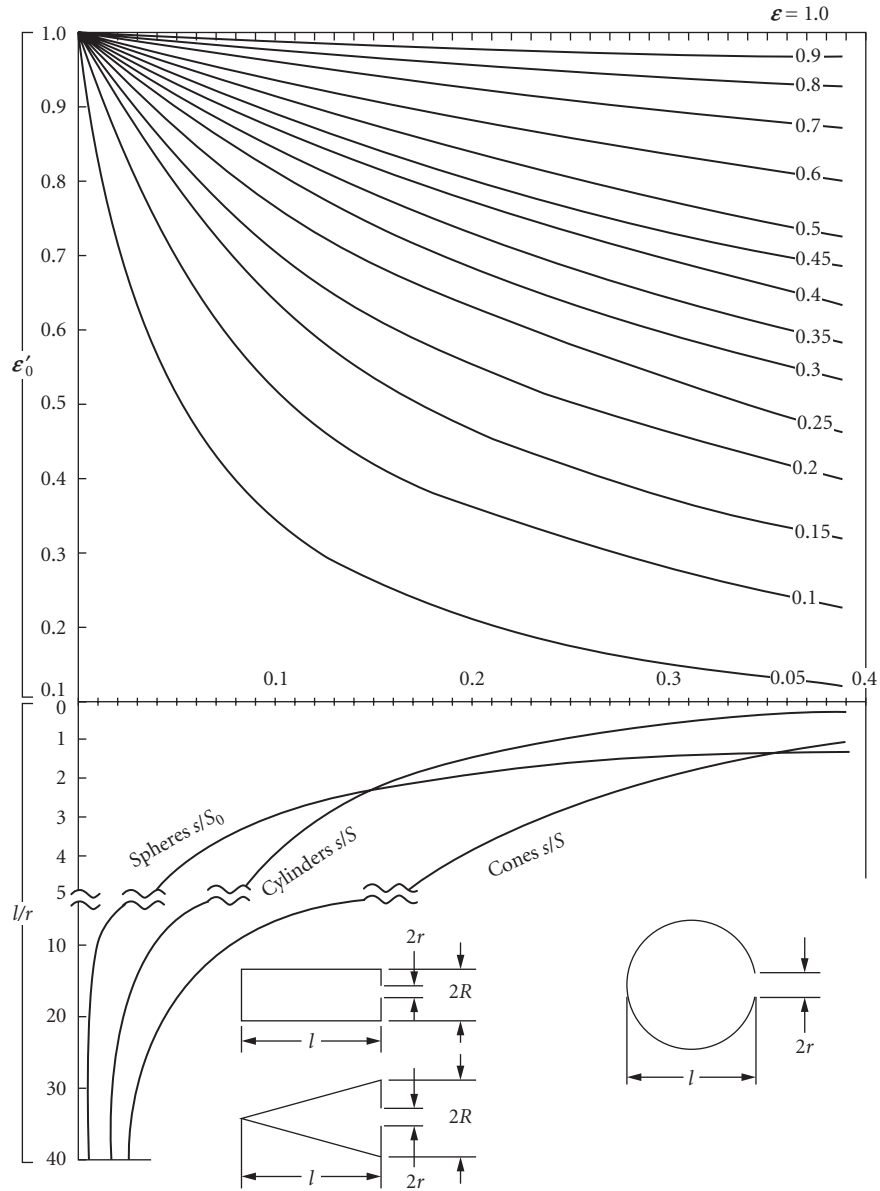


FIGURE 3 Emissivities of conical, spherical, and cylindrical cavities.

It is important to be aware of the effect of temperature gradients in a cavity. This factor is the most important in determining the quality of a blackbody, since it is not very difficult to achieve emissivities as near to unity as desired.

Manufacturers of blackbody simulators strive to achieve uniform heating of the cavity because it is only under this condition that the radiation is Planckian. The ultimate determination of a radiator that is to be used as the standard is the quality of the radiation that it emits.

A recent investigation on comparison of IR radiators is presented by Leupin et al.⁸ There has been, incidentally, a division historically between the standards of photometry and those used to establish thermal radiation and the thermodynamic temperature scale. Thus, in photometry the standard has changed from the use of candles, the Carcel lamp, the Harcourt pentane lamp, and the Hefner lamp⁹ to more modern radiators.

Baseline Standard of Radiation Although there is no internationally accepted standard of radiation, the National Institute of Standards and Technology (NIST) uses as its substitute standard the goldpoint blackbody (see Fig. 4),¹⁰ which fixes one point on the international temperature scale, now reported to be 1337.33 ± 0.34 K. Starting from this point, NIST is able to transfer fixed radiation values to working standards of radiation through an accurately constructed variable-temperature radiator as shown in Fig. 5.¹¹

The goldpoint blackbody is shown mainly for information. It is quite feasible to build a replica of the variable-temperature radiator, especially in the laboratory equipped to do fundamental radiation measurements.

Working Standards of Radiation For the calibration of instruments in the ordinary laboratory, the user is likely to use a source which is traceable to NIST, and generally supplied by NIST or one of the recognized vendors of calibrated sources, mainly in the form of a heated filament, a gaseous arc enclosed in an envelope of glass or quartz (or fused silica), or in glass with a quartz or sapphire window.

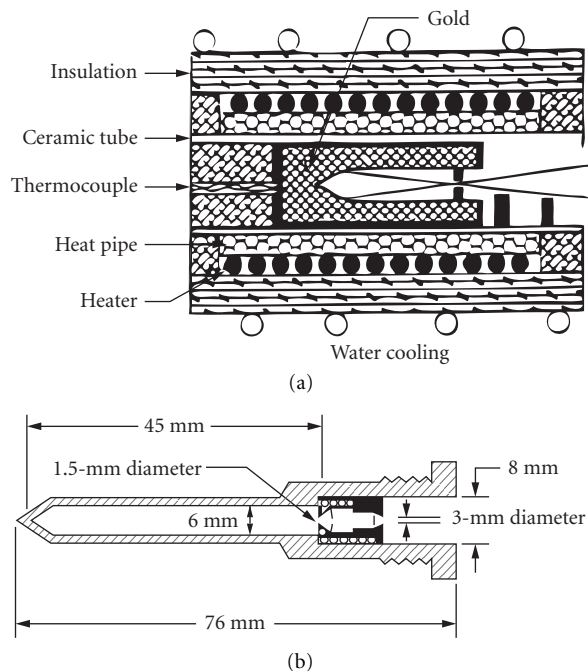


FIGURE 4 (a) Cross section of heat-pipe blackbody furnace. (b) Blackbody inner cavity dimensions.

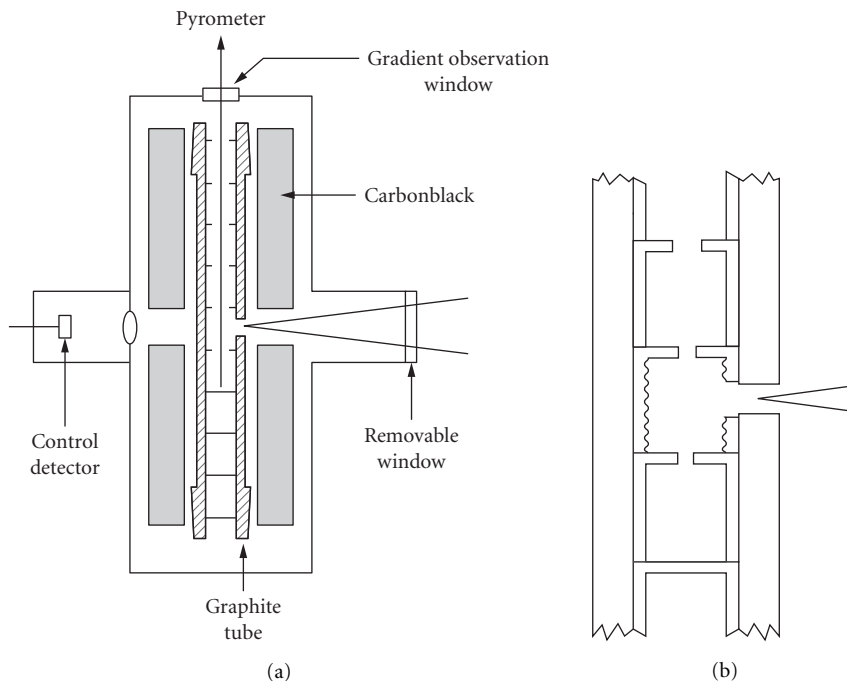


FIGURE 5 (a) Variable-temperature blackbody schematic. (b) Central section of variable-temperature blackbody.

Any source whose radiation deviates from that described by Planck's law is nonblackbody. Even the sources previously described are not strictly blackbodies, but can come as close as the user desires within the constraints of bulk and price. Any other source has an emissivity less than unity and can, and usually does, have a highly variable spectral emissivity. The lamps used by NIST, for example (see the following), fit into this category, but they differ in one large respect. They are transfer standards which have been carefully determined to emit specified radiation within certain specific spectral regions.

The following discussion of these types of sources is reproduced (in some cases with slight modifications), with permission, from the NIST Special Publication 250.¹² For specific details of calibration, and for the exact source designations, the user should contact NIST at

*U.S. Department of Commerce
National Institute of Standards and Technology
Office of Physical Measurement Services
Rm. B362, Physics Bldg.
Gaithersburg, MD 20899
Photometric Standards*

The following text on working standards is left untouched from the presentation of the earlier edition of the *Handbook* because there do not appear to be significant changes since the publication of that material. However, to the extent that information on the Internet is current, a reasonable complement to the information published in this chapter would be a search of the Internet at the location designated www.physics.nist.gov. The reader will find numerous features of the National Institute of Science and Technology from which to choose the service or other information required.

1. Sources/Lamps

Luminous Intensity Standard (100-W Frosted Tungsten Lamp, 90 cd)
 Luminous Intensity Standard (100-W Frosted Tungsten Lamp, color temp., 2700 K)
 Luminous Intensity Standard (100-W Frosted Tungsten Lamp, color temp., 2856 K)
 Luminous Intensity Standard (500-W Frosted Tungsten Lamp, 700 cd)
 Luminous Intensity Standard (1000-W Frosted Tungsten Lamp, 1400 cd)
 Luminous Intensity Standard (1000-W Frosted Tungsten Lamp, color temp., 2856 K)
 Luminous Flux Standard (25-W Vacuum Lamp about 270 lm)
 Luminous Flux Standard (60-W Gas-filled Lamp about 870 lm)
 Luminous Flux Standard (100-W Gas-Filled Lamp about 1600 lm)
 Luminous Flux Standard (200-W Gas-Filled Lamp about 3300 lm)
 Luminous Flux Standard (500-W Gas-Filled Lamp about 10,000 lm)
 Luminous Flux Standard (Miniature Lamps 7 sizes 6 to 400 lm)
 Airway Beacon Lamps for Color Temperature (500-W, 1 point in range, 2000 to 3000 K)

2. General Information

Calibration services provide access to the photometric scales realized and maintained at NIST. Lamp standards of luminous intensity, luminous flux, and color temperature, as described next, are calibrated on a routine basis.

a. Luminous Intensity Standards

Luminous intensity standard lamps supplied by NIST [100-W (90–140 cd), 500-W (approximately 700 cd), and 1000-W (approximately 1400 cd) tungsten filament lamps with C-13B filaments in inside-frosted bulbs and having medium bipost bases] are calibrated at either a set current or a specified color temperature in the range 2700 to 3000 K. Approximate 3-sigma uncertainties are 1 percent relative to the SI unit of luminous intensity and 0.8 percent relative to NIST standards.

b. Luminous Flux Standards

Vacuum tungsten lamps of 25 W and 60-, 100-, 200-, and 500-W gas-filled tungsten lamps that are submitted by customers are calibrated. Lamps must be base-up burning and rated at 120 V. Approximate 3-sigma uncertainties are 1.4 percent relative to SI units and 1.2 percent relative to NIST standards. Luminous flux standards for miniature lamps producing 6 to 400 lm are calibrated with uncertainties of about 2 percent.

c. Airway Beacon Lamps

Color temperature standard lamps supplied by NIST (airway beacon 500-W medium bipost lamps) are calibrated for color temperature in the range 2000 to 3000 K with 3-sigma uncertainties ranging from 10 to 15°.

IR Radiometric Standards

General Information

a. Spectral Radiance Ribbon Filament Lamps

These spectral radiance standards are supplied by NIST. Tungsten, ribbon filament lamps (30A/T24/13) are provided as lamp standards of spectral radiance. The lamps are calibrated at 34 wavelengths from 225 to 2400 nm, with a target area 0.6 mm wide by 0.8 mm high. Radiance temperature ranges from 2650 K at 225 nm and 2475 K at 650 nm to 1610 K at 2400 nm, with corresponding uncertainties of 2, 0.6, and 0.4 percent. For spectral radiance lamps, errors are stated as the quadrature sum of individual uncertainties at the three standard deviation level.

Figure 6 summarizes the measurement uncertainty for NIST spectral radiance calibrations.

b. Spectral Irradiance Lamps

These spectral irradiance standards are supplied by NIST. Lamp standards of spectral irradiance are provided in two forms. Tungsten filament, 1000 W quartz halogen-type FEL lamps are calibrated at 31 wavelengths in the range 250 to 2400 nm. At the working distance of 50 cm, the

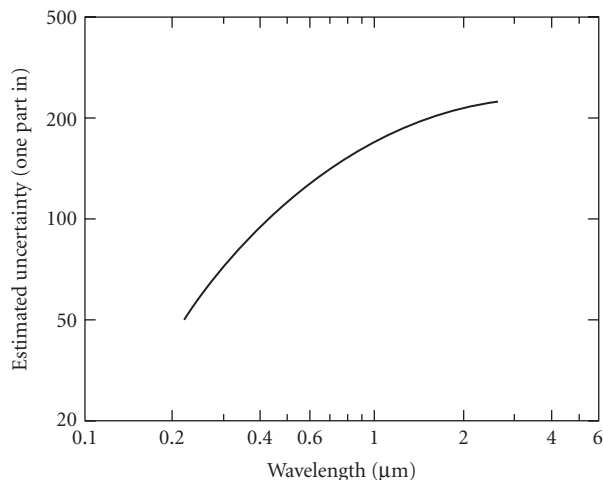


FIGURE 6 Uncertainties for NIST spectral radiance calibrations.

lamps produce 0.2 W/cm²/cm at 250 nm, 220 W/cm²/cm at 900 nm, 115 W/cm²/cm at 1600 nm, and 40 W/cm²/cm at 2400 nm, with corresponding uncertainties of 2.2, 1.3, 1.9, and 6.5 percent. For spectral irradiance lamps, errors are stated as the quadrature sum of individual uncertainties at the three standard deviation level. Deuterium lamp standards of spectral irradiance are also provided and are calibrated at 16 wavelengths from 200 to 350 nm. At the working distance of 50 cm, the spectral irradiance produced by the lamp ranges from about 0.5 W/cm²/cm at 200 nm and 0.3 W/cm²/cm at 250 nm to 0.07 W/cm²/cm at 350 nm. The deuterium lamps are intended primarily for the spectral region 200 to 250 nm. The approximate uncertainty relative to SI units is 7.5 percent at 200 nm and 5 percent at 250 nm. The approximate uncertainty in relative spectral distribution is 3 percent. It is strongly recommended that the deuterium standards be compared to an FEL tungsten standard over the range 250 to 300 nm each time the deuterium lamp is lighted to take advantage of the accuracy of the relative spectral distribution.

Figure 7 summarizes the measurement uncertainty for NIST spectral irradiance calibrations of type FEL lamps.

Radiometric Sources in the Far Ultraviolet

1. Sources

Spectral Irradiance Standard, Argon Mini-Arc (140 to 330 nm)

Spectral Radiance Standard, Argon Mini-Arc (115 to 330 nm)

Spectral Irradiance Standard, Deuterium Arc Lamp (165 to 200 nm)

2. General Information

a. Source Calibrations in the Ultraviolet

NIST maintains a collection of secondary sources such as argon maxi-arcs, argon mini-arcs, and deuterium arc lamps in the near and vacuum ultraviolet radiometric standards program to provide calibrations for user-supplied sources. The calibrations of these sources are traceable to a hydrogen arc whose radiance is calculable and which NIST maintains as a primary standard. The collection also includes tungsten strip lamps and tungsten halogen lamps whose calibrations are based on a blackbody rather than a hydrogen arc. Customer-supplied sources are calibrated in both radiance and irradiance by comparing them with NIST secondary standards.

Argon arcs are used to calibrate other sources in the wavelength range 115 to 330 nm for radiance and 140 to 330 nm for irradiance. The lower wavelength limit is determined in radiance by the cutoff of the magnesium fluoride windows used in the arcs, and in irradiance by the decrease

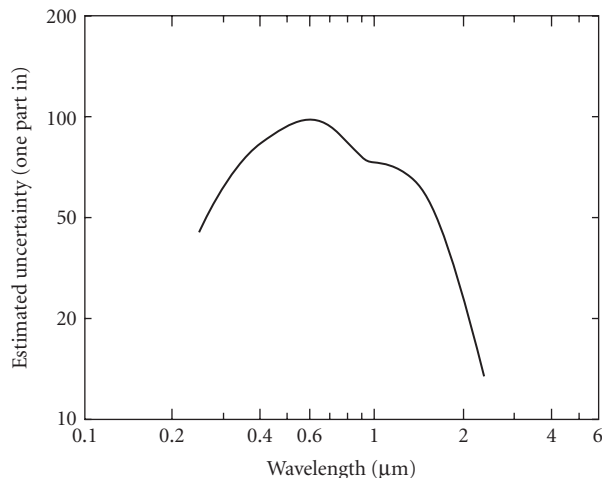


FIGURE 7 Uncertainties for NIST spectral irradiance calibrations of type FEL lamps.

in signal produced by the addition of a diffuser. Deuterium arc lamps are used in the range 165 to 200 nm, with the low wavelength cutoff due to the onset of blended molecular lines.

The high wavelength limit is the starting point of the range for the Radiometric Standards group. The tungsten lamps are used at 250 nm and above, since their signals are too weak at shorter wavelengths. It should be noted that the wavelength range of the NIST arcs partially overlap the range of tungsten lamps, thus providing an independent check on calibrations.

An argon mini-arc lamp supplied by the customer is calibrated for spectral irradiance at 10-nm intervals in the wavelength region 140 to 300 nm. Absolute values are obtained by comparison of the radiative output with laboratory standards of both spectral irradiance and spectral radiance. The spectral irradiance measurement is made at a distance of 50 cm from the field stop. Uncertainties are estimated to be less than ± 10 percent in the wavelength region 140 to 200 nm and within ± 5 percent in the wavelength region 200 to 330 nm. A measurement of the spectral transmission of the lamp window is included in order that the calibration be independent of possible window deterioration or damage. The uncertainties are taken to be two standard deviations.

The spectral radiance of argon mini-arc radiation sources is determined to within an uncertainty of less than 7 percent over the wavelength range 140 to 330 nm and 20 percent over the wavelength range 115 to 140 nm. The calibrated area of the 4-mm diameter radiation source is the central 0.3-mm diameter region. Typical values of the spectral radiance are: at 250 nm, $L(\lambda) = 30 \text{ mW/cm}^2/\text{nm/sr}$; and at 150 nm, $L(\lambda) = 3 \text{ mW/cm}^2/\text{nm/sr}$. The transmission of the demountable lamp window and that of an additional MgF_2 window are determined individually so that the user may check periodically for possible long-term variations.

The deuterium arc lamp is calibrated at 10 wavelengths from 165 to 200 nm, at a distance of 50 cm, at a spectral irradiance of about $0.5 \text{ W/cm}^2/\text{cm}$ at 165 nm, $0.5 \text{ W/cm}^2/\text{cm}$ at 170 nm, and $0.5 \text{ W/cm}^2/\text{cm}$ at 200 nm. The approximate uncertainty relative to SI units is estimated to be less than 10 percent. The lamp is normally supplied by NIST and requires 300 mA at about 100 V.

15.5 COMMERCIAL SOURCES

The commercial sources described here are derived from the 1995 edition of *The Handbook of Optics*, which were taken from catalogs available at the time and from the literature of the day and prior thereto, providing choices that have obviously been available for years. Evidently changes in

the makeup of these products are slower than those of other areas of technology, making it reasonable to retain the same sources in this chapter, mainly as examples of the types that are available. With the universality of the Internet, accessibility of information on various sources of radiation, far beyond what is attainable from a limited collection of company catalogs, is at one's fingertips. Thus, it is recommended that, in seeking information on various sources, one use the examples in the text as a reference to what can be found currently on the Internet. Experience demonstrates that, in many cases, due to the stability of the lamp industry, there will be few changes between what is found currently on the Internet and what appears in the text of this chapter.

Obviously the choice of a source is dependent on the application. Many, if not most of the sources described are multipurpose ones, although most of them have been selected specifically for scientific study, tailored in a way to produce an image amenable to different optical systems. For basic research it is usually essential to have a blackbody source, especially for infrared research, that is traceable to a Standards Laboratory, along with traceable secondary standards for calibrating research instrumentation. Many of the sources can be used to produce spectra, which are capable of calibrating spectral measuring instrumentation. Other sources are included mostly to provide the user with an array of choices.

Blackbody Simulators

Virtually any cavity can be used to produce radiation of high quality, but practicality limits the shapes to a few. The most popular shapes are cones and cylinders, the former being more popular. Spheres, combinations of shapes, and even flat-plate radiators are used occasionally. Blackbodies can be bought rather inexpensively, but there is a fairly direct correlation between cost and quality (i.e., the higher the cost the better the quality).

Few manufacturers specialize in blackbody construction. Some, whose products are specifically described here, have been specializing in blackbody construction for many years. Other companies of this description may be found, for example, in the latest *Lasers and Optronics Buying Guide*¹³ or the latest *Photonics Directory of Optical Industries*.¹⁴ These references are the latest as of the writing of this work. It is expected that they will continue in succeeding years.

A large selection of standard (or blackbody) radiators is offered by Electro-Optical Industries, Inc. (EOI), Santa Barbara, California.* Most blackbodies can be characterized as one of the following: primary, secondary, or working standard. The output of the primary must, of course, be checked with those standards retained at NIST. Figure 8 pictures an EOI blackbody and its controller. Figure 9 pictures a similar blackbody from Mikron, Inc. and its controller. All of the companies sell separate apertures (some of which are water cooled) for controlling the radiation output of the radiators. Another piece of auxiliary equipment which can be purchased is a multispeed chopper. It is impossible to cite all of the companies that sell these kinds of sources; therefore, the reader is referred to one of the buyers' guides already referenced for a relatively complete list. It is prudent to shop around for the source that suits one's own purpose.

Figure 10 demonstrates a less conventional working standard manufactured by EOI. Its grooves-and-honeycomb structure is designed to improve the absorptance of such a large and open structure. A coating with a good absorbing paint increases its absorptance further.

Incandescent Nongaseous Sources (Exclusive of High-Temperature Blackbodies)

Nernst Glower[†] The Nernst glower is usually constructed in the form of a cylindrical rod or tube from refractory materials (usually zirconia, yttria, beria, and thoria) in various sizes. Platinum leads at the ends of the tube conduct power to the glower from the source. Since the resistivity of

*Many of the sources in the text are portrayed using certain specific company products, only for the sake of demonstration. This does not necessarily imply an endorsement of these products by the author. The reader is encouraged in all cases to consult the *Photonics Directory of Optical Industries*¹⁴ or a similar directory for competitive products.

[†]Since Nernst glower is probably obsolete, this section is retained only for historical purpose.

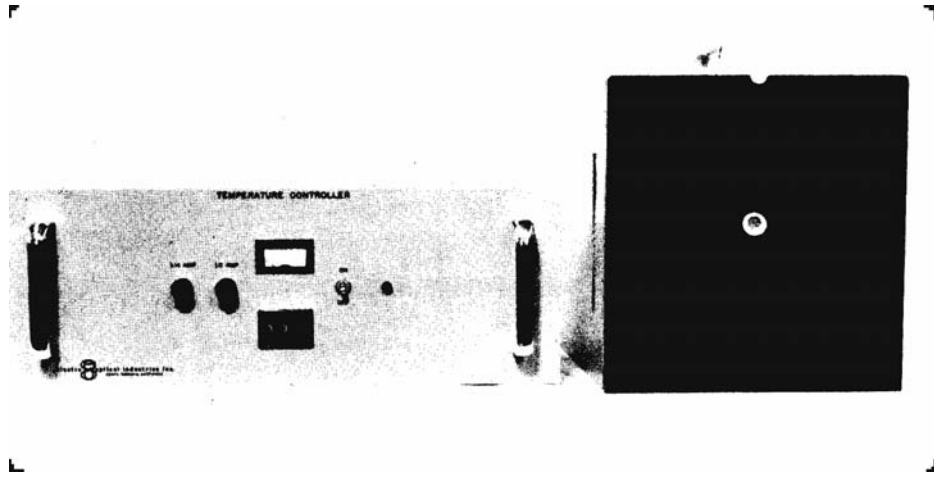


FIGURE 8 EOI blackbody.

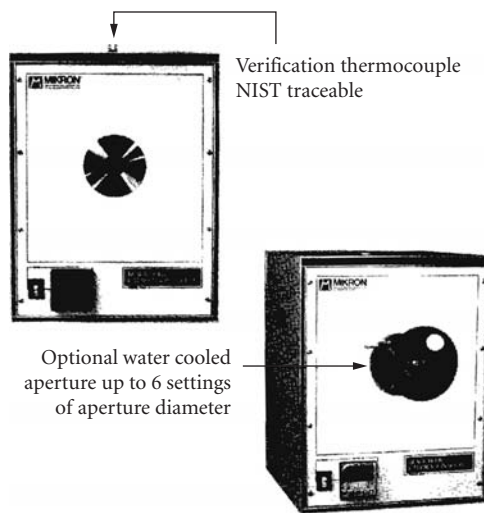


FIGURE 9 Mikron blackbody.

the material at room temperature is quite high, the working voltage is insufficient to get the glower started. Once started, its negative temperature coefficient-of-resistance tends to increase current, which would cause its destruction, so that a ballast is required in the circuit. Starting is effected by applying external heat, either with a flame or an adjacent electrically heated wire, until the glower begins to radiate.

Data from a typical glower are as follows:

1. Power requirements: 117 V, 50 to 60 A, 200 W
2. Color temperature range: 1500 to 1950 K
3. Dimension: 0.05-in. diameter by 0.3 in.

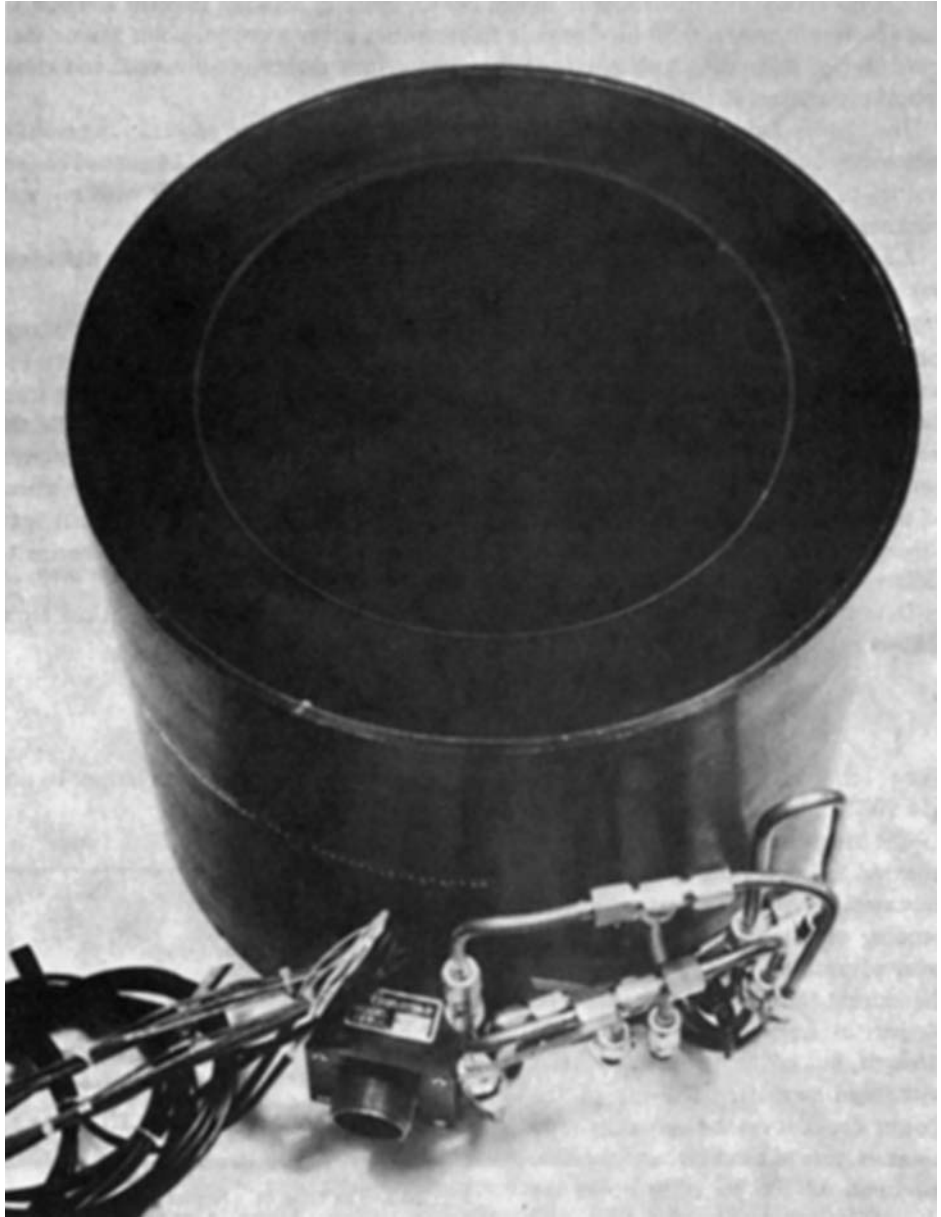


FIGURE 10 EOI model 1965. This model is 12 in. in diameter and 9 in. deep. The base is an array of intersecting conical cavities. The walls are hex-honeycomb and the temperature range is 175 to 340 K.

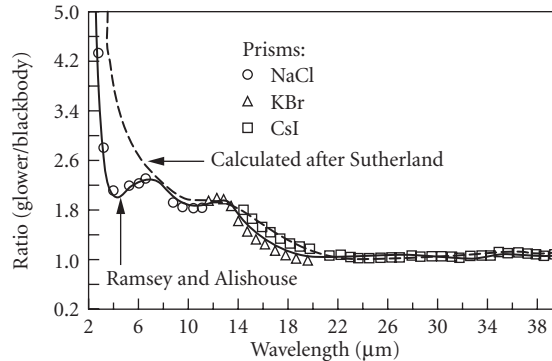


FIGURE 11 The ratio of a Nernst glower to a 900°C blackbody versus wavelength.

The spectral characteristics of a Nernst glower in terms of the ratio of its output to that of a 900°C blackbody are shown in Fig. 11.

The life of the Nernst glower diminishes as the operating temperature is increased. Beyond a certain point, depending on the particular glower, no great advantage is gained by increasing the current through the element. The glower is fragile, with low tensile strength, but can be maintained intact with rigid support. The life of the glower depends on the operating temperature, care in handling, and the like. Lifetimes of 200 to 1000 hours are claimed by various manufacturers.

Since the Nernst glower is made in the form of a long thin cylinder, it is particularly useful for illuminating spectrometer slits. Its useful spectral range is from the visible region to almost 30 μm , although its usefulness compared with other sources diminishes beyond about 15 μm . As a rough estimate, the radiance of a glower is nearly that of a graybody at the operating temperature with an emissivity in excess of 75 percent, especially below about 15 μm . The relatively low cost of the glower makes it a desirable source of moderate radiant power for optical uses in the laboratory. The makers of spectroscopic equipment constitute the usual source of supply of glowers (or of information about suppliers).

Globar The globar is a rod of bonded silicon carbide usually capped with metallic caps which serve as electrodes for the conduction of current through the globar from the power source. The passage of current causes the globar to heat, yielding radiation at a temperature above 1000°C. A flow of water through the housing that contains the rod is needed to cool the electrodes (usually silver). This complexity makes the globar less convenient to use than the Nernst glower and necessarily more expensive. This source can be obtained already mounted, from a number of manufacturers of spectroscopic equipment. Feedback in the controlled power source makes it possible to obtain high radiation output.

Ramsey and Alishouse¹⁵ provide information on a particular sample globar as follows:

1. Power consumption: 200 W, 6 A
2. Color temperature: 1470 K

They also provide the spectral characteristics of the globar in terms of the ratio of its output to that of a 900°C blackbody. This ratio is plotted as a function of wavelength in Fig. 12. Figure 13¹⁶ is a representation of the spectral emissivity of a globar as a function of wavelength. The emissivity values are only representative and can be expected to change considerably with use.

Gas Mantle The Welsbach mantle is typified by the kind found in high-intensity gasoline lamps used where electricity is not available. The mantle is composed of thorium oxide with some additive

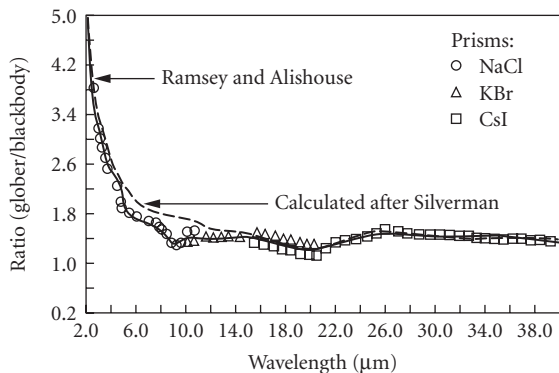


FIGURE 12 The ratio of a globar to a 900°C blackbody versus wavelength.

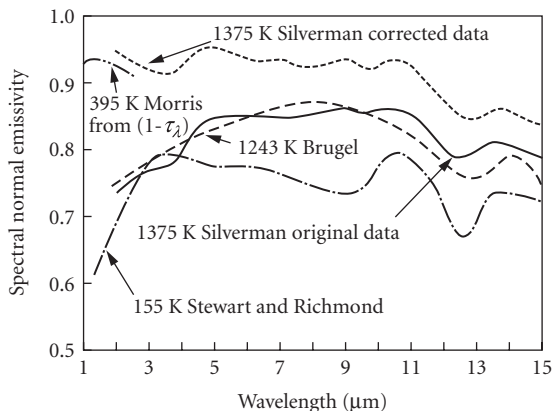


FIGURE 13 The spectral emissivity of a globar.

to increase its efficiency in the visible region. Its near-infrared emissivity is quite small, except for regions exemplified by gaseous emission, but increases considerably beyond 10 μm .

Ramsey and Alishouse¹⁵ provide information on a propane-heated sample from an experiment in which a comparison of several sources is made:

1. Color temperature: 1670 K
2. Dimensions: 25.4 by 38.1 mm

The spectral characteristics of the mantle in terms of the ratio of its output to that of a 900°C blackbody are shown in Fig. 14.

Pfund modified the gas mantle so that it became more a laboratory experimental source than an ordinary radiator. By playing a gas flame on an electrically heated mantle, he was able to increase its radiation over that from the gas mantle itself.¹⁷ Figure 15 shows a comparison of the gas mantle and the electrically heated gas mantle, with a Nernst glower. Strong¹⁸ points out that playing a flame against the mantle at an angle produces an elongated area of intense radiation useful for illuminating the slits of a spectrometer.

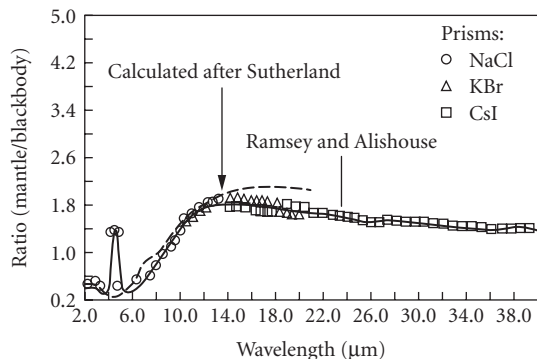


FIGURE 14 The ratio of the gas mantle to a 900°C blackbody versus wavelength.

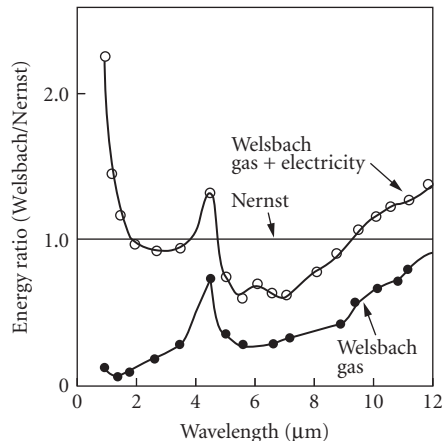


FIGURE 15 Emission relative to that of a Nernst glower (2240 K) of the gas-heated mantle (lower curve) and that of the mantle heated by gap plus electricity (upper curve).

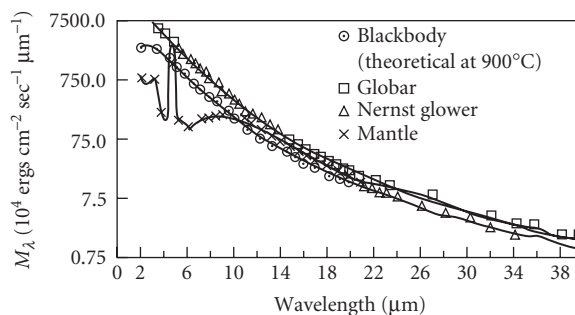


FIGURE 16 The spectral radiant emittances of a globar, Nernst glower, 900°C blackbody, and gas mantle versus wavelength.

Comparison of Nernst Glower, Globar, and Gas Mantle Figure 16 compares these three types of sources, omitting a consideration of differences in the instrumentation used in making measurements of the radiation from the sources.

Availability, convenience, and cost usually influence a choice of sources. At the very long wavelength regions in the infrared, the gas mantle and the globar have a slight edge over the Nernst glower because the Nernst glower (a convenient, small, and inexpensive source) does not have the power of the gas mantle and globar.

Tungsten-Filament Lamps A comprehensive discussion of tungsten-filament lamps is given by Carlson and Clark.¹⁹ Figures 17 to 19 show the configurations of lamp housings and filaments. The types and variations of lamps are too numerous to be meaningfully included in this chapter. The reader is referred to one of the buyer's guides for a comprehensive delineation of manufacturers from whom unlimited literature can be obtained.

Tungsten lamps have been designed for a variety of applications; few lamps are directed toward scientific research, but some bear directly or indirectly on scientific pursuits insofar as they can

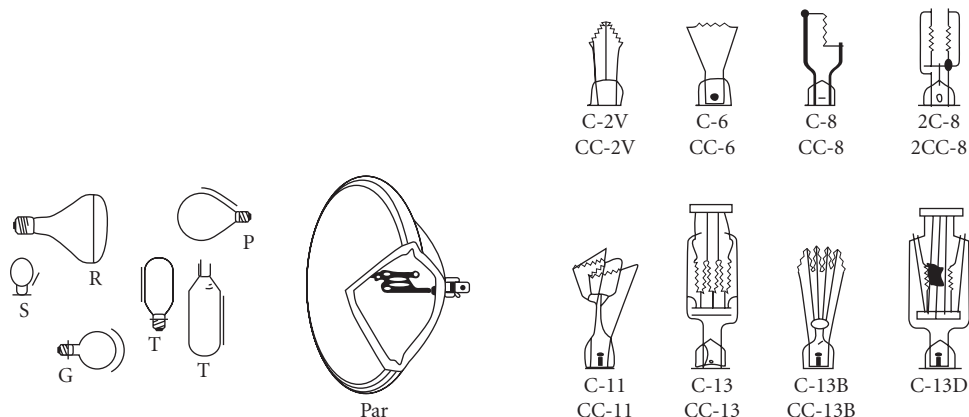


FIGURE 17 Bulk shapes most frequently used for lamps in optical devices. Letter designations are for particular shapes.

FIGURE 18 Most commonly used filament forms. Letters designate the type of filament.

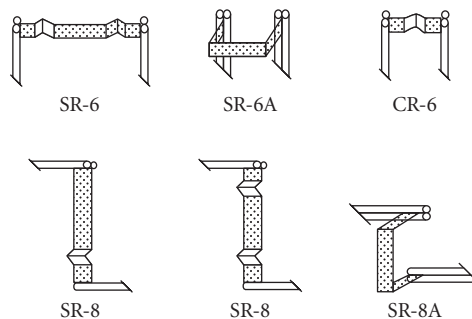
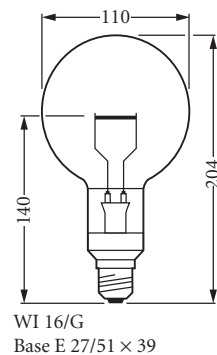
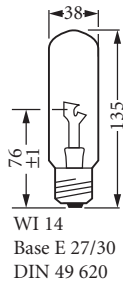
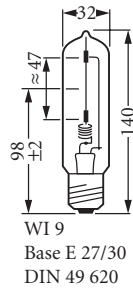


FIGURE 19 Ribbon-type tungsten filaments. Type designations are by number.

provide steady sources of numerous types of radiation. One set of sources cited here, particularly for what the manufacturer calls their scientific usefulness, is described in *Lamps for Scientific Purposes*.²⁰ Their filament structures are similar to those already described, but their designs reduce extraneous radiation and ensure the quality and stability of the desired radiation. The lamps can be obtained with a certification of their calibration values.

The physical descriptions of some of these sources are given in Fig. 20. Applications (according to the manufacturer, Osram) are photometry, pyrometry, optical radiometry, sensitometry, spectroscopy, spectrometry, polarimetry, saccharimetry, spectrophotometry, colorimetry, microscopy, microphotography, microprojection, and stroboscopy.

Quartz Envelope Lamps These are particularly useful as standards because they are longer lasting (due to action of iodine in the quartz-iodine series), can be heated to higher temperatures, are sturdier, and can transmit radiation to longer wavelengths in the infrared than glass-envelope lamps. Studer and Van Beers²¹ have shown the spectral deviation to be expected of lamps containing no iodine. The deviation, when known, is readily acceptable in lieu of the degradation in the lamp caused by the absence of iodine. The particular tungsten-quartz-iodine lamps used in accordance



Order reference	Upper limits for electric data (V) (A)		Color temperature T_f max.*	Luminance T_s max.	Dimensions of luminous width (mm)	Area height (mm)	Burning position†	Base
-----------------	--	--	-------------------------------	----------------------	-----------------------------------	------------------	-------------------	------

Lamps for scientific purposes

WI 9	8.5	6	2856 K	—	0.2	47	s	E 27
WI 14	5	16	—	2400 K	1.6	8	s	E 27
WI 16/G	9	16	—	2600 K	21	1.6	s+h	E 27
WI 17/G	9	16	—	2600 K	1.6	20	s	E 27
WI 40/G	31	6	2856 K	—	18	18	s+h	E 27
WI 41/G	31	6	2856 K	—	18	18	s+h	E 27

Lamps for scientific purposes are gas-filled, incandescent lamps for calibration of luminous intensity, luminous flux, luminance (spectral radiant temperature), color temperature (luminance temperature and spectral radiance distribution). A test certificate can be issued for these types of lamps.

Also for other types of lamp with sufficiently constant electric and photometric data, a test certificate can be issued. To order a test certificate, the order reference of the lamp, the type of measurement, and the desired burning position have to be given. Example: Lamp 41/G, measurement of the electric data and the luminous intensity for $T_f = 2856$ K (light type A), burning position vertical, base up.

Variables for which test certificates can be issued are shown in the following table by +. The sign (+) indicates that certificates can be issued for variables although the lamps were not designed for such measurements.

Type of lamp	Light intensity	Luminous flux	Luminance	Color temperature	Spectral radiance distribution‡
WI 9	+	—	—	+	—
WI 14	(+)	—	+	+	+ 300–800 mm
WI 16/G	(+)	—	+	+	+ 300–800 mm
WI 17/G	—	—	+	(+)	+ 250–800 mm 250–2500 mm
WI 40/G	+	+	—	(+)	—
WI 42/G	+	—	—	+	—

Description

WI 9:
Lamp with uncoiled straight filament.

WI 14:
Tungsten ribbon lamp with tubular bulb. The portion of the tungsten ribbon to be utilized for measurement is mounted parallel to the lamp axis and positioned approx. 8 mm off-axis in the measuring direction.

* The color temperature of 2856 K corresponds with light-type A (DIN 5035).

† s = vertical (base down); h = vertical (base up).

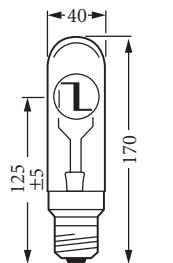
‡ Only for additional measurement of luminance or color temperature.

FIGURE 20a Lamps for scientific purposes. (Note dimensions in mm.)

with the NIST are described earlier in this chapter. Others can be obtained in a variety of sizes and wattages from General Electric, Sylvania, and a variety of other lamp manufacturers and secondary sources.

Carbon Arc

The carbon arc has been passed down from early lighting applications in three forms: low-intensity arc, flame, and high-intensity arc. The low- and high-intensity arcs are usually operated on direct



WI 16/G:

Tungsten ribbon lamp with spherical bulb. Horizontal tungsten ribbon with a small notch to indicate the measuring point. The ribbon is positioned approx. 3 mm off-axis.

WI 17/G:

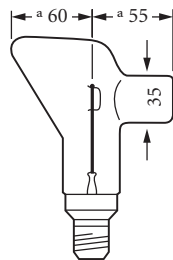
Tungsten ribbon lamp with horn-shaped bulb. The bulb has a tubular extension with a sealed-on quartz glass window (homogenized ultrasil). Vertical tungsten ribbon with a small notch to indicate the measuring point.

WI 40/G:

Standard lamp for total radiation, luminous flux, and color temperatures with conic bulb. The bulb shape prevents reflections in the direction of the plane normal of the luminous area, which is formed by the meandrous-shaped filament.

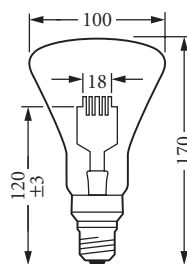
WI 41/G:

Standard lamp for light intensity and color temperature with conic bulb. Differs from the WI 40/G lamp by a black, opaque coating which covers one side of the bulb. A window is left open in the coating opposite the filament, through which over an angle of approx. $\pm 3^\circ$ a constant light intensity is emitted. The black coating prevents stray light being reflected in the measuring direction.



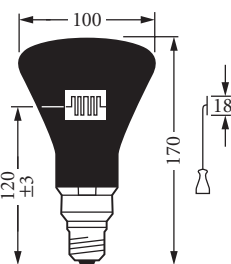
WI 17/G

Base E 27/51 × 39



WI 40/G

Base E 27/51 × 39



WI 41/G

Base E 27/51 × 39

FIGURE 20b (Continued)

current; the flame type adapts to either direct or alternating current. In all cases, a ballast must be used. In the alternating current arc, the combined radiation from the two terminals is less than that from the positive crater of the direct-current arc of the same wattage.²²

Spatial variation in the amount of light energy across the crater of dc arcs for different currents is shown in Fig. 21.

The carbon arc is a good example of an open arc, widely used because of its very high radiation and color temperatures (from approximately 3800 to 6500 K, or higher). The rate at which the material is consumed and expended during burning (5 to 30 cm/h) depends on the intensity of the arc. The arc is discharged between two electrodes that are moved to compensate for the rate of consumption of the material. The anode forms a crater of decomposing material which provides a center of very high luminosity. Some electrodes are hollowed out and filled with a softer carbon material which helps keep the arc fixed in the anode and prevents it from wandering on the anode surface.

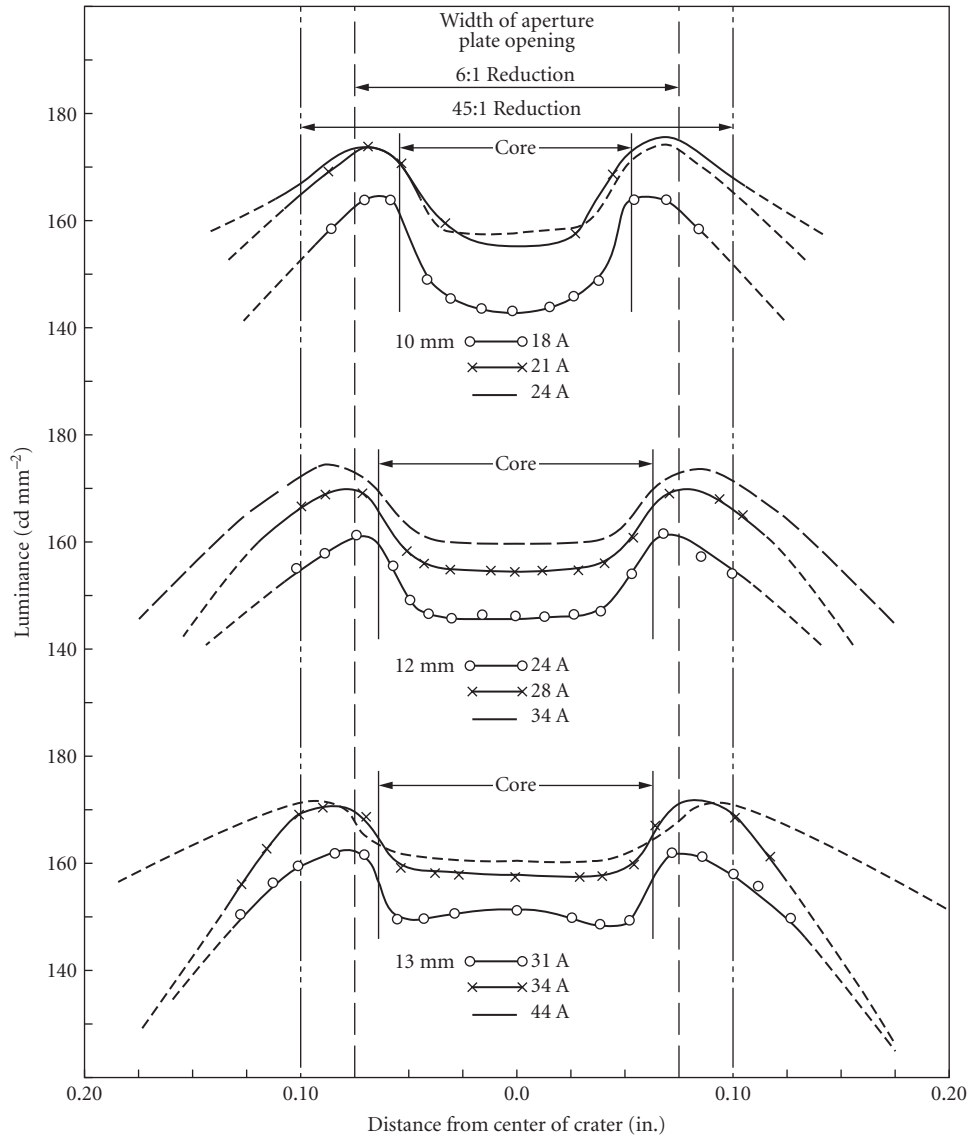


FIGURE 21 Variations in brightness across the craters of 10-, 12-, and 13-mm positive carbons of dc plain arcs operated at different currents in the regions of recommended operation.

In some cored electrodes, the center is filled with whatever material is needed to produce desired spectral characteristics in the arc. In such devices, the flame between the electrodes becomes the important center of luminosity, and color temperatures reach values as high as 8000 K.²² An example of this so-called flaming arc is shown in Fig. 22a. Figure 22b and c shows the low-intensity dc carbon arc and the high-intensity dc carbon arc with rotating positive electrodes. Tables 2 and 3 give characteristics of dc high-intensity and flame carbon arcs.

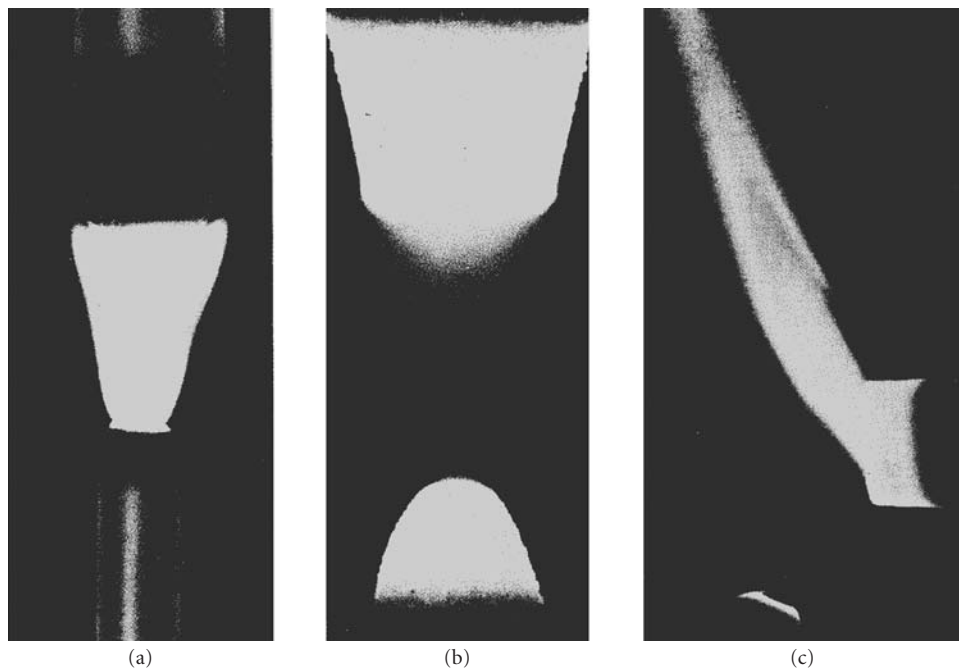


FIGURE 22 Various types of carbon arc: (a) flame type; (b) low-intensity dc arc; and (c) high-intensity dc arc with rotating positive carbon.

A spectrum of low-intensity arc (Fig. 23) shows the similarity between the radiation from it and a 3800-K blackbody, except for the band structure at 0.25 and 0.39 μm . In Koller²² an assortment of spectra are given for cored carbons containing different materials. Those for a core of soft carbon and for a polymetallic core are shown in Figs. 24 and 25. Because radiation emitted from the carbon arc is very intense, this arc supplants, for many applications, sources which radiate at lower temperatures. Among the disadvantages in using the carbon arc are its inconvenience relative to the use of other sources (e.g., lamps) and its relative instability. However, Null and Lozier²³ have studied the properties of the low-intensity carbon arc extensively and have found that under the proper operating conditions the carbon arc can be made quite stable; in fact, in their treatise they recommend its use as a standard of radiation at high temperatures.

Enclosed Arc and Discharge Sources (High-Pressure)

Koller²² states that the carbon arc is generally desired if a high intensity is required from a single unit but that it is less efficient than the mercury arc. Other disadvantages are the short life of the carbon with respect to mercury, and combustion products which may be undesirable. Worthing² describes a number of the older, enclosed, metallic arc sources, many of which can be built in the laboratory for laboratory use. Today, however, it is rarely necessary to build one's own source unless it is highly specialized.

*The Infrared Handbook*¹ compiles a large number of these types of sources, some of which will be repeated here, in case that publication would not be currently available to the reader. However, the reader should take caution that many changes might have occurred in the characteristics of these sources and in the supplier whose product is preferred. Consultation with the Photonics Directory (see preceding) is usually a good procedure. In some cases a certain type of source described previously

TABLE 2 DC Carbon Arcs

	Low Intensity			Nonrotating High Intensity			Rotating High Intensity			
	1	2	3	4	5	6	7	8	9	10
Type of carbon	Microscope	Projector	Projector	Projector	Projector	Projector	Projector	Searchlight	Studio	
Positive carbon:										
Diameter (mm)	5	7	8	10	11	13.6	13.6	16	16	16
Length (in.)	8	12-14	12-14	20	20	22	22	22	22	22-30
Negative carbon:										
Diameter	6 mm	6 mm	7 mm	11/32 in.	3/8 in.	0.5 in.	0.5 in.	11 mm	17/32 in.	7/16 in.
Length (in.)	4.5	9	9	9	9	9	9	12	9	12-48
Arc current (A)	5	50	70	105	120	160	180	150	225	400
Arc volts (dc)	59	40	42	59	57	66	74	78	70	80
Arc power (W)	295	2,000	2,940	6,200	6,840	10,600	13,300	11,700	15,800	32,000
Burning rate (in. h ⁻¹)										
Positive carbon	4.5	11.6	13.6	21.5	16.5	17	21.5	8.9	20.2	55
Negative carbon	2.1	4.3	4.3	2.9	2.4	2.2	2.5	3.9	2.2	3.5
Approximate crater diameter (in.)	0.12	0.23	0.28	0.36	0.39	0.5	0.5	0.55	0.59	0.59
Maximum luminance of crater (cd cm ⁻²)	15,000	55,000	83,000	90,000	85,000	96,000	95,000	65,000	68,000	45,000
Forward crater candlepower	975	10,500	22,000	36,000	44,000	63,000	78,000	68,000	99,000	185,000
Crater lumens [†]	3,100	36,800	77,000	126,000	154,000	221,000	273,000	250,000	347,000	660,000
Total lumens [‡]	3,100	55,000	115,000	189,000	231,000	368,000	410,000	374,000	521,000	999,000
Total lumens per arc watt	10.4	29.7	39.1	30.5	33.8	34.7	30.8	32	33	30.9
Color temperature (K) [§]	3,600	5,950	5,500-6,500	5,500-6,500	5,500-6,500	5,500-6,500	5,500-6,500	5,400	4,100	5,800-6,100

[†]Typical applications: 1, microscope illumination and projection; 2 to 7, motion-picture projection; 8, searchlight projection; 9, motion-picture-set lighting and motion-picture and television background projection.

[‡]Includes light radiated in forward hemisphere.

[§]Includes light from crater and arc flame in forward hemisphere.

[¶]Crater radiation only.

TABLE 3 Flame-Type Carbon Arcs

Type of carbon Flame materials Burning position ^g Upper carbon ^d	Application Number ^e									
	1	3	3	4	5	6	7 ^b	8 ^{c,d}	9 ^{d,e}	10 ^f
C		E	Sunshine Rare earth Vertical	Sunshine Rare earth Vertical	W	Enclosed arc None Vertical	Photo Rare earth Vertical	Sunshine Rare earth Horizontal	Photo Rare earth Horizontal	Studio Rare earth Vertical
Polymetallic Vertical	22 mm	Strontium Vertical	22 mm	22 mm	Polymetallic Vertical	1/2 in. 3-16 Vertical	1/2 in. 12 Vertical	6 mm 6.5 Horizontal	9 mm 8 Horizontal	8 mm 12 Vertical
Diameter	12	22 mm	12	12	22 mm	1/2 in.	1/2 in.	6 mm	9 mm	7 mm
Length (in.)	13 mm	13 mm	13 mm	13 mm	13 mm	3-16	3-16	6 mm	9 mm	9
Lower carbon ^d	12	12	12	12	12	16	16	6.5	8	8
Diameter	60	60	60	80	80	138	50	40	95	40
Length (in.)	50	50	50	50	50	2.2	1.9	24	30	37 dc
Arc current (A)	3	3	3	4	4	5.9	1.9	1	2.85	1.5
Arc voltage (ac) ^h	2,100	6,300	9,100	10,000	8,400	1,170	6,700	4,830	14,200	11,000
Arc power (kW)	23,000	69,000	100,000	110,000	92,000	13,000	74,000	53,000	156,000	110,000
Candlepower ⁱ	7.6	23	33.3	27.5	23	5.9	39.8	53	54.8	73.5
Lumens per arc watt			12,800 ^j	24,000 ^j			7,420 ^j	6,590	8,150	4,700
Color temperature (K)										
Spectral intensity ($\mu\text{W cm}^{-2}$)										
1 m from arc axis:										
Below 270 nm	540.0	180.0	102	140	1,020		95	11	100	12
270-320 nm	540.0	150.0	186	244	1,860		76	49	100	48
320-400 nm	1,800.0	1,200.0	2,046	2,816	3,120	1,700	684	415	1,590	464
400-450 nm	300.0	1,100.0	1,704	2,306	1,480	177	722	405	844	726
450-700 nm	600.0	4,050.0	3,210	3,520	2,600	442	2,223	1,602	3,671	3,965
700-1125 nm	1,580.0	2,480.0	3,032	3,500	3,220	1,681	1,264	1,368	5,632	2,123
Above 1125 nm	9,480.0	10,290.0	9,820	11,420	14,500	6,600	5,189	3,290	8,763	4,593
Total	14,930	19,460	20,100	24,000	27,800	10,600	10,253	7,140	20,600	11,930

Spectral radiation (percent of input power):										
Below 270 nm	1.8	0.6	0.34	0.35	2.55	0.5	0.11	0.32	0.08	
270–320 nm	1.8	0.5	0.62	0.61	4.65	0.4	0.49	0.35		
320–400 nm	6.0	4.0	6.82	7.04	7.80	7.7	4.15	5.59	3.09	
400–450 nm	1.3	3.7	5.68	5.90	3.70	0.8	4.05	2.96	4.84	
450–700 nm	2.0	13.5	10.7	8.80	6.50	11.7	16.02	12.86	26.43	
700–1125 nm	5.27	8.27	10.1	8.75	8.05	7.6	13.68	10.75	14.15	
Above 1125 nm	31.6	34.3	32.7	28.55	36.25	27.3	32.90	30.60	30.62	
Total	49.77	64.87	67.00	60.00	69.50	54.0	71.40	72.20	79.53	

^a Typical applications: 1 to 5 and 8, photochemical, therapeutic, accelerated exposure testing, or accelerated plant growth; 6, 7, and 9 blue-printing diazo printing, photo copying, and graphic arts; 10, motion-picture and television studio lighting.

^b Photographic white-flame carbons.

^c High-intensity copper-coated sunshine carbons.

^d Both carbons are same in horizontal, coaxial arcs.

^e High-intensity photo carbons.

^f Motion-picture-studio carbons

^g All combinations shown are operated coaxially.

^h All operated on alternating current except item 10.

ⁱ Horizontal candlepower, transverse to arc axis.

^j Deviates enough from blackbody colors to make color temperature of doubtful meaning.

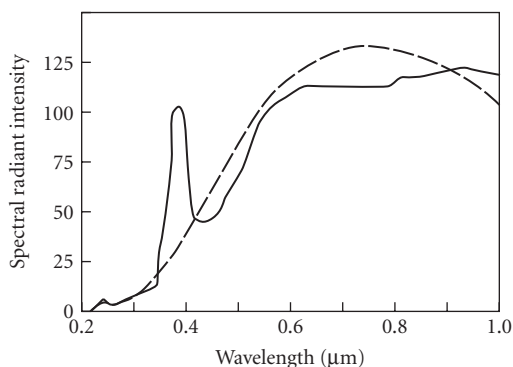


FIGURE 23 Spectral distribution of radiant flux from 30-A, 55-V dc low-intensity arc with 12-mm positive carbon (solid line) and a 3800-K blackbody radiator (broken line).

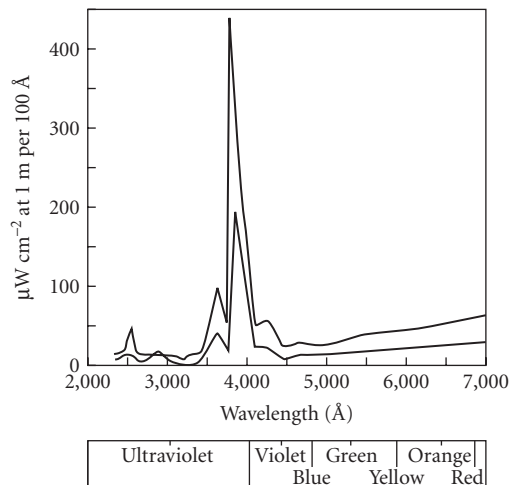


FIGURE 24 Spectral energy distribution of carbon arc with core of soft carbon. Upper curve: 60-A ac 50-V across the arc; lower curve: 30-A ac 50-V across the arc.

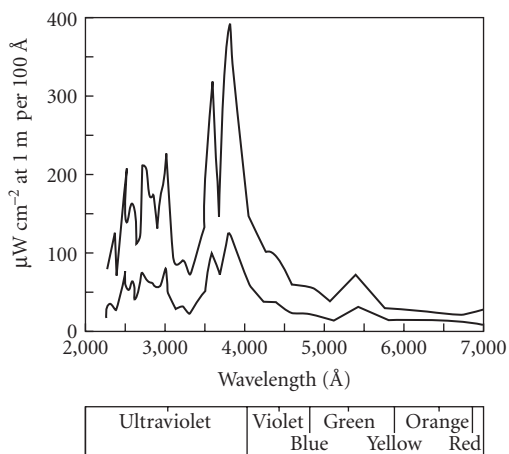


FIGURE 25 Spectral energy distribution of carbon arc with polymetallic-cored carbons. Upper curve: 60-A ac 50-V across the arc; lower curve: 30-A ac 50-V across the arc.

may not still exist. Thus, whereas some manufacturers were less compliant in providing data, they should be expected to respond more readily to a potential customer.

*Uviarc** This lamp is an efficient radiator of ultraviolet radiation. The energy distribution of one type is given in Fig. 26. Since the pressure of this mercury-vapor lamp is intermediate between the usual high- and the low-pressure lamps, little background (or continuum) radiation

*Registered trademark of General Electric.

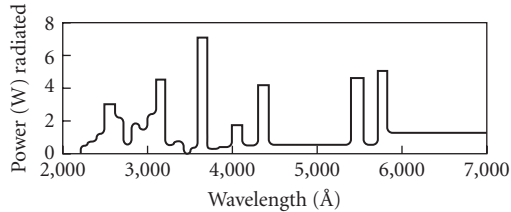


FIGURE 26 Intensity distribution of UA-2 intermediate-pressure lamp.

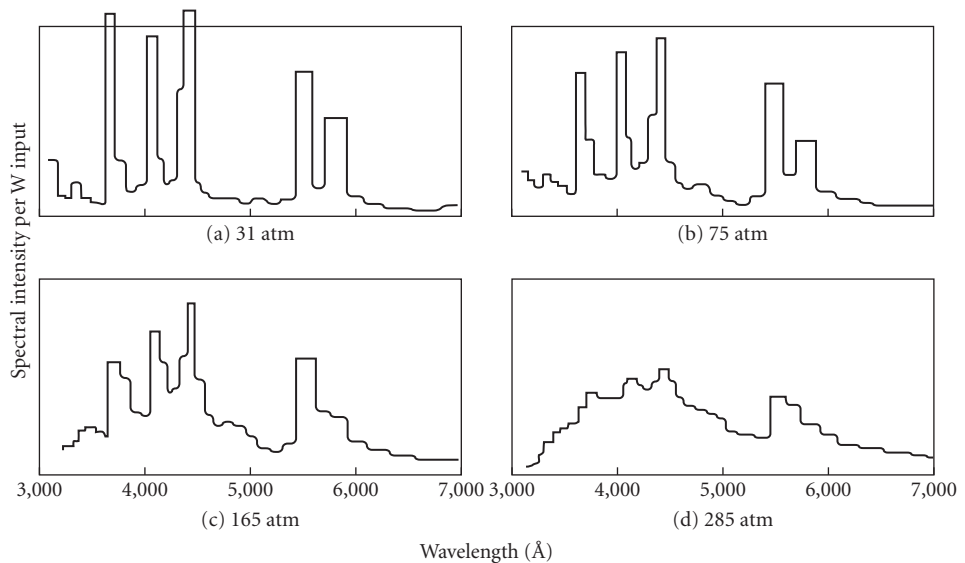


FIGURE 27 Emission spectrum of high-pressure mercury-arc lamps showing continuum background.

is present. In the truly high-pressure lamp, considerable continuum radiation results from greater molecular interaction. Figure 27²⁴ shows the dependence on pressure of the amount of continuum in mercury lamps of differing pressure. Bulb shapes and sizes are shown in Fig. 28.

Mercury Arcs A widely used type of high-pressure, mercury-arc lamp and the components necessary for its successful operation are shown in Fig. 29. The coiled tungsten cathode is coated with a rare-earth material (e.g., thorium). The auxiliary electrode is used to help in starting. A high resistance limits the starting current. Once the arc is started, the operating current is limited by ballast supplied by the high reactance of the power transformer. Spectral data for clear, 400-W mercury lamps of this type are given in Fig. 30.

Multivapor Arcs In these lamps, argon and mercury provide the starting action. Then sodium iodide, thallium iodide, and indium iodide vaporize and dissociate to yield the bulk of the lamp radiation. The physical appearance is like that of mercury lamps of the same general nature. Ballasts are similar to their counterparts for the mercury lamp. Up-to-date information on these sources should be obtained from the General Electric Corporation Lamp Division in Nela Park near Cleveland, Ohio. Spectral features of these sources are given in Fig. 31.

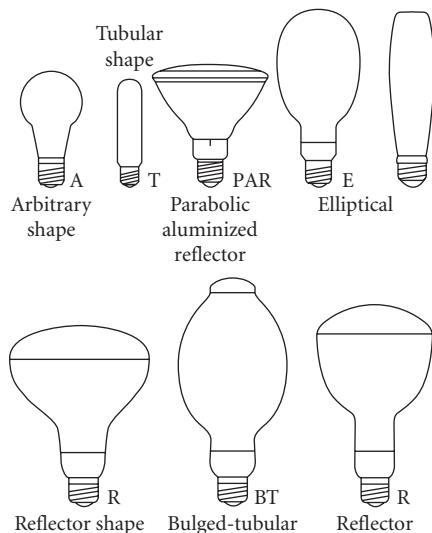


FIGURE 28 Bulb shapes and sizes (not to scale).

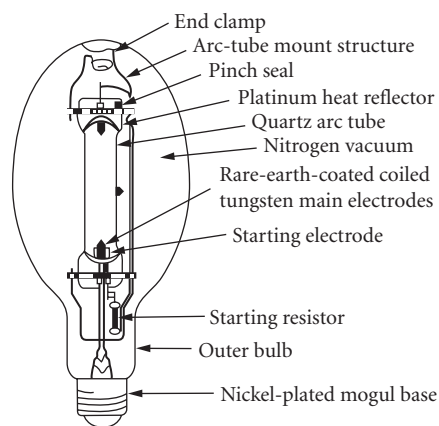


FIGURE 29 High-pressure mercury lamp showing various components.

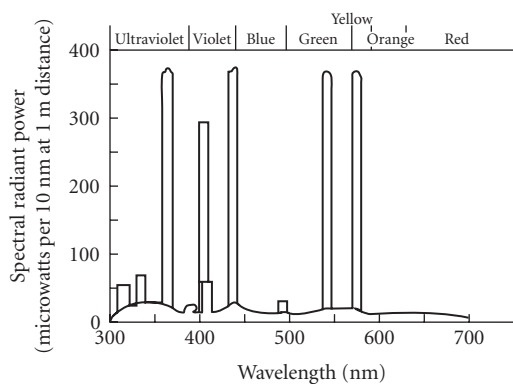


FIGURE 30 Spectral energy distribution for clear mercury-arc lamp.

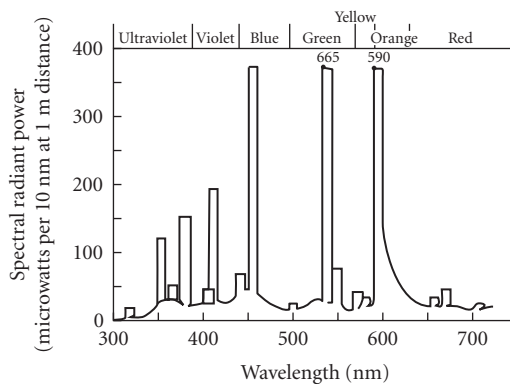


FIGURE 31 Spectral energy distribution of multivapor-arc lamp.

Lucalox[®] Lamps The chief characteristics of this lamp are high-pressure sodium discharge and a high temperature withstanding ceramic, Lucalox (translucent aluminum oxide), to yield performance typified in the spectral output of the 400-W Lucalox lamp shown in Fig. 32. Ballasts for this lamp are described in the General Electric *Bulletin TP-109R*.²⁵

Capillary Mercury-Arc Lamps²² As the pressure of the arc increases, cooling is required to avoid catastrophic effects on the tube. The AH6 tube (Fig. 33) is constructed with a quartz bulb wall and a quartz outer jacket, to allow 2800 K radiation to pass, or a Pyrex^{®†} outer jacket to eliminate ultraviolet. Pure water is forced through at a rapid rate, while the tube is maintained at a potential of 840 V.

[®]Registered trademark of General Electric.

[†]Registered trademark of Corning Glass Works.

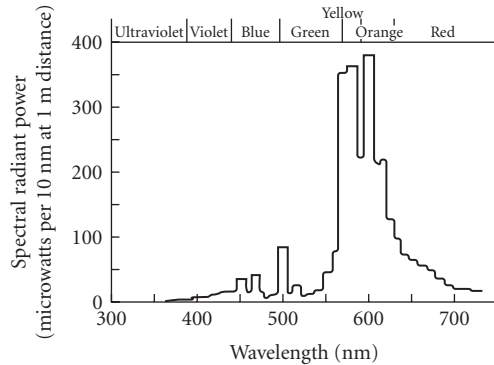


FIGURE 32 Spectral output of 400-W Lucalox lamp.

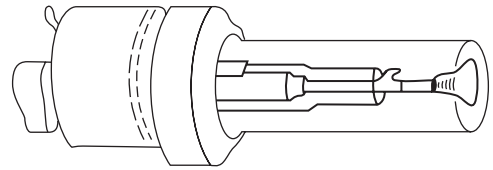


FIGURE 33 Water-cooled high-pressure (110 atm) mercury-arc lamp showing lamp in water jacket.

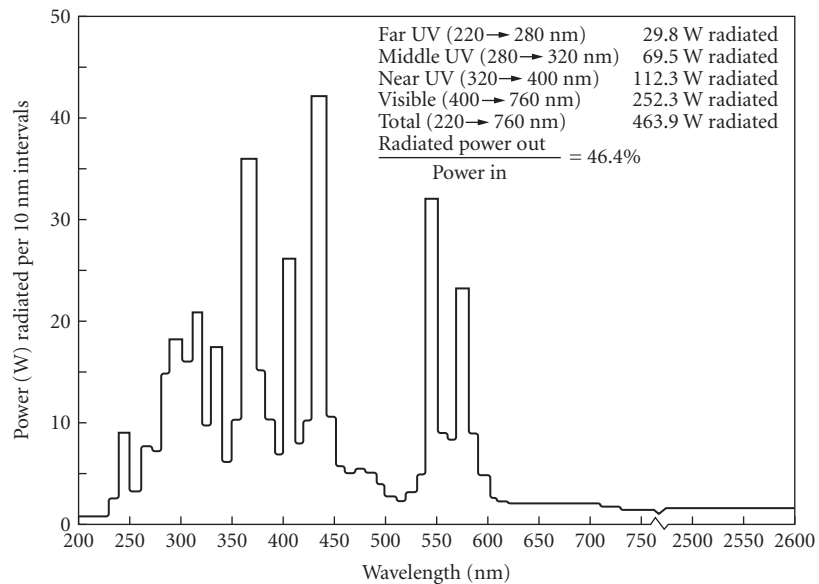


FIGURE 34 Spectral energy distribution of type BH6-1 mercury capillary lamp.

The spectral characteristics of certain tubes²⁶ are shown in Fig. 34. This company does not appear in the *Photonics Guide of 1989*, so the catalog referenced in the figure may not be current.

Compact-Source Arcs^{19,27} Some common characteristics of currently available compact-source arc lamps are as follows:

1. A clear quartz bulb of roughly spherical shape with extensions at opposite ends constituting the electrode terminals. In some cases, the quartz bulb is then sealed within a larger glass bulb, which is filled with an inert gas.
2. A pair of electrodes with relatively close spacing (from less than 1 mm to about 1 cm); hence the sometimes-used term short-arc lamps.

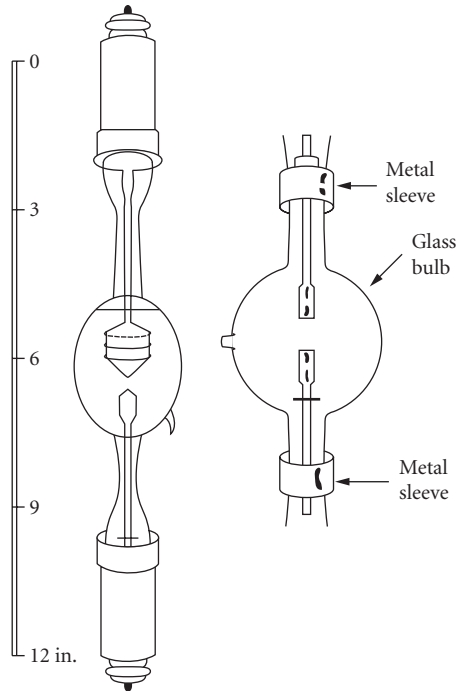


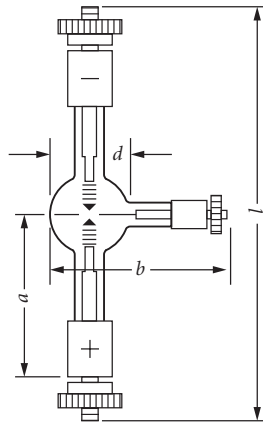
FIGURE 35 Construction of different lamps showing differences in relative sizes of electrodes for dc (*left*) and ac (*right*) operation.

3. A filling of gas or vapor through which the arc discharge takes place.
4. Extreme electrical loading of the arc gap, which results in very high luminance, internal pressures of many atmospheres, and bulb temperatures as high as 900°C . Precautions are necessary to protect people and equipment in case the lamps should fail violently.
5. The need for a momentary high-voltage ignition pulse, and a ballast or other auxiliary equipment to limit current during operation.
6. Clean, attention-free operation for long periods of time.

These lamps are designated by the chief radiating gases enclosed as mercury, mercury xenon, and xenon lamps.

Figure 35 shows a compact-source construction for a 1000-W lamp. Since starting may be a problem, some lamps (Fig. 36) are constructed with a third (i.e., a starting) electrode, to which a momentary high voltage is applied for starting (and especially restarting) while hot. The usual ballast is required for compact-source arcs. For stability, these arcs, particularly mercury and mercury-xenon, should be operated near rated power on a well-regulated power supply.²⁷

The spatial distribution of luminance from these lamps is reported in the literature already cited, and typical contours are shown in Fig. 37. Polar distributions are similar to those shown in Fig. 38. Spectral distributions are given in Figs. 39 through 41 for a 1000-W ac mercury lamp, a 5-kW dc xenon lamp, and 1000-W dc mercury-xenon lamp. Lamps are available at considerably less wattage.



Lamp (order reference)	HBO 200	
Type of current	DC	AC
Lamp supply voltage	V	>105 220
Operating voltage of lamp	$V \frac{L_1}{L_2}$	65...47 $\frac{61 \pm 4}{53 \pm 4}$
Operating current at operating voltage range	$A \frac{L_1}{L_2}$	3.1...4.2 $\frac{3.6}{4.2}$
Rated power of lamp	W	200
Luminous flux	lm	9,500
Luminous efficacy	lm/W	47.5
Light intensity	cd	1,000
Average luminance	cd / cm ²	40,000
Arc (width $w \times$ height h^4)	mm	0.6 \times 2.2
Average lamp life	h	200
Diameter d	mm	18
Length l_{max}	mm	108
Distance a	mm	41 \pm 2
Width b	mm	45
Burning position with stamped base down		s^{45}

FIGURE 36 Construction of a lamp with a third, starting electrode.

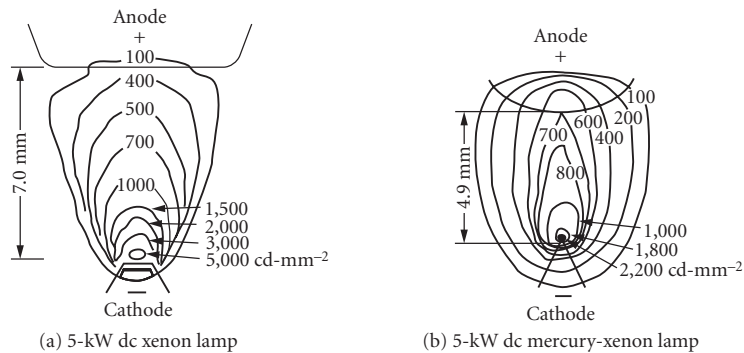


FIGURE 37 Spatial luminance distribution of compact-arc lamps.

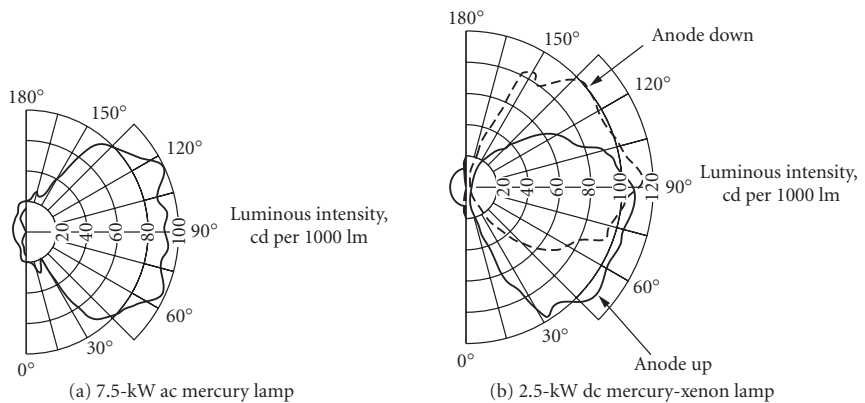


FIGURE 38 Polar distribution of radiation in planes that include arc axis. Asymmetry in (b) is due to unequal size of electrodes.

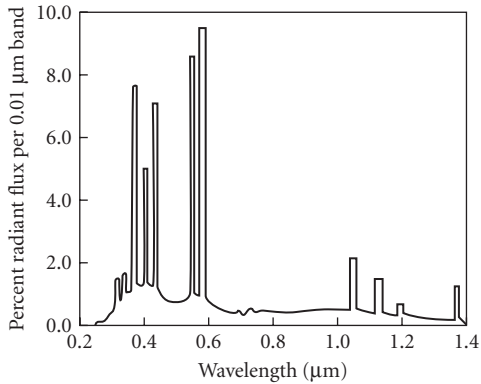


FIGURE 39 Spectral distribution of radiant intensity from a 1000-W ac mercury lamp perpendicular to the lamp axis.

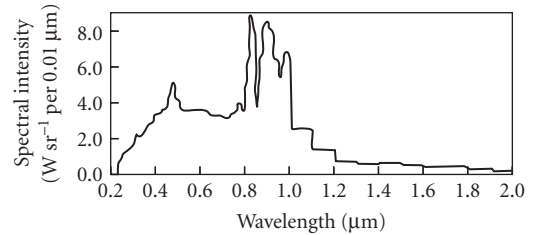


FIGURE 40 Spectral distribution of radiant intensity from a 5-kW dc xenon lamp perpendicular to the lamp axis with electrode and bulb radiation excluded.

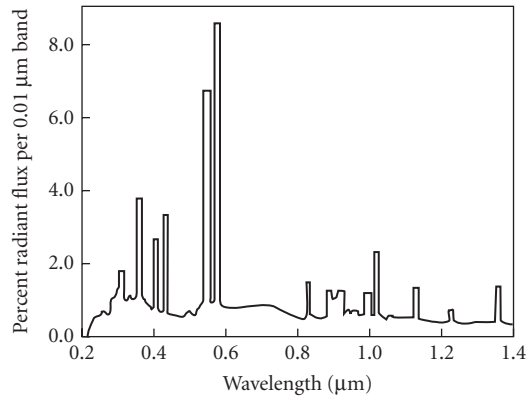


FIGURE 41 Spectral distribution of radiant flux from a 1000-W mercury-xenon lamp.

Cann²⁷ reports on some interesting special lamps tested by Jet Propulsion Laboratories for the purpose of obtaining a good spectral match to the solar distribution. The types of lamps tested were Xe, Xe-Zn, Xe-Cd, Hg-Xe-Zn, Hg-Xe-Cd, Kr, Kr-Zn, Kr-Cd, Hg-Kr-Zn, Hg-Kr-Cd, Ar, Ne, and Hg-Xe with variable mercury-vapor pressure. For details, the reader should consult the literature.

A special design of a short-arc lamp manufactured by Varian²⁸ is shown in Fig. 42. Aside from its compactness and parabolic sector, it has a sapphire window which allows a greater amount of IR energy to be emitted. It is operated either dc or pulsed, but the user should obtain complete specifications, because the reflector can become contaminated, with a resultant decrease in output.



FIGURE 42 High-pressure, short-arc xenon illuminators with sapphire windows. Low starting voltage, 150 through 800 W; VIX150, VIX300, VIX500, VIX800.

Enclosed Arc and Discharge Sources (Low-Pressure)²²

With pressure reduction in a tube filled with mercury vapor, the 2537 Å line becomes predominant so that low-pressure mercury tubes are usually selected for their ability to emit ultraviolet radiation.

Germicidal Lamps These are hot-cathode lamps which operate at relatively low voltages. They differ from ordinary fluorescent lamps which are used in lighting in that they are designed to transmit ultraviolet, whereas the wall of the fluorescent lamp is coated with a material that absorbs ultraviolet and reemits visible light. The germicidal lamp is constructed of glass of 1-mm thickness which transmits about 65 percent of the 2537 Å radiation and virtually cuts off shorter wavelength ultraviolet radiation.

Sterilamp[®] Types These cold cathode lamps start and operate at higher voltages than the hot-cathode type and can be obtained in relatively small sizes as shown in Fig. 43. Operating characteristics of the Sterilamps should be obtained from the manufacturer.

Black-Light Fluorescent Lamps This fluorescent lamp is coated with a phosphor efficient in the absorption of 2537Å radiation, emitting ultraviolet radiation in a broadband around 3650 Å. The phosphor is a cerium-activated calcium phosphate, and the glass bulb is impervious to shorter wavelength ultraviolet radiation. Characteristics of one type are given in Table 4.

Hollow Cathode Lamps A device described early in this century and used for many years by spectroscopists is the hollow-cathode tube. The one used by Paschen² consisted of a hollow metal cylinder and contained a small quantity of inert gas, yielding an intense cathode-glow characteristic of the cathode constituents. Materials that vaporize easily can be incorporated into the tube so that their spectral characteristics predominate.²

Several companies sell hollow-cathode lamps which do not differ significantly from those constructed in early laboratories. The external appearance of these modern tubes shows the marks of mass production and emphasis on convenience. They come with a large number of vaporizable elements, singly or in multiples, and with Pyrex[®] or quartz windows. A partial list of the characteristics or the lamps available from two manufacturers is given in Table 5. Their physical appearance is shown in Fig. 44a. A schematic of the different elements obtainable in various lamps is shown in Fig. 44b.²⁹

²⁹Registered trademark of Westinghouse Electric.

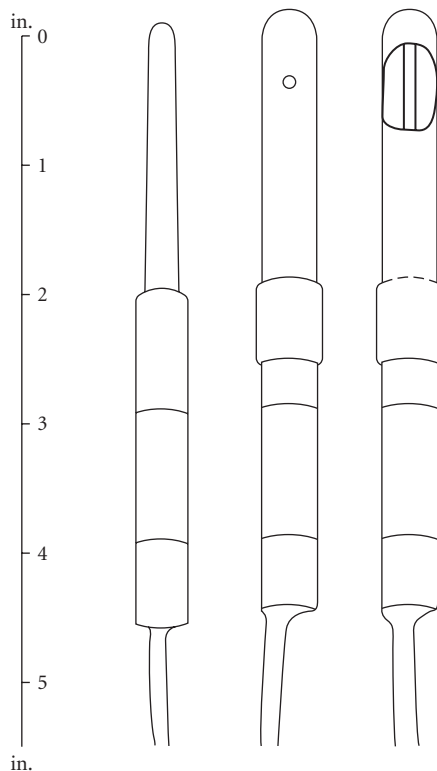


FIGURE 43 Pen-Ray low-pressure lamp. Pen-Ray is a registered trademark of Ultraviolet Products, Inc.

TABLE 4 Spectral Energy Distribution for Black-Light (360 BL) Lamps

(W)	Length (in.)	3200–3800 Å		Total Ultraviolet below 3800 Å		Total Visible (W)	3800–7600 Å %*	Erythematel Flux
		(W)	%*	(W)	%*			
6	9	0.55	9.1	0.56	9.4	0.1	1.7	250
15	18	2.10	14.0	2.20	14.6	0.4	2.7	950
30	36	4.60	15.3	4.70	15.8	0.9	3.0	2100
40	48	6.70	16.8	6.90	17.3	1.5	3.8	3000

* Percentage of input power.

Electrodeless Discharge Lamps^{30–32} The electrodeless lamp gained popularity when Meggers used it in his attempt to produce a highly precise standard of radiation. Simplicity of design makes laboratory construction of this type of lamp easy. Some of the simplest lamps consist of a tube, containing the radiation-producing element, and a microwave generator, for producing the electric field (within the tube) which in turn excites the elemental spectra. Lamps of this type can be purchased with specially designed microwave cavities for greater efficiency in coupling. Those made of fused quartz can transmit from ultraviolet to near infrared. The electrodeless lamp is better able than the

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps*

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Aluminum	Q	A	B	3092	WL22804	Copper	P	N	A	3247	JA-45-458
	P	N	A	3092	JA-45-452		Q	A	B	3247	WL22606
	Q	A	A	3092	WL22870		Q	A	A	3247	WL22879
Antimony	Q	N	B	3092	WL22929	Dysprosium	Q	N	B	3247	JA-45-490
	Q	N	A	3092	WL22954		Q	N	A	3247	WL23042
	Q	A	B	2311	WL22840		Q	N	B	4212	JA-45-595
	Q	A	A	2311	WL22872		Q	N	A	4212	WL22880
Arsenic	Q	N	B	2311	JA-45-461	Europium	Q	N	B	4008	JA-45-571
	Q	N	A	2311	WL22956		Q	N	A	4008	WL22881
	Q	N	B	1937	JA-45-315		Q	N	A	4594	JA-45-572
Barium	Q	N	A	1937	WL22873	Gadolinium	P	N	A	4079	WL22975
	Q	N	B	1937	JA-45-315		Q	N	B	4079	JA-45-573
	P	N	A	5536	JA-45-480		Q	N	A	4079	WL22986
Beryllium	Q	N	B	2349	WL23407	Gallium	Q	N	B	4172	JA-45-470
	Q	A	B	3068	WL22841		Q	N	A	4172	WL22884
Bismuth	Q	A	A	3068	WL22874	Germanium	Q	N	B	2651	JA-45-575
	Q	N	B	3068	JA-45-469		Q	N	B	2651	JA-45-313
	Q	N	A	3068	WL22957		Q	A	B	2676	WL22839
Boron	Q	A	B	2497	JA-45-568	Gold	Q	A	B	2676	WL22883
	Q	A	A	2497	WL22917		Q	N	B	2676	JA-45-467
Cadmium	Q	A	B	3261	WL22816	Hafnium	Q	N	B	3072	JA-45-303
	Q	A	A	3261	WL22875		Q	N	A	4104	WL22885
	Q	N	B	3261	JA-45-462		Q	N	B	4104	JA-45-576
Calcium	Q	N	A	3261	WL22958	Indium	Q	A	B	3040	WL22867
	P	N	A	4227	JA-45-440		Q	A	B	3040	WL22915
Cerium	Q	N	B	-	JA-45-569	Iridium	Q	A	A	3040	JA-45-471
	Q	N	A	4556	WL22978		Q	N	B	2850	JA-45-577
Cesium	P	A	A	4566	WL22817	Iron	P	A	A	3270	WL22602
	P	N	A	4566	JA-45-141		Q	A	B	3720	WL22611

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (Continued)

Element	Window†	Gas Fill‡	Size§	Analytical Line (Å)	Catalog Number¶	Element	Window†	Gas Fill‡	Size§	Analytical Line (Å)	Catalog Number¶
Chromium	P	A	A	3579	WL22812	Iron, high-purity	Q	N	B	3720	JA-45-155
	Q	A	B	3579	WL22521		P	N	A	3720	WL22820
	Q	A	A	3579	WL22877		Q	A	A	3720	WL22886
Cobalt	Q	N	B	3579	JA-45-454	Lanthanum	Q	N	A	3720	WL22887
	Q	N	A	3579	WL22959		Q	N	B	3720	WL22837
	P	A	A	3454	WL22813		Q	N	A	3720	WL22888
	Q	A	B	3454	WL22814		Q	A	B	5501	WL22846
	Q	A	A	3454	WL22878		Q	A	A	5501	WL22889
Lead	Q	N	B	3454	JA-45-456	Palladium	Q	A	A	5501	JA-45-495
	Q	N	A	3454	WL22953		Q	N	B	5501	JA-45-495
	Q	A	B	2833	WL22838		Q	A	B	3404	WL22857
	Q	A	A	2833	WL22890		Q	A	A	3404	WL22911
	Q	N	B	2833	JA-45-468		Q	N	B	3404	JA-45-475
Lithium 6	Q	N	A	2833	WL22952	Phosphorus	Q	N	A	3404	WL22970
	P	N	A	6708	JA-45-579		Q	N	B	2136	JA-45-449
Lithium 7	P	N	A	6708	WL22925	Platinum	Q	N	A	2136	WL22990
	P	A	A	6708	JA-45-580		Q	A	B	2659	WL22851
Lithium, natural	P	A	A	6708	WL22926	Potassium	Q	A	A	2659	WL22896
	P	A	A	6708	WL22825		Q	N	B	2659	JA-45-466
	P	N	A	6708	JA-45-444		Q	N	A	4044	JA-45-484
Lutetium	Q	A	B	6708	WL23115	Praseodymium	Q	N	B	4951	JA-45-585
	Q	N	A	3282	JA-45-581		Q	N	A	4951	WL22982
Magnesium	Q	N	A	3282	WL23010	Rhodium	Q	N	A	3460	JA-45-489
	Q	A	B	2852	WL22609		Q	N	A	3460	WL22967
	Q	A	A	2852	WL22891		Q	A	B	3435	WL22850
	Q	N	A	2852	WL22951		Q	A	A	3435	WL22897
	Q	N	B	2852	JA-45-451		Q	N	B	3435	JA-45-476
Manganese	Q	A	B	2795	WL22608	Rubidium	P	N	A	7800	JA-45-443
	Q	A	A	2795	WL22815		Q	N	B	7800	WL23046
	Q	N	B	2795	JA-45-472		Q	N	B	3499	JA-45-586
	Q	N	A	2795	WL22961		Q	N	B	4760	JA-45-587
	Q	A	A	2795	WL22876		Q	N	A	4760	WL22899

Mercury	Q	A	B	2537	JA-45-493	Scandium	Q	N	B	3912	JA-45-309
	Q	A	A	2557	WL22892	Selenium	Q	A	B	1960	WL22843
Molybdenum	Q	A	B	3133	WL22805		Q	A	A	1960	WL22898
	Q	A	A	3133	WL22893		Q	N	B	1960	JA-45-477
	Q	N	B	3133	JA-45-460		Q	N	A	1960	WL22963
Neodymium	Q	N	A	3133	WL22962	Silicon	Q	A	B	2516	WL22832
	Q	N	B	4925	JA-45-582		Q	A	A	2516	WL22900
	Q	N	A	4925	WL22980		Q	N	B	2516	JA-45-479
Nickel	P	A	A	3415	WL22605		Q	N	A	2516	WL22964
	Q	A	B	3415	WL22663	Silver	Q	A	B	3281	JA-45-483
	Q	N	B	3415	JA-45-457		Q	A	A	3281	WL22901
	Q	A	A	3415	WL22894	Sodium	P	A	A	5890	WL22864
	Q	N	A	3415	WL22895		P	N	A	5890	JA-45-485
Niobium	Q	N	B	4059	JA-45-486	Strontium	P	N	A	4607	JA-45-481
	Q	N	A	4059	WL22912	Sulphur	Q	N	B	-	JA-45-588
Osmium	Q	A	B	2909	JA-45-584						
Tantalum	Q	A	B	2714	JA-45-488	Zirconium	Q	A	B	3601	JA-45-482
	Q	A	A	2714	WL22913		Q	A	A	3601	WL22914
	Q	N	B	2714	WL22971		Q	N	B	3601	WL22998
	Q	N	A	2714	WL22972						
Tellurium	Q	A	B	2143	WL22842						
	Q	A	A	2143	WL22902						
	Q	N	B	2143	JA-45-473	Aluminum	P	N	C	3092	JA-45-36009
	Q	N	A	2143	WL22965	Antimony	Q	N	C	2311	JA-4-36010
Terbium	Q	N	B	4326	JA-45-589	Arsenic	Q	A	C	1937	JA-45-36011
	Q	N	A	4326	WL22903	Barium	P	N	C	5536	JA-45-36012
Thallium	Q	N	B	3776	WL23408	Beryllium	Q	N	C	2349	JA-45-36013
Thorium	Q	N	A	3245	WL23028	Bismuth	Q	N	C	3068	JA-45-36014
	Q	N	B	3245	JA-45-590	Boron	Q	A	C	2497	JA-45-36015
Thulium	Q	N	B	4105	JA-45-591	Cadmium	Q	N	C	3261	JA-45-36016
	Q	N	A	4105	WL23008	Calcium	P	N	C	4227	JA-45-36017
Tin	Q	A	B	2863	WL22822	Cerium	Q	N	C	-	JA-45-36019
	Q	A	A	2863	WL22904	Cesium	P	N	C	4556	JA-45-36020
	Q	N	B	2863	JA-45-463	Chromium	P	N	C	3579	JA-45-36021
	Q	N	A	2863	WL22966	Cobalt	Q	N	C	3454	JA-45-36022
Titanium	Q	N	B	3643	JA-45-592	Copper	P	N	C	3247	JA-45-36024
	Q	N	A	3643	WL22992	Dysprosium	P	N	C	4212	JA-45-36025

Single-Element 36000 Series

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (*Continued*)

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Tungsten	Q	N	B	4009	JA-45-465	Erbium	P	N	C	4008	JA-45-36026
	Q	A	B	4009	WL22849	Europium	P	N	C	4594	JA-45-36027
	Q	N	A	4009	WL22905	Gadolinium	P	N	C	4079	JA-45-36028
Uranium	Q	A	A	4009	WL22906	Gallium	Q	N	C	4172	JA-45-36029
	Q	N	B	5027	JA-45-447	Germanium	Q	N	C	2651	JA-45-36030
	Q	N	A	5027	WL22907	Gold	Q	N	C	2676	JA-45-36031
Vanadium	Q	A	B	3184	WL22856	Hafnium	Q	N	C	3072	JA-45-36032
	Q	A	A	3184	WL22910	Holmium	P	N	C	4104	JA-45-36033
	Q	N	B	3184	JA-45-453	Indium	Q	N	C	3040	JA-45-36034
Ytterbium	Q	N	A	3184	WL22974	Iridium	Q	N	C	2850	JA-45-36036
	Q	A	B	3988	JA-45-593	Iron	Q	N	C	3720	JA-45-36037
	Q	A	A	3988	WL22984	Lanthanum	P	N	C	5501	JA-45-36038
Yttrium	P	N	A	4102	WL22976	Lead	Q	N	C	2833	JA-45-36039
	Q	N	B	4102	JA-45-594	Lithium 6	P	N	C	6708	JA-45-36090
	Q	N	A	4102	WL22988	Lithium 7	P	N	C	6708	JA-45-36091
Zinc	Q	A	B	2139	WL22607	Lithium, natural	P	N	C	6708	JA-45-36040
	Q	N	B	2139	JA-45-459	Lutetium	P	N	C	3282	JA-45-36041
	Q	A	A	2139	WL22908	Magnesium	Q	N	C	2852	JA-45-36042
Mercury	Q	N	A	2139	WL22909	Manganese	Q	N	C	2795	JA-45-36043
	Q	A	C	2537	JA-45-36044		Multiple-Element 22000 Series				
	Q	N	C	3133	JA-45-36045	Aluminum, calcium	Q	N	B	-	WL23246
Neodymium	P	N	C	4925	JA-45-36046						
Nickel	Q	N	C	3415	JA-45-36047	Aluminum, calcium, magnesium	Q	A	B	-	WL22604
	P	N	C	4059	JA-45-36023						
Osmium	Q	A	C	2909	JA-45-36048	Aluminum, calcium, magnesium	Q	A	A	-	WL22871
Palladium	Q	N	C	3404	JA-45-36049						
Phosphorus	Q	N	C	2136	JA-45-36050	Aluminum, calcium, magnesium	Q	N	B	-	JA-45-450
	Q	N	C	2659	JA-45-36051						
Potassium	P	N	C	4044	JA-45-36052	Aluminum, calcium, magnesium	Q	N	A	-	WL22955

Praseodymium	P	N	C	4951	JA-45-36053	Aluminum, calcium, magnesium, iron	Q	N	B	-	JA-45-310
Rhenium	P	N	C	3460	JA-45-36056						
Rhodium	P	N	C	3435	JA-45-36057						
Rubidium	P	N	C	7800	JA-45-36058	Aluminum, calcium, magnesium, lithium	Q	N	B	-	JA-45-436
Ruthenium	P	A	C	3499	JA-45-36059						
Samarium	P	N	C	4760	JA-45-36060	Aluminum, calcium, magnesium, lithium	Q	A	A	-	WL23036
Scandium	P	N	C	3912	JA-45-36061						
Selenium	Q	N	C	1960	JA-45-36062	Aluminum, calcium, strontium	P	N	A	-	WL23403
Silicon	Q	N	C	2516	JA-45-36063						
Silver	P	A	C	3281	JA-45-36064	Antimony, arsenic, bismuth	Q	N	B	-	WL23147
Sodium	P	N	C	5890	JA-45-36065						
Strontium	P	N	C	4607	JA-45-36066	Arsenic, nickel	Q	N	B	-	JA-45-434
Sulphur	Q	N	C	-	JA-45-36067	Arsenic, selenium, tellurium	Q	N	B	-	JA-45-598
Tantalum	Q	A	C	2714	JA-45-36068	Barium, calcium, strontium	P	N	A	-	JA-45-437
Tellurium	Q	N	C	2143	JA-45-36069	Barium, calcium, silicon, magnesium	Q	N	B	-	JA-45-478
Terbium	P	N	C	4326	JA-45-36070						
Thallium	Q	N	C	3776	JA-45-36071	Cadmium, copper					
Thorium	Q	N	C	3245	JA-45-36072,	zinc, lead	Q	N	B	-	JA-45-597
Thulium	P	N	C	4105	JA-45-36073	Cadmium, silver, zinc, lead	Q	N	B	-	JA-45-308
Tin	Q	N	C	2863	JA-45-36074	Calcium, magnesium, strontium	Q	N	B	-	WL23605
Titanium	P	N	C	3643	JA-45-36075	Calcium, magnesium, zinc	Q	N	B	-	JA-45-311
Tungsten	Q	N	C	4009	JA-45-36076	Calcium, magnesium, aluminum, lithium	Q	N	B	-	WL23158
Uranium	P	N	C	5027	JA-45-36077	Calcium, zinc	Q	N	B	-	JA-45-304
Vanadium	Q	N	C	3184	JA-45-36078	Chromium, iron, manganese, nickel	Q	N	B	-	JA-45-442
Ytterbium	P	A	C	3988	JA-45-36079	Chromium, cobalt, nickel	Q	N	B	-	WL23174

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (Continued)

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Yttrium	P	N	C	4102	JA-45-36080	Chromium, copper	Q	N	B	-	JA-45-306
Zinc	Q	N	C	2139	JA-45-36081						
Zirconium	P	A	C	3601	JA-45-36082						
Chromium, manganese	Q	N	B	-	WL23499	Antimony, arsenic, bismuth	Q	N	C	-	JA-45-36203
Chromium, cobalt, copper, manganese, nickel	Q	N	B	-	WL23601	Barium, calcium, strontium, magnesium	Q	N	C	-	JA-45-36228
Chromium, cobalt, copper, iron,	Q	N	B	-	JA-45-599	Cadmium, silver, zinc, lead	Q	N	C	-	JA-45-36205
manganese, nickel	Q	N	B	-	JA-45-599						
Cobalt, copper,	Q	N	B	-	JA-45-305	Cadmium, copper, zinc, lead	Q	N	C	-	JA-45-36227
gold, nickel	Q	N	B	-	WL23295	Calcium, magnesium	Q	N	C	-	JA-45-36092
Cobalt, copper,	Q	N	B	-	WL23291						
zinc, molybdenum	Q	N	B	-	WL23291	Calcium, magnesium, zinc	Q	N	C	-	JA-45-36097
Cobalt, iron	Q	N	B	-	WL23426	Calcium, zinc	Q	N	C	-	JA-45-36093
Cobalt, nickel	Q	N	B	-	JA-45-431	Chromium, iron, manganese, nickel	Q	N	C	-	JA-45-36201
Copper, gallium	Q	N	B	-	JA-45-312	Chromium, cobalt, copper	Q	N	C	-	JA-45-36094
Copper, iron	Q	N	B	-	JA-45-435	manganese, nickel	Q	N	C	-	JA-45-36103
Copper, iron, manganese	Q	N	B	-	JA-45-301	Chromium, cobalt, copper, manganese, nickel	Q	N	C	-	JA-45-36096
Copper, iron, molybdenum	Q	N	B	-	JA-45-307	Chromium, copper, iron, nickel, silver	Q	N	C	-	JA-45-36108
Copper, iron, gold, nickel	Q	N	B	-	JA-45-492	Cobalt, copper, iron, manganese, molybdenum	Q	N	C	-	JA-45-36102
Copper, iron, manganese, zinc	Q	N	B	-	JA-45-491	Copper, zinc, lead, tin	Q	N	C	-	JA-45-36202
Copper, manganese	Q	N	B	-							

Copper, nickel	Q	N	B	-	WL23441A	Copper, iron	Q	N	C	JA-45-36200
Copper, nickel, zinc	Q	N	B	-	WL23405	Copper, iron, nickel	Q	N	C	JA-45-36101
Copper, zinc, molybdenum	Q	N	B	-	JA-45-496	Copper, iron, lead, nickel, zinc	Q	N	C	JA-45-36204
Copper, zinc, lead, silver	Q	N	B	-	JA-45-448	Copper, iron, manganese, zinc	Q	N	C	JA-45-36105
Copper, zinc, lead, tin	Q	N	B	-	JA-45-438	Sodium, potassium	P	A	C	JA-45-36095
Gold, nickel	Q	N	B	-	JA-4S-433					
Gold, silver	Q	N	B	-	WL23269					
Indium, silver	Q	N	B	-	WL23294					
Lead, silver, zinc	Q	N	B	-	WL23171					
Magnesium, zinc	Q	N	B	-	WL23455					
Sodium, potassium	P	N	A	-	JA-45-439					
Sodium, potassium	P	A	A	-	WL23230					
Zinc, lead, tin	Q	N	B	-	WL23404					
Multiple-Element 36000 Series										
Aluminum, calcium, magnesium	Q	N	C		JA-45-36099					
Aluminum, calcium, magnesium, lithium	Q	N	C		JA-45-36250					

*Tubes listed in this table are issued by Fisher Scientific and produced by Westinghouse Electric.

¹P = Pyrex, Q = quartz, ²N = neon, A = argon, ³A = 1 1/2-in. diameter, B = 1-in. diameter, C = 2-in. diameter, ⁴WL = Westinghouse, JA = Jarrell-Ash.



FIGURE 44a Hollow-cathode spectral tubes described in Table 5.

		Transition elements															
Li	Be											B					
Na	Mg	Group 8										Al	Si	P	S		
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Gn	Ge	As	Se		
Rb	Sr	Y	Zr	Nb	Mo		Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te		
Cs	Ba	La	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi			
Lanthanides	Ce	Pr	Nd		Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb	Lu			
Actinides	Th		U														

FIGURE 44b Periodic table showing the prevalence of elements obtainable in hollow-cathode tubes.

arc lamp to produce stable radiation of sharp spectral lines; this makes it useful in spectroscopy and interferometry. The Hg 198 lamp makes a suitable secondary standard of radiation.

Spectral Lamps³¹ Some manufacturers produce groups of arc sources, which are similar in construction and filled with different elements and rare gases, and which yield discontinuous or monochromatic radiation throughout most of the ultraviolet and visible spectrum. They are called

spectral lamps. The envelopes of these lamps are constructed of glass or quartz, depending on the part of the spectrum desired. Thus, discrete radiation can be obtained from around 2300 Å into the near infrared. Figure 45³¹ represents the various atomic lines observable from Osram spectral lamps. Figure 46³³ gives a physical description of various spectral lamps obtainable from Philips. Table 6 lists the characteristics of the various types of lamps obtainable from Philips.

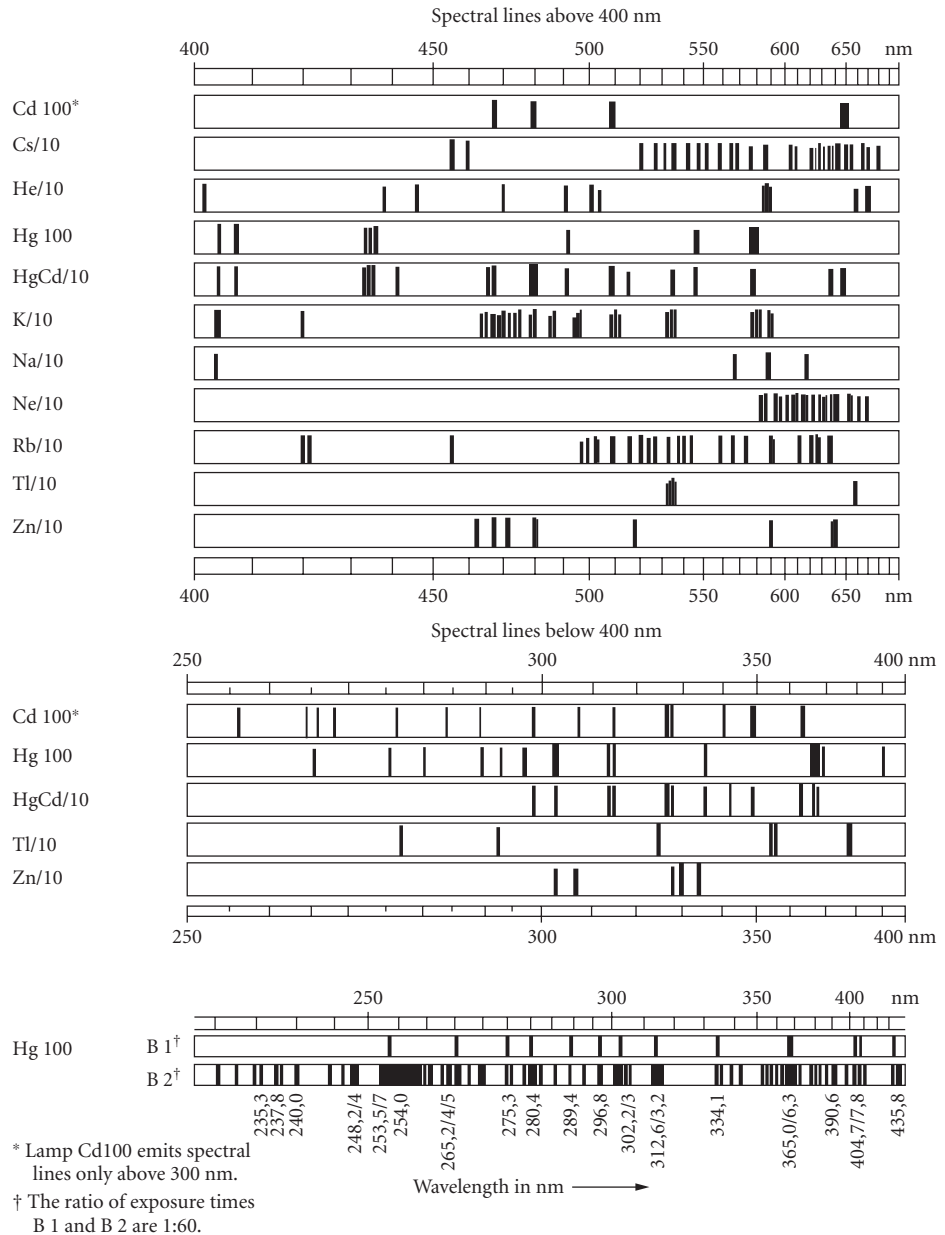


FIGURE 45 Spectral lamps.

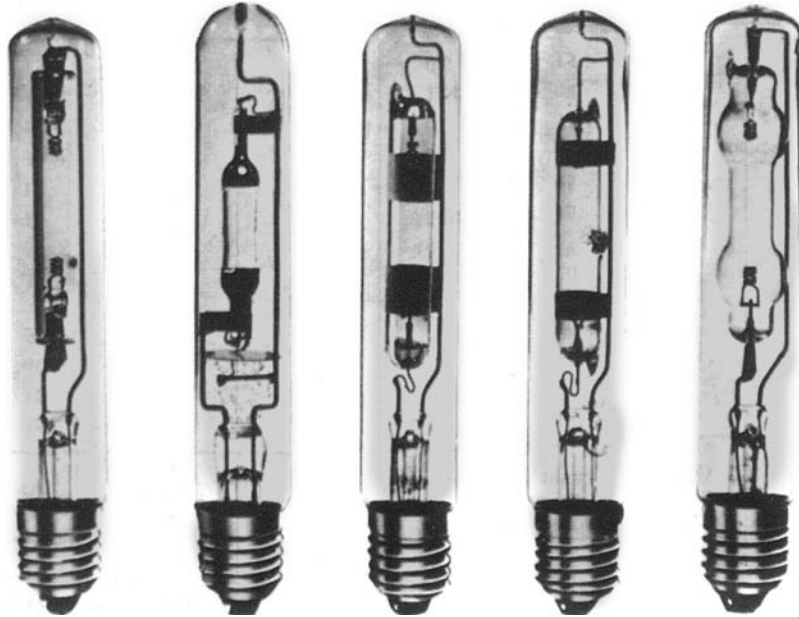


FIGURE 46 Examples of Philips spectral lamps.

TABLE 6 Specifications of Philips Spectral Lamps

Catalog Number	Symbols	Type	Material		Operating Current (A)	Wattage	Arc Length (mm)
			Burner	Envelope			
26-2709	Hg	Mercury (low-pressure)	Quartz	Glass	0.9	15	40
26-2717	Hg	Mercury (high-pressure)	Quartz	Glass	0.9	90	30
26-2725	Cd	Cadmium	Quartz	Glass	0.9	25	30
26-2733	Zn	Zinc	Quartz	Glass	0.9	25	30
26-2741	Hg, Cd, Zn	Mercury, cadmium, and zinc	Quartz	Glass	0.9	90	30
26-2758	He	Helium	Glass	Glass	0.9	45	32
26-2766	Ne	Neon	Glass	Glass	0.9	25	40
26-2774	A	Argon	Glass	Glass	0.9	15	40
26-2782	Kr	Krypton	Glass	Glass	0.9	15	40
26-2790	Xe	Xenon	Glass	Glass	0.9	10	40
26-2808	Na	Sodium	Glass	Glass	0.9	15	40
26-2816	Rb	Rubidium	Glass	Glass	0.9	15	40
26-2824	Cs	Caesium	Glass	Glass	0.9	10	40
26-2832	K	Potassium	Glass	Glass	0.9	10	40
26-2857	Hg	Mercury (low-pressure)	Quartz	Quartz	0.9	15	40
26-2865	Hg	Mercury (high-pressure)	Quartz	Quartz	0.9	90	30
26-2873	Cd	Cadmium	Quartz	Quartz	0.9	25	30
26-2881	Zn	Zinc	Quartz	Quartz	0.9	25	30
26-2899	Hg, Cd, Zn	Mercury, cadmium, and zinc	Quartz	Quartz	0.9	90	30
26-2907	In	Indium*	Quartz	Quartz	0.9	25	25
26-2915	Tl	Thallium	Quartz	Quartz	0.9	20	30
26-2923	Ga	Gallium	Quartz	Quartz	0.9	20	30

* Requires a Tesla coil to cause it to strike initially.

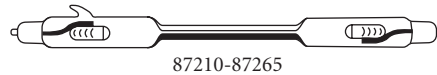


FIGURE 47 Physical construction of Pluecker spectrum tubes.

TABLE 7 Gas Fills in Plueker Tubes*

Cenco Number	Type
87210	Argon Gas
87215	Helium Gas
87220	Neon Gas
87225	Carbonic Acid Gas
87230	Chlorine Gas
87235	Hydrogen Gas
87240	Nitrogen Gas
87242	Air
87245	Oxygen Gas
87255	Iodine Vapor
87256	Krypton Gas
87258	Xenon
87260	Mercury Vapor
87265	Water Vapor

*Consists of glass tube with overall length of 25 cm with capillary portion about 8.5 to 10 cm long. Glass-to-metal seal wires are welded in metal caps with loops for wire connection are firmly sealed to the ends. Power supply no. 87208 is recommended as a source of excitation.

Pluecker Spectrum Tubes³⁴ These are inexpensive tubes made of glass (Fig. 47) with an overall length of 25 cm and capillary portion of 8.5 to 10 cm long. They operate from an ordinary supply with a special transformer which supports the tubes in a vertical position and maintains the voltage and current values adequate to operate the discharge and regulate the spectral intensity. Table 7 lists the various gases in available tubes.

Concentrated Arc Lamps

Zirconium Arc²⁵ The cathodes of these lamps are made of a hollow refractory metal containing zirconium oxide. The anode, a disk of metal with an aperture, resides directly above the cathode with the normal to the aperture coincident with the longitudinal axis of the cathode. Argon gas fills the tube. The arc discharge causes the zirconium to heat (to about 3000 K) and produce an intense, very small source of light. These lamps have been demonstrated in older catalogs from the Cenco Company in a number of wattages (from 2 to 300). The end of the bulk through which the radiation passes comes with ordinary curvature or (for a slight increase in price) flat. Examples are shown in Fig. 48.

Tungsten-Arc (Photomicrographic) Lamp²⁵ The essential elements of this discharge-type lamp (see Fig. 49) are a ring electrode and a pellet electrode, both made of tungsten. The arc forms between these electrodes, causing the pellet to heat incandescently. The ring also incandesces, but to a lesser extent. Thus, the hot pellet (approximately 3100 K) provides an intense source of small-area

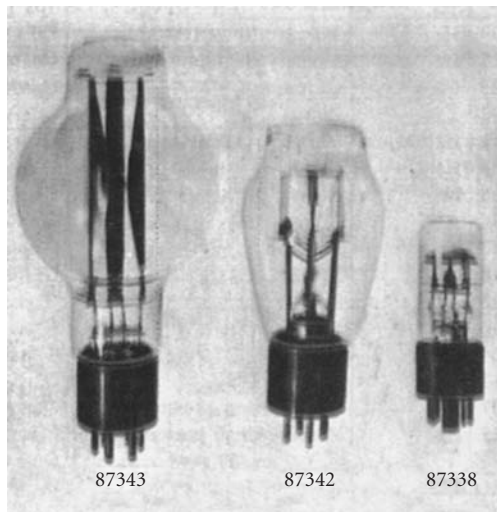


FIGURE 48 Physical construction of some zirconium arc lamps. Two 2-W lamps are available but not shown here.

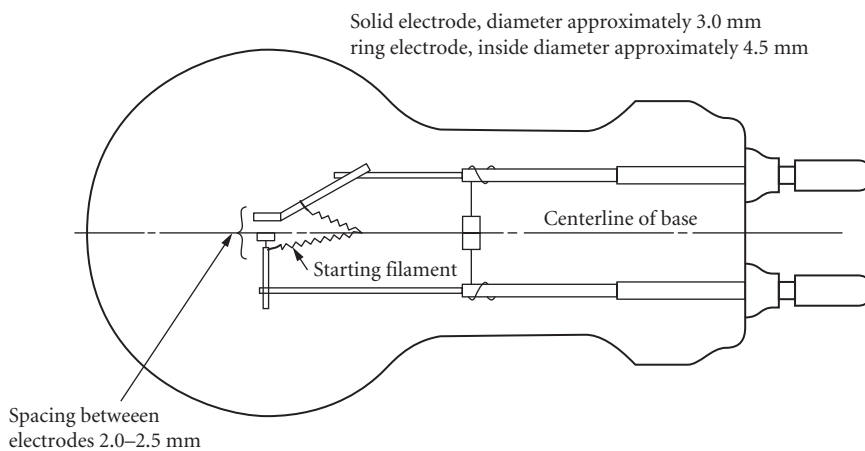


FIGURE 49 Construction of tungsten-arc lamp. The lamp must be operated baseup on a well-ventilated housing and using a special high-current socket which does not distort the position of the posts.

radiation. A plot of the spectral variation of this radiation is given in Fig. 50. As with all tungsten sources, evaporation causes a steady erosion of the pellet surface with the introduction of gradients, which is not serious if the pellet is used as a point source.

General Electric, manufacturer of the 30A/PS22 photomicrographic lamp, which uses a 30 Å operating current, states that this lamp requires a special heavy-duty socket obtainable through certain manufacturers suggested in its brochure, which may now be out of print in the original, but obtainable presumably as a copy from GE.

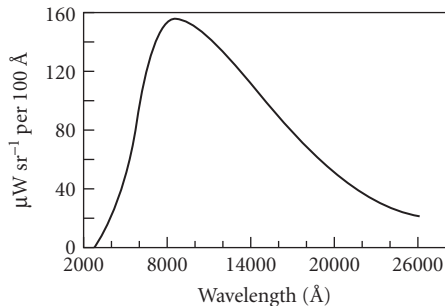


FIGURE 50 Spectral distribution of a 30/PS-22 photomicrographic lamp (superceded by the 330 watt, 30 amp PS-70).

Glow Modulator Tubes³⁵

According to technical data supplied by Sylvania, these are cold-cathode light sources uniquely adaptable to high-frequency modulation. (These tubes are now manufactured by The English Electric Valve Company, Elmsford, New York.) Pictures of two types are shown in Fig. 51. The cathode is a small hollow cylinder, and the high ionization density in the region of the cathode provides an intense source of radiation. Figure 52 is a graph of the light output as a function of tube current. Figure 53 is a graph depicting the response of the tube to a modulating input. The spectral outputs of a variety of tubes are shown in Fig. 54. Table 8 gives some of the glow-modulator specifications.

Hydrogen and Deuterium Arcs

For applications requiring a strong continuum in the ultraviolet region, the hydrogen arc at a few millimeters pressure provides a useful source. It can be operated with a cold or hot cathode. One hot-cathode type is shown in Fig. 55. Koller²² plots a distribution for this lamp down to about 200 Å.

Deuterium lamps (Fig. 56) provide a continuum in the ultraviolet with increased intensity over the hydrogen arc. Both lamps have quartz envelopes. The one on the left is designed for operation down to 2000 Å; the one on the right is provided with a Suprasil[®] window to increase the ultraviolet range down to 1650 Å. NIST is offering a deuterium lamp standard of spectral irradiance between 200 and 350 nm. The lamp output at 50 cm from its medium bipost base is about 0.7 W cm^{-3} at 200 nm and drops off smoothly to 0.3 W cm^{-3} at 250 nm and 0.07 W cm^{-3} at 350 nm. A working standard of the deuterium lamp can be obtained also, for example, from Optronic Laboratories, Incorporated, Orlando, Florida.

Other Commercial Sources

Activated-Phosphor Sources Of particular importance and convenience in the use of photometers are sources composed of a phosphor activated by radioactive substances. Readily available, and not subject to licensing with small quantities of radioactive material, are the ^{14}C -activated phosphor light sources. These are relatively stable sources of low intensity, losing about 0.02 percent per year due to the half-life of ^{14}C and the destruction of phosphor centers.

[®]Registered trademark of Heraeus-Amersil.

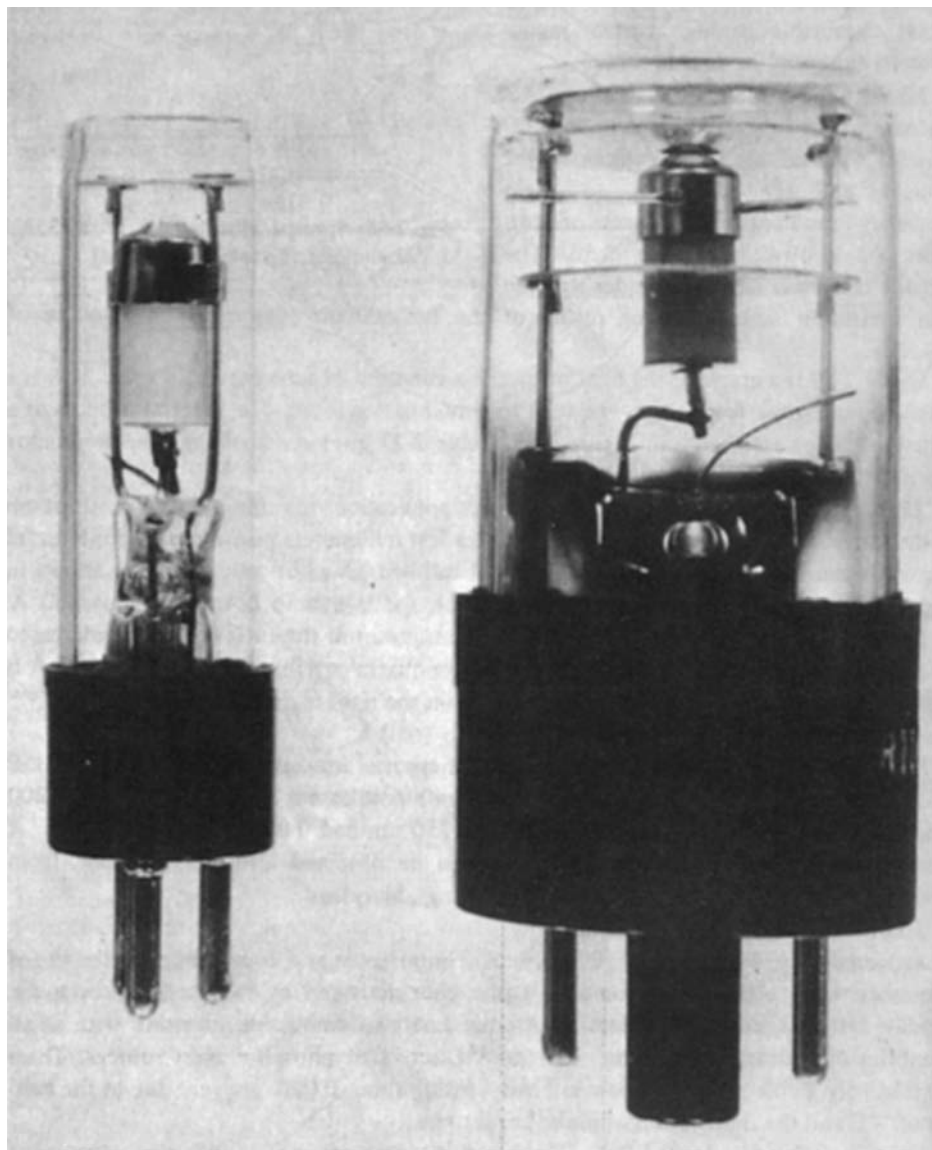


FIGURE 51 Construction of two glow modulator tubes.

Other (High-Energy) Sources Radiation at very high powers can be produced. Sources are synchrotrons, plasmotrons, arcs, sparks, exploding wires, shock tubes, and atomic and molecular beams, to name but a few. Among these, one can purchase in convenient, usable form precisely controlled spark-sources for yielding many joules of energy in a time interval of the order of microseconds. The number of vendors will be few, but check the directories.

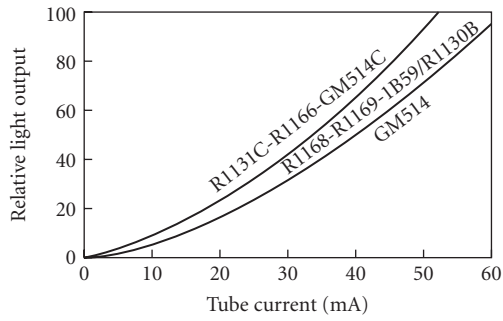


FIGURE 52 Variations of the light output from a glow modulator tube as a function of tube current.

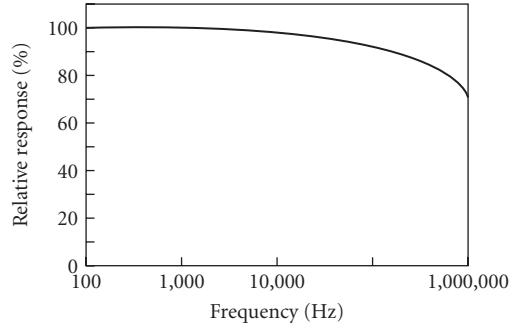
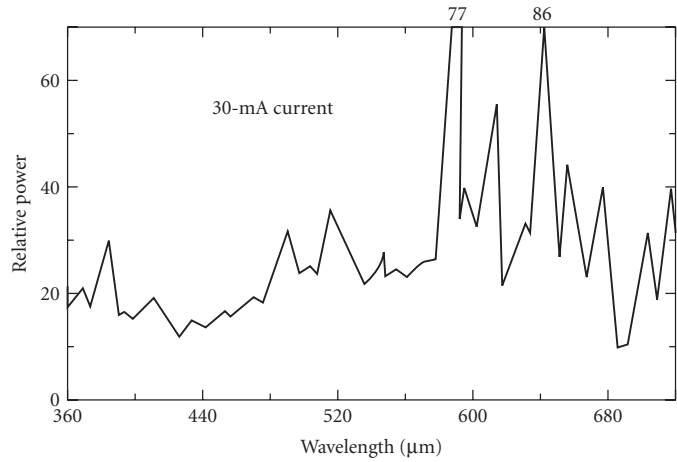
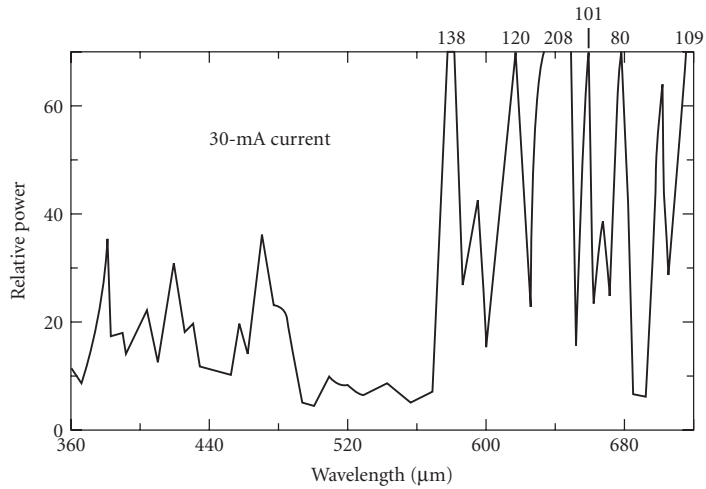


FIGURE 53 Response of the glow modulator tube to a modulating input.



(a) GM514C-R1166-R1131C



(b) R1168-R1169-1B59/R1130D-GM514

FIGURE 54 Spectral variation of the output of glow modulator tubes.

TABLE 8 Glow-Modulator Specifications

No.*	Maximum Operating Voltage	Current (mA)		Minimum Starting Voltage (V)	Crater Diameter (in.)	Approximate Light Center Length (in.)	Light Output (cd)	Brightness (cd in. ⁻²)	Rated Life (h)	Base Type	Bulb Type	Maximum Overall Length (in.)	Maximum Diameter (in.)	Color of Discharge
		Average	Peak											
GM-514	160	5-25	55	240	0.056	1-3/4	0.1 at 25 mA	41 at 25 mA	100 at 15 mA	3-pin miniature [†]	T-4 1/2	2-5/8	4 1/64	Blue-red
GM-514C	160	5-15	35	240	0.093	1-3/4	0.1 at 15 mA	15 at 15 mA	25 at 10 mA	3-pin miniature [‡]	T-4 1/2	2-5/8	4 1/64	White
IB59/ R-1130D	150	5-35	75	225	0.056	2	0.13 at 30 mA	43 at 30 mA	250 at 20 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red
R-1131C	150	3-25	55	225	0.093	2	0.2 at 25 mA	29 at 25 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	White
R-1166	150	3-25	55	225	0.093	2	0.2 at 25 mA	29 at 25 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	White
R-1168	150	5-15	30	225	0.015	2	0.023 at 15 mA	132 at 15 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red
R-1169	150	5-25	45	225	0.025	2	0.036 at 15 mA	72 at 15 mA	250 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red

[†]Type R-1166 is opaque-coated with the exception of a circle 3/8 in. in diameter at end of lamp. All other types have clear-finish bulb.

[‡]Pins 1 and 3 arc anode; pin 2 cathode.

[§]Pin 7 anode; pin 3 cathode.

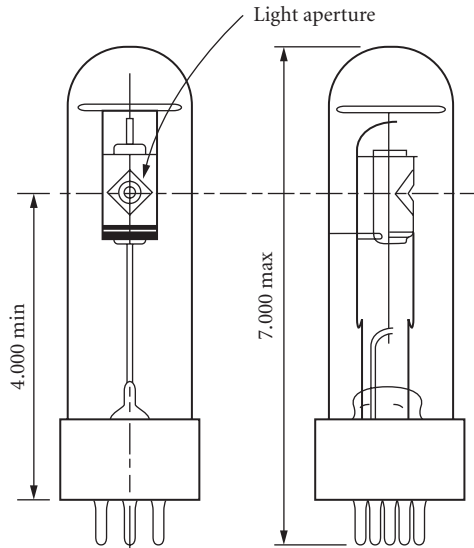


FIGURE 55 Hydrogen-arc lamp.

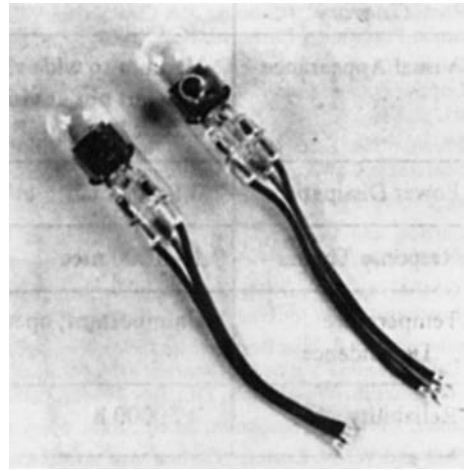


FIGURE 56 Two types of deuterium arc lamps.

Other Special Sources An enormous number of special-purpose sources are obtainable from manufacturers and scientific instrument suppliers. One source that remains to be mentioned is the so-called miniature, sub- and microminiature lamps. These are small, even tiny, incandescent bulbs of glass or quartz, containing tungsten filaments. They serve excellently in certain applications where small, intense radiators of visible and near-infrared radiation are needed. Second-source vendors advertise in the trade magazines.

15.6 REFERENCES

1. A. J. LaRocca, "Artificial Sources," in *The Infrared Handbook*, rev. ed., ONR, Washington, D.C., 1985, chap. 2.
2. A. G. Worthing, "Sources of Radiant Energy," in W. E. Forsythe (ed.), *Measurement of Radiant Energy*, McGraw-Hill, New York, 1937, chap. 2.
3. F. K. Richtmeyer and E. H. Kennard, *Introduction to Modern Physics*, McGraw-Hill Book Company, New York, 1947, p. 159.
4. T. J. Quinn, "The Calculation of the Emissivity of Cylindrical Cavities Giving Near Back-body Radiation," *British Journal of Applied Physics*, vol. 18, 1967, p. 1105.
5. J. C. DeVos, "Evaluation of the Quality of a Conical Blackbody," *Physica*, North Holland Publishing, Amsterdam, Netherlands, vol. 20, 1954, p. 669.
6. K. Irani, "Theory and Construction of Blackbody Calibration Sources," *Proceedings of SPIE*, vol. 4360, March 2001, p. 347.
7. Andre Gouffé, "Corrections d'Ouverture des Corps-noir Artificiels Compte Tenu des Diffusions Multiples Internes," in *Revue d'Optique*, vol. 24, Masson, Paris, 1945, p. 1.
8. A. Leupin, H. Vetsch, and F. K. Kneibuhl, "Investigation, Comparison, and Improvement of Technical Infrared Radiators," *Infrared Physics*, vol. 30, no. 3, 1990, pp. 199–258.
9. J. W. T. Walsh, *Photometry*, 3d ed., Dover, New York, 1965.
10. K. D. Mielenz, R. D. Saunders, and J. B. Shumaker, "Spectroradiometric Determination of the Freezing Temperature of Gold," *Journal of Research of the National Institute of Standards and Technology*, vol. 95, no. 1, Jan.–Feb. 1990, pp. 49–67.

11. *Spectral Radiance Calibrations*. NBS Special Publication, Jan. 1987, p. 250–251.*
12. *NIST Calibration Services Users Guide 1989*, NIST Special Publication, 1989, p. 250.*
13. *Lasers and Optronics 1990 Buying Guide*, Elsevier Communications, Morris Plains, N.J., 1990.*
14. *Photonics Directory of Optical Industries*, Laurin Publishing Co., Pittsfield, Mass., 1994.*
15. W. Y. Ramsey and J. C. Alishouse, “A Comparison of Infrared Sources,” *Infrared Physics*, vol. 8, Pergamon Publishing, Elmsford, N.Y., 1968, p. 143.
16. J. C. Morris, “Comments on the Measurements of Emittance of the Globar Radiation Source,” *Journal of the Optical Society of America*, vol. 51, Optical Society of America, Washington, D.C., July 1961, p. 798.
17. A. H. Pfund, “The Electric Welsbach Lamp,” *Journal of the Optical Society of America*, vol. 26, Optical Society of America, Washington, D.C., Dec. 1936, p. 439.
18. J. Strong, *Procedures in Experimental Physics*, Prentice-Hall, New York, 1938, p. 346.
19. F. E. Carlson and C. N. Clark, “Light Sources for Optical Devices,” in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. 1, Academic Press, New York, 1965, p. 80.
20. G. M. B. H. Osram, *Lamps for Scientific Purposes*, Munchen, West Germany, 1966. (See also the later Osram brochure, *Light for Cine Projection, Technology and Science*, 1987.)*
21. F. J. Studer and R. F. Van Beers, “Modification of Spectrum of Tungsten Filament Quartz-Iodine Lamps due to Iodine Vapor,” *Journal of the Optical Society of America*, vol. 54, no. 7, Optical Society of America, Washington, D.C., July 1964, p. 945.
22. L. R. Koller, *Ultraviolet Radiation*, 2d ed., John Wiley and Sons, New York, 1965.
23. M. R. Null and W. W. Lozier, “Carbon Arc as a Radiation Standard,” *Journal of the Optical Society of America*, vol. 52, no. 10, Optical Society of America, Washington, D.C., Oct. 1962, pp. 1156–1162.
24. E. B. Noel, “Radiation from High Pressure Mercury Arcs,” *Illuminating Engineering*, vol. 36, 1941, p. 243.
25. General Electric, *Bulletin TP-109R*, Cleveland, 1975.*
26. Illumination Industries, *Catalog No. 108-672-3M*, Sunnyvale, Calif., 1972.*
27. M. W. P. Cann, *Light Sources for Remote Sensing Systems*, in NASA-CR-854, Aug. 1967.*
28. Varian Associates, Palo Alto, Calif., 1969.*
29. Fisher Scientific, *Special Catalog to Spectrophotometer Users*, Pittsburgh, Mass., 1972.*
30. E. F. Worden, R. G. Gutmacher, and J. F. Conway, “Use of Electrodeless Discharge Lamps in the Analysis of Atomic Spectra,” *Applied Optics*, vol. 2, no. 7, Optical Society of America, Washington, D.C., July 1963, pp. 707–713.
31. W. F. Meggers and F. O. Westfall, “Lamps and Wavelengths of Mercury 198,” *NBS Journal of Research*, vol. 44, National Bureau of Standards, Washington, D.C., 1950, pp. 447–55.
32. W. F. Meggers, “Present Experimental Status of Rare Earth Spectra,” *Journal of the Optical Society of America*, vol. 50, Optical Society of America, Washington, D.C., 1960, p. 405.
33. Ealing Corporation, *Ealing Catalog, Optical Components Section*, South Natick, Mass., 1976–77. (See also *Optical Services Supplement*, no. 1, 1969–70, p. 26; See also *Ealing Electro-optics Product Guide*, 1990.)*
34. Central Scientific, *Cenco Scientific Education Catalog*, Physics-Light Section, Chicago, 1975, p. 602.*
35. GTE Sylvania, *Special Purpose Lamps*, Lighting Products Group, Danvers, Mass., in TR-29R, May 1966.*

*Many of these references are likely to be inaccessible as shown, either because they are out-of-date or perhaps otherwise obsolete. They are retained here because there are up-to-date versions of many of them, and the user is advised to use them as starting points in a search of the Internet for current information. Most of the companies still exist and have more recent catalogs that are available.

William T. Silfvast

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

16.1 GLOSSARY

A_2	radiative transition probability from level 2 to all other possible lower-lying levels
A_{21}	radiative transition probability from level 2 to level 1
a_L	scattering losses within a laser cavity for a single pass through the cavity
B_{12}	Einstein B coefficient associated with absorption
B_{21}	Einstein B coefficient associated with stimulated emission
E_1, E_2	energies of levels 1 and 2 above the ground state energy for that species
g_1, g_2	stability parameters for laser modes when describing the laser optical cavity
g_1, g_2	statistical weights of energy levels 1 and 2 that indicate the degeneracy of the levels
g_{21}	gain coefficient for amplification of radiation within a medium at a wavelength of λ_{21}
I_{sat}	saturation intensity of a beam in a medium; intensity at which exponential growth will cease to occur even though the medium has uniform gain (energy/time-area)
N_1, N_2	population densities (number of species per unit volume) in energy levels 1 and 2
r_c	radius of curvature of the expanding wavefront of a gaussian beam
R_1, R_2	reflectivities of mirrors 1 and 2 at the desired wavelength
T_1	lifetime of a level when dominated by collisional decay
T_2	average time between phase-interrupting collisions of a species in a specific excited state
t_{opt}	optimum mirror transmission for a laser of a given gain and loss
$w(z)$	beam waist radius at a distance z from the minimum beam waist for a gaussian beam
w_0	minimum beam waist radius for a gaussian mode
α_{12}	absorption coefficient for absorption of radiation within a medium at wavelength λ_{21}
γ_{21}	angular frequency bandwidth of an emission or absorption line
Δt_p	pulse duration of a mode-locked laser pulse
$\Delta\nu$	frequency bandwidth over which emission, absorption, or amplification can occur

$\Delta\nu_D$	frequency bandwidth (FWHM) when the dominant broadening process is Doppler or motional broadening
η	index of refraction of the laser medium at the desired wavelength
λ_{21}	wavelength of a radiative transition occurring between energy levels 2 and 1
ν_{21}	frequency of a radiative transition occurring between energy levels 2 and 1
σ_{21}	stimulated emission cross section (area) associated with levels 2 and 1
σ_{21}^D	stimulated emission cross section at line center when Doppler broadening dominates (area)
σ_{21}^H	stimulated emission cross section at line center when homogeneous broadening dominates (area)
τ_2	lifetime of energy level 2
τ_{21}	lifetime of energy level 2 if it can only decay to level 1

16.2 INTRODUCTION

A laser is a device that amplifies light and produces a highly directional, high-intensity beam that typically has a very pure frequency or wavelength. It comes in sizes ranging from approximately one-tenth the diameter of a human hair to the size of a very large building, in powers ranging from 10^{-9} to 10^{20} W and in wavelengths ranging from the microwave to the soft-x-ray spectral regions with corresponding frequencies from 10^{11} to 10^{17} Hz. Lasers have pulse energies as high as 10^4 J and pulse durations as short as 6×10^{-15} seconds. They can easily drill holes in the most durable of materials and can weld detached retinas within the human eye.

Lasers are a key component of some of our most modern communication systems and are the “phonograph needle” of compact disc players. They are used for heat treatment of high-strength materials, such as the pistons of automobile engines, and provide a special surgical knife for many types of medical procedures. They act as target designators for military weapons and are used in the checkout scanners we see everyday at the supermarket.

The word *laser* is an acronym for *Light Amplification by Stimulated Emission of Radiation*. The laser makes use of processes that increase or amplify light signals after those signals have been generated by other means. These processes include (1) stimulated emission, a natural effect that arises out of considerations relating to thermodynamic equilibrium, and (2) optical feedback (present in most lasers) that is usually provided by mirrors. Thus, in its simplest form, a laser consists of a gain or amplifying medium (where stimulated emission occurs) and a set of mirrors to feed the light back into the amplifier for continued growth of the developing beam (Fig. 1).

The entire spectrum of electromagnetic radiation is shown in Fig. 2, along with the region covered by currently existing lasers. Such lasers span the wavelength range from the far infrared part of

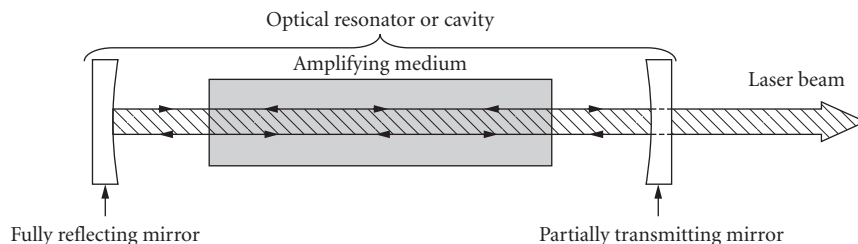


FIGURE 1 Simplified diagram of a laser, including the amplifying medium and the optical resonator.

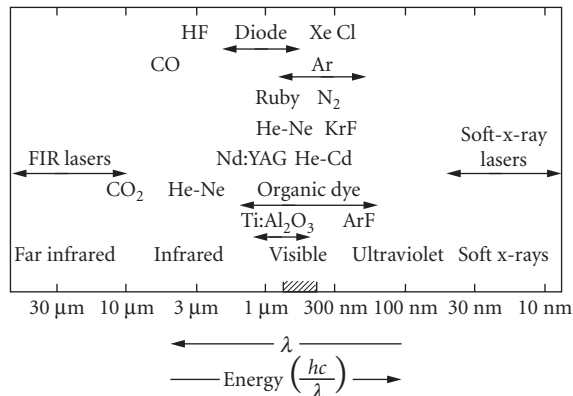


FIGURE 2 The portion of the electromagnetic spectrum that involves lasers, along with the general wavelengths of operation of most of the common lasers.

the spectrum ($\lambda = 1000 \mu\text{m}$) to the soft-x-ray region ($\lambda = 3 \text{ nm}$), thereby covering a range of almost six orders of magnitude! There are several types of units that are used to define laser wavelengths. These range from micrometers (μm) in the infrared to nanometers (nm) and angstroms (\AA) in the visible, ultraviolet (UV), vacuum ultraviolet (VUV), extreme ultraviolet (EUV or XUV), and soft-x-ray (SXR) spectral regions.

This chapter provides a brief overview of how a laser operates. It considers a laser as having two primary components: (1) a region where light amplification occurs which is referred to as a *gain medium* or an *amplifier*, and (2) a *cavity*, which generally consists of two mirrors placed at either end of the amplifier.

Properties of the amplifier include the concept of discrete excited energy levels and their associated finite lifetimes. The broadening of these energy levels will be associated with the emission linewidth which is related to decay of the population in these levels.

Stimulated emission will be described, and the formulas for calculating the amount of gain that can occur via stimulated emission will be given in terms of the radiative properties of the medium. The concept of the saturation intensity will be introduced and related to the amount of gain that is necessary for laser output. The addition of mirrors at the ends of the amplifier will be used to increase the gain length and to reduce the divergence of the amplified beam. The threshold conditions for laser output will be described in terms of the amplifier properties and the mirror reflectivities. This section will conclude with a review of excitation or pumping processes that are used to produce the necessary population density in the upper laser level.

Cavity properties will begin with a discussion of both longitudinal and transverse cavity modes which provide the laser beam with a gaussian-shaped transverse profile. The properties of those gaussian beams will be reviewed. The types of optical cavities that allow stable operation of laser modes will then be described. A number of special types of laser cavity arrangements and techniques will be reviewed, including unstable resonators, Q-switching, mode-locking, and ring lasers. A brief review will then be given of the various common types of gaseous, liquid, and solid-state lasers.

Additional information related to spectral lineshape and the mechanisms of spectral broadening can be found in Chap. 10 (Vol. I), "Optical Spectroscopy and Spectroscopic Lineshapes." Other related material can be found in Chap. 8 (Vol. IV), "Fundamental Optical Properties of Solids," and Chap. 33 (Vol. I), "Holography and Holographic Instruments." As lasers are widely used in many of the devices and techniques discussed in other chapters in this *Handbook*, the reader is directed to those topics for information on specific lasers.

16.3 LASER PROPERTIES ASSOCIATED WITH THE LASER GAIN MEDIUM

Energy Levels and Radiation^{1,2}

Nearly all lasers involve electronic charge distributions of atoms, molecules, organic dye solutions, or solids that make transitions from one energy state or level E_2 to another lower-lying level E_1 . The loss of energy resulting from this transition is given off in the form of electromagnetic radiation. The relationship between the energy difference between the levels, $E_2 - E_1$ or ΔE_{21} , and the frequency ν_{21} of radiation occurring as a result of the transition, is determined by the Einstein relationship $E_{21} = h\nu_{21}$ where h is Planck's constant. It was first shown by Bohr in 1913 that the discrete set of emission wavelengths from a hydrogen discharge could be explained by the occurrence of discrete energy levels in the hydrogen atom that have a fixed relationship. This discrete arrangement of energy levels was later shown to occur in other atoms, in molecules, and also in liquids and solids. In atoms these energy levels are very precisely defined and narrow in width ($\approx 10^9$ Hz) and can be accurately calculated with sophisticated atomic physics codes. In molecules and high-density materials the locations of the levels are more difficult to calculate and they tend to be much broader in width, the largest widths occurring in liquids and solids (up to 5×10^{13} Hz).

The lowest energy level of a species is referred to as the *ground state* and is usually the most stable state of the species. There are some exceptions to this, for example, ground states of ionized species or unstable ground states of some molecular species such as excimer molecules. Energy levels above the ground state are inherently unstable and have lifetimes that are precisely determined by the arrangement of the atoms and electrons associated with any particular level as well as to the particular species or material. Thus, when an excited state is produced by applying energy to the system, that state will eventually decay by emitting radiation over a time period ranging from 10^{-15} seconds or less to times as long as seconds or more, depending upon the particular state or level involved. For *strongly allowed* transitions that involve the electron charge cloud changing from an atomic energy level of energy E_2 to a lower-lying level of energy E_1 , the radiative decay time τ_{21} can be approximated by $\tau_{21} \approx 10^4 \lambda_{21}^2$ where τ_{21} is in seconds and λ_{21} is the wavelength of the emitted radiation in meters. λ_{21} is related to ν_{21} by the relationship $\lambda_{21}\nu_{21} = c/\eta$ where c is the velocity of light (3×10^8 m/s) and η is the index of refraction of the material. For most gases, η is near unity and for solids and liquids it ranges between 1 and 10, with most values ranging from ≈ 1.3 to 2.0.

Using the expression suggested above for the approximate value of the lifetime of an excited energy level, one obtains a decay time of several ns for green light ($\lambda_{21} = 5 \times 10^{-7}$ m). This represents a minimum radiative lifetime since most excited energy levels have a weaker radiative decay probability than mentioned above and would therefore have radiative lifetimes one or two orders of magnitude longer. Other laser materials such as molecules, organic dye solutions, and semiconductor lasers have similar radiative lifetimes. The one exception is the class of dielectric solid-state laser materials (both crystalline and glass) such as ruby and Nd:YAG, in which the lifetimes are of the order of 1 μ s to 3 ms. This much longer radiative lifetime in solid-state laser materials is due to the nature of the particular state of the laser species and to the crystal matrix in which it is contained. This is a very desirable property for a laser medium since it allows excitation and energy storage within the laser medium over a relatively long period of time.

Emission Linewidth and Line Broadening of Radiating Species^{1,3}

Assume that population in energy level 2 decays to energy level 1 with an exponential decay time of τ_{21} and emits radiation at frequency ν_{21} during that decay. It can be shown by Fourier analysis that the exponential decay of that radiation requires the frequency width of the emission to be of the order of $\Delta\nu \approx 1/2\pi\tau_{21}$. This suggests that the energy width ΔE_2 of level 2 is of the order of $\Delta E_2 = h/2\pi\tau_{21}$. If the energy level 2 can decay to more levels than level 1, with a corresponding decay time of τ_2 , then its energy is broadened by an amount $\Delta E_2 = h/2\pi\tau_2$. If the decay is due primarily to

radiation at a rate A_{2i} to one or more individual lower-lying levels i , then $1/\tau_2 = A_2 = \sum_i A_{2i}$. A_2 represents the total radiative decay rate of level 2, whereas A_{2i} is the specific radiative decay rate from level 2 to a lower-lying level i .

If population in level 2 decays radiatively at a radiative rate A_2 and population in level 1 decays radiatively at a rate A_1 , then the emission linewidth of radiation from level 2 to 1 is given by

$$\Delta\nu_{21} = \frac{\sum_i A_{2i} + \sum_j A_{1j}}{2\pi} \quad (1)$$

which is referred to as the *natural linewidth* of the transition and represents the sum of the widths of levels 2 and 1 in frequency units. If, in the above example, level 1 is a ground state with infinite lifetime or a long-lived metastable level, then the natural linewidth of the emission from level 2 to level 1 would be represented by

$$\Delta\nu_{21} = \frac{\sum_i A_{2i}}{2\pi} \quad (2)$$

since the ground state would have an infinite lifetime and would therefore not contribute to the broadening. This type of linewidth or broadening is known as *natural broadening* since it results specifically from the radiative decay of a species. Thus the natural linewidth associated with a specific transition between two levels has an inherent value determined only by the factors associated with specific atomic and electronic characteristics of those levels.

The emission-line broadening or natural broadening described above is the minimum line broadening that can occur for a specific radiative transition. There are a number of mechanisms that can increase the emission linewidth. These include collisional broadening, phase-interruption broadening, Doppler broadening, and isotope broadening. The first two of these, along with natural broadening, are all referred to as *homogeneous broadening*. Homogeneous broadening is a type of emission broadening in which all of the atoms radiating from the specific level under consideration participate in the same way. In other words, all of the atoms have the identical opportunity to radiate with equal probability.

The type of broadening associated with either Doppler or isotope broadening is referred to as *inhomogeneous broadening*. For this type of broadening, only certain atoms radiating from that level that have a specific property such as a specific velocity, or are of a specific isotope, participate in radiation at a certain frequency within the emission bandwidth.

Collisional broadening is a type of broadening that is produced when surrounding atoms, molecules, solvents (in the case of dye lasers), or crystal structures interact with the radiating level and cause the population to decay before it has a chance to decay by its normal radiative processes. The emission broadening is then associated with the faster decay time T_1 , or $\Delta\nu = 1/2\pi T_1$.

Phase-interruption broadening or *phonon broadening* is a type of broadening that does not increase the decay rate of the level, but it does interrupt the phase of the rotating electron cloud on average over a time interval T_2 which is much shorter than the radiative decay time τ_2 which includes all possible radiative decay channels from level 2. The result of this phase interruption is to increase the emission linewidth beyond that of both natural broadening and T_1 broadening (if it exists) to an amount $\Delta\nu = 1/2\pi T_2$.

Doppler broadening is a type of inhomogeneous broadening in which the Doppler effect shifts the frequencies of radiating atoms moving toward the observer to a higher value and the frequencies of atoms moving away from the observer to a lower value. This effect occurs only in gases since they are the only species that are moving fast enough to produce such broadening. Doppler broadening is the dominant broadening process in most visible gas lasers. The expression for the Doppler linewidth (FWHM) is given by

$$\Delta\nu_D = 7.16 \cdot 10^{-7} \nu_o \sqrt{\frac{T}{M}} \quad (3)$$

in which ν_o is the center frequency associated with atoms that are not moving either toward or away from the observer, T is the gas temperature in kelvin and M is the atomic or molecular weight (number of nucleons/atom or molecule) of the gas atoms or molecules.

Isotope broadening also occurs in some gas lasers. It becomes the dominant broadening process if the specific gas consists of several isotopes of the species and if the isotope shifts for the specific radiative transition are broader than the Doppler width of the transition. The helium-cadmium laser is dominated more by this effect than any other laser since the naturally occurring cadmium isotopic mixture contains eight different isotopes and the isotope shift between adjacent isotopes (adjacent neutron numbers) is approximately equal to the Doppler width of the individual radiating isotopes. This broadening effect can be eliminated by the use of isotopically pure individual isotopes, but the cost for such isotopes is often prohibitive.

All homogeneous broadening processes have a frequency distribution that is described by a Lorentzian mathematical function

$$I_{21}(\nu) = I_o \frac{\gamma_{21}/4\pi^2}{(\nu - \nu_o)^2 + (\gamma_{21}/4\pi)^2} \quad (4)$$

in which γ_{21} represents the decay rate of level 2, I_o is the total emission intensity of the transition over the entire linewidth, and ν_o is the center frequency of the emission line. In Eq. (4), γ_{21} is determined by the relationship $\gamma_{21} = 2\pi\Delta\nu_{21}$. For natural broadening, $\Delta\nu_{21}$ is given by either Eq. (1) or Eq. (2), whichever is applicable. For T_1 -dominated broadening, $\Delta\nu_{21} = 1/2\pi T_1$, and for T_2 -dominated broadening, $\gamma_{21} = 1/2\pi T_2$.

The frequency distribution for Doppler broadening is described by a gaussian function

$$I(\nu) = \frac{2(\ln 2)^{1/2}}{\pi^{1/2}\Delta\nu_D} I_o \exp\left[-\frac{4\ln 2(\nu - \nu_o)^2}{\Delta\nu_D^2}\right] \quad (5)$$

Both of these lineshape functions are indicated in Fig. 3. In this figure, both the total emission intensity integrated over all frequencies and the emission linewidth (full width at half maximum or FWHM) for both functions are identical.

Isotope broadening involves the superposition of a series of either lorentzian shapes or gaussian shapes for each isotope of the species, separated by the frequencies associated with the isotope shifts of that particular transition.

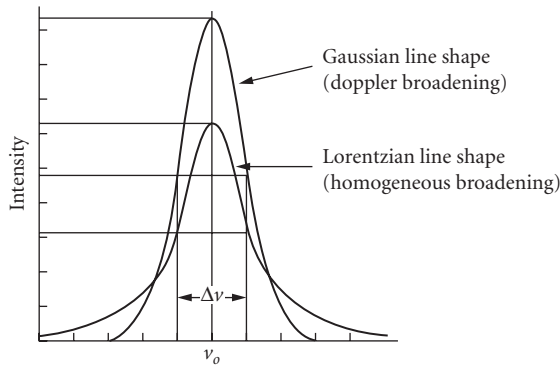


FIGURE 3 Lineshape functions for both homogeneous broadening, with a lorentzian shape, and Doppler broadening (inhomogeneous), with a gaussian shape. Both lines are arranged with equal linewidths (FWHM) and equal total intensities.

TABLE 1 Amplifier Parameters for a Wide Range of Lasers

Type of Laser	λ_{21} (nm)	τ_2 (s)	$\Delta\nu_{21}$ (Hz)	$\Delta\lambda_{21}$ (nm)	σ_{21} (m ²)	ΔN_{21} (m ⁻³)	L (m)	g_{21} (m ⁻¹)
Helium-Neon	632.8	3×10^{-7}	2×10^9	2×10^{-3}	3×10^{-17}	5×10^{15}	0.2	0.15
Argon	488.0	1×10^{-8}	2×10^9	1.6×10^{-3}	5×10^{-16}	1×10^{15}	0.2–1.0	0.5
He-Cadmium	441.6	7×10^{-7}	2×10^9	1.3×10^{-3}	8×10^{-18}	4×10^{16}	0.2–1.0	0.3
Copper vapor	510.5	5×10^{-7}	2×10^9	1.3×10^{-3}	8×10^{-18}	6×10^{17}	1.0–2.0	5
CO ₂	10,600	4	6×10^7	2.2×10^{-2}	1.6×10^{-20}	5×10^{19}	0.2–2.0	0.8
Excimer	248.0	9×10^{-9}	1×10^{13}	2	2.6×10^{-20}	1×10^{20}	0.5–1.0	2.6
Dye (Rh6G)	577.0	5×10^{-9}	5×10^{13}	60	1.2×10^{-20}	2×10^{22}	0.01	240
Semiconductor	800	1×10^{-9}	1×10^{13}	20	1×10^{-19}	10^{24}	0.00025	100,000
Nd:Yag	1064.1	2.3×10^{-4}	1.2×10^{11}	0.4	6.5×10^{-23}	1.6×10^{23}	0.1	10
Nd:Glass	1054	3.0×10^{-4}	7.5×10^{12}	26	4.0×10^{-24}	8×10^{23}	0.1	3
Cr:LiSAF	840	6.7×10^{-5}	9.0×10^{13}	250	5.0×10^{-24}	2×10^{24}	0.1	10
Ti:Al ₂ O ₃	760	3.2×10^{-6}	1.5×10^{14}	400	4.1×10^{-23}	5×10^{23}	0.1	20

Table 1 gives examples of the dominant broadening process and the value of the broadening for most of the common commercial lasers.

Bandwidths of laser transitions in semiconductors are actually made narrower by making the active regions of the material extremely thin in one or more dimensions. In doing so the energy levels become quantized and thus behave more like single atom atomic levels. Quantizing in one dimension, by making the thickness of the order of 50 to 100 nm, leads to a quantum well laser. Narrowing and thereby quantizing in two dimensions is referred to as a quantum wire and in three dimensions, a quantum dot. The advantages of quantizing the dimensions is that the reduced thickness leads to a significantly reduced heat loss during the excitation of the semiconductor as well as a narrower laser emission linewidth because of the smaller size of the electron energy distribution in the upper laser level. In the case of the quantum dot, the material takes on atom-like properties because the energies are quantized in all three dimensions and the lasing threshold is reduced much more so than even with the quantum well laser. Of course there is less laser gain medium produced in such materials per unit volume, because of the reduced gain volume, and thereby less laser power per unit volume. This can be made up by having many such gain media in parallel and/or in series, taking into account the heat removal requirements of the pumping process.

Stimulated Radiative Processes— Absorption and Emission^{1,2}

Two types of stimulated radiative processes, absorption and stimulated emission, occur between energy levels 1 and 2 of a gain medium when light of frequency ν_{21} corresponding to an energy difference $\Delta E_{21} = (E_2 - E_1) = h\nu_{21}$ passes through the medium. These processes are proportional to the light intensity I as indicated in Fig. 4 for a two-level system as well as to the stimulated absorption and emission coefficients B_{12} and B_{21} , respectively. These coefficients are related to the frequency ν_{21} and the spontaneous emission probability A_{21} associated with the two levels. A_{21} has units of (1/s).

Absorption results in the loss of light of intensity I when the light interacts with the medium. The energy is transferred from the beam to the medium by raising population from level 1 to the higher-energy level 2. In this situation, the species within the medium can either reradiate the energy and return to its initial level 1, it can reradiate a different energy and decay to a different level, or it can lose the energy to the surrounding medium via collisions, which results in the heating of the medium, and return to the lower level. The absorption probability is proportional to the intensity I which has units of energy/s-m² times B_{12} , which is the absorption probability coefficient for that transition. B_{12} is one of the Einstein B coefficients and has the units of m³/energy-s².

Stimulated emission results in the increase in the light intensity I when light of the appropriate frequency ν_{21} interacts with population occupying level 2 of the gain medium. The energy is given

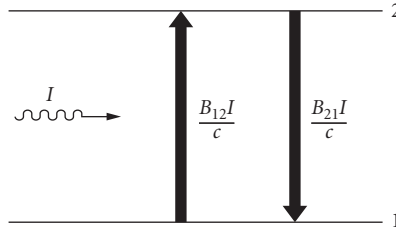


FIGURE 4 Stimulated emission and absorption processes that can occur between two energy levels, 1 and 2, and can significantly alter the population densities of the levels compared to when a beam of intensity I is not present.

up by the species to the radiation field. In the case of stimulated emission, the emitted photons or bundles of light have exactly the same frequency ν_{21} and direction as the incident photons of intensity I that produce the stimulation. B_{21} is the associated stimulated emission coefficient and has the same units as B_{12} . It is known as the other Einstein B coefficient.

Einstein showed the relationship between B_{12} , B_{21} , and A_{21} as

$$g_2 B_{21} = g_1 B_{12} \quad (6)$$

and

$$A_{21} = \frac{8\pi h \nu_{21}^3}{c^3} B_{21} \quad (7)$$

where g_2 and g_1 are the statistical weights of levels 2 and 1 and h is Planck's constant. Since $A_{21} = 1/\tau_{21}$ for the case where radiative decay dominates and where there is only one decay path from level 2, B_{21} can be determined from lifetime measurements or from absorption measurements on the transition at frequency ν_{21} .

Population Inversions^{1,2}

The two processes of absorption and stimulated emission are the principal interactions involved in a laser amplifier. Assume that a collection of atoms of a particular species is energized to populate two excited states 1 and 2 with population densities N_1 and N_2 (number of species/m³) and state 2 is at a higher energy than state 1 by an amount ΔE_{21} as described in the previous section. If a photon beam of energy ΔE_{21} , with an intensity I_0 and a corresponding wavelength $\lambda_{21} = c/\nu_{21} = hc/\Delta E_{21}$, passes through this collection of atoms, then the intensity I after the beam emerges from the medium can be expressed as

$$I = I_0 e^{\sigma_{21}(N_2 - (g_2/g_1)N_1)L} = I_0 e^{\sigma_{21}\Delta N_{21}L} \quad (8)$$

where σ_{21} is referred to as the *stimulated emission cross section* with dimensions of m², L is the thickness of the medium (in meters) through which the beam passes, and $N_2 - (g_2/g_1)N_1 = \Delta N_{21}$ is known as the *population inversion density*. The exponents in Eq. (8) are dimensionless quantities that can be either greater or less than unity, depending upon whether N_2 is greater than or less than $(g_2/g_1)N_1$.

The general form of the stimulated emission cross section per unit frequency is given as

$$\sigma_{21} = \frac{\lambda_{21}^2 A_{21}}{8\pi \Delta\nu} \quad (9)$$

in which $\Delta\nu$ represents the linewidth over which the stimulated emission or absorption occurs.

For the case of homogeneous broadening, at the *center* of the emission line, σ_{21} is expressed as

$$\sigma_{21}^H = \frac{\lambda_{21}^2 A_{21}}{4\pi^2 \Delta\nu_{21}^H} \quad (10)$$

where $\Delta\nu_{21}^H$ is the homogeneous emission linewidth (FWHM) which was described earlier for several different situations.

For the case of Doppler broadening, σ_{21}^D can be expressed as

$$\sigma_{21}^D = \sqrt{\frac{\ln 2}{16\pi^3}} \frac{\lambda_{21}^2 A_{21}}{\Delta\nu_{21}^D} \quad (11)$$

at the *center* of the emission line and $\Delta\nu_{21}^D$ is the Doppler emission linewidth expressed earlier in Eq. (3).

For all types of matter, the population density ratio of levels 1 and 2 would normally be such that $N_2 \ll N_1$. This can be shown by the Boltzmann relationship for the population ratio in thermal equilibrium which provides the ratio of N_2/N_1 to be

$$\frac{N_2}{N_1} = \frac{g_2}{g_1} e^{-\Delta E_{21}/kT} \quad (12)$$

in which T is the temperature and k is Boltzmann's constant. Thus, for energy levels separated by energies corresponding to visible transitions, in a medium at or near room temperature, the ratio of $N_2/N_1 \cong e^{-100} = 10^{-44}$. For most situations in everyday life we can ignore the population N_2 in the upper level and rewrite Eq. (8) as

$$I = I_o e^{-\sigma_{12} N_1 L} = I_o e^{-\alpha_{12} L} \quad (13)$$

where $\sigma_{21} = (g_2/g_1)\sigma_{12}$ and $\alpha_{12} = \sigma_{12} N_1$. Equation (13) is known as Beer's law, which is used to describe the absorption of light within a medium. α_{12} is referred to as the absorption coefficient with units of m^{-1} .

In laser amplifiers we cannot ignore the population in level 2. In fact, the condition for amplification and laser action is

$$N_2 - (g_2/g_1)N_1 > 0 \quad \text{or} \quad g_1 N_2 / g_2 N_1 > 1 \quad (14)$$

since, if Eq. (14) is satisfied, the exponent of Eq. (8) will be greater than 1 and I will emerge from the medium with a greater value than I_o , or amplification will occur. The condition of Eq. (14) is a necessary condition for laser amplification and is known as a population inversion since N_2 is greater than $(g_2/g_1)N_1$. In most cases, (g_2/g_1) is either unity or close to unity.

Considering that the ratio of $g_1 N_2 / g_2 N_1$ [Eq. (14)] could be greater than unity does not follow from normal thermodynamic equilibrium considerations [Eq. (12)] since it would represent a population ratio that could never exist in thermal equilibrium. When Eq. (14) is satisfied and amplification of the beam occurs, the medium is said to have *gain* or *amplification*. The factor $\sigma_{21}(N_2 - (g_2/g_1)N_1)$ or $\sigma_{21}\Delta N_{21}$ is often referred to as the *gain coefficient* g_{21} and is given in units of m^{-1} such that $g_{21} = \sigma_{21}\Delta N_{21}$. Typical values of σ_{21} , ΔN_{21} , and g_{21} are given in Table 1 for a variety

of lasers. The term g_{21} is also referred to as the *small-signal gain coefficient* since it is the gain coefficient determined when the laser beam intensity within the laser gain medium is small enough that stimulated emission does not significantly alter the populations in the laser levels.

Gain Saturation²

It was stated in the previous section that Eq. (14) is a necessary condition for making a laser but it is not a sufficient condition. For example, a medium might satisfy Eq. (14) by having a gain of $e^{g_{21}L} \approx 10^{-10}$, but this would not be sufficient to allow any reasonable beam to develop. Lasers generally start by having a pumping process that produces enough population in level 2 to create a population inversion with respect to level 1. As the population decays from level 2, radiation occurs spontaneously on the transition from level 2 to level 1 equally in all directions within the gain medium. In most of the directions very little gain or enhancement of the spontaneous emission occurs, since the length is not sufficient to cause significant growth according to Eq. (8). It is only in the elongated direction of the amplifier, with a much greater length, that significant gain exists and, consequently, the spontaneous emission is significantly enhanced. The requirement for a laser beam to develop in the elongated direction is that the exponent of Eq. (8) be large enough for the beam to grow to the point where it begins to significantly reduce the population in level 2 by stimulated emission. The beam will eventually grow, according to Eq. (8), to an intensity such that the stimulated emission rate is equal to the spontaneous emission rate. At that point the beam is said to reach its saturation intensity I_{sat} , which is given by

$$I_{\text{sat}} = \frac{h\nu_{21}}{\sigma_{21}^H \tau_{21}} \quad (15)$$

The saturation intensity is that value at which the beam can no longer grow exponentially according to Eq. (8) because there are no longer enough atoms in level 2 to provide the additional gain. When the beam grows above I_{sat} it begins to extract significant energy since at this point the stimulated emission rate exceeds the spontaneous emission rate for that transition. The beam essentially takes energy that would normally be radiated in all directions spontaneously and redirects it via stimulated emission, thereby increasing the beam intensity.

Threshold Conditions with No Mirrors

I_{sat} can be achieved by having any combination of values of the three parameters σ_{21}^H , ΔN_{21} , and L large enough that their product provides sufficient gain. It turns out that the requirement to reach I_{sat} can be given by making the exponent of Eq. (8) have the following range of values:

$$\sigma_{21}^H \Delta N_{21} L \approx 10-20 \quad \text{or} \quad g_{21}^H L \approx 10-20 \quad (16)$$

where the specific value between 10 and 20 is determined by the geometry of the laser cavity. Equation (16) suggests that the beam grows to a value of $I/I_o = e^{10-20} = 2 \times 10^4 - 5 \times 10^8$, which is a very large amplification.

Threshold Conditions with Mirrors

One could conceivably make L sufficiently long to always satisfy Eq. (16), but this is not practical. Some lasers can reach the saturation intensity over a length L of a few centimeters, but most require much longer lengths. Since one cannot readily extend the lasing medium to be long enough to achieve I_{sat} , the same result is obtained by putting mirrors around the gain medium. This effectively increases the path length by having the beam pass many times through the amplifier.

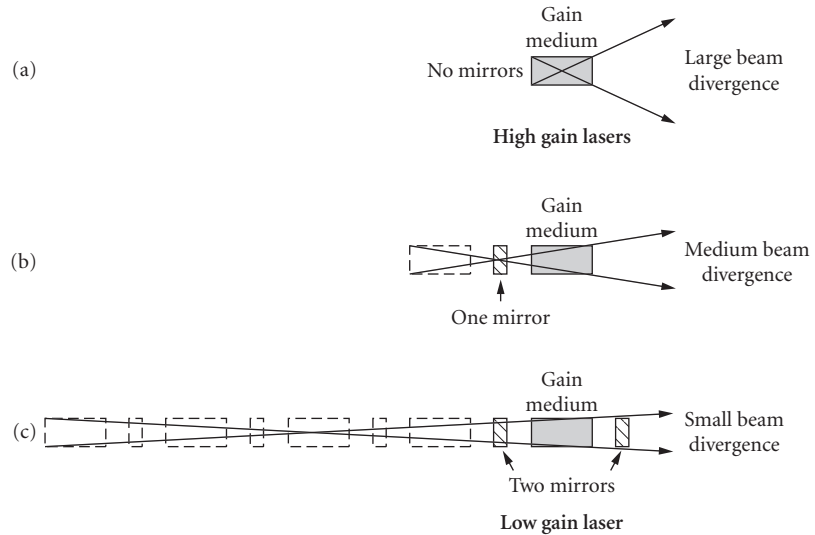


FIGURE 5 Laser beam divergence for amplifier configurations having high gain and (a) no mirrors or (b) one mirror, and high or low gain and (c) two mirrors.

A simple understanding of how the mirrors affect the beam is shown in Fig. 5. The diagram effectively shows how multiple passes through the amplifier can be considered for the situation where flat mirrors are used at the ends of the amplifier. The dashed outlines that are the same size as the gain media represent the images of those gain media produced by their reflections in the mirrors. It can be seen that as the beam makes multiple passes, its divergence narrows up significantly in addition to the large increase in intensity that occurs due to amplification. For high-gain lasers, such as excimer or organic dye lasers, the beam only need pass through the amplifier a few times to reach saturation. For low-gain lasers, such as the helium-neon laser, it might take 500 passes through the amplifier in order to reach I_{sat} .

The laser beam that emerges from the laser is usually coupled out of the amplifier by having a partially transmitting mirror at one end of the amplifier which typically reflects most of the beam back into the medium for more growth. To ensure that the beam develops, the transmission of the output coupling mirror (as it is usually referred to) must be lower than the gain incurred by the beam during a round-trip pass through the amplifier. If the transmission is higher than the round-trip gain, the beam undergoes no net amplification. It simply never develops. Thus a relationship that describes the threshold for laser oscillation balances the laser gain with the cavity losses. In the most simplified form, those losses are due to the mirrors having a reflectivity less than unity. Thus, for a round-trip pass-through the laser cavity, the threshold for inversion can be expressed as

$$R_1 R_2 e^{\sigma_{21} \Delta N_{21} 2L} > 1 \quad \text{or} \quad R_1 R_2 e^{g_{21} 2L} > 1 \quad (17a)$$

which is a similar requirement to that of Eq. (16) since it defines the minimum gain requirements for a laser.

A more general version of Eq. (17a) can be expressed as

$$R_1 R_2 (1 - a_L)^2 e^{(\sigma_{21} \Delta N_{21} - \alpha) 2L} = 1 \quad \text{or} \quad R_1 R_2 (1 - a_L)^2 e^{(g_{21} - \alpha) 2L} = 1 \quad (17b)$$

in which we have included a distributed absorption α throughout the length of the gain medium at the laser wavelength, as well as the total scattering losses a_L per pass through the cavity (excluding the gain medium and the mirror surfaces). The absorption loss α is essentially a separate absorbing transition within the gain medium that could be a separate molecule as in an excimer laser,

absorption from either the ground state or from the triplet state in a dye laser, or absorption from the ground state in the broadband tunable solid-state lasers and in semiconductor lasers or from the upper laser state in most solid-state lasers. The scattering losses a_L , per pass, would include scattering at the windows of the gain medium, such as Brewster angle windows, or scattering losses from other elements that are inserted within the cavity. These are typically of the order of one or two percent or less.

Laser Operation above Threshold

Significant power output is achieved by operating the laser at a gain greater than the threshold value defined above in Eqs. (16) and (17). For such a situation, the higher gain that would normally be produced by increased pumping is reduced to the threshold value by stimulated emission. The additional energy obtained from the reduced population inversion is transferred to the laser beam in the form of increased laser power. If the laser has low gain, as most cw (continuously operating) gas lasers do, the gain and also the power output tend to stabilize rather readily.

For solid-state lasers, which tend to have higher gain and also a long upper-laser-level lifetime, a phenomenon known as *relaxation oscillations*⁷ occurs in the laser output. For pulsed (non-Q-switched) lasers in which the gain lasts for many microseconds, these oscillations occur in the form of a regularly repeated spiked laser output superimposed on a lower steady-state value. For cw lasers it takes the form of a sinusoidal oscillation of the output. The phenomenon is caused by an oscillation of the gain due to the interchange of pumped energy between the upper laser level and the laser field in the cavity. This effect can be controlled by using an active feedback mechanism, in the form of an intensity-dependent loss, in the laser cavity.

Laser mirrors not only provide the additional length required for the laser beam to reach I_{sat} , but they also provide very important resonant cavity effects that will be discussed in a later section. Using mirrors at the ends of the laser gain medium (or amplifier) is referred to as having the gain medium located within an optical cavity.

How Population Inversions Are Achieved⁴

It was mentioned earlier that population inversions are not easily achieved in normal situations. All types of matter tend to be driven toward thermal equilibrium. From an energy-level standpoint, to be in thermal equilibrium implies that the ratio of the populations of two excited states of a particular material, whether it be a gaseous, liquid, or solid material, is described by Eq. (12). For any finite value of the temperature this leads to a value of $N_2/(g_2/g_1)N_1$ that is always less than unity and therefore Eq. (14) can never be satisfied under conditions of thermal equilibrium. Population inversions are therefore produced in either one of two ways: (1) selective excitation (pumping) of the upper-laser-level 2, or (2) more rapid decay of the population of the lower-laser-level 1 than of the upper-laser-level 2, even if they are both populated by the same pumping process.

The first requirement mentioned above was met in producing the very first laser, the ruby laser.⁵ In this laser the flash lamp selectively pumped chromium atoms to the upper laser level (through an intermediate level) until the ground state (lower laser level) was depleted enough to produce the inversion (Fig. 6). Another laser that uses this selective pumping process is the copper vapor laser⁶ (CVL). In this case, electrons in a gaseous discharge containing the copper vapor have a much preferred probability of pumping the upper laser level than the lower laser level (Fig. 7). Both of these lasers involve essentially three levels.

The second type of excitation is used for most solid-state lasers, such as the Nd³⁺ doped yttrium aluminum garnet laser⁷ (commonly referred to as the Nd:YAG laser), for organic dye lasers,⁸ and many others. It is probably the most common mechanism used to achieve the necessary population inversion. This process involves four level² (although it can include more) and generally occurs via excitation from the ground state 0 to an excited state 3 which energetically lies above the upper-laser-level 2. The population then decays from level 3 to level 2 by nonradiative processes (such as

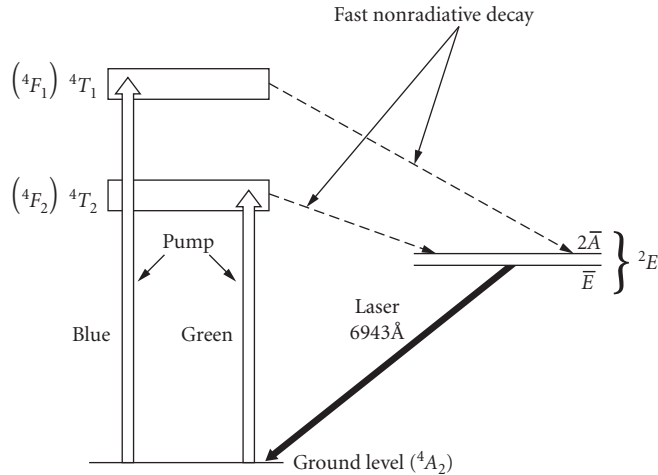


FIGURE 6 Energy-level diagram for a ruby laser showing the pump wavelength bands and the laser transition. The symbols 4T_1 and 4T_2 are shown as the appropriate designations for the pumping levels in ruby along with the more traditional designations in parenthesis.

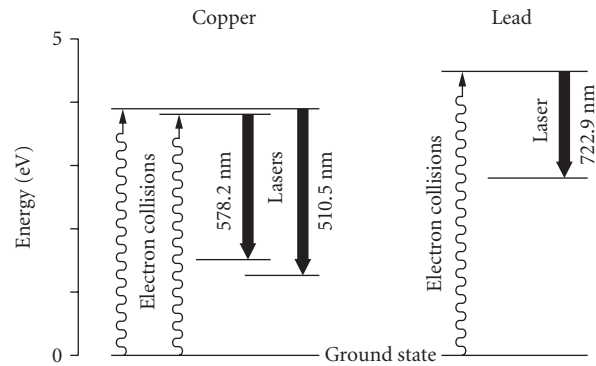


FIGURE 7 Energy-level diagrams of the transient three-level copper and lead vapor lasers showing the pump processes as well as the laser transitions.

collisional processes in gases), but can also decay radiatively to level 2. If the proper choice of materials has been made, the lower-laser-level 1 in some systems will decay rapidly to the ground state (level 0) which allows the condition $N_2 > (g_2/g_1)N_1$ to be satisfied. This situation is shown in Fig. 8 for an Nd:YAG laser crystal.

Optimization of the Output Coupling from a Laser Cavity⁹

A laser will operate with any combination of mirror reflectivities subject to the constraints of the threshold condition of Eq. (17a or b). However, since lasers are devices that are designed to use the

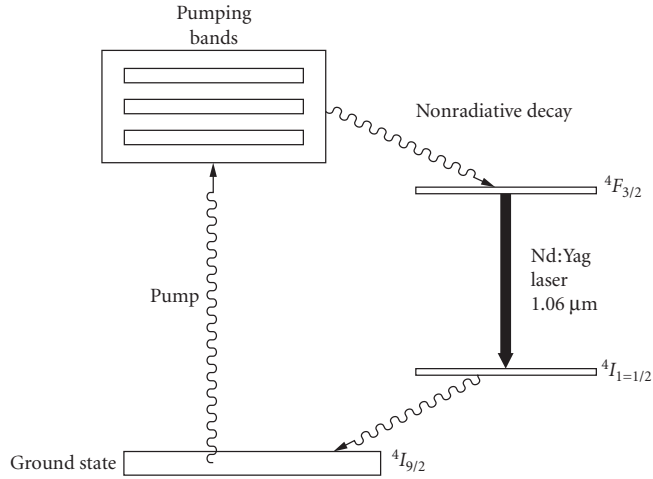


FIGURE 8 Energy-level diagram of the Nd:YAG laser indicating the four-level laser excitation process.

laser power in various applications, it is desirable to extract the most power from the laser in the most efficient manner. A simple expression for the optimum laser output coupling is given as

$$t_{\text{opt}} = (a_L g_{21} L)^{1/2} - a_L \tag{18}$$

in which a_L is the absorption and scattering losses per pass through the amplifier (the same as in Eq. (17a and b), g_{21} is the small signal gain per pass through the amplifier, and L is the length of the gain medium. This value of t_{opt} is obtained by assuming that equal output couplings are used for both mirrors at the ends of the cavity. To obtain all of the power from one end of the laser, the output transmission must be doubled for one mirror and the other mirror is made to be a high reflector.

The intensity of the beam that would be emitted from the output mirror can be estimated to be

$$I_{t_{\text{max}}} = \frac{I_{\text{sat}} t_{\text{opt}}^2}{2a_L} \tag{19}$$

in which I_{sat} is the saturation intensity as obtained from Eq. (15). If all of the power is desired from one end of the laser, as discussed above, then $I_{t_{\text{max}}}$ would be doubled in the above expression.

Pumping Techniques to Produce Inversions

Excitation or pumping of the upper laser level generally occurs by two techniques: (1) particle pumping and (2) optical or photon pumping. No matter which process is used, the goal is to achieve sufficient pumping flux and, consequently enough, population in the upper-laser-level 2 to exceed the requirements of either Eq. (16) or Eq. (17).

Particle Pumping²⁰ Particle pumping occurs when a high-speed particle collides with a laser species and converts its kinetic energy to internal energy of the laser species. Particle pumping occurs mostly with electrons as the pumping particles. This is especially common in a gas discharge where a voltage is applied across a low-pressure gas and the electrons flow through the tube in the form of a discharge current that can range from a few milliamps to tens of amperes, depending upon the particular laser and the power level desired. This type of excitation process is used for lasers such as the argon (Fig. 9) and krypton ion lasers, the copper vapor laser, excimer lasers, and the molecular nitrogen laser.

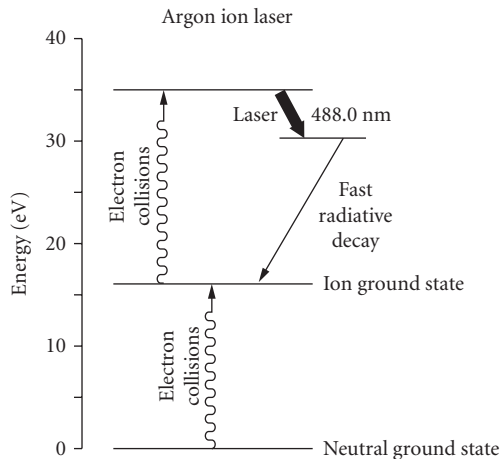


FIGURE 9 Energy-level diagram of the argon ion laser indicating the two-step excitation process.

Two other well-known gas lasers, the helium-neon laser and the helium-cadmium laser, operate in a gas discharge containing a mixture of helium gas and the laser species (neon gas or cadmium vapor). When an electric current is produced within the discharge, the high-speed electrons first pump an excited metastable state in helium (essentially a storage reservoir). The energy is then transferred from this reservoir to the upper laser levels of neon or cadmium by collisions of the helium metastable level with the neon or cadmium ground-state atoms as shown in Fig. 10. Electron collisions with the cadmium ion ground state have also been shown to produce excitation in the case of the helium-cadmium laser.

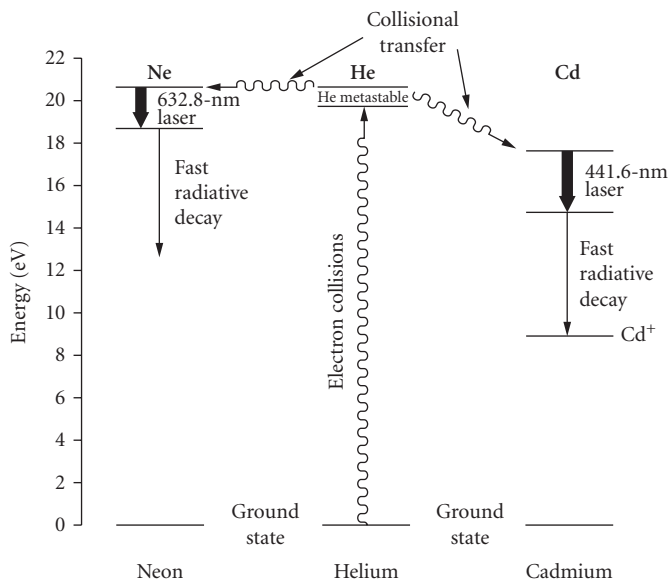


FIGURE 10 Energy-level diagrams of the helium-neon (He-Ne) laser and the helium-cadmium (He-Cd) laser that also include the helium metastable energy levels that transfer their energy to the upper laser levels by collisions.

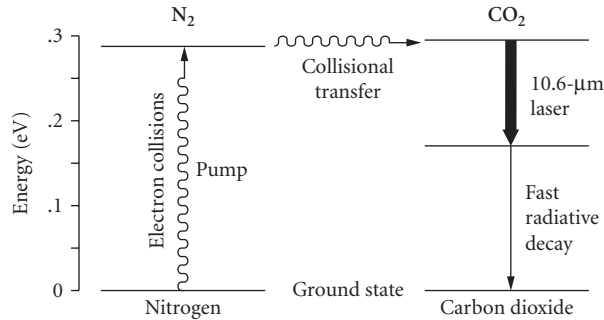


FIGURE 11 Energy-level diagram of the carbon dioxide (CO_2) laser along with the energy level of molecular nitrogen that collisionally transfers its energy to the CO_2 upper laser level.

In an energy-transfer process similar to that of helium with neon or cadmium, the CO_2 laser operates by using electrons of the gas discharge to produce excitation of molecular nitrogen vibrational levels that subsequently transfer their energy to the vibrational upper laser levels of CO_2 as indicated in Fig. 11. Helium is used in the CO_2 laser to control the electron temperature and also to cool (reduce) the population of the lower laser level via collisions of helium atoms with CO_2 atoms in the lower laser level.

High-energy electron beams and even nuclear reactor particles have also been used for particle pumping of lasers, but such techniques are not normally used in commercial laser devices.

Optical Pumping⁷ Optical pumping involves the process of focusing light into the gain medium at the appropriate wavelength such that the gain medium will absorb most (or all) of the light and thereby pump that energy into the upper laser level as shown in Fig. 12. The selectivity in pumping the laser level with an optical pumping process is determined by choosing a gain medium having significant absorption at a wavelength at which a suitable pump light source is available. This of course implies that the absorbing wavelengths provide efficient pumping pathways to the upper laser level. Optical pumping requires very intense pumping light sources, including flash lamps and other lasers. Lasers that are produced by optical pumping include organic dye lasers and solid-state lasers. The two types of energy level arrangements for producing lasers via optical pumping were described in detail in the section “How Population Inversions Are Achieved” and shown in Figs. 6 and 8.

Flash lamps used in optically pumped laser systems are typically long, cylindrically shaped, fused quartz structures of a few millimeters to a few centimeters in diameter and 10 to 50 cm in

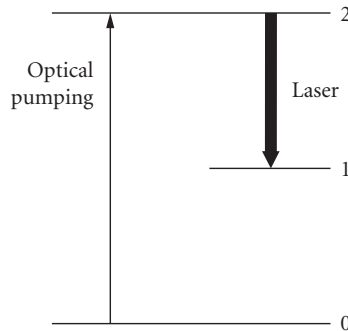


FIGURE 12 A general diagram showing optical pumping of the upper laser level.

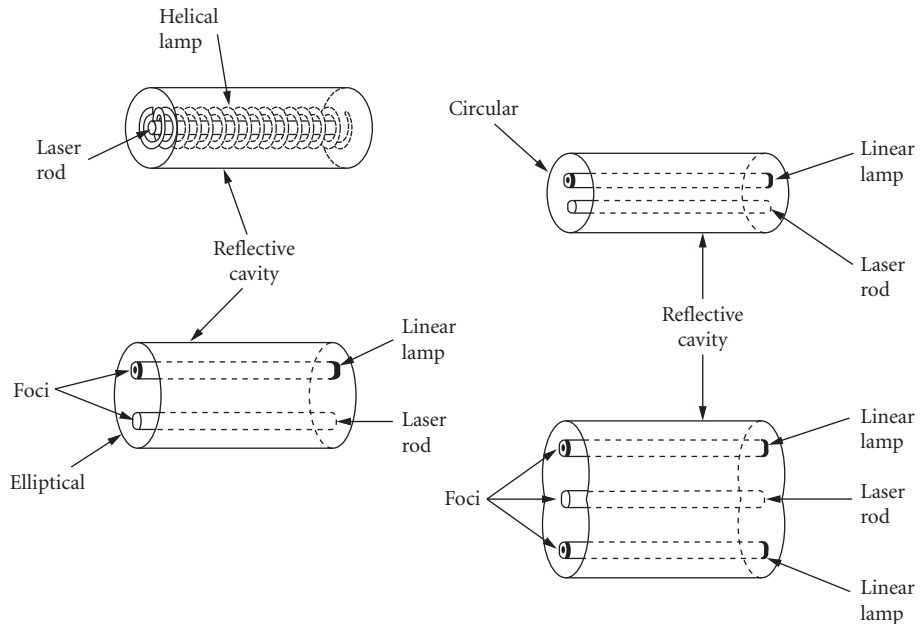


FIGURE 13 Flash lamp pumping arrangements for solid-state laser rods showing the use of helical lamps as well as linear lamps in a circular cavity, a single elliptical cavity, and a double elliptical cavity.

length. The lamps are filled with gases such as xenon and are initiated by running an electrical current through the gas. The light is concentrated into the lasing medium by using elliptically shaped reflecting cavities that surround both the laser medium and the flash lamps as shown in Fig. 13. These cavities efficiently collect and transfer to the laser rod most of the lamp energy in wavelengths within the pump absorbing band of the rod. The most common flash lamp-pumped lasers are Nd:YAG, Nd:glass, and organic dye lasers. New crystals such as Cr:LiSAF and HoTm:YAG are also amenable to flash lamp pumping.

Solid-state lasers also use energy-transfer processes as part of the pumping sequence in a way similar to that of the He-Ne and He-Cd gas lasers. For example, Cr^{3+} ions are added into neodymium-doped crystals to improve the absorption of the pumping light. The energy is subsequently transferred to the Nd^{3+} laser species. Such a process of adding desirable impurities is known as *sensitizing*.

Lasers are used as optical pumping sources in situations where (1) it is desirable to be able to concentrate the pump energy into a small-gain region or (2) it is useful to have a narrow spectral output of the pump source in contrast with the broadband spectral distribution of a flash lamp.

Laser pumping is achieved by either transverse pumping (a direction perpendicular to the direction of the laser beam) or longitudinal pumping (a direction in the same direction as the emerging laser beam). Frequency doubled and tripled pulsed Nd:YAG lasers are used to transversely pump organic dye lasers⁸ that provide continuously tunable laser radiation over the near-ultraviolet, visible, and near-infrared spectral regions (by changing dyes at appropriate wavelength intervals). For transverse pumping, the pump lasers are typically focused into the dye medium with a cylindrical lens to provide a 1- to 2-cm-long (but very narrow) gain medium in the liquid dye solution. The dye concentration is adjusted to absorb the pump light within a millimeter or so into the dye cell to provide the very high concentration of gain near the surface of the cell.

Both cw and mode-locked argon ion lasers and Nd:YAG lasers are typically used for longitudinal or end pumping of cw and/or mode-locked organic dye lasers and also of solid-state gain media. In this pumping arrangement the pump laser is focused into a very thin gain region, which is provided by either a thin jet stream of flowing dye solution (Fig. 14) or a solid-state crystal such as $\text{Ti:Al}_2\text{O}_3$,

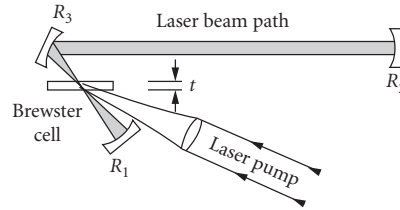


FIGURE 14 End-pumping cavity arrangement for either organic dye lasers or solid-state lasers.

The thin gain medium is used, in the case of the generation of ultrashort mode-locked pulses, so as to allow precise timing of the short-duration pump pulses with the ultrashort laser pulses that develop within the gain medium as they travel within the optical cavity.

Thin-disk lasers are diode-pumped solid-state lasers that efficiently produce high output power with good beam quality. Such lasers have gain media with a very short axial dimension of several hundred microns and wide transverse dimensions of several centimeters. The disks are antireflection coated for both the lasing and pumping wavelengths on the front side and a high reflection coating on the rear side. The laser crystal is mounted on a heat sink to efficiently remove the wasted heat of the pumping process. An output coupling mirror is mounted in front of the disk to provide multiple passes of the beam through the gain medium, for maximum power extraction, and to provide good mode quality. Yb:YAG is the most successful laser material for this type of laser.

Gallium arsenide semiconductor diode lasers, operating at wavelengths around $0.8 \mu\text{m}$, can be effectively used to pump Nd:YAG lasers because the laser wavelength is near that of the strongest absorption feature of the pump band of the Nd:YAG laser crystal, thereby minimizing excess heating of the laser medium. Also, the diode lasers are very efficient light sources that can be precisely focused into the desired mode volume of a small Nd:YAG crystal (Fig. 15a) which results in minimal waste of the pump light. Figure 15b shows how close-coupling of the pump laser and the Nd:YAG crystal can be used to provide compact efficient diode laser pumping. The infrared output of the

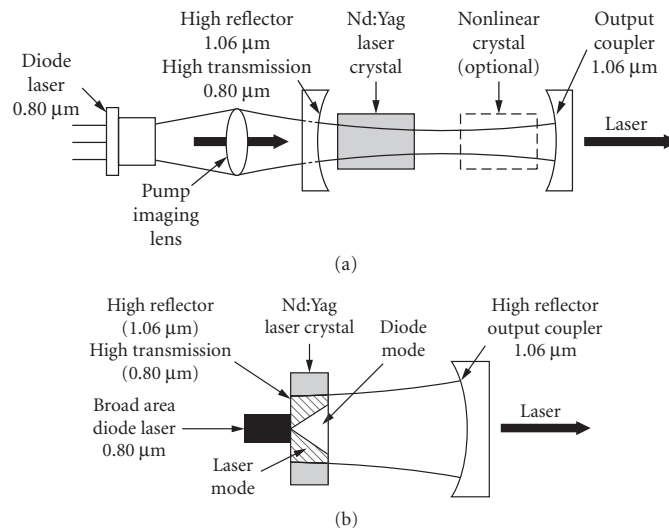


FIGURE 15 Pumping arrangements for diode-pumped Nd:YAG lasers showing both (a) standard pumping using an imaging lens and (b) close-coupled pumping in which the diode is located adjacent to the laser crystal.

Nd:YAG laser can then also be frequency doubled as indicated in Fig. 15a, using nonlinear optical techniques to produce a green laser beam in a relatively compact package.

Semiconductor Diode Laser Pumping Semiconductor laser pumping occurs when electrons are made to flow from an n -type semiconductor to a p -type semiconductor. In this case, as opposed to particle pumping described previously, it is not the kinetic energy of the electrons that does the excitation. Instead, it is the electrons themselves flowing into a p -doped material that produces the inversion. An analogy might be the water in a mountain region approaching a waterfall. The water is already at the upper energy site and loses its energy when it cascades down the waterfall. In the same sense, the electrons already have sufficient potential energy when they are pulled into the p -type material via an external electric field. Once they arrive, a population inversion exists where they recombine with the holes and cascade downward to produce the recombination radiation.

16.4 LASER PROPERTIES ASSOCIATED WITH OPTICAL CAVITIES OR RESONATORS

Longitudinal Laser Modes^{1,2}

When a collimated optical beam of infinite lateral extent (a plane wave) passes through two reflecting surfaces of reflectivity R and also of infinite extent that are placed normal (or nearly normal) to the beam and separated by a distance d , as shown in Fig. 16a, the plot of transmission versus wavelength for the light as it emerges from the second reflecting surface is shown in Fig. 16b. The transmission reaches a maximum of 100 percent (if there are no absorption losses at the reflecting surfaces) at frequency spacings of $\Delta\nu = c/2\eta d$ where c is the speed of light in a vacuum, and η is the index of refraction of the medium between the mirrors. This frequency-selective optical device is known as a Fabry-Perot interferometer and has many useful applications in optics.

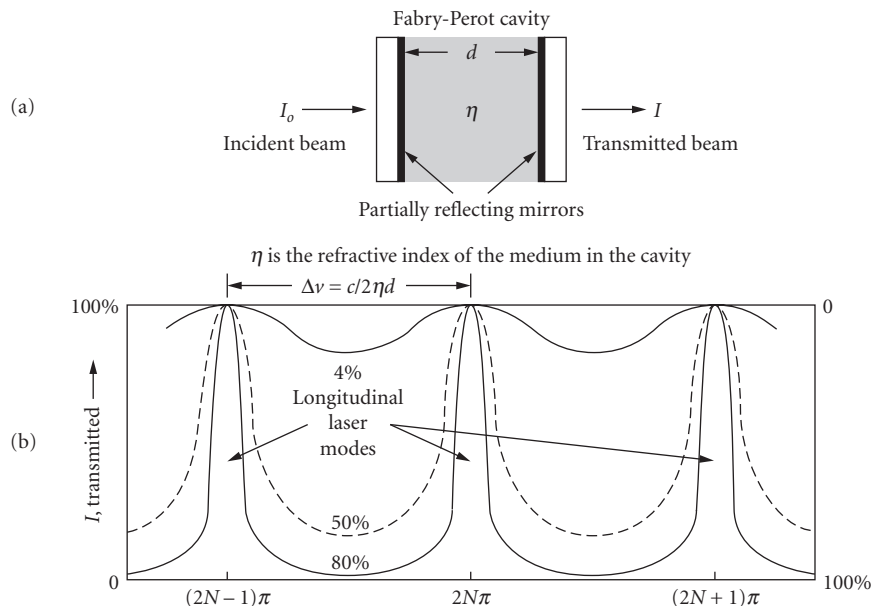


FIGURE 16 Fabry-Perot optical cavity consisting of two plane-parallel mirrors with a specific reflectivity separated by a distance d indicating the frequency spacing of longitudinal modes.

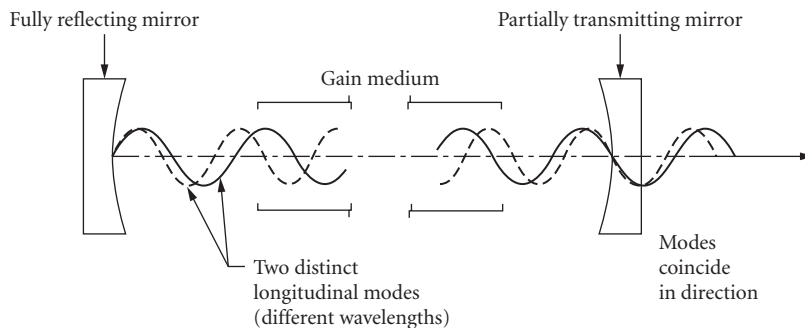


FIGURE 17 Laser resonator showing two distinct longitudinal modes traveling in the same direction but with slightly different frequencies.

The transmission through this device, as shown in Fig. 16*b* is enhanced at regular frequency or wavelength intervals due to the development of standing waves that resonate within the optical cavity. The enhancement occurs at frequencies (or wavelengths) at which complete sinusoidal half-cycles of the electromagnetic wave exactly “fit” between the mirrors such that the value of the electric field of the wave is zero at the mirror surfaces.

If a laser amplifier is placed between two mirrors in the same arrangement as described above, the same standing waves tend to be enhanced at frequency intervals of

$$\nu = n(c/2\eta d) \quad (20)$$

where n is an integer that indicates the number of half wavelengths of the laser light that fit within the spacing d of the two mirrors. In a typical laser operating in the visible spectral region, n would be of the order of 30,000 to 40,000. In such a laser, the output of the laser beam emerging from the cavity is very strongly enhanced at the resonant wavelengths as shown in Fig. 17, since these are the wavelengths that have the lowest loss within the cavity. The widths of the resonances shown in Fig. 17 are those of a passive Fabry-Perot cavity. When an active gain medium is placed within the cavity, the linewidth of the beam that is continually amplified as it reflects back and forth between the mirrors is narrowed even further.

These enhanced regions of very narrow frequency output are known as *longitudinal modes* of the laser. They are referred to as modes since they represent discrete values of frequency associated with the integral values of n at which laser output occurs. Lasers operating on a single longitudinal mode with ultrastable cavities and ultrahigh reflectivity mirrors have generated linewidths as narrow as a few hundred hertz or less. Since the longitudinal or temporal coherence length of a beam of light is determined by $c/\Delta\nu$, a very narrow laser linewidth can provide an extremely long coherence length and thus a very coherent beam of light.

For a typical gas laser (not including excimer lasers), the laser gain bandwidth is of the order of 10^9 to 10^{10} Hz. Thus, for a laser mirror cavity length of 0.5 m, the mode spacing would be of the order of 300 MHz and there would be anywhere from 3 to 30 longitudinal modes operating within the laser cavity. For an organic dye laser or a broadband solid-state laser, such as a $\text{Ti:Al}_2\text{O}_3$ laser, there could be as many as one million distinct longitudinal modes, each of a slightly different frequency than the next one, oscillating at the same time. However, if mode-locking is not present, typically only one or a few modes will dominate the laser output of a homogeneously broadened laser gain medium.

Transverse Laser Modes¹¹⁻¹⁴

The previous section considered the implications of having a collimated or parallel beam of light of infinite lateral extent pass through two infinite reflecting surfaces that are arranged normal to the direction of propagation of the beam and separated by a specific distance d . We must now

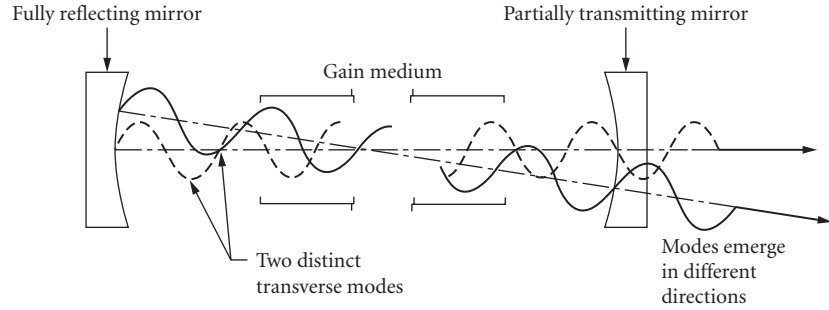


FIGURE 18 Laser resonator showing two distinct transverse modes traveling in different directions with slightly different frequencies.

consider the consequences of having the light originate within the space between the two mirrors in an amplifier that has only a very narrow lateral extent limited by either the diameter of the mirrors or by the diameter of the amplifying medium. The beam evolves from the spontaneous emission within the gain medium and eventually becomes a nearly collimated beam when it reaches I_{sat} since only rays traveling in a very limited range of directions normal to the laser mirrors will experience enough reflections to reach I_{sat} . The fact that the beam has a restricted aperture in the direction transverse to the direction of propagation causes it to evolve with a slight transverse component due to diffraction, which effectively causes the beam to diverge. This slight divergence actually consists of one or more distinctly separate beams that can operate individually or in combination.

These separate beams that propagate in the z direction are referred to as *transverse modes* as shown in Fig. 18. They are characterized by the various lateral spatial distributions of the electric field vector $\mathbf{E}(x, y)$ in the x - y directions as they emerge from the laser. These transverse amplitude distributions for the waves can be described by the relationship.

$$E_{pq}(x, y) = H_p\left(\frac{\sqrt{2}x}{w}\right) H_q\left(\frac{\sqrt{2}y}{w}\right) e^{-(x^2+y^2)/w^2} \quad (21)$$

In this solution, p and q are positive integers ranging from zero to infinity that designate the different modes which are associated with the order of the Hermite polynomials. Thus every set of p, q represents a specific distribution of wave amplitude at one of the mirrors, or a specific transverse mode of the open-walled cavity. We can list several Hermite polynomials as follows:

$$\begin{aligned} H_0(u) &= 1 & H_1(u) &= 2u \\ H_2(u) &= 2(2u^2 - 1) \\ H_m(u) &= (-1)^m e^{u^2} \frac{d^m(e^{-u^2})}{du^m} \end{aligned} \quad (22)$$

The spatial intensity distribution would be obtained by squaring the amplitude distribution function of Eq. (21). The transverse modes are designated TEM for transverse electromagnetic. The lowest order mode is given by TEM_{00} . It could also be written as TEM_{n00} in which n would designate the longitudinal mode number [Eq. (20)]. Since this number is generally very large for optical frequencies, it is not normally given.

The lowest order TEM_{00} mode has a circular distribution with a gaussian shape (often referred to as the *gaussian mode*) and has the smallest divergence of any of the transverse modes. Such a mode can be focused to a spot size with dimensions of the order of the wavelength of the beam. It has a minimum width or waist $2w_0$ that is typically located between the laser mirrors (determined

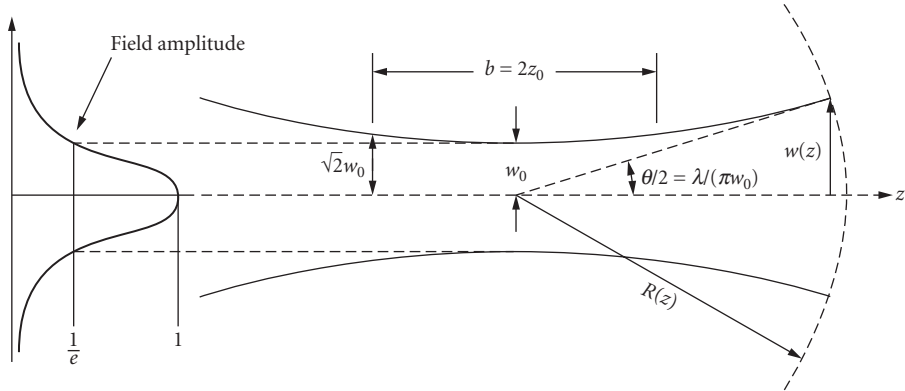


FIGURE 19 Parameters of a gaussian-shaped beam which are the features of a TEM_{00} transverse laser mode.

by the mirror curvatures and separation) and expands symmetrically in opposite directions from that minimum waist according to the following equation:

$$w(z) = w_0 \left[1 + \left(\frac{\lambda z}{\eta \pi w_0^2} \right)^2 \right]^{1/2} = w_0 \left[1 + \left(\frac{z}{z_0} \right)^2 \right]^{1/2} \quad (23)$$

where $w(z)$ is the beam waist at any location z measured from w_0 , η is the index of refraction of the medium, and $z_0 = \eta \pi w_0^2 / \lambda$ is the distance over which the beam waist expands to a value of $\sqrt{2} w_0$.

The waist $w(z)$ at any location, as shown in Fig. 19, describes the transverse dimension within which the electric field distribution of the beam decreases to a value of 37 percent ($1/e$) of its maximum on the beam axis and within which 86.5 percent of the beam energy is contained. The TEM_{00} mode would have an intensity distribution that is proportional to the square of Eq. (21) for $p=q=0$ and, since it is symmetrical around the axis of propagation, it would have a cylindrically symmetric distribution of the form

$$I(r, z) = I_0 e^{-2r^2/w^2(z)} \quad (24)$$

where I_0 is the intensity on the beam axis.

The beam also has a wavefront curvature given by

$$R(z) = z \left[1 + \left(\frac{\eta \pi w_0^2}{\lambda z} \right)^2 \right] \quad (25)$$

which is indicated in Fig. 19, and a far-field angular divergence given by

$$\theta = \lim_{z \rightarrow \infty} \frac{2w(z)}{z} = \frac{2\lambda}{\pi w_0} = 0.64 \frac{\lambda}{w_0} \quad (26)$$

For a symmetrical cavity formed by two mirrors, each of radius of curvature R , separated by a distance d and in a medium in which $\eta=1$ the minimum beam waist w_0 is given by

$$w_0^2 = \frac{\lambda}{2\pi} [d(2R-d)]^{1/2} \quad (27)$$

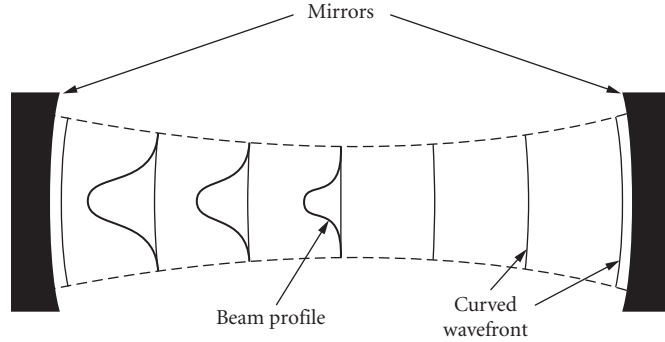


FIGURE 20 A stable laser resonator indicating the beam profile at various locations along the beam axis as well as the wavefront at the mirrors that matches the curvature of the mirrors.

and the radius of curvature r_c of the wavefront is

$$r_c = z + \frac{d(2R-d)}{4z} \quad (28)$$

For a confocal resonator in which $R=d$, w_o is given by

$$w_o = \sqrt{\frac{\lambda d}{2\pi}} \quad (29)$$

and the beam waist (spot size) at each mirror located a distance $d/2$ from the minimum is

$$w = \sqrt{\frac{\lambda d}{\pi}} \quad (30)$$

Thus, for a confocal resonator, $w(d/2)$ at each of the mirrors is equal to $\sqrt{2}w_o$ and thus at that location, $z=z_o$. The distance between mirrors for such a cavity configuration is referred to as the *confocal parameter* b such that $b=2z_o$. In a stable resonator, the curvature of the wavefront at the mirrors, according to Eqs. (25) and (28), exactly matches the curvature of the mirrors as shown in Fig. 20.

Each individual transverse mode of the beam is produced by traveling a specific path between the laser mirrors such that, as it passes from one mirror to the other and returns, the gain it receives from the amplifier is at least as great as the total losses of the mirror, as indicated from Eq. (17), plus the additional diffraction losses produced by either the finite lateral extent of the laser mirrors or the finite diameter of the laser amplifier or some other optical aperture placed in the system, whichever is smaller. Thus the TEM_{00} mode is produced by a beam passing straight down the axis of the resonator as indicated in Fig. 19.

Laser Resonator Configurations and Cavity Stability^{4,14}

There are a variety of resonator configurations that can be used for lasers. The use of slightly curved mirrors leads to much lower diffraction losses of the transverse modes than do plane parallel mirrors, and they also have much less stringent alignment tolerances. Therefore, most lasers use curved mirrors for the optical cavity. For a cavity with two mirrors of curvature R_1 and R_2 , and a separation distance d , a number of possible cavity configurations are shown in Fig. 21.

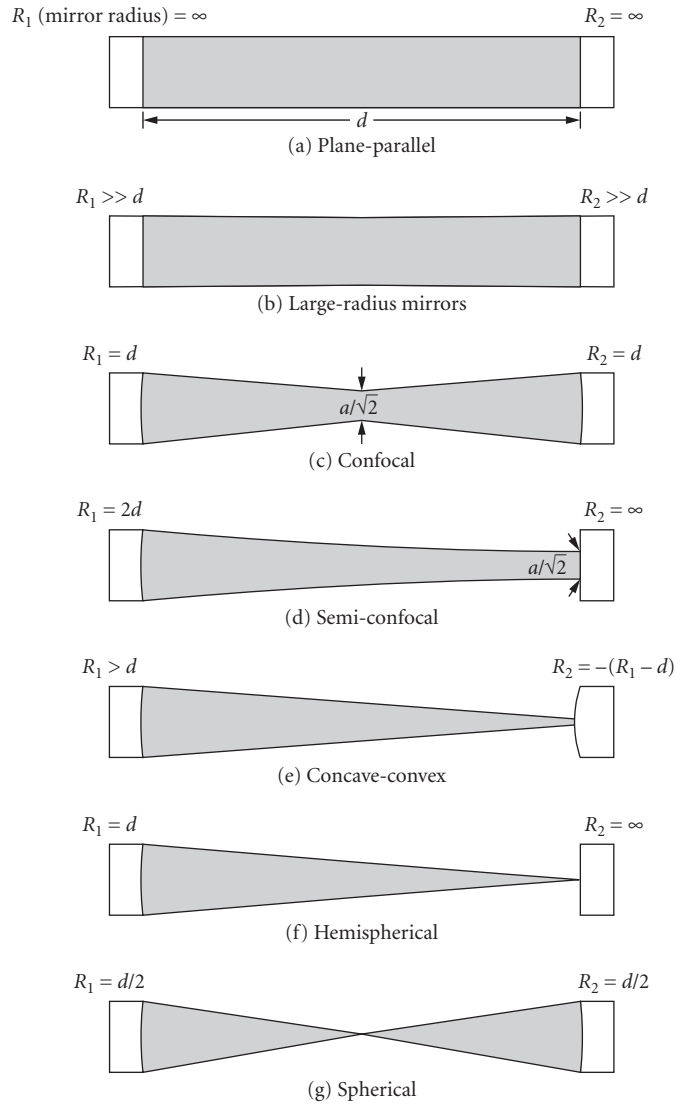


FIGURE 21 Possible two-mirror laser cavity configurations indicating the relationship of the radii of curvature of the mirrors with respect to the separation between mirrors.

A relationship between the radii of curvature and the separation between mirrors can be defined as

$$g_1 = 1 - \frac{d}{R_1} \quad \text{and} \quad g_2 = 1 - \frac{d}{R_2} \quad (31)$$

such that the condition for stable transverse modes is given by

$$0 < g_1 g_2 < 1 \quad (32)$$

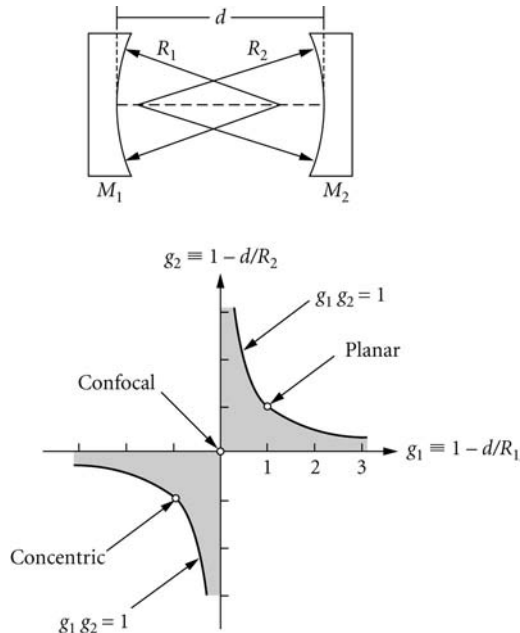


FIGURE 22 Stability diagram for two-mirror laser cavities indicating the shaded regions where stable cavities exist.

A stable mode is a beam that can be maintained with a steady output and profile over a relatively long period of time. It results from a cavity configuration that concentrates the beam toward the resonator axis in a regular pattern as it traverses back and forth within the cavity, rather than allowing it to diverge and escape from the resonator. In considering the various possible combinations of curved mirror cavities, one must keep the relation between the curvatures and mirror separation d within the stable regions of the graph, as shown in Fig. 22, in order to produce stable modes. Thus it can be seen from Fig. 22 that not all configurations shown in Fig. 21 are stable. For example, the planar, confocal, and concentric arrangements are just on the edge of stability according to Fig. 22.

There may be several transverse modes oscillating simultaneously “within” a single longitudinal mode. Each transverse mode can have the same value of n [Eq. (20)] but will have a slightly different value of d as it travels a different optical path between the resonator mirrors, thereby generating a slight frequency shift from that of an adjacent transverse mode. For most optical cavities, the mode that operates most easily is the TEM_{00} mode, since it travels a direct path along the axis of the gain medium.

16.5 SPECIAL LASER CAVITIES

Unstable Resonators¹⁴

A laser that is operating in a TEM_{00} mode, as outlined above, typically has a beam within the laser cavity that is relatively narrow in width compared to the cavity length. Thus, if a laser with a relatively wide gain region is used, to obtain more energy in the output beam, it is not possible to extract the energy from that entire region from a typically very narrow low-order gaussian mode.

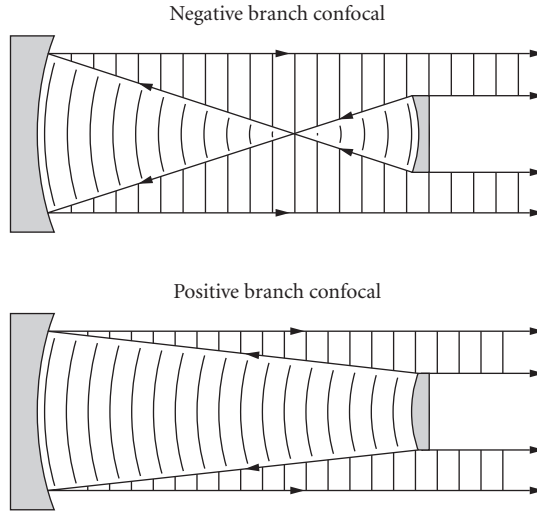


FIGURE 23 Unstable resonator cavity configurations showing both the negative and positive branch confocal cavities.

A class of resonators has been developed that can extract the energy from such wide laser volumes and also produce a beam with a nearly gaussian profile that makes it easily focusable. These resonators do not meet the criteria for stability, as outlined above, but still provide a good beam quality for some types of lasers. This class of resonators is referred to as *unstable resonators*.

Unstable resonators are typically used with high-gain laser media under conditions such that only a few passes through the amplifier will allow the beam to reach the saturation intensity and thus extract useful energy. A diagram of two unstable resonator cavity configurations is shown in Fig. 23. Figure 23b is the positive branch confocal geometry and is one of the most common unstable resonator configurations. In this arrangement, the small mirror has a convex shape and the large mirror a concave shape with a separation of length d such that $R_2 - R_1 = 2d$. With this configuration, any ray traveling parallel to the axis from left to right that intercepts the small convex mirror will diverge toward the large convex mirror as though it came from the focus of that mirror. The beam then reflects off of the larger mirror and continues to the right as a beam parallel to the axis. It emerges as a reasonably well-collimated beam with a hole in the center (due to the obscuration of the small mirror). The beam is designed to reach the saturation intensity when it arrives at the large mirror and will therefore proceed to extract energy as it makes its final pass through the amplifier. In the far field, the beam is near gaussian in shape, which allows it to be propagated and focused according to the equations described in the previous section. A number of different unstable resonator configurations can be found in the literature for specialized applications.

Q-Switching¹⁵

A typical laser, after the pumping or excitation is first applied, will reach the saturation intensity in a time period ranging from approximately 10 ns to 1 μ s, depending upon the value of the gain in the medium. For lasers such as solid-state lasers the upper-laser-level lifetime is considerably longer than this time (typically 50 to 200 μ s). It is possible to store and accumulate energy in the form of population in the upper laser level for a time duration of the order of that upper-level lifetime. If the laser cavity could be obscured during this pumping time and then suddenly switched into the system at a time order of the upper-laser-level lifetime, it would be possible for the gain, as well as

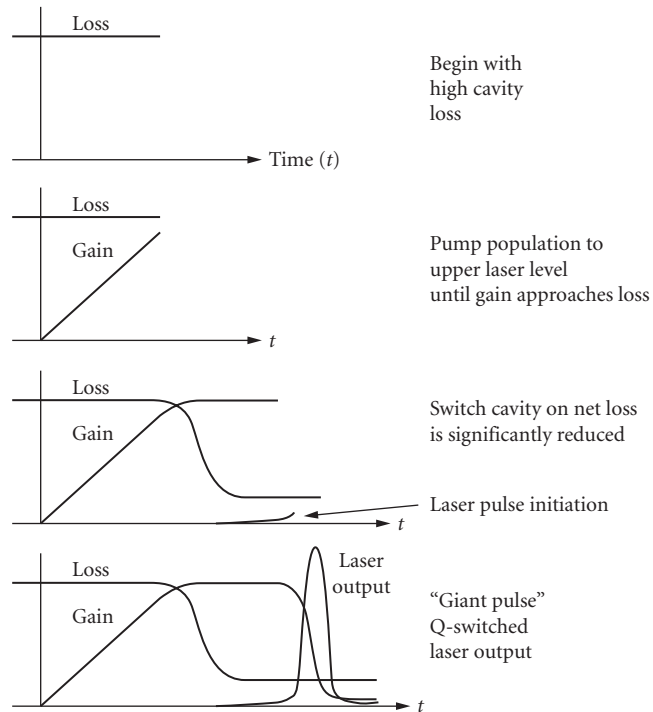


FIGURE 24 Schematic diagram of a Q-switched laser indicating how the loss is switched out of the cavity and a giant laser pulse is produced as the gain builds up.

the laser energy, to reach a much larger value than it normally would under steady-state conditions. This would produce a "giant" laser pulse of energy from the system.

Such a technique can in fact be realized and is referred to as *Q-switching* to suggest that the cavity *Q* is changed or switched into place. The cavity is switched on by using either a rapidly rotating mirror or an electro-optic shutter such as a Pockel cell or a Kerr cell. Nd:YAG and Nd:glass lasers are the most common Q-switched lasers. A diagram of the sequence of events involved in Q-switching is shown in Fig. 24.

Another technique that is similar to Q-switching is referred to as *cavity dumping*. With this technique, the intense laser beam inside of a normal laser cavity is rapidly switched out of the cavity by a device such as an acousto-optic modulator. Such a device is inserted at Brewster's angle inside the cavity and is normally transparent to the laser beam. When the device is activated, it rapidly inserts a high-reflecting surface into the cavity and reflects the beam out of the cavity. Since the beam can be as much as two orders of magnitude higher in intensity within the cavity than that which leaks through an output mirror, it is possible to extract high power on a pulsed basis with such a technique.

Mode-Locking¹⁴

In the discussion of longitudinal modes it was indicated that such modes of intense laser output occur at regularly spaced frequency intervals of $\Delta\nu = c/2\eta d$ over the gain bandwidth of the laser medium. For laser cavity lengths of the order of 10 to 100 cm these frequency intervals range from approximately 10^8 to 10^9 Hz. Under normal laser operation, specific modes with the highest gain

tend to dominate and quench other modes (especially if the gain medium is homogeneously broadened). However, under certain conditions it is possible to obtain all of the longitudinal modes lasing simultaneously, as shown in Fig. 25a. If this occurs, and the modes are all phased together so that they can act in concert by constructively and destructively interfering with each other, it is possible to produce a series of giant pulses separated in a time Δt of

$$\Delta t = 2\eta d/c \tag{33}$$

or approximately 1 to 10 ns for the cavity lengths mentioned above (see Fig. 25b). The pulse duration is approximately the reciprocal of the separation between the two extreme longitudinal laser modes or

$$\Delta t_p = \frac{1}{n\Delta\nu} \tag{34}$$

as can also be seen in Fig. 25b. This pulse duration is approximately the reciprocal of the laser gain bandwidth. However, if the index of refraction varies significantly over the gain bandwidth, then all of the frequencies are not equally spaced and the mode-locked pulse duration will occur only within the frequency width over which the frequency separations are approximately the same.

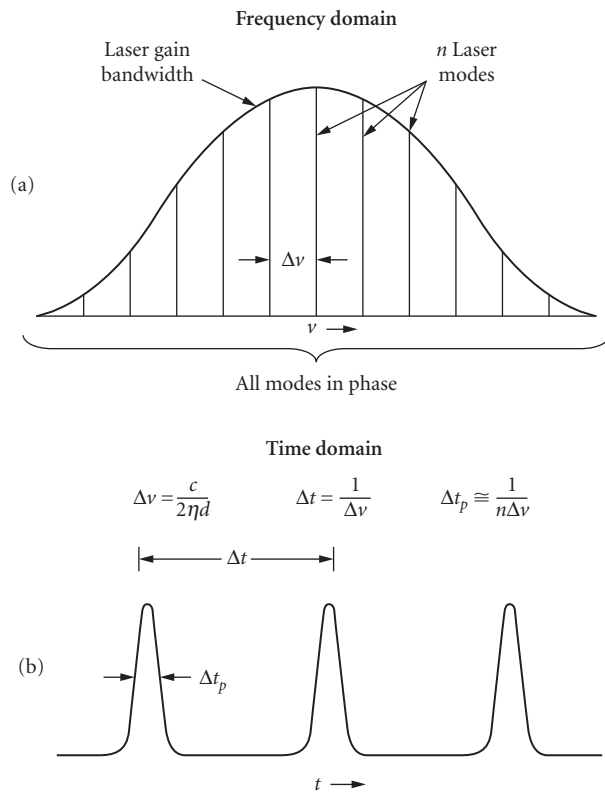


FIGURE 25 Diagrams in both the frequency and time domains of how mode-locking is produced by phasing n longitudinal modes together to produce an ultrashort laser pulse.

The narrowest pulses are produced from lasers having the largest gain bandwidth such as organic dye lasers, and solid-state lasers including Ti:Al₂O₃, Cr:Al₂O₃, and Cr:LiSAIF. The shortest mode-locked pulses to date, 4.4×10^{-15} seconds, have been produced in Ti:sapphire at a center wavelength of 820 nm at a rep rate of 80 MHz. Using pulse compression techniques, such pulses have been made as short as 3.5×10^{-15} seconds.

Such short-pulse operation is referred to as *mode-locking* and is achieved by inserting a very fast shutter within the cavity which is opened and closed at the intervals of the round-trip time of the short laser pulse within the cavity. This shutter coordinates the time at which all of the modes arrive at the mirror and thus brings them all into phase at that location. Electro-optic shutters, short duration gain pumping by another mode-locked laser, or passive saturable absorbers are techniques that can serve as the fast shutter. The second technique, short-duration gain pumping, is referred to as *synchronous pumping*. Three fast saturable absorber shutter techniques for solid-state lasers include colliding pulse mode-locking,¹⁶ additive pulse mode-locking,¹⁷ and Kerr lens mode-locking.¹⁸

Extremely short soft-x-ray pulses have also been produced via the interaction of intense laser beams with atoms. These pulses have been made as short as 70 as (70×10^{-18} seconds).

Distributed Feedback Lasers¹⁹

The typical method of obtaining feedback into the laser gain medium is to have the mirrors located at the ends of the amplifier as discussed previously. It is also possible, however, to provide the reflectivity within the amplifying medium in the form of a periodic variation of either the gain or the index of refraction throughout the medium. This process is referred to as *distributed feedback* (DFB). Such feedback methods are particularly effective in semiconductor lasers in which the gain is high and the fabrication of periodic variations is not difficult. The reader is referred to the reference section at the end of this chapter for further information concerning this type of feedback.

Ring Lasers¹⁴

Ring lasers are lasers that have an optical path within the cavity that involves the beam circulating in a loop rather than passing back and forth over the same path. This requires optical cavities that have more than two mirrors. The laser beam within the cavity consists of two waves traveling in opposite directions with separate and independent resonances within the cavity. In some instances an optical device is placed within the cavity that provides a unidirectional loss. This loss suppresses one of the beams, allowing the beam propagating in the other direction to become dominant. The laser output then consists of a traveling wave instead of a standing wave and therefore there are no longitudinal modes. Such an arrangement also eliminates the variation of the gain due to the standing waves in the cavity (spatial hole burning), and thus the beam tends to be more homogeneous than that of a normal standing-wave cavity. Ring lasers are useful for producing ultrashort mode-locked pulses and also for use in laser gyroscopes as stable reference sources.

16.6 SPECIFIC TYPES OF LASERS

Lasers can be categorized in several different ways including wavelength, material type, and applications. In this section we will summarize them by material type such as gas, liquid, solid-state, and semiconductor lasers. We will include only lasers that are available commercially since such lasers now provide a very wide range of available wavelengths and powers without having to consider special laboratory lasers.

Gaseous Laser Gain Media

Helium-Neon Laser¹⁰ The helium-neon laser was the first gas laser. The most widely used laser output wavelength is a red beam at 632.8 nm with a cw output ranging from 1 to 100 mW and sizes ranging from 10 to 100 cm in length. It can also be operated at a wavelength of 543.5 nm for some specialized applications. The gain medium is produced by passing a relatively low electrical current (10 mA) through a low pressure gaseous discharge tube containing a mixture of helium and neon. In this mixture, helium metastable atoms are first excited by electron collisions as shown in Fig. 10. The energy is then collisionally transferred to neon atom excited states which serve as upper laser levels.

Argon Ion Laser¹⁰ The argon ion laser and the similar krypton ion laser operate over a wide range of wavelengths in the visible and near-ultraviolet regions of the spectrum. The wavelengths in argon that have the highest power are at 488.0 and 514.5 nm. Power outputs on these laser transitions are available up to 20 W cw in sizes ranging from 50 to 200 cm in length. The gain medium is produced by running a high electric current (many amperes) through a very low pressure argon or krypton gas. The argon atoms must be ionized to the second and third ionization stages (Fig. 9) in order to produce the population inversions. As a result, these lasers are inherently inefficient devices.

Helium-Cadmium Laser¹⁰ The helium-cadmium laser operates cw in the blue at 441.6 nm, and in the ultraviolet at 353.6 and 325.0 nm with powers ranging from 20 to 200 mW in lasers ranging from 40 to 100 cm in length. The gain medium is produced by heating cadmium metal and evaporating it into a gaseous discharge of helium where the laser gain is produced. The excitation mechanisms include Penning ionization (helium metastables collide with cadmium atoms) and electron collisional ionization within the discharge as indicated in Fig. 10. The laser uses an effect known as *cataphoresis* to transport the cadmium through the discharge and provide the uniform gain medium.

Copper Vapor Laser¹⁰ This pulsed laser provides high-average powers of up to 100 W at wavelengths of 510.5 and 578.2 nm. The copper laser and other metal vapor lasers of this class, including gold and lead lasers, typically operate at a repetition rate of up to 20 kHz with a current pulse duration of 10 to 50 ns and a laser output of 1 to 10 mJ/pulse. The copper lasers operate at temperatures in the range of 1600°C in 2- to 10-cm-diameter temperature-resistant tubes typically 100 to 150 cm in length. The lasers are self-heated such that all of the energy losses from the discharge current provide heat to bring the plasma tube to the required operating temperature. Excitation occurs by electron collisions with copper atoms vaporized in the plasma tube as indicated in Fig. 7.

Carbon-Dioxide Laser¹⁰ The CO₂ laser, operating primarily at a wavelength of 10.6 μm, is one of the most powerful lasers in the world, producing cw powers of over 100 kW and pulsed energies of up to 10 kJ. It is also available in small versions with powers of up to 100 W from a laser the size of a shoe box. CO₂ lasers typically operate in a mixture of carbon dioxide, nitrogen, and helium gases. Electron collisions excite the metastable levels in nitrogen molecules with subsequent transfer of that energy to carbon dioxide laser levels as shown in Fig. 11. The helium gas acts to keep the average electron energy high in the gas discharge region and to cool or depopulate the lower laser level. This laser is one of the most efficient lasers, with conversion from electrical energy to laser energy of up to 30 percent.

Excimer Laser²⁰ The rare gas-halide excimer lasers operate with a pulsed output primarily in the ultraviolet spectral region at 351 nm in xenon fluoride, 308 nm in xenon chloride, 248 nm in krypton fluoride, and 193 nm in argon fluoride. The laser output, with pulse durations of 10 to 50 ns, is typically of the order of 0.2 to 1.0 J/pulse at repetition rates up to several hundred hertz. The lasers are relatively efficient (1 to 5 percent) and are of a size that would fit on a desktop. The excitation occurs via electrons within the discharge colliding with and ionizing the rare gas molecules and at the same time disassociating the halogen molecules to form negative halogen ions. These two species

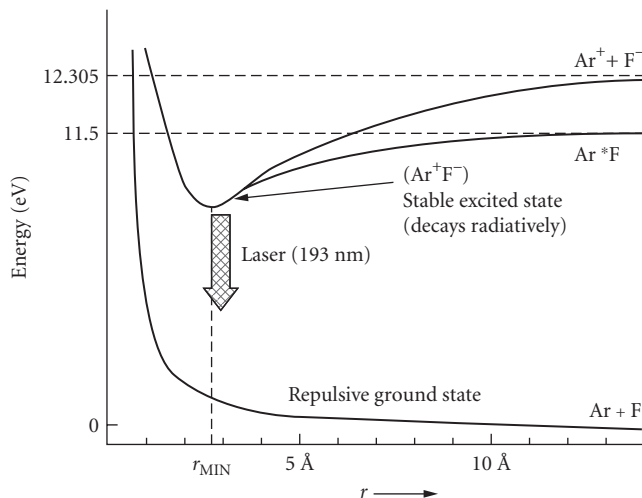


FIGURE 26 Energy-level diagram of an argon fluoride (ArF) excimer laser showing the stable excited state (upper laser level) and the unstable or repulsive ground state.

then combine to form an excited molecular state of the rare gas-halogen molecule which serves as the upper laser state. The molecule then radiates at the laser transition and the lower level advantageously disassociates since it is unstable, as shown in Fig. 26, for the ArF excimer molecule. The excited laser state is an excited state dimer which is referred to as an *excimer state*.

X-Ray Laser²¹ Laser output in the soft-x-ray spectral region has been produced in plasmas of highly ionized ions of a number of atomic species. The highly ionized ions are produced by the absorption of powerful solid-state lasers focused onto solid material of the desired atomic species. Since mirrors are not available for most of the soft-x-ray laser wavelengths (4 to 30 nm), the gain has to be high enough to obtain laser output in a single pass through the laser amplifier [Eq. (16)]. The lasers with the highest gain are the selenium laser (Se^{24+}) at 20.6 and 20.9 nm and the germanium laser (Ge^{22+}) at 23.2, 23.6, and 28.6 nm.

Liquid Laser Gain Media

Organic Dye Lasers⁸ A dye laser consists of a host or solvent material, such as alcohol or water, into which is mixed a laser species in the form of an organic dye molecule, typically in the proportion of one part in ten thousand. A large number of different dye molecules are used to make lasers covering a wavelength range of from 320 to 1500 nm with each dye having a laser bandwidth of the order of 30 to 50 nm. The wide, homogeneously broadened gain spectrum for each dye allows laser tunability over a wide spectrum in the ultraviolet, visible, and near-infrared. Combining the broad gain spectrum (Fig. 27) with a diffraction-grating or prism-tuning element allows tunable laser output to be obtained over the entire dye emission spectrum with a laser linewidth of 10 GHz or less. Dye lasers are available in either pulsed (up to 50 to 100 mJ/pulse) or continuous output (up to a few watts) in tabletop systems that are pumped by either flash lamps or by other lasers such as frequency doubled or tripled YAG lasers or argon ion lasers. Most dye lasers are arranged to have the dye mixture circulated by a pump into the gain region from a much larger reservoir since the dyes degrade at a slow rate during the excitation process.

Dye lasers, with their broad gain spectrum, are particularly attractive for producing ultrashort mode-locked pulses. Some of the shortest light pulses ever generated, of the order of 6×10^{-15} seconds,

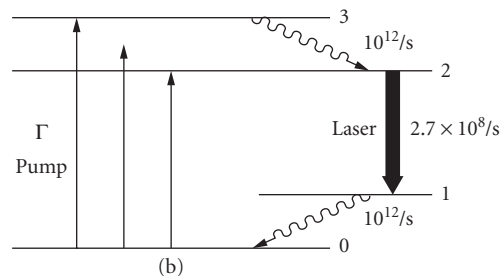
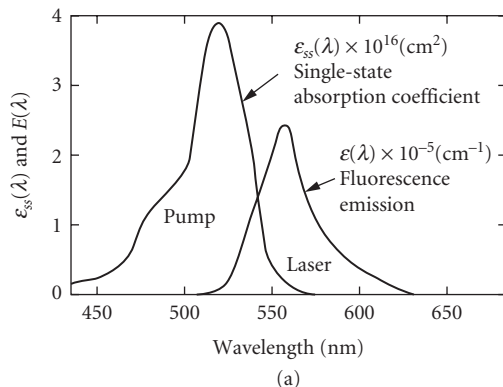


FIGURE 27 Absorption and emission spectra along with the energy-level diagram of an Rh6G organic dye laser showing both (a) the broad pump and emission bandwidths and (b) the fast decay of the lower laser level.

were produced with mode-locked dye lasers. A mode-locked dye laser cavity is shown in Fig. 14 with a thin dye gain region located within an astigmatically compensated laser cavity.

Dielectric Solid-State Laser Gain Media

Ruby Laser^{5,7} The ruby laser, with an output at 694.3 nm, was the first laser ever developed. It consisted of a sapphire (Al_2O_3) host material into which was implanted a chromium laser species in the form of Cr^{3+} ions at a concentration of 0.05 percent as the amplifying medium. The ruby laser involves a three-level optical pumping scheme, with the excitation provided by flash lamps, and operates either in a pulsed or cw mode. The three-level scheme for the ruby laser (Fig. 6) requires a large fraction of the population to be pumped out of the ground state before an inversion occurs. Therefore, the ruby laser is not as efficient as other solid-state lasers such as the Nd:YAG laser which employs a four-level pumping scheme (Fig. 8). It is therefore no longer used as much as it was in the early days.

Nd:YAG and Nd:Glass Lasers⁷ Neodymium atoms, in the form of Nd^{3+} ions, are doped into host materials, including crystals such as yttrium-aluminum-garnet (YAG) and various forms of glass, in concentrations of approximately one part per hundred. The pumping scheme is a four-level system as shown in Fig. 8. When implanted in YAG crystals to produce what is referred to as an Nd:YAG laser, the laser emits primarily at 1.06 μm with continuous powers of up to 250 W and with pulsed

powers as high as several mega watts. Difficulties in growing Nd:YAG crystals limit the size of commonly used laser rods to a maximum of 1 cm in diameter and 10 cm in length. Although this size limitation is somewhat restrictive, a YAG crystal has the advantage of high thermal conductivity allowing the rapid removal of wasted heat due to inefficient excitation. Efforts are being made to use Nd:YAG powders to be able to construct larger Nd:YAG laser rods, thereby providing the advantages of Nd:YAG laser without their usual size and laser power limitations. Slab geometries have recently been developed to compensate for focusing due to thermal gradients and gradient-induced stresses within the amplifier medium, thereby allowing higher average powers to be achieved. Efficient GaAs laser diodes are also being used to pump Nd:YAG amplifiers (see Fig. 15) since the diode pump wavelength matches a very strong pump absorption wavelength for the Nd^{3+} ion. This absorption wavelength is near the threshold pump energy, thereby minimizing the generation of wasted heat. Also, the diode pump laser can be made to more efficiently match the Nd laser mode volume than can a flash lamp.

Under long-pulse, flash lamp-pumped operation, Nd:YAG lasers can produce 5 J/pulse at 100 Hz from a single rod. The pulses are trains of relaxation oscillation spikes lasting 3 to 4 ms which is obtained by powerful, long-pulse flash lamp pumping. Q-switched Nd:YAG lasers, using a single Nd:YAG laser rod, can provide pulse energies of up to 1.5 J/pulse at repetition rates up to 200 Hz (75 W of average power).

Currently power levels of solid state Nd:YAG lasers are approaching 2 kW with repetition rates of nearly 10 kHz. These lasers are presently being developed for use in EUV microlithography as well as in applications such as liquid crystal display production, micromachining, and silicon processing (where carbon dioxide lasers are presently being used.)

Nd:glass lasers have several advantages over Nd:YAG lasers. They can be made in much larger sizes, allowing the construction of very large amplifiers. Nd:glass has a wider gain bandwidth which makes possible the production of shorter mode-locked pulses and a lower stimulated emission cross section. The latter property allows a larger population inversion to be created and thus more energy to be stored before energy is extracted from the laser. Nd:glass amplifiers operate at a wavelength near that of Nd:YAG, at 1.054 μm , for example, with phosphate glass, with a gain bandwidth 10 to 15 times broader than that of Nd:YAG. This larger gain bandwidth lowers the stimulated emission cross section compared to that of Nd:YAG, which allows increased energy storage when used as an amplifier, as mentioned above.

The largest laser system in the world is located at the National Ignition Facility (NIF) at Lawrence Livermore National Laboratories in Livermore, California. It consists of 192 separate laser beams traveling about 1000 ft from their origination (one of two master laser oscillators), to the center of a target chamber. The amplifiers are Nd-doped phosphate glass disc-shaped amplifiers with diameters of over a meter, installed at Brewster's angle within the laser beam path to reduce reflection losses. The beam of each of the 192 amplifiers is capable of producing an energy of 20,000 J with pulse durations ranging from 100 fs to 25 ps. One of the objectives of the laser facility is to test the concept of laser fusion as a source of useful commercial power production.

Neodymium-YLF Lasers Nd:YLF has relatively recently become an attractive laser with the successful implementation of diode laser pumping of solid state laser materials. It has the advantage over Nd:YAG in that the upper-laser-level lifetime is twice as long, thereby allowing twice the energy storage. Its lower stimulated emission cross section than that of Nd:YAG also increases the energy storage and thus the power output per unit volume of the crystal. It has a relatively large thermal conductivity, similar to Nd:YAG. Also the output is polarized and the crystal exhibits low thermal birefringence. Because its emission wavelength matches that of phosphate glass, the laser makes an ideal laser oscillator for seeding large Nd:glass amplifiers for laser fusion studies.

Neodymium:Yttrium Vanadate Lasers The Nd:VO_4 laser has also come into prominence with the use of diode pumping. Its advantages over Nd:YAG include a 5 times larger stimulated emission cross section (and hence higher gain) as well as a four times stronger absorption cross section with a 6 times wider pump absorption cross section centered at 809 nm. Therefore it can be effectively pumped in crystals of only a few millimeters in length and is therefore attractive for use in producing

small diode-pumped lasers. Typically this laser is frequency doubled or tripled intracavity to produce several watts of power at either 532 or 355 nm. It can be operated either cw or Q-switched at repetition rates of up to 100 KHz. It can also be operated at 1.342 μm .

Fiber Lasers The first major application of fiber lasers was in their use as amplifiers in the field of optical communications. The erbium-doped fiber is installed directly in the fiber-optic transmission line and pumped through the fiber itself. The useful wavelengths are in the 1.53 μm region where optical fibers have their lowest loss. A more recent application is in the production of high power output from large mode-area fibers. These lasers utilize single-emitter semiconductor diodes as the light source to pump the cladding of rare-earth doped optical fibers. Pulsed Ytterbium fiber lasers have demonstrated high average powers and peak powers of up to 25 kW with variable repetition rates from a few kilohertz to up to 400 kHz. Fiber lasers are being considered as the replacement technology for conventional solid state lasers due to their compactness, high wall-plug efficiency, high average and peak powers, stability, close to diffraction-limited beam quality and lack of thermal effects.

Ti:Al₂O₃ Laser and Other Broad Bandwidth Solid-State Lasers⁷

Another class of solid-state lasers provides emission and gain over a bandwidth of the order of 100 to 400 nm in the near-infrared, as indicated in Fig. 28. The pump absorption band is in the visible spectrum, thus allowing such pump sources as flash lamps and other lasers. The Ti:Al₂O₃ laser is perhaps the most well-known laser of this category in that it has the widest bandwidth, covering a wavelength range from 0.67 μm to greater than 1.07 μm . Other lasers of this type include alexandrite (Cr:BeAl₂O₄), lasing from 0.7 to 0.8 μm , and Cr:LiSAF which lases from 0.8 to 1.05 μm . These lasers are used in applications where either wide tunability or short-pulse production are desired. Pulses as short as 4 fs have been produced with mode-locked versions of these lasers. Ti:Al₂O₃ lasers, although offering very wide gain bandwidth, have relatively short upper-laser-level lifetimes (3 μs), thereby making them less efficient when pumping with conventional flash lamps. Cr:LiSAF lasers

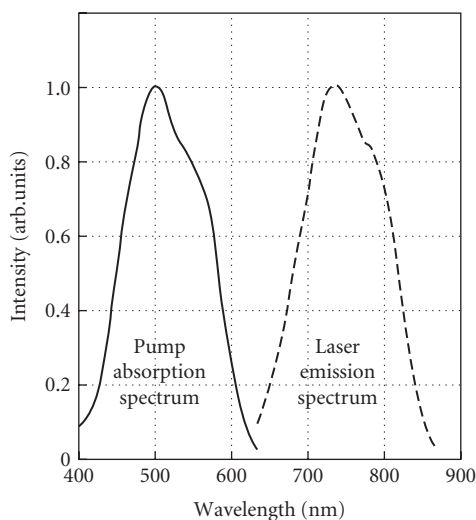


FIGURE 28 Absorption and emission spectra of a Ti:Al₂O₃ laser crystal showing the extremely broad emission (gain) spectrum for this material.

have longer upper-level lifetimes (67 μs), closer to that of Nd:YAG or Nd:glass and have demonstrated efficient laser operation with pumping technologies developed for Nd:YAG lasers.

Color-Center Laser²² Color-center laser gain media are produced by a different form of impurity species than most solid-state lasers. Special defect centers (F-centers) are produced within alkali-halide crystals at a density of 1 part in 10,000 by irradiation with x rays. These defect centers have absorbing regions in the visible portion of the spectrum and emission (and gain) in the near-infrared. A variety of crystals are used to span the laser wavelength spectrum from 0.8 to 4.0 μm . Disadvantages of color-center lasers include operation at temperatures well below room temperature and the necessity to re-form the color centers at intervals of weeks or months in most cases.

Semiconductor Laser Gain Media

Semiconductor Lasers²³ Semiconductor or diode lasers, typically about the size of a grain of salt, are the smallest lasers yet devised. They consist of a p - n junction formed in semiconductor crystal such as GaAs or InP in which the p -type material has an excess of holes (vacancies due to missing electrons) and the n -type material has an excess of electrons. When these two types of materials are brought together to form a junction, and an electric field in the form of a voltage is applied across the junction in the appropriate direction, the electrons and holes are brought together and recombine to produce recombination radiation at or near the wavelength associated with the bandgap energy of the material. The population of electrons and holes within the junction provides the upper-laser-level population, and the recombination radiation spectrum is the gain bandwidth $\Delta\nu$ of the laser, typically of the order of 0.5 to 1.0 nm.

The extended gain length required for these lasers is generally provided by partially reflecting cleaved parallel faces at the ends of the crystals which serve as an optical cavity. Because the cavity is so short, the longitudinal modes are spaced far apart in frequency ($\Delta\nu \cong 1-5 \times 10^{11}$ Hz or several tenths of a nanometer), and thus it is possible to obtain single longitudinal mode operation in such lasers. They require a few volts to operate with milliamperes of current.

Heterostructure semiconductor lasers include additional layers of different materials of similar electronic configurations, such as aluminum, indium, and phosphorous, grown adjacent to the junction to help confine the electron current to the junction region in order to minimize current and heat dissipation requirements (see Fig. 29). The laser mode in the transverse direction is either

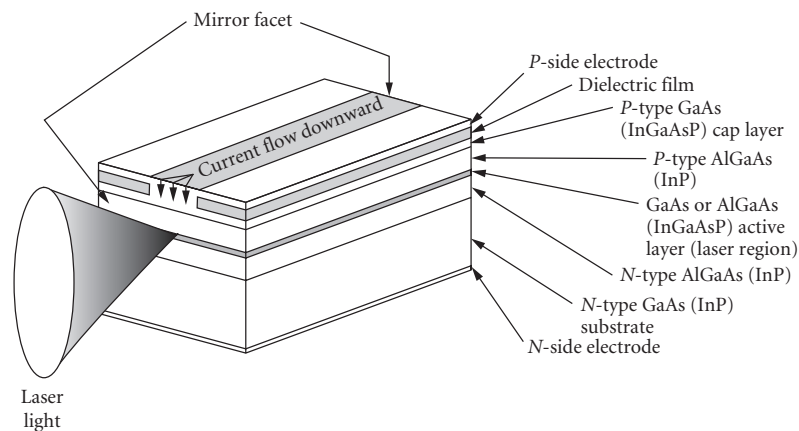


FIGURE 29 A typical heterostructure semiconductor laser showing the various layers of differential materials and the narrow region where current flows to produce the gain region in the active layer.

controlled by gain guiding, in which the gain is produced over a specific narrow lateral extent determined by fabrication techniques, or by index guiding, in which the index of refraction in the transverse direction is varied to provide total internal reflection of a guided mode. Quantum-well lasers have a smaller gain region (cross section), which confines the excitation current and thus the laser mode to an even smaller lateral region, thereby significantly reducing the threshold current and also the heat dissipation requirements. Because of these low threshold requirements, quantum-well semiconductor lasers are used almost exclusively for most semiconductor laser applications.

Multiple semiconductor lasers fabricated within the same bulk structure, known as *semiconductor arrays*, can be operated simultaneously to produce combined cw power outputs of up to 10 W from a laser crystal of dimensions of the order of 1 mm or less. Semiconductor lasers have also been fabricated in multiple arrays mounted vertically on a chip with the mirrors in the plane of the chip. They are known as *vertical cavity semiconductor lasers* (see below). Each of the individual lasers has dimensions of the order of 5 μm and can be separately accessed and excited.

Semiconductor lasers operate over wavelengths ranging from 400 nm to 2.2 μm by using special doping materials to provide the expanded or contracted bandgap energies that provide the varied wavelengths. The newest additions to this class of lasers are based upon the GaN laser materials with the active region consisting of various combinations of InGaN dopings that provide laser wavelengths in the green, blue, and violet portions of the spectrum.

Quantum Cascade Lasers Laser action in semiconductor lasers typically occurs when recombination radiation occurs across the band gap of the semiconductor. Quantum cascade lasers are different in that the radiation occurs from transitions between quantized conduction band states. Such transitions are inherently low-energy transitions and hence the laser output occurs in the middle infrared at wavelengths ranging from 3 to 24 μm . In the laser process, electrons are fed into the injector region of each stage via an electric field and transition across a mini-band from $n = 3$ to $n = 2$ quantum well levels thereby emitting the near infrared radiation. Several of these stages are stacked together in series to produce high power output. The active region is typically made up of aluminum indium arsenide and gallium indium arsenide. The lasers are particularly useful for operation in the two atmospheric windows at 3 to 5 μm and 8 to 13 μm .

Vertical Cavity Surface-Emitting Lasers (VCSELs) The vertical cavity laser is a different type of semiconductor laser than the typical edge-emitting lasers in that the emission occurs from the top surface of the laser and the cavity mirrors are comprised of multilayer dielectric coatings on the top and bottom surfaces of the very thin gain medium. These lasers can be made, using lithographic techniques, in large arrays on a microchip and can also be tested on the chip before being cleaved into individual lasers. Applications include optical fiber data transmission, absorption spectroscopy, and laser printers.

Laser Gain Media in Vacuum

Free Electron Laser²⁴ Free-electron lasers are significantly different than any other type of laser in that the laser output does not result from transitions between discrete energy levels in specific materials. Instead, a high-energy beam of electrons, such as that produced by a synchrotron, traveling in a vacuum with kinetic energies of the order of 1 MeV, are directed to pass through a spatially varying magnetic field produced by two regular arrays of alternating magnet poles located on opposite sides of the beam as shown in Fig. 30. The alternating magnetic field causes the electrons to oscillate back and forth in a direction transverse to the beam direction. The frequency of oscillation is determined by the electron beam energy, the longitudinal spacing of the alternating poles (the magnet period), and the separation between the magnet arrays on opposite sides of the beam. The transverse oscillation of the electrons causes them to radiate at the oscillation frequency and to thereby stimulate other electrons to oscillate and thereby radiate at the same frequency, in phase with the originally oscillating electrons. The result is an intense tunable beam of light emerging from the end

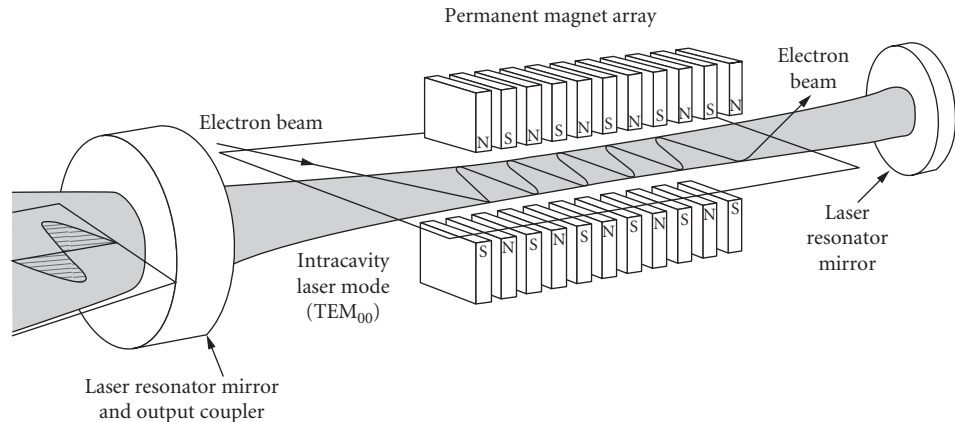


FIGURE 30 A general diagram of a free-electron laser showing how the electron beam is introduced into the cavity and how the alternating magnets cause the beam to oscillate to produce laser radiation.

of the device. Mirrors can be placed at the ends of the radiated beam to produce feedback and extra amplification.

Free-electron lasers have operated at wavelengths ranging from the near-ultraviolet ($0.25\ \mu\text{m}$) to the far-infrared (6 mm) spectral regions. They are efficient devices and also offer the potential of very high average output powers.

16.7 REFERENCES

1. A. Corney, *Atomic and Laser Spectroscopy*, Clarendon Press, Oxford, 1977.
2. P. W. Milonni and J. H. Everly, *Lasers*, John Wiley & Sons, New York, 1988.
3. H. G. Kuhn, *Atomic Spectra*, 2d ed., Academic Press, New York, 1969.
4. J. T. Verdeyen, *Laser Electronics*, 2d ed., Prentice-Hall, Englewood Cliffs, New Jersey, 1989.
5. T. H. Maiman, *Nature* **187**:493–494 (1960).
6. W. T. Walter, N. Solimene, M. Piltch, and G. Gould, *IEEE J. of Quant. Elect.* **QE-2**:474–479 (1966).
7. W. Koehner, *Solid-State Laser Engineering*, 3d ed., Springer-Verlag, New York, 1992.
8. F. J. Duarte and L. W. Hillman (eds.), *Dye Laser Principles*, Academic Press, New York, 1990.
9. W. W. Rigrod, *J. Appl. Phys.* **34**:2602–2609 (1963); **36**:2487–2490 (1965).
10. C. S. Willett, *An Introduction to Gas Lasers: Population Inversion Mechanisms*, Pergamon Press, Oxford, 1974.
11. A. G. Fox and T. Li, *Bell Syst. Tech. J.* **40**:453–488 (1961).
12. G. D. Boyd and J. P. Gordon, *Bell Syst. Tech. J.* **40**:489–508 (1961).
13. H. Kogelnik and T. Li, *Proc. IEEE* **54**:1312–1329 (1966).
14. A. E. Siegman, *Lasers*, University Science Books, Mill Valley, Calif., 1986.
15. W. Wagner and B. Lengyel, *J. Appl. Phys.* **34**:2040–2046 (1963).
16. R. L. Fork, B. I. Greene, and C. V. Shank, *Appl. Phys. Lett.* **38**:671–672 (1981).
17. J. Mark, L. Y. Liu, K. L. Hall, H. A. Haus, and E. P. Ippen, *Opt. Lett.* **14**:48–50 (1989).
18. D. K. Negus, L. Spinelli, N. Goldblatt, and G. Feuget, *Proc. on Advanced Solid State Lasers* (Optical Soc. of Am.) **10**:120–124 (1991).
19. H. Kogelnik and C. V. Shank, *J. Appl. Phys.* **43**:2327–2335 (1972).

20. M. Rokni, J. A. Mangano, J. H. Jacobs, and J. C Hsia, *IEEE Journ. of Quant. Elect.* **QE-14**:464–481 (1978).
21. R. C. Elton, *X-Ray Lasers*, Academic Press, New York, 1990.
22. L. F. Mollenauer and D. H. Olson, *J. Appl. Phys.* **46**:3109–3118 (1975).
23. A. Yariv, *Quantum Electronics*, John Wiley & Sons, New York, 1989.
24. L. R. Elias, W. M. Fairbank, J. M. Madey, H. A. Schewttman, and T. J. Smith, *Phys. Rev. Lett.* **36**:717–720 (1976).

LIGHT-EMITTING DIODES

Roland H. Haitz, M. George Craford, and Robert H. Weissman

Hewlett-Packard Co.

San Jose, California

17.1 GLOSSARY

c	velocity of light
E_g	semiconductor energy bandgap
h	Planck's constant
I_T	total LED current
J	LED current density
k	Boltzmann's constant
M	magnification
n_0	low index of refraction medium
n_1	high index of refraction medium
q	electron charge
T	temperature
V	applied voltage
η_i	internal quantum efficiency
θ_c	critical angle
λ	emission wavelength
τ	total minority carrier lifetime
τ_n	nonradiative minority carrier lifetime
τ_r	radiative minority carrier lifetime

17.2 INTRODUCTION

Over the past 25 years the light-emitting diode (LED) has grown from a laboratory curiosity to a broadly used light source for signaling applications. In 1992 LED production reached a level of approximately 25 billion chips, and \$2.5 billion worth of LED-based components were shipped to original equipment manufacturers.

This chapter covers light-emitting diodes from the basic light-generation processes to descriptions of LED products. First, we will deal with light-generation mechanisms and light extraction. Four major types of device structures—from simple grown or diffused homojunctions to complex double heterojunction devices are discussed next, followed by a description of the commercially important semiconductors used for LEDs, from the pioneering GaAsP system to the AlGaInP system that is currently revolutionizing LED technology. Then processes used to fabricate LED chips are explained—the growth of GaAs and GaP substrates, the major techniques used for growing the epitaxial material in which the light-generation processes occur, and the steps required to create LED chips up to the point of assembly. Next the important topics of quality and reliability—in particular, chip degradation and package-related failure mechanisms—will be addressed. Finally, LED-based products, such as indicator lamps, numeric and alphanumeric displays, optocouplers, fiber-optic transmitters, and sensors, are described.

This chapter covers the mainstream structures, materials, processes, and applications in use today. It does not cover certain advanced structures, such as quantum well or strained layer devices, a discussion of which can be found in Chap. 19, “Semiconductor Lasers.” The reader is also referred to Chap. 19 for current information on edge-emitting LEDs, whose fabrication and use are similar to lasers.

For further information on the physics of light generation, the reader should consult Refs. 1 to 11. Semiconductor material systems for LEDs are discussed in Refs. 13 to 24. Crystal growth, epitaxial, and wafer fabrication processes are discussed in detail in Refs. 25 to 29.

17.3 LIGHT-GENERATION PROCESSES

When a p - n junction is biased in the forward direction, the resulting current flow across the boundary layer between the p and n regions has two components: holes are injected from the p region into the n region and electrons are injected from the n region into the p region. This so-called minority-carrier injection disturbs the carrier distribution from its equilibrium condition. The injected minority carriers recombine with majority carriers until thermal equilibrium is reestablished. As long as the current continues to flow, minority-carrier injection continues. On both sides of the junction, a new steady-state carrier distribution is established such that the recombination rate equals the injection rate.^{1,2}

Minority-carrier recombination is not instantaneous. The injected minority carriers have to find proper conditions before the recombination process can take place. Both energy and momentum conservation have to be met. Energy conservation can be readily met since a photon can take up the energy of the electron-hole pair, but the photon doesn't contribute much to the conservation of momentum. Therefore, an electron can only combine with a hole of practically identical and opposite momentum. Such proper conditions are not readily met, resulting in a delay. In other words, the injected minority carrier has a finite lifetime τ_r before it combines radiatively through the emission of a photon.² This average time to recombine radiatively through the emission of light can be visualized as the average time it takes an injected minority carrier to find a majority carrier with the right momentum to allow radiative recombination without violating momentum conservation.

Unfortunately, radiative recombination is not the only recombination path. There are also crystalline defects, such as impurities, dislocations, surfaces, etc., that can trap the injected minority carriers. This type of recombination process may or may not generate light. Energy and momentum conservation are met through the successive emission of phonons. Again, the recombination process is not instantaneous because the minority carrier first has to diffuse to a recombination site. This nonradiative recombination process is characterized by a lifetime τ_n .²

Of primary interest in design of light-emitting diodes is the maximization of the radiative recombination relative to the nonradiative recombination. In other words, it is of interest to develop conditions where radiative recombination occurs fairly rapidly compared with nonradiative recombination. The effectiveness of the light-generation process is described by the fraction of the injected minority carriers that recombine radiatively compared to the total injection. The internal quantum efficiency η_i can be calculated from τ_r and τ . The combined recombination processes lead to a total minority-carrier lifetime τ given by Eq. (1):

$$\frac{1}{\tau} = \frac{1}{\tau_r} + \frac{1}{\tau_n} \quad (1)$$

η_i is simply computed from Eq. (1) as the fraction of carriers recombining radiatively:²

$$\eta_i = \frac{\tau_n}{\tau_r + \tau_n} \quad (2)$$

Of interest are two simple cases: in the case of excellent material quality (large τ_n) or efficient radiative recombination conditions (small τ_r), the internal quantum efficiency approaches 100 percent. For the opposite case ($\tau_n \ll \tau_r$), we find $\eta_i \approx \tau_n/\tau_r \ll 1$. As discussed under “Material Systems,” there are several families of III–V compounds with internal quantum efficiencies approaching 100 percent. There are also other useful semiconductor materials with internal quantum efficiencies in the 1- to 10-percent range.

To find material systems for LEDs with a high quantum efficiency, one has to understand the band structure of semiconductors. The band structure describes the allowed distribution of energy and momentum states for electrons and holes (see Fig. 1 and Ref. 2). In practically all semiconductors the lower band, also known as the *valence band*, has a fairly simple structure, a

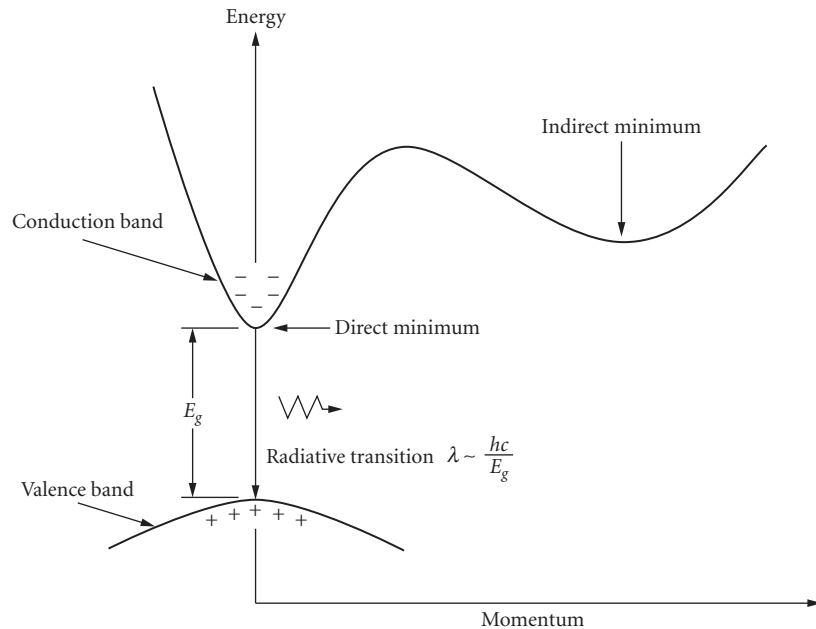


FIGURE 1 Energy band structure of a direct semiconductor showing radiative recombination of electrons in the conduction band with holes in the valence band.

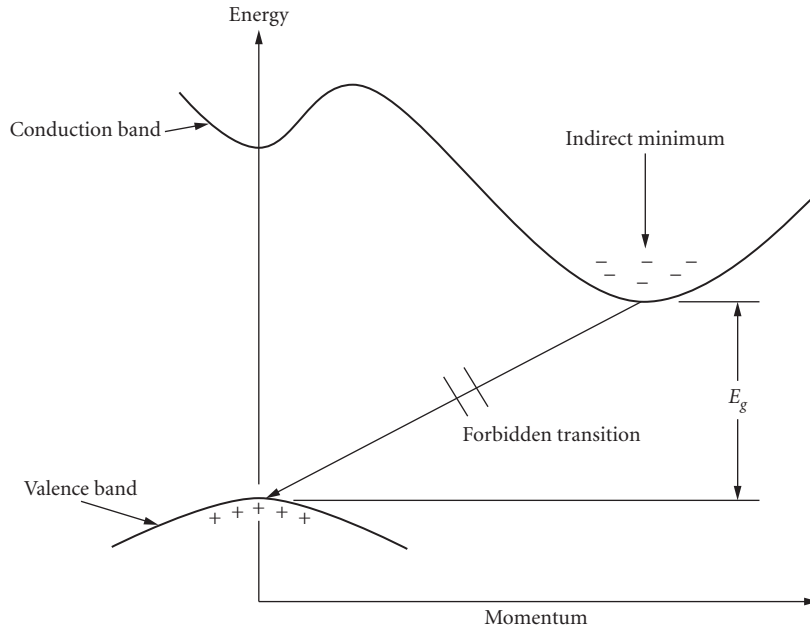


FIGURE 2 Energy band structure of an indirect semiconductor showing the conduction-band minima and valence-band maximum at different positions in momentum space. Radiative recombination of conduction-band electrons and valence-band holes is generally forbidden.

paraboloid around the $\langle 0, 0, 0 \rangle$ crystalline direction. Holes will take up a position near the apex of the paraboloid and have very small momentum. The upper band, also known as the *conduction band*, is different for various semiconductor materials. All semiconductors have multiple valleys in the conduction band. Of practical interest are the valleys with the lowest energy. Semiconductor materials are classified as either *direct* or *indirect*.¹⁻³ In a direct semiconductor, the lowest valley in the conduction band is directly above the apex of the valence-band paraboloid. In an indirect semiconductor, the lowest valleys are not at $\langle 0, 0, 0 \rangle$, but at different positions in momentum per energy space (see Fig. 2). Majority or minority carriers mostly occupy the lowest energy states, i.e., holes near the top of the valence band paraboloid and electrons near the bottom of the lowest conduction-band valley.

In the case of a direct semiconductor the electrons are positioned directly above the holes at the same momentum coordinates. It is relatively easy to match up electrons and holes with proper momentum-conserving conditions. Thus, the resulting radiative lifetime τ_r is short. On the other hand, electrons in an indirect valley will find it practically impossible to find momentum-matching holes and the resulting radiative lifetime will be long. Injected carriers in indirect material generally recombine nonradiatively through defects.

In a direct semiconductor, such as GaAs, the radiative lifetime τ_r is in the range of 1 to 100 ns, depending on doping, temperature, and other factors. It is relatively easy to grow crystals with sufficiently low defect density such that τ_n is in the same range as τ_r .

For indirect semiconductors, such as germanium or silicon, the radiative recombination process is extremely unlikely, and τ_r is in the range of seconds.¹ In this case, $\tau_r \gg \tau_n$, and practically all injected carriers recombine nonradiatively.

The wavelength of the photons emitted in a radiative recombination event is determined by the energy difference between the recombining electron-hole pair. Since carriers relax quickly to an energy level near the top of the valence band (holes), or the bottom of the conduction band

(electrons), we have the following approximation for the wavelength λ of the emitted photon (see Fig. 1):

$$\lambda \approx hc/E_g \quad (3)$$

where h = Planck's constant, c = velocity of light, and E_g = bandgap energy.

This relation is only an approximation since holes and electrons are thermally distributed at levels slightly below the valence-band maximum and above the conduction-band minimum, resulting in a finite linewidth in the energy or wavelength of the emitted light. Another modification results from a recombination between a free electron and a hole trapped in a deep acceptor state. See Ref. 2 for a discussion of various recombination processes.

To change the wavelength or energy of the emitted light, one has to change the bandgap of the semiconductor material. For example, GaAs with a bandgap of 1.4 eV has an infrared emission wavelength of 900 nm. To achieve emission in the visible red region, the bandgap has to be raised to around 1.9 eV. This increase in E_g can be achieved by mixing GaAs with another material with a wider bandgap, for instance GaP with $E_g = 2.3$ eV. By adjusting the ratio of arsenic to phosphorous the bandgap of the resulting ternary compound, GaAsP, can be tailored to any value between 1.4 and 2.3 eV.³

The resulting band structure with varying As to P ratio is illustrated in Fig. 3. Note, the two conduction-band valleys do not move upward in energy space at the same rate. The direct valley moves up faster than the indirect valley with increasing phosphorous composition. At a composition of around 40 percent GaP and 60 percent GaAs, the direct and indirect valleys are about equal in energy. When the valleys are approximately equal in energy, electrons in the conduction band can scatter from the direct valley into the indirect valley. While the direct valley electrons still undergo rapid radiative recombination, the indirect valley electrons have a long radiative lifetime and either have to be scattered back to the direct valley or they will recombine nonradiatively. In other words, near this crossover

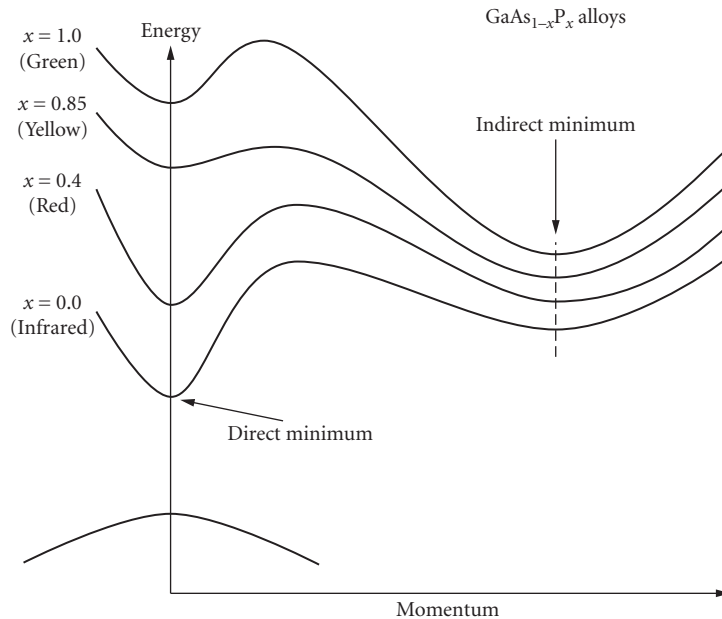


FIGURE 3 Energy band diagram for various alloys of the GaAs_{1-x}P_x material system showing the direct and indirect conduction-band minima for various alloy compositions.

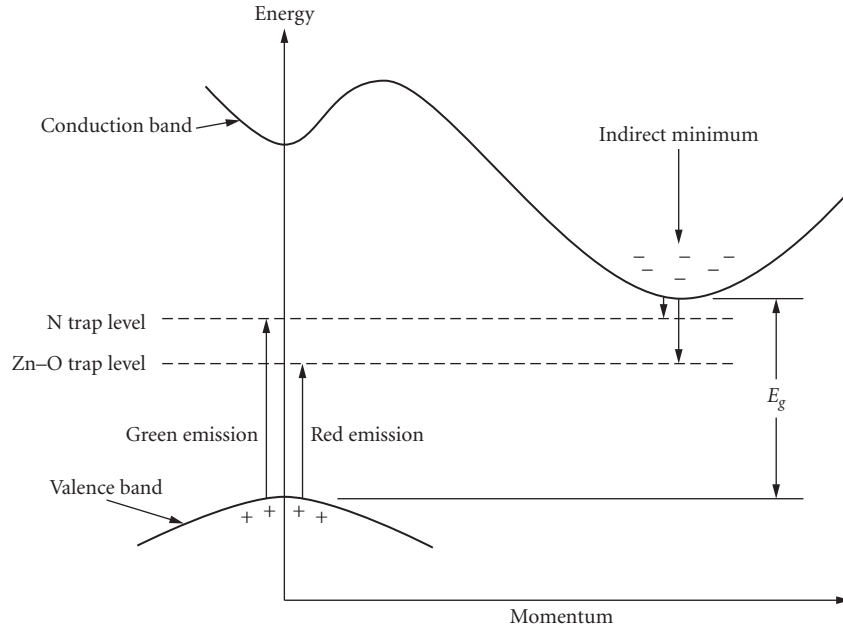


FIGURE 4 Formation of excitons (electron-hole pairs) by the addition of isoelectronic dopants N and ZnO to an indirect semiconductor. The excitons have a high probability to recombine radiatively.

between direct and indirect valleys, the radiative efficiency drops off dramatically, and for compositions with greater than 40 percent phosphorous the direct radiative recombination is practically nonexistent.^{3,4}

The above discussion indicates that indirect semiconductors are not suitable for efficient generation of light through minority-carrier recombination. Fortunately, the introduction of so-called isoelectronic impurities can circumvent this limitation and introduces a new radiative recombination process.⁵ A commonly used isoelectronic trap is generated by substituting a nitrogen atom for phosphorous in the GaAsP system^{6–11} (see Fig. 4). Since N and P both have five electrons in their outer shell, the trap is electrically neutral. However, the stronger electronegativity of N relative to P can result in capture of an electron from the conduction band. Since the electron is very tightly bound to the impurity atom, its wave function in momentum space is spread out and has reasonable magnitude at $\langle k=0 \rangle$ in momentum space.¹⁰ The negatively charged defect can attract a free hole to form a loosely bound electron-hole pair or “exciton.” This electron-hole pair has a high probability to recombine radiatively. The energy of the emitted light is less than E_g . Another isoelectronic trap in GaP is formed by ZnO pairs (Zn on a Ga site and O on a P site) (see Fig. 4). The ZnO trap is deeper than the N trap, resulting in longer wavelength emission in the red region of the spectrum.¹

The recombination process for exciton recombination is quite complex. For a detailed analysis, the reader is referred to Ref. 10. One result of this analysis is the recognition that the bound exciton has a relatively long lifetime in the range of 100 to 1000 ns. Light emission by exciton recombination is generally slower than emission due to direct band-to-band recombination.

17.4 LIGHT EXTRACTION

Generating light efficiently within a semiconductor material is only one part of the problem to build an efficient light source. The next challenge is the extraction of light from within the LED chip to the outside. The designer must consider total internal reflection.¹ According to Snell’s law, light can

escape from a medium of high index of refraction n_1 into a medium of low index refraction n_0 only if it intersects the surface between the two media at an angle from normal less than the critical angle θ_c with θ_c being defined by Eq. (4):

$$\theta_c = \arcsin n_0/n_1 \quad (4)$$

Most semiconductor LEDs have an isotropic emission pattern as seen from within the light-generating material. Assuming a cubic shape for the LED chip, because of internal reflections, only a small fraction of the isotropically emitted light can escape any of the six surfaces. As a case in point, let us calculate the emission through the top surface. For typical light-emitting semiconductors, n_1 is in the range of 2.9 to 3.6. If $n_1 = 3.3$ and $n_0 = 1.0$ (air), we find $\theta_c = 17.6^\circ$. The emission from an isotropic source into a cone with a half angle of θ_c is given by $(1 - \cos \theta_c)/2$. After correcting for Fresnel reflections, only 1.6 percent of the light generated escapes through the LED top surface into air. Depending on chip and p - n junction geometry, virtually all of the remaining light (98.4 percent) is reflected and absorbed within the LED chip.

The fraction of light coupled from chip to air is a function of the number of surfaces through which the chip can transmit light effectively. Most LED chips are called “absorbing substrate” (AS) chips. In such a chip, the starting substrate material (discussed later under “Substrate Technology”) has a narrow bandgap and absorbs all the light with energy greater than the bandgap of the substrate. Consider the case of a GaAsP LED grown on a GaAs substrate. The emitted light ($E_g > 1.9\text{eV}$) is absorbed by the GaAs substrate ($E_g = 1.4\text{ eV}$). Thus, a GaAsP-emitting layer on a GaAs substrate can transmit only through its top surface. Light transmitted toward the side surfaces or downward is absorbed.

To increase light extraction, the substrate or part of the epitaxial layers near the top of the chip has to be made of a material transparent to the emitted light. The “transparent substrate” (TS) chip is designed such that light transmitted toward the side surfaces within θ_c half-angle cones can escape. Assuming that there is negligible absorption between the point of light generation and the side walls, this increases the extraction efficiency by a factor of 5 (5 instead of 1 escape cones).

In a TS chip, additional light can be extracted if the side walls are nonplanar, i.e., if light from outside an escape cone can be scattered into an escape cone. This process increases the optical path within the chip and is very dependent on residual absorption. In a chip with low absorption and randomizing side surfaces, most of the light should escape. Unfortunately, in practical LED structures, there are several absorption mechanisms left such as front and back contacts, crystal defects, and absorption in areas where secondary radiative recombination is inefficient.

A common approach is to use a hybrid chip with properties between AS and TS chips. These chips utilize a thick, transparent window layer above the light-emitting layer. If this layer is sufficiently thick, then most of the light in the top half of the cones transmitted toward the side surfaces will reach the side of the chip before hitting the substrate. In this case of hybrid chips, the efficiency is between that of AS and TS chips as shown in Table 1.

Another important way to increase extraction efficiency is derived from a stepwise reduction in the index of refraction from chip to air. If the chip is first imbedded in a material with an intermediate index, i.e., plastic with $n_2 = 1.5$, then the critical angle θ_c between chip and plastic is increased to 27° . The extraction efficiency relative to air increases by the ratio of $(n_2/n_0)^2$ plus some additional correction for Fresnel-reflection losses. The gain from plastic encapsulation is usually around

TABLE 1 Extraction Efficiency into Air or Plastic for Three Types of Commonly Used LED Chips

Chip Type	No. Cones	Typical extraction efficiency	
		Air (%)	Plastic (%)
AS	1	1.5	4
Thick window	3	4.5	12
TS	5	7.5	20

2.7 times compared to air. Chips with multipath internal reflection will result in lower gains. It is important to note that this gain can be achieved only if the plastic/air interface can accommodate the increased angular distribution through proper lenslike surface shaping or efficient scattering optics. Table 1 illustrates the approximate extraction efficiencies achieved by the three dominant chip structures in air and in plastic. The numbers assume only first-pass extraction, limited absorption, and no multiple reflections within the chip.

17.5 DEVICE STRUCTURES

LED devices come in a broad range of structures. Each material system (see following section) requires a different optimization. The only common feature for all LED structures is the placement of the p - n junction where the light is generated. The p - n junction is practically never placed in the bulk-grown substrate material for the following reasons:

- The bulk-grown materials such as GaAs, GaP, and InP usually do not have the right energy gap for the desired wavelength of the emitted light.
- The light-generating region requires moderately low doping that is inconsistent with the need for a low series resistance.
- Bulk-grown material often has a relatively high defect density, making it difficult to achieve high efficiency.

Because of these reasons, practically all commercially important LED structures utilize a secondary growth step on top of a single-crystal bulk-grown substrate material. The secondary growth step consists of a single-crystal layer lattice matched to the substrate. This growth process is known as *epitaxial growth* and is described in a later section of this chapter.

The commonly used epitaxial structures can be classified into the following categories:

- Homojunctions
 - grown
 - diffused
- Heterojunctions
 - single confinement
 - double confinement

Grown Homojunctions

Figure 5 illustrates one of the simplest design approaches to an LED chip. An n -type GaAs layer with low to moderate doping density is grown on top of a highly doped n -type substrate by a vapor or liquid-phase epitaxial process (see “Epitaxial Technology”). After a growth of 5 to 10 μm , the doping is changed to p type for another 5 to 10 μm . A critical dimension is the thickness of the epitaxial p layer. The thickness should be larger than the diffusion length of electrons. In other words, the electrons should recombine radiatively in the epitaxially grown p layer before reaching the surface. The p layer should be of sufficiently high quality to meet the condition for efficient recombination, i.e., $\tau_n \gg \tau_r$. In addition, the side surfaces may have to be etched to remove damage. Damage and other defects where the p - n junction intercepts the chip surface can lead to a substantial leakage current that reduces efficiency, especially at low drive levels.

The structure of Fig. 5 was used in some of the earliest infrared emitters (wavelength 900 nm). Efficiency was low, typically 1 percent. Modern infrared emitters use Si-doped GaAs (see Fig. 6). The detailed recombination mechanism in GaAs:Si even today is quite controversial and goes beyond the scope of this publication. The recombination process has two important characteristics: (1) the radiative lifetime τ_r is relatively slow, i.e., in the range of 1 μs and (2) the wavelength is shifted to 940 nm.

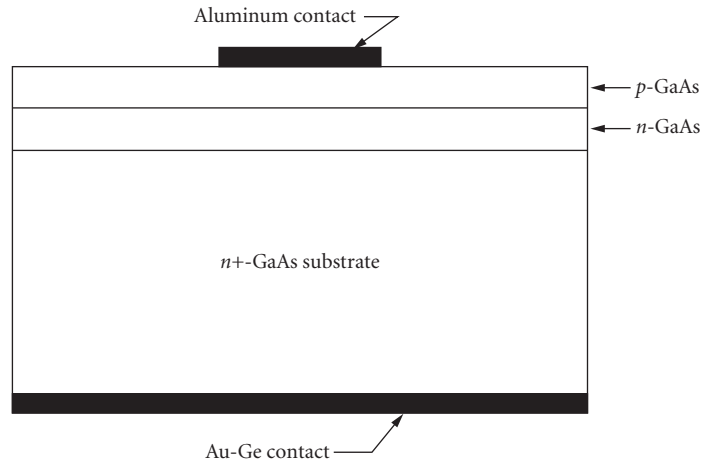


FIGURE 5 Cross section of an infrared LED chip. A p - n junction is formed by epitaxially growing n - and p -doped GaAs onto an n^+ -doped GaAs substrate.

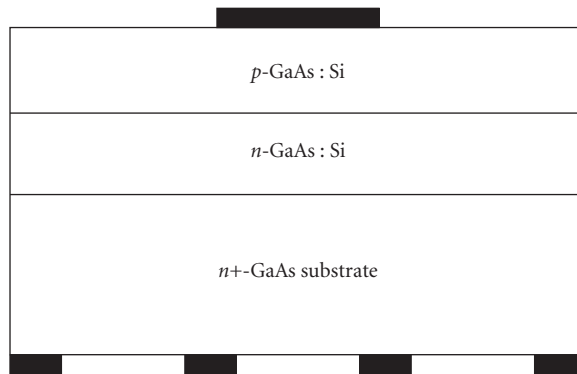


FIGURE 6 A high-efficiency IR LED made by LPE growth of GaAs that is doped with silicon on both the p and n sides of the junction. To increase external quantum efficiency, a partly reflective back contact is employed.

At this wavelength the GaAs substrate is partly transparent, making this device a quasi-TS structure with efficiencies into plastic of 5 to 10 percent.

Diffused Homojunction

The chip structure of Fig. 5 can also be produced by a zinc (Zn) diffusion into a thick n layer. The commercially most significant structure of this type is shown in Fig. 7. By replacing 40 percent of the As atoms with P atoms, the bandgap is increased to 1.92 eV to make a GaAsP LED that emits visible red light. In this case, the p - n junction is diffused selectively by using a deposited layer of silicon-nitride as a diffusion mask. This structure of Fig. 7 has several advantages over the structure of Fig. 5. Lateral diffusion of Zn moves the intersection of the junction with the chip surface

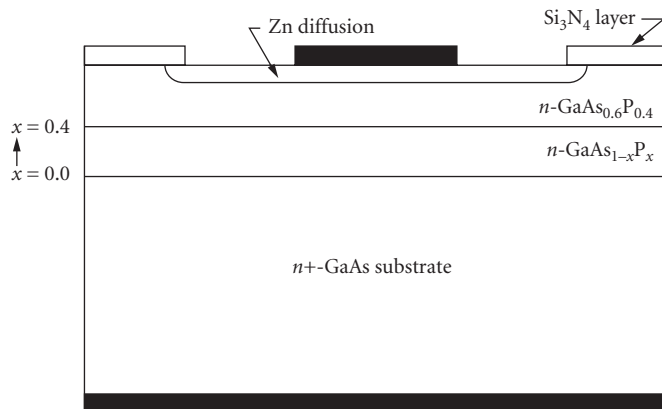


FIGURE 7 GaAsP LED which emits at 650 nm. On a GaAs substrate, a layer is grown whose composition varies linearly from GaAs to GaAs $_{0.6}$ P $_{0.4}$, followed by a layer of constant composition. Zinc is selectively diffused using an Si_3N_4 mask to form the light-emitting junction.

underneath the protecting Si_3N_4 layer. This layer protects the junction from contamination and adds to the long-term stability of the device. In addition, it is important for applications requiring more than one clearly separated light-emitting area. For instance, seven-segment displays, such as those used in the early handheld calculators, are made by diffusing seven long and narrow stripes into a single chip of GaAsP material. (See Fig. 8.) This chip consists of eight (seven segments plus decimal point) individually addressable p -regions (anodes) with a common n -type cathode. Such a chip is feasible only in an AS-type structure because the individual segments have to be optically isolated from each other. A TS-type structure results in unacceptable levels of crosstalk.

Figure 7 shows another feature of practical LED devices. The composition with 40 percent P and 60 percent As has a lattice constant (atomic spacing) that is different from the GaAs substrate. Such a lattice mismatch between adjacent layers would result in a very high density of dislocations. To reduce this problem to an acceptable level, one has to slowly increase the phosphorous composition from 0 percent at the GaAs interface to 40 percent over a 10- to 20- μm -thick buffer layer. Typically, the buffer layer is graded linearly. The phosphorous composition is increased linearly from bottom to top. The thicker the buffer layer, the lower is the resulting dislocation density. Cost constraints keep the layer in the 10- to 15- μm range. The layer of constant composition (40 percent P) has to be thick enough to accommodate the Zn diffusion plus the diffusion length of minority carriers. A thickness range of 5 to 10 μm is typical.

Another variation of a homojunction is shown in the TS-chip structure of Fig. 9. Instead of an absorbing GaAs substrate, one starts with a transparent GaP substrate. The graded layer has an inverse gradient relative to the chip shown in Fig. 7. The initial growth is 100 percent GaP phasing in As linearly over 10 to 15 μm . At 15 percent As, the emission is in the yellow range (585 nm), at 25 percent in the orange range (605 nm), and at 35 percent in the red range (635 nm). Figure 3 shows the approximate band structure of this material system. The composition range mentioned above has an indirect band structure. To obtain efficient light emission, the region of minority-carrier injection is doped with nitrogen forming an isoelectronic recombination center (see exciton recombination in the first section).

Figure 9 shows an important technique to increase extraction efficiency. In a TS-chip, the major light loss is due to free carrier absorption at the alloyed contacts. Rather than covering the entire bottom surface of the chip with contact metal, one can reduce the contact area either by depositing small contact islands (see Fig. 6) or by placing a dielectric mirror (deposited SiO_2) between the substrate and the unused areas of the back contact (Fig. 9). This dielectric mirror increases the efficiency by 20 to 50 percent at the expense of higher manufacturing cost.

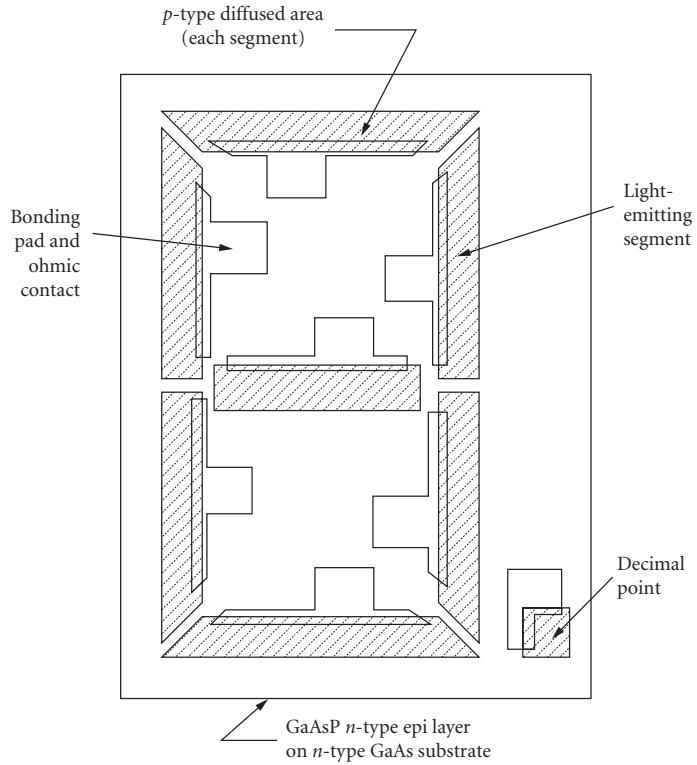


FIGURE 8 Monolithic seven-segment display chip with eight separate diffused regions (anodes) and a common cathode (the GaAs substrate).

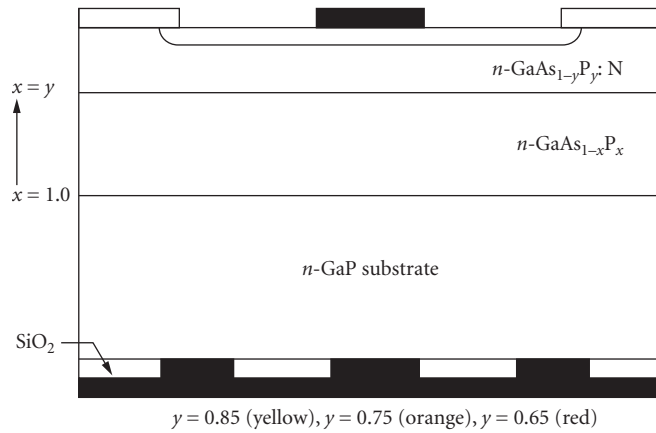


FIGURE 9 Cross section of a $\text{GaAs}_{1-x}\text{P}_x$ LED which, by changing the composition “ x ,” produces red, orange, or yellow light. The top layer is doped with nitrogen to increase the quantum efficiency. The GaP substrate is transparent to the emitted light. Also shown is a reflective back contact made by depositing the contact metallization on top of SiO_2 .

Single Heterojunctions

Heterojunctions introduce a new variable: local variation of the energy bandgap resulting in carrier confinement. Figure 10 shows a popular structure for a red LED emitter chip. A p -type layer of GaAlAs with 38 percent Al is grown on a GaAs substrate. The GaAlAs alloy system can be lattice-matched to GaAs; therefore, no graded layer is required as in the GaAsP system of Fig. 7. Next, an n -type layer is grown with 75 percent Al. The variation of E_g from substrate to top is illustrated in Fig. 11. Holes accumulate in the GaAlAs p layer with the narrower bandgap. Electrons are injected from the n layer into this p layer. The holes have insufficient energy to climb the potential barrier to the wide-bandgap material. Holes are confined to the p layer. In the p layer, the radiative recombination time is very short because of the high concentration of holes. As a result, the internal quantum efficiency is quite high. A variation of this structure is a widely used infrared emitter that emits at 880 nm.

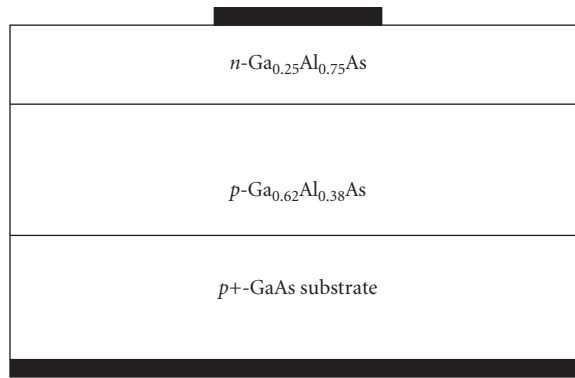


FIGURE 10 Cross section of a single heterostructure LED.

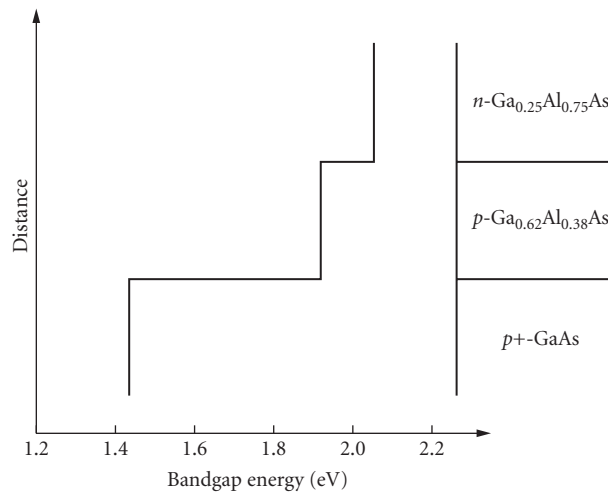


FIGURE 11 Variation in energy bandgap for the various layers in the GaAlAs LED shown in Fig. 10.

Double Heterostructures

The double heterostructure shown in Fig. 12 repeats the p -side confinement of Fig. 10 on the n side. An n -type buffer layer is grown on the GaAs substrate to create a high-quality surface onto which the first n -type GaAlAs confinement layer with 75 percent Al is grown. The active or light-generation layer is a 3- μm -thick p -type layer with 38 percent Al. The top p -type confinement layer again uses 75 percent Al. This structure with the energy-band diagram of Fig. 13 has two advantages: (1) There is no hole injection into the n -type layer with reduced efficiency and a slow hole recombination in the lowly doped n layer. (2) The high electron and hole density in the active layer reduces τ_r , thus increasing device speed and efficiency. The increased speed is quite important for LED sources in fiber-optic communication applications. (See “Fiber Optics” subsection later in this chapter.)

The double heterostructure of Fig. 12 represents a one-dimensional containment of injected carriers. Injection and light emission occurs across the entire lateral dimension of the chip. For fiber-optic applications, the light generated over such a large area cannot be effectively coupled into small-diameter fiber. A rule of thumb for fiber coupling requires that the light-emitting area be equal to or, preferably, smaller than the fiber core diameter. This rule requires lateral constraint of carrier injection. The localized diffusion of Fig. 7 is not applicable to grown structures such as those in Fig. 12. The preferred solution inserts an n layer between buffer and lower confinement layer (see Fig. 14). A hole etched into the n layer allows current flow. Outside of this hole, the p - n junction between n layer and lower confinement layer is reverse-biased, thus blocking any current flow. The disadvantage of this approach is a complication of the growth process. In a first growth process, the n layer is grown. Then a hole is etched into the n layer using standard photolithographic etching techniques. Finally, a second epitaxial growth is used for the remaining layers.

Another technique to constrain current injection utilizes a small ohmic contact.¹² It is used frequently in conjunction with InP-based fiber-optic emitters (see Fig. 15). An SiO_2 layer limits contact to a small-diameter (typically 25 μm) hole that results in a relatively small light-emitting area. The etched lens shown for this structure helps to collimate the light for more efficient coupling onto the fiber.¹²

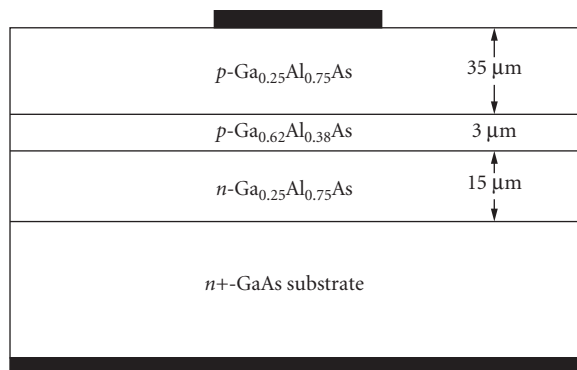


FIGURE 12 Structure of a double heterostructure (DH) GaAlAs LED. The DH is composed of a $\text{Ga}_{0.62}\text{Al}_{0.38}\text{As}$ layer surrounded on either side by a $\text{Ga}_{0.25}\text{Al}_{0.75}\text{As}$ layer. The thick top layer acts as window to increase light extraction through the side walls of the chip.

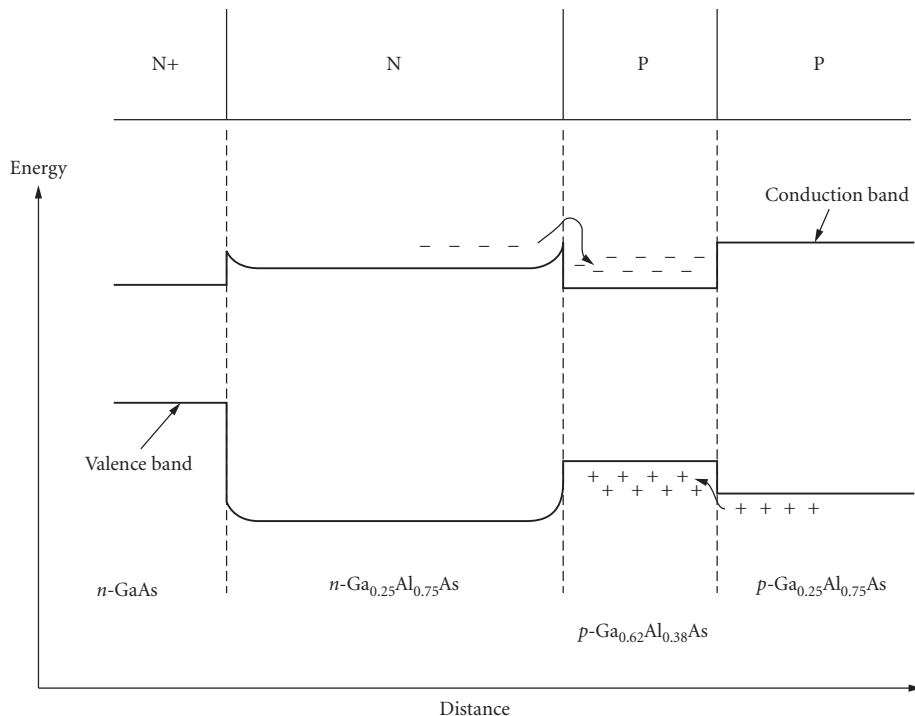


FIGURE 13 Energy band diagram of the GaAlAs LED shown in Fig. 12. The LED is forward biased. Electrons and holes are confined in the p -doped $\text{Ga}_{0.62}\text{Al}_{0.38}\text{As}$ layer, which increases the radiative recombination efficiency.

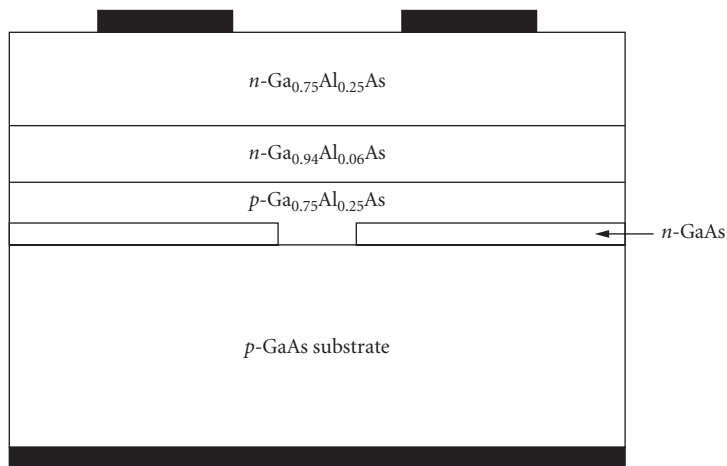


FIGURE 14 Cross section of an LED with three-dimensional carrier confinement. A DH structure is used to confine injected carriers in the $\text{Ga}_{0.94}\text{Al}_{0.06}\text{As}$ layer (direction perpendicular to the junction). The patterned, n -type GaAs layer is used to limit current flow in the lateral direction. The small emitting area and the 820-nm emission of this LED makes it ideal for fiber-optic applications.

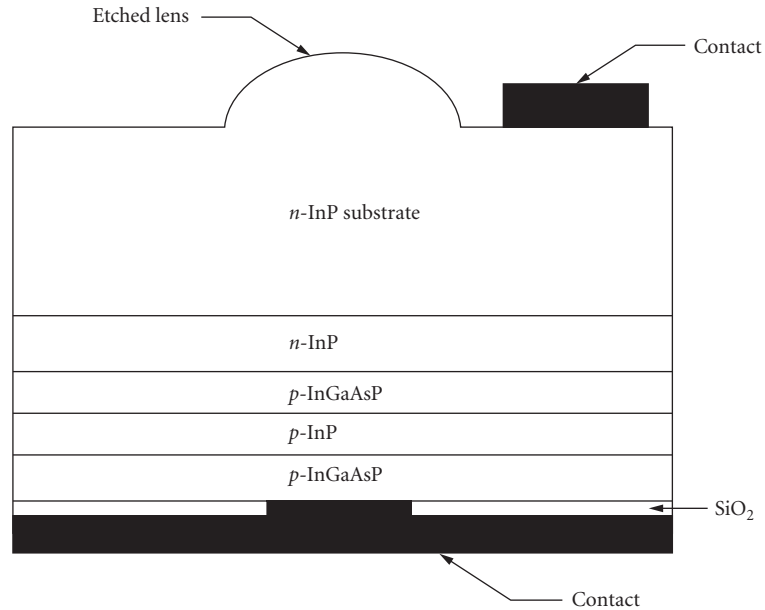


FIGURE 15 Structure of a 1300-nm LED used for optical fiber communications. The cross section shows the DH-layer configuration, limited area back contact for emission-size control, and an etched lens at the top of the chip which, magnifies ($M = 2$) the source area and collimates the light for effective coupling into a fiber.

17.6 MATERIAL SYSTEMS

The GaAs_{1-x}P_x System

The most widely used alloy for LEDs is the ternary GaAs_{1-x}P_x system, including its two binary components GaAs and GaP. This system is best described by the composition parameter x with $0 \leq x \leq 1$. For $x = 0$, we have GaAs and for $x = 1$ the composition is GaP. For $x \leq 0.4$, the alloy has a direct bandgap. GaAs was developed in the early 1960s as an infrared emitter with a wavelength of 910 nm and an efficiency in the range of 1 percent. This emitter was soon followed by a Si-doped variety. As discussed earlier, this leads to an emission wavelength of 940 nm, a wavelength at which the GaAs substrate is partly transparent. The resulting efficiency is increased substantially and, depending on configuration, is in the 5- to 10-percent range. However, the recombination process is quite slow, resulting in rise and fall times in the 50-ns to 1.0- μ s range (see Table 2). One other drawback is caused by the low absorption coefficient of Si detectors at 940 nm. To absorb 90 percent of the light requires a detector thickness of 60 to 70 μ m. Conventional photo transistors are quite suitable as detectors. Integrated photo ICs with their 5- to 7- μ m-thick epitaxial layers are very inefficient as detectors at 940 nm.

To shift the wavelength toward the near-infrared or into the visible spectrum, one has to grow a ternary alloy, a mixture between GaAs and GaP. Of commercial interest are two alloys with $x = 0.3$ and $x = 0.4$ grown on a GaAs substrate (see Table 2). The $x = 0.4$ alloy was the first commercially produced material with a wavelength in the visible range of the spectrum. Grown on an absorbing substrate, it has a modest luminous efficacy of around 0.2 lm/A. [Luminous efficacy is the luminous (visible) flux output measured in lumens divided by the electrical current input.] The absorbing substrate allows the integration of multiple light sources into a single chip without crosstalk. Such

TABLE 2 Performance Summary of LED Chips in the GaAs_{1-x}P_x System

x	Growth Process*	Isol. Dopant	Substrate†	Dominant Wavelength (nm)	Colon‡	Luminous Efficacy§ (lm/A)	Quantum Efficiency§ (%)	Speed (ns)
0	LPE	—	AS	910	IR		1	50
0	LPE	—	TS	940	IR		10	1000
0.3	VPE	—	AS	700	IR		0.5	50
0.4	VPE	—	AS	650	Red	0.2		50
0.65	VPE	N	TS	635	Red	2.5		300
0.75	VPE	N	TS	605	Orange	2.5		300
0.85	VPE	N	TS	585	Yellow	2.5		300
1.0	LPE	N	TS	572	Y/green	6		300
1.0	LPE	—	TS	565	Green	1		
1.0	LPE	ZnO	TS	640	Red	1		

*LPE = liquid phase epitaxy; VPE = vapor phase epitaxy.

†AS = absorbing substrate; TS = transparent substrate.

‡IR = infrared.

§Into plastic (index = 1.5).

monolithic seven-segment displays became the workhorse display technology for handheld calculators from 1972 to 1976.^{9,13} Today this alloy is used in LED arrays for printers.¹⁴

The $x = 0.3$ alloy with a wavelength of 700 nm became important in the mid 1970s as a light source in applications using integrated photodetectors. It has 3 to 5 times the quantum efficiency of the $x = 0.4$ alloy (see Table 2), but has a lower luminous efficacy because of the much-reduced eye sensitivity at 700 nm.

For $x > 0.4$, the GaAsP material system becomes indirect (see earlier under “Light-Generation Processes” and Fig. 3). The quantum efficiency decreases faster than the increase in eye sensitivity.⁴ The only way to achieve a meaningful efficiency is through the use of isoelectronic dopants as described earlier. The choices for isoelectronic dopants that have been successful are N for GaAsP¹¹ and N or ZnO for GaP. Nitrogen doping is used widely for alloys with $x = 0.65$ to $x = 1.0$.⁶⁻¹⁰ The resulting light sources cover the wavelength range from 635 nm to approximately 565 nm (see Table 2). Since these alloys are either GaP or very close in composition to GaP, they are all grown on GaP substrates. The resulting transparent-substrate chip structure increases the luminous efficacy.

In the case of the binary GaP compound, the dominant wavelength depends on N concentration. With low concentrations, practically all N atoms are isolated in single sites. With increasing concentration, N atoms can arrange themselves as pairs or triplets. The resulting electron traps have lower energy states, which shift the emitted light toward longer wavelength. Phonon coupling can also reduce the emission energy. Commercially significant are two compositions: (1) undoped GaP which emits at 565 nm (dominant wavelength) with a reasonably green appearance and a low efficiency of around 1 lm/A, and (2) GaP with a nitrogen concentration in the range of 10^{19} cm⁻³ with a substantially higher luminous efficiency of around 6 lm/A at 572 nm. At this wavelength, the color appearance is yellow-green, often described as chartreuse.

Three nitrogen-doped ternary alloys of GaAsP are commercially important for red, orange, and yellow. The red source with $x = 0.65$ has an efficiency in the range of 2 to 3 lm/A. With increasing bandgap or decreasing wavelength, the drop in quantum efficiency is compensated by an increase in eye sensitivity, resulting in a practically wavelength-independent luminous efficiency for the range of 635 to 585 nm.^{8,15}

ZnO-doped GaP is an interesting material. The quantum efficiency of such chips is relatively high, around 3 percent. However, the linewidth is quite broad. The quantum efficiency peaks at 700 nm, but the luminous efficiency peaks at 640 nm (dominant wavelength). In other words, most of the photons are emitted at wavelengths with low eye sensitivity. Another problem of GaP:ZnO is saturation. The deep ZnO electron trap causes very slow exciton recombination. At high injection currents, all traps are saturated and most of the injected carriers recombine

nonradiatively. At low injection levels ($\leq 1 \text{ A/cm}^2$), the efficacy is relatively high, 3 to 5 lm/A. At a more useful density of 10 to 30 A/cm^2 , the emission saturates, resulting in an efficacy of around 1 lm/A.

The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ System

The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ material system has a direct bandgap for $0 \leq x \leq 0.38$. This system has one very significant advantage over the GaAsP system described earlier, the entire alloy range from $x = 0$ to $x = 1$ can be lattice-matched to GaAs. In other words, every alloy composition can be directly grown on any other alloy composition without the need for transition layers. This feature allows the growth of very abrupt heterojunctions, i.e., abrupt transition in composition and bandgap. These heterojunctions add one important property not available in the GaAsP system: carrier containment (see earlier under "Device Structures"). Carrier containment reduces the movement of injected carriers in a direction perpendicular to the junction. Thus, carrier density can be increased beyond the diffusion-limited levels. This results in increased internal quantum efficiency and higher speed. Another benefit is reduced absorption and improved extraction efficiency (under "Light Extraction").

Of practical significance are two compositions: $x = 0.06$ and $x = 0.38$ (see Table 3). Both compositions exist in single and double heterojunction variations (see under "Device Structures"). The double heterojunctions usually have a 1.5 to 2.0 times advantage in efficiency and speed. In all cases, the efficiency strongly depends on the thickness of the window layer and, to a lesser degree, on the thickness of the transparent layer between active layer and absorbing substrate (see Fig. 12). Chips with a transparent substrate have an additional efficiency improvement of 1.5 to 3.0 times again, depending on layer thickness and contact area. The efficiency variation is best understood by counting exit cones as described in the text in conjunction with Table 1. For $x = 0.06$, the internal quantum efficiency of a double heterojunction approaches 100 percent. For $x = 0.38$, the direct and indirect valleys are practically at the same level and the internal quantum efficiency is reduced to the range of 50 percent, again dependent upon the quality of the manufacturing process.

The best compromise for efficiency and speed is the $x = 0.06$ alloy as a double heterostructure. Depending on layer thickness, substrate, and contact area, these devices have efficiencies of 5 to 20 percent and rise/fall times of 20 to 50 ns. This alloy is becoming the workhorse for all infrared applications demanding power and speed. A structural variation as shown in Fig. 14 is an important light source for fiber-optic communication.

The $x = 0.38$ alloy is optimized for applications in the visible spectrum. The highest product of quantum efficiency and eye response is achieved at $x = 0.38$ and $\lambda = 650 \text{ nm}$. The single heterostructure on an absorbing substrate has an efficacy of around 4 lm/A. The equivalent double heterostructure is in the 6- to 8-lm/A range. On a transparent GaAlAs substrate, the efficacy is typically in the 15- to 20-lm/A range and results of as high as 30 lm/A have been reported in the literature.¹⁶ The major application for these red LEDs is in light-flux-intensive applications, such as message panels and automotive stoplights. A variation optimized for speed is widely used for optical communication using plastic fiber.

TABLE 3 Performance Summary of LED Chips in the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ System

x	λ (nm)	Substrate*	Structure [†]	Efficiency or Efficacy	Speed (ns)
0.06	820	AS, TW	DH	8%	30
0.06	820	TS	DH	15%	30
0.38	650	AS, TW	SH	4 lm/A	
0.38	650	AS, TW	DH	8 lm/A	
0.38	650	TS	DH	16 lm/A	

*AS = absorbing substrate; TW = thick window layer; TS = transport substrate.

[†]DH = double heterostructure; SH = single heterostructure.

The AlInGaP System

The AlInGaP system has most of the advantages of the AlGaAs system with the additional advantage that it has a higher-energy direct energy gap of 2.3 eV that corresponds to green emission at 540 nm. AlInGaP can be lattice-matched to GaAs substrates. Indium occupies about half of the Group III atomic sites. The ratio of aluminum to gallium can be changed without affecting the lattice match, just as it can in the AlGaAs material system, since AlP and GaP have nearly the same lattice spacing. This enables the growth of heterostructures that have the efficiency advantages described in the previous section.

Various AlInGaP device structures have been grown. A simple DH structure with an AlInGaP active layer surrounded by higher bandgap AlInGaP confining layers has been effective for injection lasers, but has not produced efficient surface-emitting LEDs.¹⁷ The main problem has been that AlInGaP is relatively resistive and the top AlInGaP layer is not effective in distributing the current uniformly over the chip. This is not a problem with lasers since the top surface is covered with metal and the light is emitted from the edge of the chip.

The top layer must also be transparent to the light that is generated. Two window layers used are AlGaAs or GaP on top of the AlInGaP heterostructure.^{18,19} AlGaAs has the advantage that it is lattice-matched and introduces a minimum number of defects at the interface, but it has the disadvantage that it is somewhat absorptive to the yellow and green light which is generated for high-aluminum compositions. The highest AlInGaP device efficiencies have been achieved using GaP window layers.^{18,20} GaP has the advantage that it is transparent to shorter wavelengths than AlGaAs and that it is easy to grow thick GaP layers, using either VPE or LPE, on top of the AlInGaP DH that was grown by MOVPE. Both VPE and LPE have substantially higher growth rates than MOCVD. The various growth techniques are discussed later under “Epitaxial Technology.”

AlInGaP devices with 45- μm -thick GaP window layers have achieved external quantum efficiencies exceeding 5 percent in the red and yellow regions of the emission spectrum.²¹ This is more than twice as bright as devices that have thinner AlGaAs window layers.

Green-emitting AlInGaP devices have also been grown which are brighter than the conventional GaP and GaP:N green emitters.^{18,20,21} Substantial further improvement in green is expected since the quantum efficiency is not as high as would be expected based on the energy position of the transition from a direct to an indirect semiconductor.

The performance of AlInGaP LEDs compared to the most important other types of visible emitters is shown in Figs. 16 and 17. GaAsP on a GaAs substrate and GaP: ZnO are not shown since their luminous efficacy are off the chart at the bottom and the lower-right-hand corner, respectively. It is clear from Fig. 17 that the luminous efficacy of AlInGaP is substantially higher than the other technologies in all color regions except for red beyond 640 nm. Since forward voltage is typically about 2 V, the lumen per watt value for a given device is about one-half the lumen per ampere value that is given in Tables 2 and 3. The quantum efficiency of AlInGaP is also better than all of the other technologies except for the highest-performance AlGaAs devices operating at about 650 nm. Because of the eye sensitivity variations (see C.I.E. curve in Fig. 16), the 620 nm (red/orange) AlInGaP devices have a higher luminous efficacy than 650 nm AlGaAs LEDs (see Fig. 17).

Blue LED Technology

Blue emitters have been commercially available for more than a decade, but have only begun to have a significant impact on the market in the last few years. SiC is the leading technology for blue emitters with a quantum efficiency of about 0.02 percent and 0.04 lm/A luminous performance. SiC devices are not much used due to their high price and relatively low performance efficiency.

Other approaches for making blue LEDs are the use of II–VI compounds such as ZnSe or the nitride system GaN, AlGaN, or AlGaInN. It has been difficult to make good *p-n* junctions in these materials. Recently improved *p-n* junctions have been demonstrated in both ZnSe^{22,23} and GaN.²⁴

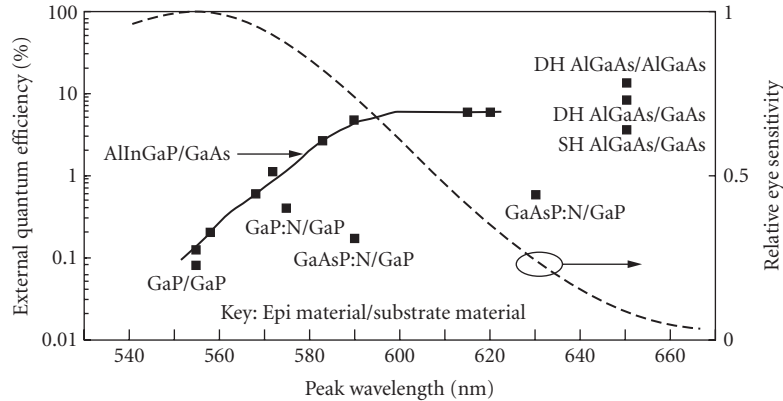


FIGURE 16 External quantum efficiency as a function of peak wavelength for various types of visible LEDs. Below 590 nm the efficiency of AlInGaP LEDs decreases due to the approaching transition from direct to indirect semiconductor. The human eye sensitivity curve is also shown. Since the eye response increases sharply from 660 to 540 nm, it partially makes up for the drop in AlInGaP LED efficiency. The resulting luminous performance is shown in Fig. 17.

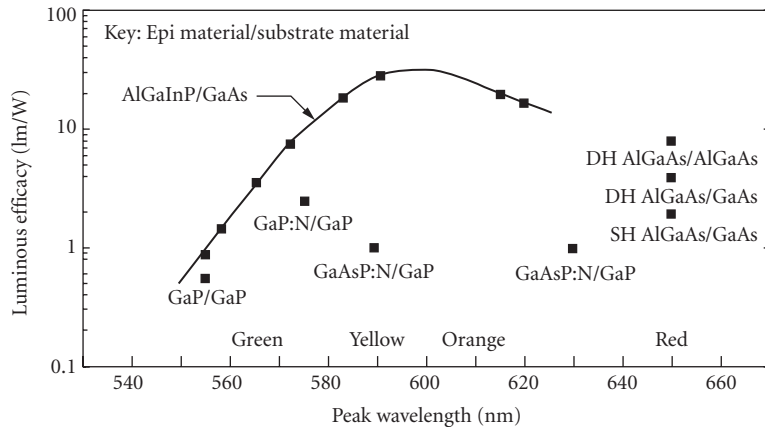


FIGURE 17 Luminous efficacy for AlInGaP LEDs versus wavelength compared to other LED technologies. AlInGaP LEDs are more than an order of magnitude brighter in the orange and yellow regions than other LEDs. AlInGaP LEDs compare favorably to the best AlGaAs red LEDs.

Device performance is still in the 0.1-lm/A range and reliability is unproven. However, this recent progress is very encouraging for blue-emission technology and could lead to a high-performance device in the next few years. Both ZnSe and the nitride system have a major advantage over SiC because they are direct bandgap semiconductors, so a much higher internal quantum efficiency is possible. However, it is difficult to find suitable lattice-matched substrates for these materials.

17.7 SUBSTRATE TECHNOLOGY

Substrate Criteria

There are several requirements for substrates for LEDs. The substrate must be as conductive as possible both thermally and electrically to minimize power loss. In order to minimize defects it should match the epitaxial layers as closely as possible in atomic lattice spacing, and in the coefficient of thermal expansion. The substrate should also have a low defect density itself. Finally, the substrate should, if possible, be transparent to the light generated by the LED structure since this will enhance the external quantum efficiency.

Substrate Choices

The substrates used for nearly all visible LEDs are GaAs and GaP. GaAs or InP is used for infrared devices, depending upon the device structure required. Substrate parameters are summarized in Table 4, along with Si and Ge for comparison. GaAs is used for the AlGaAs and AlInGaP material systems since they can be lattice-matched to it. GaAs is also used for the GaAs_{1-x}P_x system for $x \leq 0.4$ because it is more nearly lattice-matched and, since it is absorbing, it is useful in multiple-junction devices where optical crosstalk must be minimized (see under “Diffused Homojunction” and Fig. 8). GaP is used for compositions of $x > 0.6$ due to its transparency and closer lattice match. However, neither GaP nor GaAs are well matched to GaAsP, and grading layers are required to grow epitaxial layers with decent quality. InP is the choice for long wavelength emitters made using the InGaAs or InGaAsP materials systems, due to the better lattice match.

Substrate Doping

Generally, substrates are *n* type and are doped with Te, S, or Si, although sometimes Se and Sn are also used. In some cases, particularly for some AlGaAs LEDs and laser structures, *p*-type substrates are required and Zn is nearly always the dopant. The doping levels are typically in the 10¹⁸-cm⁻³ range. Basically, the substrates are doped as heavily as possible to maximize conductivity. However, the doping must be below the solubility limit to eliminate precipitates and other defects. In the case of substrates that are transparent to the light which is generated, such as GaP, with a GaAs_{1-x}P_x epitaxial layer, the doping should also remain below the level at which substantial free carrier absorption occurs.

Growth Techniques

Substrates can be grown by either the Bridgeman or Liquid Encapsulated Czochralski (LEC) technique. The LEC technique is the most widely used. Both techniques are described in detail elsewhere and will be only briefly summarized here.²⁵

TABLE 4 Properties of Common Semiconductor Substrates

Substrate	Lattice Parameter	Energy Gap @ 300 K (eV)		Melting Point (°C)
GaAs	5.653	1.428	Direct	1238
GaP	5.451	2.268	Indirect	1467
InP	5.868	1.34	Direct	1062
Si	5.431	1.11	Indirect	1415
Ge	5.646	0.664	Indirect	937

The LEC technique for GaAs consists of a crucible containing a molten GaAs solution, into which a single crystal “seed” is dipped. The temperature is carefully controlled so that the molten GaAs slowly freezes on the seed as the seed is rotated and raised out of the molten solution. By properly controlling the temperature, rotation rate, seed lift rate, etc., the seed can be grown into a single crystal weighing several kilograms and having a diameter of typically 2 to 4 in for GaAs and for 2 to 2.5 in for GaP. At the GaAs and GaP melting points As and P would rapidly evaporate from the growth crucible if they were not contained with a molten boric oxide layer covering the growth solution. This layer is the reason for the name “liquid encapsulated.” The seed is dipped through the boric oxide to grow the crystal. The growth chamber must be pressurized to keep the phosphorus and arsenic from bubbling through the boric oxide. The growth pressure for GaP is 80 atm, and for GaAs is 20 atm or less, depending upon the approach for synthesis and growth. The LEC technique, used for GaAs, GaP, and InP, is similar to the Czochralski technique used for silicon, but the silicon process is much simpler since encapsulation is not required and the growth can be done at atmospheric pressure.

The Bridgeman and the gradient-freeze technique, which is a variation of the Bridgeman technique, can also be used to grow compound semiconductors. In this technique the growth solution and seed are contained in a sealed chamber so liquid encapsulation is not required. Growth is accomplished by having a temperature gradient in the solution, with the lowest temperature at the melting point in the vicinity of the seed. Growth can be accomplished by lowering the temperature of the entire chamber (gradient freeze), or by physically moving the growth chamber relative to the furnace (Bridgeman technique) to sweep the temperature gradient through the molten solution. The growing crystal can be in either a vertical or horizontal position.

GaAs for LEDs is commonly grown using either LEC or horizontal gradient freeze, also called “boat grown,” but sometimes a vertical Bridgeman approach is used. GaP and InP are almost always grown using LEC but sometimes a vertical Bridgeman approach has also been used.

17.8 EPITAXIAL TECHNOLOGY

Growth Techniques Available

Epitaxial layers are grown using one of several techniques, depending on the material system. The most common techniques are liquid phase epitaxy (LPE), which is primary used to grow GaAs, GaP, and AlGaAs and vapor phase epitaxy (VPE), which is used to grow GaAsP. Metalorganic vapor phase epitaxy (MOVPE) is also used to grow AlGaAs, GaInAsP, and AlInGaP. Molecular beam epitaxy (MBE) is used for lasers and high-speed devices but is not used for high-volume commercial LEDs at this time. It has been used to grow blue ZnSe-based lasers and LEDs and could be important in the future. All of these epitaxial techniques have been discussed extensively elsewhere and will be only briefly described here.

LPE

LPE growth consists of a liquid growth solution, generally gallium, which is saturated with the compound to be grown.²⁶ The saturated solution is placed in contact with the substrate at the desired growth temperature, and cooled. As the solution cools, an epitaxial film is grown on the substrate. The technique has the advantage that it is relatively easy to grow high-quality epitaxial layers, and materials containing aluminum (such as AlGaAs) can be readily grown. The disadvantages are that composition control can be difficult. Also, the growth of epitaxial structures involving more than two or three layers, particularly thin layers, can be mechanically complicated since each layer requires a separate growth solution that must be carefully saturated and sequentially brought into contact with the substrate.

One important use of the LPE technique has been for the growth of GaP:ZnO and GaP:N for red and green LEDs, respectively. These devices each consist of two relatively thick layers: an *n*-type layer,

followed by the growth of a *p*-type layer. While other growth techniques can be used to grow GaP LEDs, the best results have been obtained using LPE. As a result of a high-volume, low-cost production technology has evolved, which produces more visible LED chips than any other technique. Another major use of LPE is for the growth of GaAs: Si for infrared emitters. LPE is the only technique with which it has been possible to produce the recombination center that gives rise to the 940-nm emission characteristic of this material. The GaAs: Si structures are generally grown from a single growth solution. At high temperatures the silicon is incorporated on Ga sites and the layers are *n* type. As the solution cools the silicon becomes preferentially incorporated on the As sites and the layer becomes *p* type.

AlGaAs devices for both visible (red) and IR devices are also generally grown by LPE. AlGaAs devices can also be grown by MOCVD, but LPE has the advantage that thick layers can be more easily grown. This is important for high extraction efficiency (see under "Device Structures" and Fig. 12). In the case of visible devices at 650 nm, the internal quantum efficiency is also higher using LPE than MOCVD. This is not understood, but the result is that virtually all of the visible AlGaAs LEDs are produced using LPE.

VPE

VPE is the other major commercial epitaxial technology for LEDs.^{9,27} VPE consists of a quartz chamber containing the substrate wafers at the appropriate growth temperature. The reactants are transported to the substrates in gaseous form. The technique is mainly used for the growth of GaAsP which, along with GaP, dominates the high-volume visible LED market. In this case HCl is passed over Ga metal to form gallium chlorides, and AsH₃ and PH₃ are used to provide the As and P compounds. Appropriate dopant gases are added to achieve the *n*- and *p*-type doping. NH₃ is used to achieve nitrogen doping for the growth of GaAsP:N. The VPE technique has the advantages that it is relatively easy to scale up the growth chamber size so large quantities of material can be grown and layer composition and thickness can be easily controlled by adjusting the flow conditions. A limitation of VPE is that it has not been possible to grow high-quality compounds containing aluminum because the aluminum-bearing reactants attack the quartz chamber resulting in contaminated films. Thus, AlGaAs and AlInGaP, the emerging high-brightness technologies, cannot be grown using this technique.

MOCVD

MOCVD growth, like conventional VPE, uses gases to transport the reactants to the substrates in a growth chamber.²⁸ However, in this case metallorganic compounds such as trimethylgallium (TMG) are used for one or more of the reactants. A major difference between VPE and MOCVD is that in the case of MOCVD the decomposition of the source gas (e.g., TMG) occurs as a reaction at or near the substrate surface, and the substrate is in the hottest area of the reactor such that the decomposition occurs on the substrate instead of the walls of the growth chamber. The walls of the growth chamber remain relatively cool. This is the key factor that makes MOCVD suitable for the growth of aluminum-bearing compounds which, unlike the VPE situation, do not react significantly with the cooler reactor walls. Thus AlGaAs and AlInGaP can be readily grown with MOCVD, and this technology is widely used for infrared AlGaAs LEDs and lasers, and for the emerging visible AlInGaP laser and LED technology.

MOCVD is also used for the growth of GaN and AlGaN that are candidates for blue emission, and for the growth of II–VI compounds, such as ZnSe, that is also a potential blue emitter. However, at this time the key limitation in obtaining blue ZnSe emitters is the growth of low-resistivity *p*-type ZnSe. For reasons that are not yet understood, low-resistivity *p*-type ZnSe has been grown using MBE only.

MBE

MBE is a high vacuum growth technique in which the reactants are essentially evaporated onto the substrates under very controlled conditions.²⁹ MBE, like MOCVD, can be used to deposit compounds

containing aluminum. The growth rates using MBE are generally slower than the other epitaxial techniques, so MBE is most suitable for structures requiring thin layers and precise control of layer thickness. MBE equipment is somewhat more expensive than the equipment used for the other types of epitaxial growth, so it has not been suitable for the high-volume, low-cost production that is required for most types of LEDs. MBE has generally been utilized for lasers and high-speed devices where control of complicated epitaxial structures is critical and where relatively low volumes of devices are required. One advantage that MBE has over the other growth technologies is that the reactants utilized are generally less hazardous. Consequently, MBE equipment is often cheaper and easier to install since there are less safety issues and safety-code restrictions to deal with.

17.9 WAFER PROCESSING

Wafer Processing Overview

Wafer processing of compound semiconductors for LED applications has many of the same general steps used to process silicon integrated-circuit wafers, namely passivation, diffusion, metallization, testing, and die fabrication. The LED device structures are much simpler, so fewer steps are required; but, due to the materials involved, the individual steps are generally different and sometimes more complicated.

Compound semiconductor processing has been described in detail elsewhere, so only a brief summary is discussed here.³⁰ Some types of LED structures require all of the processing steps listed here, but in many cases fewer process steps are required. An example is a GaP or AlGaAs device with a grown *p-n* junction. For these devices no passivation or diffusion is required.

Passivation

Some types of LED structures, particularly multijunction structures, require a passivation layer prior to diffusion, as shown in Fig. 7. This layer must be deposited relatively free of pinholes, be patternable with standard photolithographic techniques, and must block the diffusing element, generally zinc. In the case of silicon, a native oxide is grown which is suitable for most diffusions. Unfortunately, the compound semiconductors do not form a coherent native oxide as readily as silicon. Silicon nitride (Si_3N_4) is the most widely used passivation layer for LEDs. Si_3N_4 is grown by reacting silane (SiH_4) and ammonia at high temperature in a furnace. Si_3N_4 blocks zinc very effectively, and is easily grown, patterned, and removed. Sometimes an SiO_2 layer is used in conjunction with Si_3N_4 for applications such as protecting the surface of the compound semiconductor during high-temperature processing. Silicon oxynitride can also be used instead of or in addition to pure Si_3N_4 . Silicon oxynitride is somewhat more complicated to deposit and control, but can have superior properties, such as a better match of coefficient of expansion, resulting in lower stress at the interface.

Diffusion

Generally, only *p*-type impurities, usually Zn, are diffused in compound semiconductors. *N*-type impurities have prohibitively small diffusion coefficients for most applications. Zn is commonly used because it diffuses rapidly in most materials and because it is nontoxic in contrast to Be, which also diffuses rapidly. Mg is another reasonable *p*-type dopant, but it diffuses more slowly than zinc. Diffusions are generally done in evacuated and sealed ampoules using metallic zinc as source material. A column V element such as As is also generally added to the ampoule to provide an overpressure that helps to prevent decomposition of the semiconductor surface during diffusion. Diffusion conditions typically range from 600 to 900°C for times ranging from minutes to days, depending upon the material and device involved. Junction depths can range from a fraction of a μm to a more than 10 μm .

Open-tube diffusions have also been employed but have generally been harder to control than the sealed ampoule approach, often because of surface decomposition problems. Open-tube diffusions have the advantage that one does not have to deal with the expense and hazard of sealing, breaking, and replacing quartz ampoules.

A third type of diffusion that has been used is a “semisealed” ampoule approach in which the ampoule can be opened and reused. The diffusion is carried out at atmospheric pressure and the pressure is controlled by having a one-way pressure relief valve on the ampoule.

Contacting

The contacts must make good ohmic contact to both the *p*- and *n*-type semiconductor, and the top surface of the top contact must be well suited for high-speed wire bonding. Generally, multilayer contacts are required to meet these conditions. Evaporation, sputtering, and *e*-beam deposition are all employed in LED fabrication. The *p*-type contact generally uses an alloy of either Zn or Be to make the ohmic contact. An Au-Zn alloy is the most common due to the toxicity of Be. The Au-Zn can be covered with a layer of Al or Au to enable high-yield wire bonding. A refractory metal barrier layer may be included between the Au-Zn and top Al or Au layer to prevent intermixing of the two layers and the out-diffusion of Ga, both of which can have a deleterious effect on the bondability of the Al or Au top layer. The *n*-type contact can be similar to the *p*-type contact except that an element which acts as an *n*-type dopant, commonly Ge, is used instead of Zn. An Au-Ge alloy is probably most frequently used to form the *n*-type ohmic contact since it has a suitably low melting point. If the *n*-type contact is the top, or bonded, contact it will be covered by one or more metallic layers to enhance bondability.

Testing

The key parameters that need to be tested are light output, optical rise and fall times, emission wavelength, forward voltage, and leakage current. The equipment used is similar to that used to test other semiconductor devices except that a detector must be added to measure light output. Rise and fall times and wavelength are generally measured on only a sample basis and not for each device on a wafer. In order to test the individual LEDs, the devices must be isolated on the wafer. This occurs automatically for LEDs that are masked and diffused, but if the LEDs are sheet diffused or have a grown junction, the top layer must be processed to isolate individual junctions. This can be accomplished by etching or sawing with a dicing saw. Generally, sawing is used, followed by an etch to remove saw damage, because the layers are so thick that etching deep (>10 μm) grooves is required. It is advisable to avoid deep groove etching because undercutting and lateral etching often occur and the process becomes hard to control. In many cases LED junctions are not 100 percent tested. This is particularly true of GaP and AlGaAs red-emitting devices in which the top layer may be 30 μm thick. Wafers of this type can be sampled by “coring” through the top layer of the wafer in one or more places with an ultrasonic tool in order to verify that the wafer is generally satisfactory. Later, when chips are fully processed, chips can be selected from several regions of the wafer and fully tested to determine if the wafer should be used or rejected.

Die Fab

Die fab is the process of separating the wafer into individual dice so they are ready for assembly. Generally, the wafer is first mounted on a piece of expandable tape. Next the wafer is either scribed or sawed to form individual dice. Mechanical diamond scribing or laser scribing were the preferred technologies in the past. Mechanical scribing has zero kerf loss, but the chips tend to have jagged edges and visual inspection is required. Laser scribing provides uniform chips but the molten waste material from between the chips damages neighboring chips, and in the case of full function chips

the edges of the junction can be damaged by the laser. As a result of the limitations of scribing, sawing (using a dicing saw with a thin diamond impregnated blade) has become the technology of choice for most LEDs.

The kerf loss for sawing has been reduced to about 40 μm and the chip uniformity is excellent such that a minimum of inspection and testing is required. For most materials a “cleanup” etch is required after sawing to remove work damage at the edges of the chips, which can both affect the electrical performance and absorb light. The wafer remains on the expandable tape during the sawing process. After sawing, the tape is expanded so that the chips are separated. The tape is clamped in a ring that keeps it expanded and the chips aligned. In this form the chips are easily individually picked off the tape by the die-attach machine that places the chips in the LED package.

17.10 LED QUALITY AND RELIABILITY

LEDs offer many advantages over other types of light sources. They have long operating life, they operate over a wide temperature range, and they are unaffected by many adverse environmental conditions. LED devices also are mechanically robust, making them suitable for applications where there is high vibration, shock, or acceleration. Excellent quality and reliability are obtainable when an LED product is properly designed, fabricated, packaged, tested, and operated.

Product quality is defined as “fitness for use” in a customer’s application. Quality is measured in units of the average number of defective parts per million shipped (i.e., ppm), and is inferred from product sampling and testing. LED product quality is assured by (1) robust chip and product design, (2) high-quality piece parts, (3) well-controlled fabrication processes, (4) use of statistical process control during manufacturing, (5) careful product testing, and (6) proper handling and storage. Most III–V LEDs are comparable in quality to the best silicon devices manufactured today. Well-designed LED products have total defect levels well below 100 ppm.

Reliability measures the probability that a product will perform its intended function under defined use conditions over the useful life of the product. Probability of survival is characterized by a failure rate, which is calculated by dividing the number of failures by the total number of operating hours (number of products tested per x hours operated). Common measures for reliability are percent failures per 1000 hours (percent/khr) and number failures per 10^9 hours (FITS). LED failure rates typically are better than 0.01 percent per khr at 50°C.

The reliability of an LED product is dependent on the reliability of the LED semiconductor chip and on the robustness of the package into which the chip is placed. Interactions between the chip and package can affect product reliability as well. Aspects of LED packaging and LED chip reliability are discussed in the following paragraphs.

LED Package Reliability

The package into which the LED chip is assembled should provide mechanical stability, electrical connection, and environmental protection. To evaluate package integrity, stress tests such as temperature cycling, thermal and mechanical shock, moisture resistance, and vibration are used to establish the worst-case conditions under which a product can survive. Generally, product data sheets contain relevant information about safe conditions for product application and operation.

Plastic materials are commonly used to package LEDs (see under “LED-Based Products”). Thermal fatigue is a limiting factor in plastic-packaged LEDs. Take the case of the plastic LED lamp shown later in Fig. 21. Because of the different materials used (epoxy plastic, copper lead frame, gold wire, III–V LED, etc.) and the different coefficients of expansion of these materials, temperature changes cause internal stresses. If the package is not well designed and properly assembled, thermal changes can cause cracking, chip-attach failure, or failure of the wire bond (open circuits). Careful design can reduce these problems to negligible levels over wide temperature ranges. Today’s high-quality plastic lamps are capable of being cycled from -55 to $+100^\circ\text{C}$ for 100 cycles without failure.

Long-term exposure to water vapor can lead to moisture penetration through the plastic, subjecting the chips to humidity. High humidity can cause chip corrosion, plastic delamination, or surface leakage problems. Plastic-packaged LEDs are typically not harmed when used under normal use conditions. Accelerated moisture resistance testing can be used to test the limits of LED packages. Plastic materials have been improved to the point that LED products can withstand 1000 hours of environmental testing at elevated temperatures and high humidity (i.e., 85°C and 85 percent RH).

The thermal stability of plastic packaging materials is another important parameter. Over normal service conditions, the expansion coefficient of plastic is relatively constant. Above the so-called glass transition temperature T_g , the coefficient increases rapidly. Reliable operation of plastic-packaged LEDs generally requires operation at ambient temperatures below T_g . Failures associated with improper soldering operations, such as too high a soldering temperature or for too long a time can cause the package to fail. Excessive storage temperatures also should be avoided. When an LED product is operated, internal ohmic heating occurs; hence, the safe operating maximum temperature is generally somewhat lower than the safe storage temperature.

LED Chip Reliability

The reliability characteristics of the LED chip determine the safe limits of operation of the product. When operated at a given temperature and drive current there is some probability that the LED will fail. In general, LED failure rates can be separated into three time periods: (1) infantile failure, (2) useful life, and (3) the wearout period (see Fig. 18). During the infant mortality period, failures occur due to weak or substandard units. Typically, the failure rate decreases during the infantile period until no weak units remain. During the “useful life” period, the failure rate is relatively low and constant. The number of failures that do occur are random in nature and cannot be eliminated by more testing. The useful life of an LED is a function of the operating temperature and drive conditions. Under normal use conditions, LEDs have useful lives exceeding 100,000 hours. The wearout period is characterized by a rapidly rising failure rate. Generally, wearout for LEDs is not a concern, as the useful life far exceeds the useful life of the product that the LEDs are designed into.

The principle failure mode for LED chips is light-output degradation. In the case of a visible lamp or display, failure is typically defined as a 50 percent decrease in light output from its initial

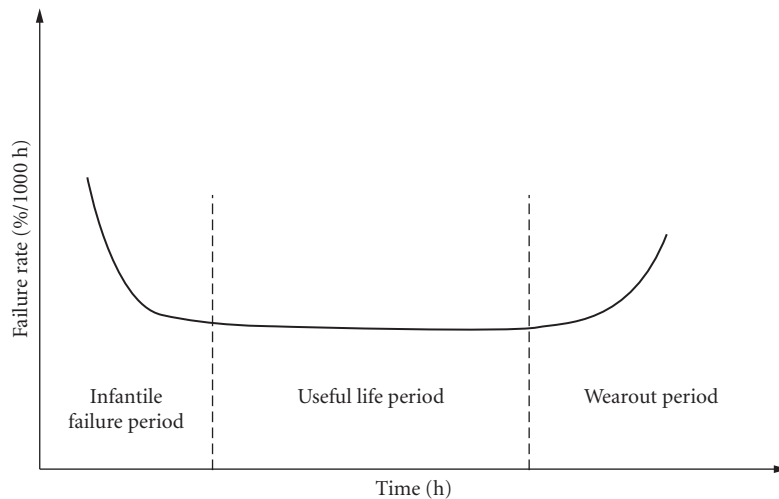


FIGURE 18 Plot of LED failure rate versus time, showing the infantile failure period (decreasing rate), the useful life period (constant, low rate), and the wear-out period (rapidly rising rate).

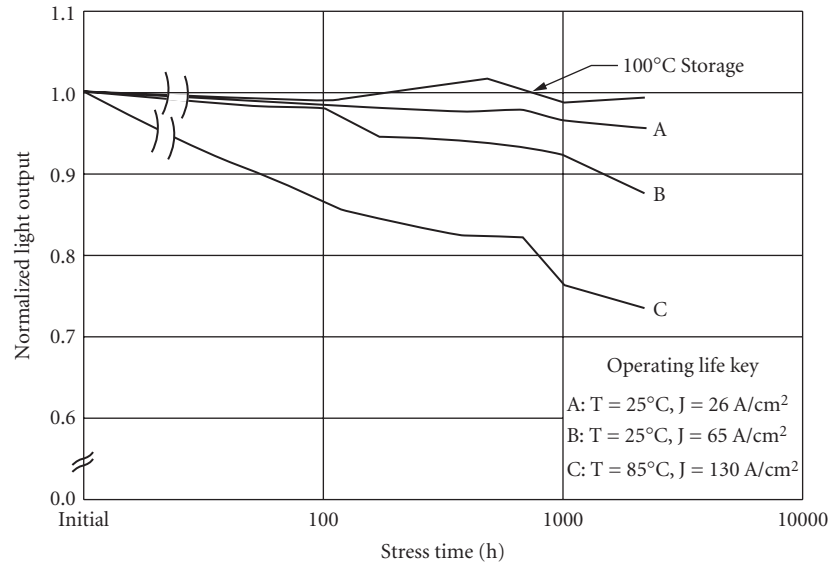


FIGURE 19 Curves of light output versus time for GaAsP indicator lamps stressed under various conditions. Light output is normalized to the initial value. Each curve shows the average degradation of 20 lamps

value, since that is the level where the human eye begins to observe a noticeable change. For infrared emitters or visible LEDs where flux is sensed by a semiconductor detector, failure is commonly defined as a 20 to 30 percent decrease in flux output.

Figure 19 shows degradation curves for a direct bandgap GaAsP LED packaged in a 5-mm plastic lamp.³¹ Current must flow for degradation to occur, as negligible change is observed after 1 khr of 100°C storage. Degradation is a function of the temperature at the p - n junction and the junction current density. As shown in Fig. 19, a larger decline in light output is observed as junction current density and/or temperature at the junction increases. The dependence of degradation on current density J is superlinear, varying as J^x with $1.5 < x < 2.5$. Hence, accelerated-aging tests typically use high currents and temperatures to shorten the time needed to observe LED degradation. The maximum stress level shown in Fig. 19 is 200 percent of the maximum allowable drive current specified in the data sheet.

Light-output degradation in GaAsP LEDs is due to an increase in the nonradiative space-charge recombination current. Total current flowing through the LED is made up of the sum of diffusion current and space-charge recombination current, as shown in the following equations:

$$I_T(V, t) = A(t)e^{qV/kT} + B(t)e^{qV/2kT} \quad (5)$$

where q is electron charge, k is Boltzmann's constant, and T is temperature.³² The first term is diffusion current that produces light output, while the second term is space-charge recombination that is nonradiative. At fixed I_T , if $B(t)$ increases, then the diffusion term must decrease and, hence, the light output decreases. The reason for the increase in space-charge recombination in GaAsP LEDs is not fully understood.

The degradation characteristics of GaAlAs LEDs differ from those of GaAsP LEDs. Typical curves of normalized light output versus time for GaAlAs LEDs are shown in Fig. 20. "Good" units show negligible decrease in light output when operated under normal service conditions. Gradual degradation may occur, but it is relatively uncommon. The predominate failure mode in GaAlAs LEDs is catastrophic degradation. The light emission decreases very rapidly over a period of less than 100 hours

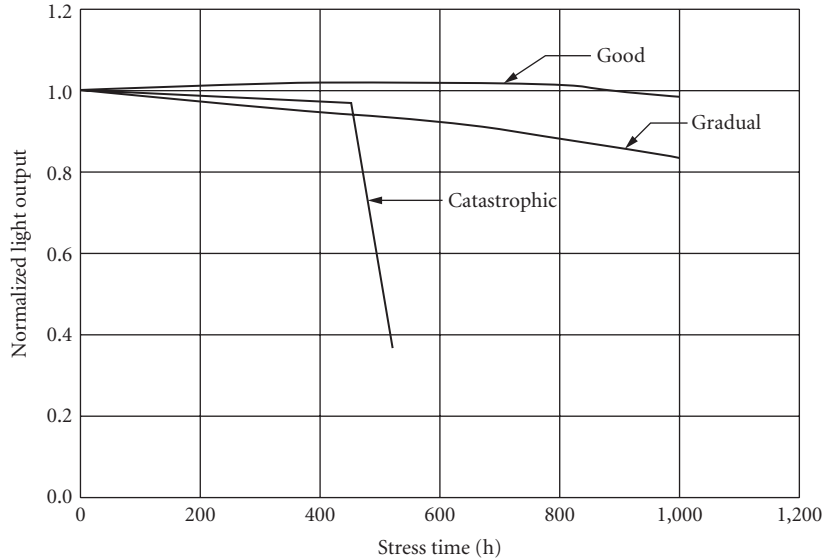


FIGURE 20 Light output degradation of AlGaAs LEDs. Three modes are shown: “good” devices with negligible light-output decrease over time, devices which degrade gradually over time, and “catastrophic” degradation devices where the flux rapidly decreases over a short time period

and simultaneously nonradiative regions (“dark-spot” or “dark-line defects”) are observed to form. The catastrophic degradation mechanism is described in detail in Refs. 33 and 34. In brief, the dark regions are caused by nonradiative recombination at dislocation networks that grow rapidly from a crystal dislocation located in the light-emitting region of the LED. Network formation depends on carrier recombination, both nonradiative, which creates mobile point defects, and radiative recombination, which enhances the movement of the point defects to the growing network.

Formation of dark-line defects is enhanced by mechanical stress either present in the LED chip or occurring during assembly. Properly designed products reduce such stress by minimizing bending caused by the different coefficients of expansion in the LED, and by stress-free die attach, wire bond, and encapsulation of the LED during assembly.

Failures due to dark-spot or dark-line defects can be effectively screened out by operating the LED at high current and temperature. Units with defects typically fail within the first few hundred hours. GaAlAs LEDs with dark-spot or dark-line defects are screened out by means of visual inspection and/or by eliminating units with large decreases in light output. GaAlAs LEDs used for fiber-optic applications have small emitting areas (see under “Fiber Optics” discussion and Fig. 14). Due to the high current densities present in such devices, high temperature and current burn-in is used extensively for these types of LEDs.

Another degradation mode in LEDs is the change in the reverse breakdown characteristics over time. The reverse characteristics become soft and the breakdown voltage may decrease to a very low value. Several mechanisms have been observed in LEDs. One cause is localized avalanche breakdown due to microplasma formation at points where electric fields are high. Microplasmas have been observed in GaAs and GaAsP LEDs.

Damage or contamination of the edges of an LED chip can cause increased surface leakage and reduced reverse breakdown. Incomplete removal of damage during die separation operations and damage induced during handling and assembly are known to cause reverse breakdown changes. Die-attach materials also can unintentionally contaminate the edges of LEDs. Copper, frequently found in LED packaging materials, can diffuse into the exposed surfaces of the LED, causing excess leakage and, in some cases, light-output degradation.³⁵ Chips whose p - n junction extends to the edges (i.e., Fig. 6) are very susceptible to damage and/or contamination.

17.11 LED-BASED PRODUCTS

Indicator Lamps

The simplest LED product is an indicator lamp or its infrared equivalent. The most frequently used lamp is shown in Fig. 21. A LED chip with a typical dimension of $250 \times 250 \mu\text{m}$ is attached with conductive silver-loaded epoxy into a reflective cavity coined into the end of a silver-plated copper or steel lead frame. The top of the LED chip is connected with a thin $25\text{-}\mu\text{m}$ gold wire to the second terminal of the lead frame. The lead frame subassembly is then embedded in epoxy. The epoxy serves several functions: (1) it holds the assembly together and protects the delicate chip and bond wire; (2) it increases the light extraction from the chip (see under "Light Extraction," discussed earlier); and (3) it determines the spatial light distribution.

There are a large number of variations of the lamp shown in Fig. 21. Besides the obvious variation of source wavelength, there are variations of size, shape, radiation pattern, etc. The cross section of the plastic body ranges from 2 to 10 mm. The radiation pattern is affected by three factors: the shape of the dome, the relative position of the chip/reflector combination, and by the presence of a diffusant in the epoxy. Figure 22*a* shows the radiation pattern of a lamp using clear plastic as an encapsulant. The rays emanating from chip and reflector are collimated into a narrow beam. For many indicator lamps a broader viewing angle is desired such as that shown in the radiation pattern of Fig. 22*b*. This effect is achieved by adding a diffusant, such as glass powder, to the clear plastic. Another variation is shape. Common shapes are round, square, rectangular, or triangular.

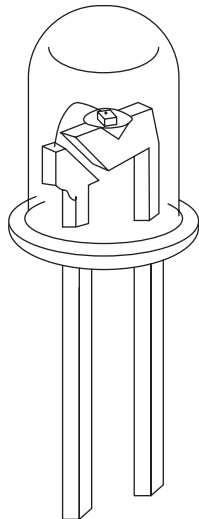


FIGURE 21 Plastic indicator lamp. The LED chip is placed in a reflector coined into the end of one electrode lead. The top of the chip is connected with a gold wire to the second electrode. The electrodes are encapsulated in plastic to form a mechanically robust package.

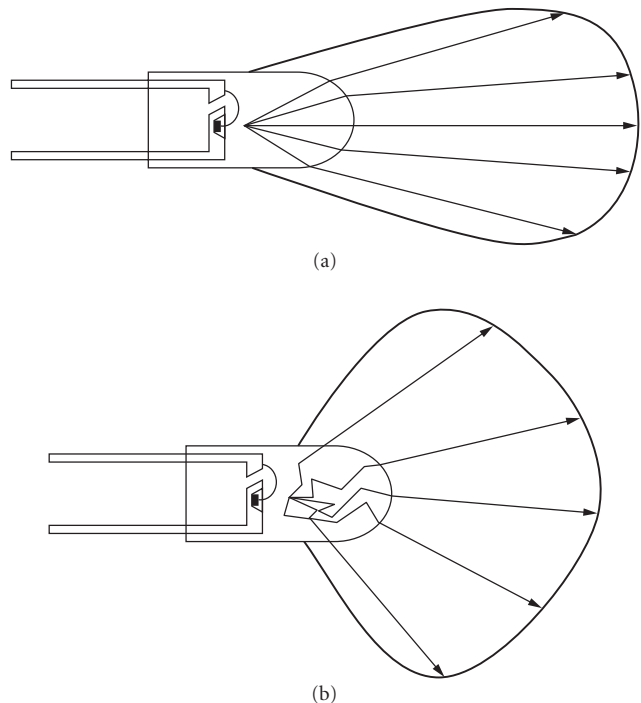


FIGURE 22 Radiation pattern of two types of LED indicator lamps: (a) lamp with clear plastic package with a narrow beam and (b) lamp with a diffusing plastic package (glass powder added) with a broader radiation pattern.

Another variation is achieved by placing two different chips into the reflector cup. For instance, a lamp with green and red chips connected to the second post in an antiparallel fashion can operate as a red indicator if the reflector post is biased positive, and as a green indicator if it is biased negative. Rapidly switching between the two polarities one can achieve any color between the two basic colors, i.e., yellow or orange, depending on current and duty cycle.

A number of chips can be combined in a single package to illuminate a rectangular area. These so-called annunciator assemblies range in size from 1 to several cm. They typically use 4 chips per cm^2 . By placing an aperture-limiting symbol or telltale in front of the lit area, these structures are cost-effective means to display a fixed message, such as warning lights in an automotive dashboard.

Numeric Displays

Numeric displays are usually made up of a nearly rectangular arrangement of seven elongated segments in a figure-eight pattern. Selectively switching these segments generates all ten digits from 0 to 9. Often decimal point, colon, comma, and other symbols are added.

There are two main types of LED numeric displays: (1) monolithic displays and (2) stretched segment displays. All monolithic displays are based on GaAs-GaAsP technology. Seven elongated *p*-doped regions and a decimal point are diffused into an *n*-type epitaxial layer of a single or monolithic chip (see Fig. 8). Electrically, this is a structure with eight anodes and one common cathode. This monolithic approach is relatively expensive. For arms-length viewing, a character height of 3 to 5 mm is required. Adding space for bonding pads, decimal point, and edge separation, such a display consumes around 10 mm^2 of expensive semiconductor material per digit. One way to reduce material and power consumption is optical magnification. Viewing-angle limitation and distortion limit the magnification M to $M \leq 2.0$. Power consumption is reduced by M^2 —an important feature for battery-powered applications.

For digits $>5 \text{ mm}$, a stretched segment display is most cost effective. The design of Fig. 23 utilizes a $250 \times 250 \text{ }\mu\text{m}$ chip to generate a segment with dimensions of up to $8 \times 2 \text{ mm}$ for a 20-mm digit height. This corresponds to a real magnification of $M^2 = 256$, or an equivalent linear magnification of $M = 16$. This magnification is achieved without a reduction in viewing angle by using scattering optics. An LED chip is placed at the bottom of a cavity having the desired rectangular exit shape and

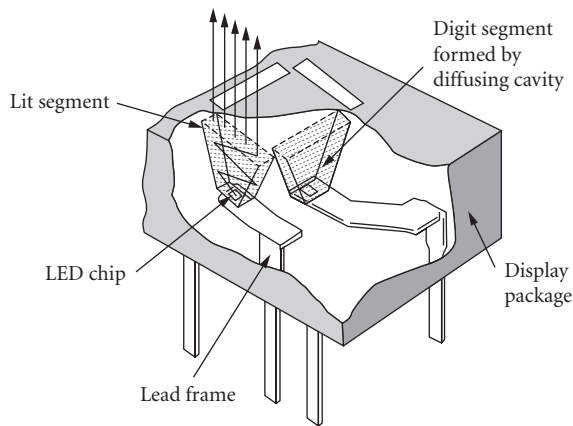


FIGURE 23 Cutaway of a seven-segment numeric LED display, showing how light from a small LED chip is stretched to a large character segment using a diffusing cavity. Characters 0 to 9 are created by turning on appropriate combinations of segments.

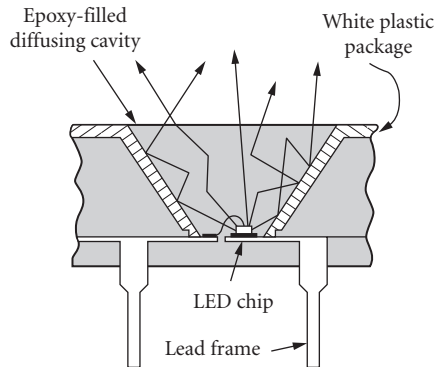


FIGURE 24 Cross section through one segment of the seven-segment numeric display shown in Fig. 23. A LED chip is placed at the bottom of a highly reflective, white plastic cavity which is filled with clear plastic containing a diffusant. Light is scattered within the cavity to produce uniform emission at the segment surface.

the cavity is filled with a diffusing plastic material (see Fig. 24). The side and bottom surfaces of the cavity are made as reflective as possible. Highly reflective white surfaces are typically used. A good white plastic surface measured in air may have a reflectivity of 94 percent compared with 98 percent for Ag and 91 percent for Al. Ag and Al achieve this reflectivity only if evaporated on a specularly smooth surface. Measured in plastic, the reflectivity of the white surface increases from 94 to 98 percent, the Ag surface remains at 98 percent, while the Al surface decreases to 86 percent. Both metallic surfaces have substantially lower reflectivities if they are evaporated onto a nonspecular surface, or if they are deposited by plating. Practically all numeric LED displays above 5-mm character heights are made using white cavity walls and diffusing epoxy within the cavity.

This case of magnification by scattering does not result in a power saving as in the case of magnification of monolithic displays. Since there is practically no reduction in emission angle, the law of energy conservation requires an increased light flux from the chip that equals the area magnification plus reflection losses in the cavity.

Alphanumeric Displays

There are two ways LEDs are used to display alphanumeric information: either by using more than seven elongated segments, i.e., 14, 16, or 24; or by using an array of LED chips in a 5×7 dot matrix format. The multiple segment products are similar in design to the monolithic or stretched segment numeric displays described above.

In the case of small monolithic characters, the number of input terminals quickly exceeds conventional pin spacing. These products are usually clusters of 4 to 16 characters combined with a decoder/driver integrated circuit within the same package. To reduce cost and power, some modest optical magnification is usually used. The segmented displays are usually larger, i.e., 12 to 25 mm and limited to 14 segments per character. At this size, there is no pin density constraint and the decoder is usually placed outside the package.

The most frequently used alphanumeric LED display is based on a 5×7 matrix per character. For small characters in the range of 5 to 8 mm, the LED chips are directly viewed and pin density limitations require an on-board decoder/driver IC. Products are offered as end-stackable clusters of 4, 8, and 16 characters.

For larger displays, the LED chip is magnified by the same optical scattering technique described earlier for numerics. Exit apertures per pixel have a diameter of 2 to 5 mm. Products are offered as 5×7 single characters or end-stackable 20×20 tiles for large message- or graphics-display panels. At this size, pin density is not a limitation.

Optocouplers

An optocoupler is a device where signal input and signal output have no galvanic connection. It is mainly used in applications as the interface between the line voltage side of a system and the low-voltage circuit functions, or in systems where the separate ground connection of interconnected subsystems causes magnetic coupling in the galvanic loop between signal and ground connections. By interrupting the galvanic loop with an optical signal path, many sources of signal interference are eliminated.

The oldest optocouplers consist of an IR LED and a photodetector facing each other in an insulating tube. The second generation utilized the so-called dual-in-line package widely used by logic ICs. In this package an IR emitter and a phototransistor are mounted face to face on two separate lead frames. The center of the package between emitter and detector is filled with a clear insulating material. The subassembly is then molded in opaque plastic to shield external light and to mechanically stabilize the assembly (Fig. 25). The second generation optocouplers have limited speed performance for two reasons: (1) the slow response time of the GaAs:Si LED and (2) the slow response of the photo-transistor detector because of the high collector-base capacitance.

The third generation of optocouplers overcomes the speed limitation. It uses an integrated photodetector and a decoupled gain element. Integration limits the thickness of the effective detection region in silicon to 5 to 7 μm . This thickness range forces a shift of the source to wavelengths shorter than the 940-nm sources used in the second-generation couplers. Third-generation couplers use GaAsP (700 nm) or AlGaAs emitters (880 nm).

Within the last decade, the optocoupler product family has seen further proliferation by adding features on the input or output side of the coupler. One proliferation resulted in couplers behaving

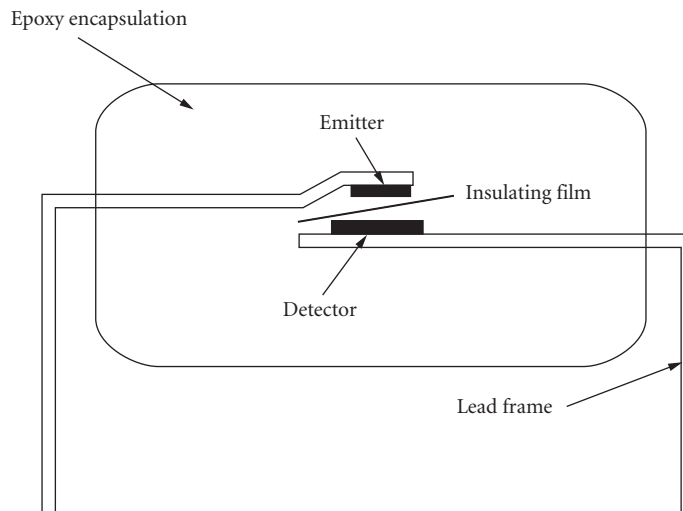


FIGURE 25 Optocoupler consisting of face-to-face emitter and detector chips. An insulating film is placed between the chips to increase the ability of the optocoupler to withstand high voltages between input and output electrodes

like logic gates. Another variation used MOS FET devices on the output side, eliminating the offset voltages of bipolar devices. These couplers are comparable to the performance of conventional relays and are classified as solid-state relays. Other types use a CMOS input driver and CMOS output circuitry to achieve data transfer rates of 50 Mb/s and CMOS interface compatibility.

Fiber Optics

LEDs are the primary light source used in fiber-optic links for speeds up to 200 Mb/s and distances up to 2 km. For higher speed and longer distances, diode lasers are the preferred source.

For fiber-optic applications, LEDs have to meet a number of requirements that go well beyond the requirements for lamps and displays. The major issues are minimum and maximum flux coupled into the fiber, optical rise and fall times, source diameter, and wavelength. Analysis of flux budget, speed, fiber dispersion, wavelength, and maximum distance is quite complex and goes far beyond the scope of this work. A simplified discussion for the popular standard, Fiber Distributed Data Interface (FDDI) will highlight the issues. For a detailed discussion, the reader is referred to Ref. 36.

Flux Budget The minimum flux that has to be coupled into the fiber is determined by receiver sensitivity (-31 dBm), fiber attenuation over the 2-km maximum distance (3 dB), connector and coupler losses (5 dB), and miscellaneous penalties for detector response variations, bandwidth limitations, jitter, etc. (3 dB). With this flux budget of 11 dB, the minimum coupled power has to be -20 dBm. Another flux constraint arises from the fact that the receiver can only handle a maximum level of power before saturation (-14 dBm). These two specifications bracket the power level coupled into the fiber at a minimum of -20 dBm (maximum fiber and connector losses) and at a maximum of -14 dBm (no fiber or connector losses for the case of very short fiber).

Speed The transmitter speed or baud rate directly translates into a maximum rise and fall time of the LED. The 125 Mdb FDDI specifications call for a maximum rise and fall time of 3.5 ns. As a rule of thumb, the sum of rise and fall time should be a little shorter than the inverse baud period (8 ns for 125 Mbd).

Source Diameter and Fiber Alignment Efficient coupling of the LED to the fiber requires a source diameter that is equal to or preferably smaller than the diameter of the fiber core. For sources smaller than the core diameter, a lens between source and fiber can magnify the source to a diameter equal to or larger than the core diameter. The magnification has two benefits: (1) It increases the coupling efficiency between source and fiber. The improvement is limited to the ratio of source area to core cross-sectional area. (2) It increases the apparent spot size to a diameter larger than the fiber core. This effect relaxes the alignment tolerance between source and core and results in substantially reduced assembly and connector costs.

Wavelength The LEDs used in fiber-optic applications operate at three narrowly defined wavelength bands 650, 820 to 870, and 1300 nm, as determined by optical fiber transmission characteristics.

650 nm This band is defined by an absorption window in acrylic plastic fiber. It is a very narrow window between two C-H resonances of the polymer material. The bottom of the window has an absorption of approximately 0.17 dB/m. However, the effective absorption is in the 0.3 to 0.4 dB/m range because the LED linewidth is comparable to the width of the absorption window and the LED wavelength changes with temperature. The 650-nm LEDs use either $\text{GaAs}_{1-x}\text{P}_x$ with $x = 0.4$ or GaAlAs_{1-x} with $x = 0.38$. Quantum efficiencies are at 0.2 and 1.5 percent, respectively. Maximum link length is in the range of 20 to 100 m, depending on source efficiency, detector sensitivity, speed, and temperature range.

820 to 870 nm This window was chosen for several reasons. GaAlAs emitters (see Fig. 14) and Si detectors are readily available at this wavelength. Early fibers had an absorption peak from water

contamination at approximately 870 nm. As fiber technology improved, the absorption peak was eliminated and the wavelength of choice moved from 820- to the 850 to 870-nm range. Fiber attenuation at 850 nm is typically 3 dB/km. Maximum link length in the 500- to 2000-m range, depending on data rate. In the 850 to 870-nm window the maximum link length is limited by chromatic dispersion. GaAlAs emitters have a half-power linewidth of approximately 35 nm. The velocity of light in the fiber is determined by the index of refraction of the fiber core. The index is wavelength-dependent resulting in dispersion of the light pulse. This dispersion grows with distance. The compounded effect of LED linewidth and fiber dispersion is a constant distance-speed product. For a typical multimode fiber and GaAlAs LED combination, this product is in the range of 100 Mbd-km.³⁶

1300 nm At this wavelength, the index of refraction as a function of wavelength reaches a minimum. At this minimum, the velocity of light is practically independent of wavelength, and chromatic dispersion is nearly eliminated. The distance-speed limitation is caused by modal dispersion. Modal dispersion can be envisioned as a different path length for rays of different entrance angles into the fiber. A ray going down the middle of the core will have a shorter path than a ray entering the fiber at the maximum acceptance angle undergoing many bends as it travels down the fiber. The resulting modal dispersion limits multimode fibers to distance bandwidth products of approximately 500 Mbd-km. LED sources used at this wavelength are GaInAsP emitters, as shown in Fig. 15. At this wavelength, fiber attenuation is typically <1.0 dB/km and maximum link length is in the range of 500 to 5000 m, depending on data rate and flux budget.

Sensors

LED/detector combinations are used in a wide range of sensor applications. They can be grouped into three classes: transmissive, reflective, and scattering sensors.

Transmissive Sensors The most widely used transmissive sensor is the slot interrupter. A U-shaped plastic holder aligns an emitter and detector face to face. It is used widely for such applications as sensing the presence or absence of paper in printers, end-of-tape in tape-recorders, erase/overwrite protection on floppy disks, and many other applications where the presence or absence of an opaque obstruction in the light path determines a system response.

A widely used slot interrupter is a two- or three-channel optical encoder. A pattern of opaque and transmissive sections moves in front of a fixed pattern with the same spatial frequency. Two optical channels positioned such that they are 90° out of phase to each other with regard to the pattern allow the measurement of both distance (number of transmissive/opaque sequences) and direction (phase of channel A with regard to channel B). Such encoders are widely used in industrial control applications, paper motion in printers, pen movement in plotters, scales, motor rotation, etc. A third channel is often used to detect an index pulse per revolution to obtain a quasi-absolute reference.

Reflective Sensors In a reflective sensor an LED, detector, and associated optical elements are positioned such that the detector senses a reflection when a reflective surface (specular or diffuse) is positioned within a narrow sensing range. A black surface or nonaligned specular surface or the absence of any reflective surface can be discriminated from a white surface or properly aligned specular surface. Applications include bar-code reading (black or white surface), object-counting on a conveyor belt (presence or no presence of a reflecting surface), and many others. Many transmissive sensor applications can be replaced by using reflective sensors and visa versa. The choice is usually determined by the optical properties of the sensing media or by cost. An emerging application for reflective sensors is blood gas analysis. The concentration of O₂ or CO₂ in blood can be determined by absorption at two different LED wavelengths, i.e., red and infrared.

Scattering Sensors One design of smoke detectors is based on light scattering. The LED light beam and the detector path are crossed. In the absence of smoke, no light from the LED can reach the detector. In the presence of smoke, light is scattered into the detector.

17.12 REFERENCES

1. A. A. Bergh and P. J. Dean, "Light-Emitting Diodes," *Proc. IEEE* vol. 60, 1972, pp. 156–224.
2. K. Gillessen and W. Shairer, *LEDs—An Introduction*, Prentice-Hall, 1987.
3. M. G. Craford, "Properties and Electroluminescence of the GaAsP Ternary System," *Progress in Solid State Chemistry*, vol. 8, 1973, pp. 127–165.
4. A. H. Herzog, W. O. Groves, and M. G. Craford, "Electroluminescence of Diffused GaAs_{1-x}P_x Diodes with Low Donor Concentrations," *J. Appl. Phys.* vol. 40, 1969, pp. 1830–1838.
5. R. A. Faulkner, "Toward a Theory of Isoelectronic Impurities in Semiconductors," *Phys. Rev.* vol. 175, 1968, pp. 991–1009.
6. W. O. Groves, A. J. Herzog, and M. G. Craford, "The Effect of Nitrogen Doping on GaAs_{1-x}P_x Electroluminescent Diodes," *Appl. Phys. Lett.* vol. 19, 1971, pp. 184–186.
7. M. G. Craford, R. W. Shaw, W. O. Groves, and A. H. Herzog, "Radiative Recombination Mechanisms in GaAsP Diodes with and without Nitrogen Doping," *J. Appl. Physics* vol. 43, 1972, pp. 4075–4083.
8. M. G. Craford, D. L. Keune, W. O. Groves, and A. H. Herzog, "The Luminescent Properties of Nitrogen Doped GaAsP Light Emitting Diodes," *J. Electron. Matls.* vol. 2, 1973, pp. 137–158.
9. M. G. Craford and W. O. Groves, "Vapor Phase Epitaxial Materials for LED Applications," *Proc. IEEE* vol. 61, 1973, pp. 862–880.
10. J. C. Campbell, N. Holonyak Jr., M. G. Craford, and D. L. Keune, "Band Structure Enhancement and Optimization of Radiative Recombination in GaAs_{1-x}P_x:N (and In_{1-x}Ga_xP:N)," *J. Appl. Phys.* vol. 45, 1974, pp. 4543–4553.
11. R. A. Logan, H. G. White, and W. Wiegmann, "Efficient Green Electroluminescence in Nitrogen-Doped GaP *p-n* Junctions," *Appl. Phys. Lett.* vol. 13, 1968, p. 139.
12. A. A. Bergh and J. A. Copeland, "Optical Sources for Fiber Transmission Systems," *Proc. IEEE*, vol. 68, 1980, pp. 1240–1247.
13. M. G. Craford, "Recent Developments in Light-Emitting Diode Technology," *IEEE Trans. Electron Devices* vol. 24, 1977, pp. 935–943.
14. H. Nather, V. Nitsche, and W. Schairer, "High Resolution Printing Capability of LED-Based Print Heads," *Proc. SPIE*, 1988, pp. 396–404.
15. M. G. Craford, "Light-Emitting Diode Displays," in *Flat-Panel Displays and CRTs*, L. E. Tannas Jr. (ed.), Van Nostrand Reinhold, 1985, pp. 289–331.
16. L. W. Cook, M. D. Camras, S. L. Rudaz, and F. M. Steranka, "High Efficiency 650 nm Aluminum Gallium Arsenide Light Emitting Diodes," *Proc. 14th International Symposium on GaAs and Related Compounds*, Institute of Physics, Bristol, 1988, pp. 777–780.
17. J. M. Dallesasse, D. W. Nam, D. G. Deppe, N. Holonyak Jr., R. M. Fletcher, C. P. Kuo, T. D. Osentowski, and M. G. Craford, "Short-Wavelength (<6400 Å) Room Temperature Continuous Operation of *p-n* InAlGaP Quantum Well Lasers," *Appl. Phys. Lett.* vol. 19, 1988, pp. 1826–1828.
18. C. P. Kuo, R. M. Fletcher, T. D. Osentowski, M. C. Lardizabel, and M. G. Craford, "High Performance AlGaInP Visible Light-Emitting Diodes," *Appl. Phys. Lett.* vol. 57, 1990, pp. 2937–2939.
19. H. Sugawara, M. Ishikawa, and G. Hatakoshi, "High-efficiency InGaAlP/GaAs Visible Light-Emitting Diodes," *Appl. Phys. Lett.* vol. 58, 1991, 1010–1012.
20. R. M. Fletcher, C. P. Kuo, T. D. Osentowski, K. H. Huang, and M. G. Craford, "The Growth and Properties of High Performance AlGaInP Emitters Using a Lattice Mismatched GaP Window Layer," *Jour. Elec. Matls.* vol. 20, 1991, pp. 1125–1130.
21. K. H. Huang, J. G. Yu, C. P. Kuo, R. M. Fletcher, T. D. Osentowski, L. J. Stinson, A. S. H. Liao, and M. G. Craford, "Twofold Efficiency Improvement in High Performance AlGaInP Light-Emitting Diodes in the 555–620 nm Spectral Region Using a Thick GaP Window Layer," *Appl. Phys. Lett.* vol. 61, 1992, pp. 1045–1047.
22. M. A. Haase, J. Qiv, J. M. DePoydt, and H. Cheng, "Blue-green Laser Diodes," *Appl. Phys. Lett.* vol. 58, 1991, pp. 1272–1275.
23. J. Jeon, J. Ding, A. V. Normikko, W. Xie, M. Kobayashi, and R. L. Gunshore, "ZnSe Based Multilayer *p/n* Junctions as Efficient Light Emitting Diodes for Display Applications," *Appl. Phys. Lett.* vol. 60, 1992, pp. 892–894.

24. S. Nakamura, M. Senoh, and T. Mukai, "Highly P-typed Mg-doped GaN Films Grown with GaN Buffer Layers," *Jpn. Jour. Appl. Phys.* vol. 30, 1991, pp. L1701–L1711.
25. A. G. Fischer, "Methods of Growing Crystals under Pressure," in *Crystal Growth*, B. R. Pamplin (ed.), Pergamon Press, Oxford, 1975, pp. 521–555.
26. R. L. Moon, "Liquid Phase Epitaxy," in *Crystal Growth*, 2d ed., B. R. Pamplin (ed.), Pergamon Press, Oxford, 1980.
27. J. W. Burd, "A Multi-Wafer Growth System for the Epitaxial Deposition of GaAs and GaAs_{1-x}P_x," *Trans. Met. Soc. AIME*, 1969, pp. 571–576.
28. G. B. Stringfellow, *Organometallic Vapor Phase Epitaxial Growth of III-V Semiconductor: Theory and Practice*, Academic Press, Oxford, 1989.
29. E. C. H. Parker (ed.), *Technology and Physics of Molecular Beam Epitaxy*, Plenum Press, New York, 1985.
30. S. K. Ghandhi, *VLSI Fabrication Principles*, John Wiley and Sons, New York, 1982.
31. Hewlett-Packard, *Optoelectronics/Fiber-Optics Application Manual*, 2d ed., McGraw-Hill, New York, 1981, p. 82.
32. A. S. Grove, *Physics and Technology of Semiconductor Devices*, sec. 6.6, John Wiley, New York, 1967.
33. O. Veda, *Material Research Society Symposium Proc.* vol. 184, 1991, p. 125.
34. M. Fukuda, *Reliability and Degradation of Semiconductor Lasers and LEDs*, Artech House, 1991.
35. A. A. Bergh, "Bulk Degradation of GaP Red LEDs," *IEEE Trans. Electron Devices* vol. 18, 1971, pp. 166–170.
36. D. C. Hanson, "Progress in Fiber Optic LAN and MAN Standards," *IEEE LCS Magazine* vol. 1, 1990, pp. 17–25.

HIGH-BRIGHTNESS VISIBLE LEDs

Winston V. Schoenfeld

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

18.1 INTRODUCTION

Over the past decade a transformation in light emitting diode (LED) application has occurred from indicator use to the more demanding solid-state lighting and illumination markets. This includes areas such as backlighting in consumer hand-held products, outdoor displays, traffic signals, and general room lighting through the replacement of incandescent and fluorescent light sources. This transformation has been driven directly by the significant increases in LED efficiency that has enabled LEDs to penetrate into such markets. These new high-efficiency LEDs are typically referred to as high-brightness LEDs (HB-LEDs), and utilize quantum well active regions to achieve their high efficiency and lumen output.

Quantum well-active regions can have either a single quantum well (SQW) or multiple quantum well (MQW) structure, reducing internal reabsorption and increasing the radiative recombination rate of the device through greater spatial overlap of electrons and holes. A discussion of quantum wells can be found in Chap. 19, “Semiconductor Lasers.” Figure 1 provides a basic structure for a modern HB-LED. It is similar to the DH structure, but utilizes a MQW active region between the n and p -type layers of the device, composed of quantum wells (QWs) that are on the order of a few nanometers in thickness. By adjusting either the composition or width of the QW, the emission wavelength of the LED can be tuned. The use of QW active regions in LEDs has not only allowed for an increase in the efficiency of the devices, it has also resulted in the ability to obtain new wavelength emissions that were not previously available due to epitaxial strain and lattice matching constraints.

18.2 THE MATERIALS SYSTEMS

There are two main semiconductor material systems currently exploited for visible HB-LEDs. These are the AlInGaP and AlInGaN systems. A detailed discussion of the AlInGaP system is in Chap. 17, “Light-Emitting Diodes.” AlInGaP is used for amber and red LEDs that fall within the 590- to 650-nm wavelength region. It is limited to this wavelength range because the AlInGaP quaternary shifts to an indirect band gap as wavelengths below 590 nm are targeted. This left a void in the blue and green spectral regions for visible HB-LEDs.

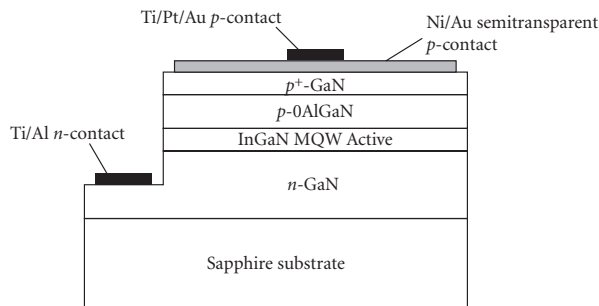


FIGURE 1 Typical structure for a modern InGaN HB-LED showing associated epitaxial layers, including the multiquantum well (MQW) active region.

In the past decade, this void has been quickly filled with MQW LEDs made from the AlInGaN material system. GaN has a wurtzite crystal structure and band gap at room temperature of about 3.2 eV (385 nm). By alloying with Al or In, the band gap can be shifted to shorter or longer wavelengths, respectively. This has led to the realization of InGaN QW active regions in HB-LEDs offering light output in the blue to green spectral regions. Although InN has a direct energy gap of 0.7 eV, phase segregation in InGaN alloys limits the ability to obtain active regions emitting at wavelengths above 530 nm. Despite this, InGaN QW active regions have successfully been used to create blue and green HB-LEDs with complete wavelength tunability across the visible 400- to 530-nm range.

18.3 SUBSTRATES AND EPITAXIAL GROWTH

Current AlInGaN HB-LEDs use one of two substrates: sapphire or silicon carbide (SiC). Despite considerable lattice mismatch between these and GaN (roughly 16 percent for sapphire and 3.5 percent for SiC), low-temperature buffer layer technologies have been developed to allow for nucleation and growth of high-quality AlInGaN HB-LEDs. Low-temperature AlN buffer layers on sapphire substrates are directly formed by MOCVD that result in dislocation entanglement just above the nucleation interface. As depicted in Fig. 2, many of the dislocations interact and annihilate. The subsequent growth of a thick (typically 3 to 4 μm) *n*-GaN buffer layer further reduces dislocation density below 10^9 per cm^2 . The use of 6H-SiC substrates offers closer lattice matching and lower defect densities, although this comes with a higher cost that has kept sapphire as a more popular substrate solution. Recent research has aimed to produce native or lattice-matched substrate solutions for GaN growth. Native GaN substrates have been produced, although the diameter of such substrates and cost remain a challenge. A sister wide band gap compound, ZnO, has strikingly similar properties

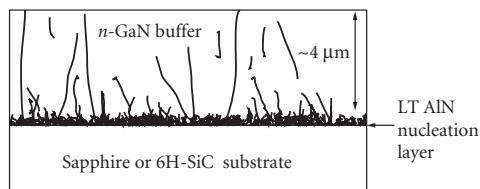


FIGURE 2 Cross-section schematic of the dislocation reduction that occurs due to use of a low temperature (LT) AlN nucleation layer. Dislocations originate from the large lattice mismatch of the nitrides with SiC or sapphire substrates, but can be reduced well below $10^9/\text{cm}^2$ when a LT-AlN nucleation layer and thick GaN buffer layer are used.

to GaN and is currently considered a good candidate for a lattice matched substrate for AlInGaN HB-LED growth. ZnO has the same wurtzite crystal structure as GaN, is nearly lattice matched to GaN, and can be easily doped n -type to form a conductive substrate. Lattice matched substrates made from ZnO have just become commercially available through hydrothermal growth.

Growth of AlInGaN HB-LEDs is accomplished commercially with MOCVD (a short discussion of MOCVD was provided in Sec. 17.8 in Chap. 17). While high quality growth of AlInGaN epitaxial films has been demonstrated by MBE, typical HB-LED structures require 4- to 5- μm -thick films that are not economically realizable in MBE systems due to the slow growth rate and limited number of substrates per growth. MOCVD reactors are able to accommodate multiple 2" substrates per growth (as many as 40 in larger systems) with the necessary uniformity and control of film thickness. One of the initial challenges in AlInGaN HB-LED growth was p -type doping of GaN. The metal-organic precursor sources of MOCVD introduce a considerable amount of hydrogen into the films that directly compensates p -type acceptors. This issue was resolved through post-growth annealing at temperatures above 800°C in which the excess hydrogen is driven out of the epitaxial films. Upon out-diffusion of the hydrogen, the p -type acceptors become active and the necessary hole injection into the structures is then possible.

18.4 PROCESSING

Many of the steps covered in Sec. 17.9 of Chap. 17 are used in the fabrication of HB-LED epitaxial wafers. Sapphire substrates are electrically insulating, imposing the need for creating both n - and p -type contacts on the front of the HB-LED surface as was shown in Fig. 1. As a result, the epitaxial structure must be etched down to the n -GaN layer in order to make electrical contact to the n -side of the HB-LED structure. AlInGaN is relatively resistant to wet chemical etching and must be dry etched using reactive ion etching (RIE) or inductively coupled plasma (ICP) methods. RIE/ICP is capable of high etch rates that are anisotropic, meaning that they are capable of etching vertically down into the structure with little to no lateral etching. SiC substrates are electrically conductive and thus allow for the n -contact to be formed on the back side of the substrate without the need for etching of the front surface.

The typical fabrication process for AlInGaN HB-LEDs utilizes four lithographic steps. These include the mesa etch, SiO₂ passivation, n -contact, and p -contact layers. A top view and associated cross section of a standard HB-LED device are given in Fig. 3, indicating these layers. Prior to the

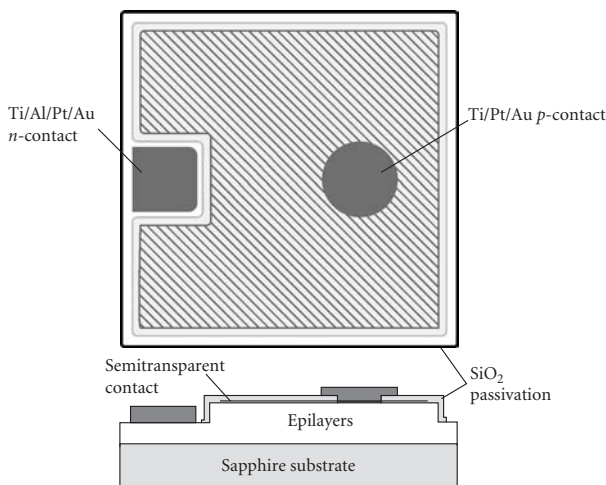


FIGURE 3 Top view and cross-section schematic of a standard nitride HB-LED. A mesa is etched to allow for contact to the n -type underlayer. A semitransparent contact on top of the mesa is used to promote uniform hole injection due to poor hole mobility in the p -GaN.

first lithographic step a semitransparent top contact is deposited on the surface (p -side) of the entire epi wafer. This is necessary in order to achieve uniform current injection across the device due to the higher resistivity of the p -GaN in comparison to the n -GaN. The most common semitransparent contacts are Ni/Au and indium tin oxide (ITO), with Ni/Au being the most widely used. Once the semitransparent contact is formed the lithographic steps are then carried out. The mesa etch is a dry etch step to access the n -GaN layers for the n -type contact as required when using sapphire substrates. The SiO_2 passivation step covers the side walls of the mesa and protects its lateral edges. Once the passivation is in place, the last two lithographic steps define the n - and p -type contacts to the device through a lift-off process. Common metallizations are Ti/Al/Pt/Au for the n -contact and Ti/Pt/Au for the p -contact to the semitransparent contact underlayer. To separate the individual LEDs on the fabricated wafer, the substrate must first be thinned from its typical 400 μm thickness to on the order of 80 to 100 μm using a multiple step wafer polishing method on automated multiwafer polishers. Once thinned, singulation of the individual die is accomplished either by a scribe and break method or by using laser separation. The sawing method typically employed for other III-V LEDs is not possible due to the substantial hardness of sapphire and SiC that greatly limits their ability to be cut using a dicing saw.

18.5 SOLID-STATE LIGHTING

The current push towards energy conservation has created a considerable interest in the replacement of conventional lighting by solid-state LED fixtures. LEDs offer the potential to considerably reduce power consumption while maintaining the necessary lumen output for lighting. Generating white light from HB-LEDs is typically accomplished by one of two methods as depicted in Fig. 4. The first method, shown in Figs. 4a and b, utilizes an AlInGaN-LED in conjunction with one or more phosphors. When a UV-LED is used, the UV light emission is absorbed and re-emitted by a mixture of red, green, and blue phosphors. As indicated in Fig. 4a the phosphors down-convert the UV light (dashed line) to visible light (solid line), and when the appropriate ratio of phosphors is used, white light is emitted. A more common LED/phosphor combination used for general illumination is the combination of a blue (~ 465 nm) AlInGaN-LED and a yttrium aluminum garnet (YAG) phosphor such as cerium doped $\text{Y}_3\text{Al}_5\text{O}_{12}$. In this approach (Fig. 4b) a portion of the blue LED emission is absorbed by the YAG phosphor and down-converted to the yellow spectral region. When the proper ratio of YAG phosphor is used, the resulting binary complimentary output of the blue LED and yellow YAG phosphor creates white light as perceived by observers. The second method for white light emission is preferred when color tuning is necessary, such as in outdoor displays. As indicated in Fig. 4c, by using red, green, and blue LEDs, one can create white light when the appropriate ratio of each is selected. This RGB approach has the added benefit of allowing the user to create any color within the associated color gamut by balancing the ratio of the individual LEDs.

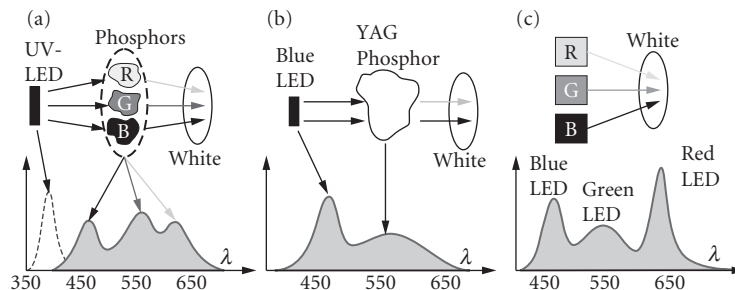


FIGURE 4 Methods for white light generation using HB-LEDs. White light is achieved by using a HB-LED and phosphors (a and b), or the combination of a red, green, and blue HB-LED (c).

This method also allows for active adjustment of the color temperature of white light emission which is not possible using the LED/phosphor approach.

The primary figure of merit for solid-state lighting is the luminous efficacy of the fixture. Luminous efficacy refers to the ratio of lumen output from the source to the power input, and has the units of lumens per Watt (lm/W). Conventional incandescent and fluorescent lights have typical luminous efficacies of 15 lm/W and 70 lm/W, respectively, varying somewhat depending on manufacturer. As of early 2008, currently available LED-based replacements have luminous efficacies as high as 80 lm/W, exceeding compact fluorescents. By comparison, lab demonstrations of white LEDs are commonly breaking 120 lm/W, with the expectation of reaching 150 lm/W in the near future. Despite these accomplishments the added cost of solid-state white lighting has slowed the market acceptance. A common cost comparison is the first cost of a light source, typically measured in terms of the cost per kilolumen (klm). Incandescent and fluorescent light sources fall in the \$0.5 to \$3.00 per klm, while current light emitting diode sources fall easily in the \$20 to \$30 per klm range. Considerable progress in efficiency continues to climb and it is expected that as the luminous efficacy and lifetime of solid-state white LEDs continues to rise that the cost per kilolumen will decrease while market acceptance increases.

18.6 PACKAGING

HB-LEDs chips for standard current use (20 to 70 mA) have lateral geometries on the order of $350 \times 350 \mu\text{m}^2$. Such LED chips are commonly packaged similarly to IR-LEDs in a 5-mm T1-3/4 format, as shown in Figs. 21 and 22 in Chap. 17. For AlInGaP HB-LEDs the substrate is conductive and silver filled epoxy is used for attaching the die to the lead frame to form one of the contacts. The other contact is formed using standard ball wire bonding to the top of the LED chip. For AlInGaN HB-LEDs, commonly grown on insulating sapphire substrates, two wire bonds are required to make the electrical connections to the lead frame. Silver-filled epoxy is still used for such LEDs since the sapphire is transparent and emission can be effectively redirected upward from the epoxy surface below the chip. In addition to the 5-mm T1-3/4 package, a considerable number of new surface mount device (SMD) packages have become available to support the introduction of HB-LEDs into the consumer hand-held device market. A schematic of a typical SMD package is provided in Fig. 5. The LED is placed on a metal lead frame that is encased in a plastic outer shell. Once the HB-LED has been die bonded into the package with silver-filled epoxy, it is then wire bonded and encapsulated with transparent epoxy for protection. Such SMD packages are typically several millimeters on a side; however, smaller versions with formats nearly identical to chip resistors are available that have

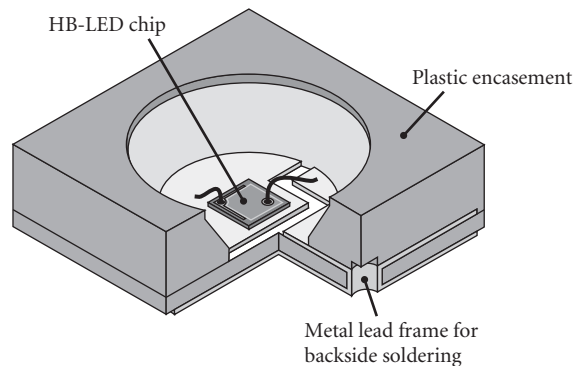


FIGURE 5 Schematic of a standard surface mount device (SMD) package for HB-LEDs.

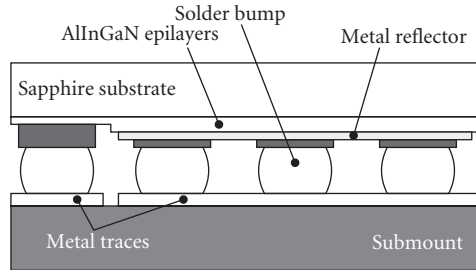


FIGURE 6 Cross-section schematic of the flip chip geometry used for high-power nitride LEDs. Flip chip packaging allows for increased heat dissipation from the LED junction, enabling them to be driven at much higher current densities.

very small form factors providing the low profile necessary for applications such as cell phone key pad backlighting.

In recent years the push for solid-state white lighting has resulted in a transition to larger HB-LED chip sizes and an increase in the operating current densities. The latter brings new constraints to packaging due to the increased need for thermal management in an effort to keep the junction temperature of the HB-LED as low as possible. While the traditional packaging formats have proven very effective for standard drive current use, they are not suitable for high drive current applications using large area HB-LEDs since they do not provide adequate thermal management. This has forced a significant change in not only the design of packaging for high drive current, large area HB-LEDs, but also the devices themselves. Most HB-LEDs use sapphire substrates that provide very poor heat dissipation from the LED junction due to the low thermal conductivity of sapphire. This has been overcome by using a flip chip approach, where the LED chip is flipped over onto a carrier and attached using solder bump methods similar to the silicon IC industry. A rough schematic of a flip chip geometry is shown in Fig. 6. The HB-LED chip is designed to support multiple solder bump attachments to the submount. A highly reflective metal layer is formed on the *p*-side of the device to redirect light emission through the backside of the device. The solder bumps affix the LED die firmly to the submount and create the electrical connections to the metal traces below, allowing for subsequent wire bonding between the submount and the external leads of the package. The solder bumps also serve as a route for efficient heat transfer from the LED junction to the submount enabling the junction temperature of the LED chip to remain low under high drive current conditions. There are a variety of external package geometries that have been developed to support flip chip HB-LEDs with no specific convention between manufacturers. Among the main considerations in the package design has been the encapsulant around the LED and providing a low thermal resistance path to the external housing. Standard clear epoxies that are used in conventional packages typically cannot withstand temperatures much above 100°C. In flip chip high drive packages this can be easily exceeded and would cause thermal damage to conventional epoxies. This has led many companies to develop a variety of new high-index encapsulants, such as silicones, that are able to withstand the higher temperature demands of high drive solid-state lighting. Low thermal resistance packaging has also followed many new routes such as packages with copper tungsten metal bases, metal core board, and high thermal conducting insulating materials (e.g., BeO and AlN).

Pamela L. Derry, Luis Figueroa, and Chi-Shain Hong

*Boeing Defense & Space Group
Seattle, Washington*

19.1 GLOSSARY

A	Constant approximating the slope of gain versus current or carrier density
C	Capacitance
c	Speed of light
D	Density of states for a transition
D_c	Density of states for the conduction band
D_v	Density of states for the valence band
d	Active layer thickness
d_{eff}	Effective beam width in the transverse direction
d_G	Guide layer thickness
dg/dN	Differential gain
E	Energy of a transition
E_c	Total energy of an electron in the conduction band
E_g	Bandgap energy
E_n	The n th quantized energy level in a quantum well
E_n^c	The n th quantized energy level in the conduction band
E_n^v	The n th quantized energy level in the valence band
E_v	Total energy of a hole in a valence band
e	Electronic charge
F_c	Quasi-Fermi level in the conduction band
F_v	Quasi-Fermi level in the valence band
f_c	Fermi occupation function for the conduction band
f_d	Damping frequency
f_o	Resonant frequency of an LRC circuit
f_p	Peak frequency

f_r	Resonance frequency
f_v	Fermi occupation function for the valence band
g	Model gain per unit length
g_{th}	Threshold modal gain per unit length
H	Heavyside function
h	Refers to heavy holes
\hbar	Plank's constant divided by 2π
I	Current
I_{off}	DC bias current before a modulation pulse
I_{on}	Bias current during a modulation pulse
I_{th}	Threshold current
J	Current density
J_o	Transparency current density
J_{th}	Threshold current density
K	Constant dependent on the distribution of spectral output function
\mathbf{k}	Wavevector
k	Boltzmann constant
L	Inductance
L	Laser cavity length
L_c	Coherence length
L_z	Quantum well thickness
l	Refers to light holes
$ M ^2$	Matrix element for a transition
m	Effective mass of a particle
m_c	Conduction band mass
m_r	Effective mass of a transition
m_v	Valence band mass
N	Carrier density
N_o	Transparency carrier density
n_{eff}	Effective index of refraction
n_r	Index of refraction
n_{sp}	Spontaneous emission factor
P	Photon density
P_{off}	Photon density before a modulation pulse
P_{on}	Photon density during a modulation pulse
R	Resistance
R_F	Front facet reflectivity
R_R	Rear facet reflectivity
T	Temperature
w	Laser stripe width
α	Absorption coefficient
α	Linewidth enhancement factor
α_i	Internal loss per unit length
β	Spontaneous emission factor
Γ	Optical confinement factor

$\Delta f_{1/2}$	Frequency spectral linewidth
$\Delta \lambda_L$	Longitudinal mode spacing
$\Delta \lambda_{1/2}$	Half-width of the spectral emission in terms of wavelength
λ	Wavelength
λ_o	Wavelength of the stimulated emission peak
τ_d	Turn-on time delay
τ_p	Photon lifetime
τ_s	Carrier lifetime

19.2 INTRODUCTION

This chapter is devoted to the performance characteristics of semiconductor lasers. In addition, some discussion is provided on fabrication and applications. In the first section we describe some of the applications being considered for semiconductor lasers. The following several sections describe the basic physics, fabrication, and operation of a variety of semiconductor laser types, including quantum well and strained layer lasers. Then we describe the operation of high-power laser diodes, including single element and arrays. A number of tables are presented which summarize the characteristics of a variety of lasers. Next we discuss the high-speed operation and provide the latest results, after which we summarize the important characteristics dealing with the spectral properties of semiconductor lasers. Finally, we discuss the properties of surface emitting lasers and summarize the latest results in this rapidly evolving field.

More than 260 references are provided for the interested reader who requires more information. In this *Handbook*, Chap. 17 (LEDs) also contains related information. For further in-depth reviews of semiconductor lasers we refer the reader to the several excellent books which have been written on the subject.¹⁻⁵

19.3 APPLICATIONS FOR SEMICONDUCTOR LASERS

The best-known application of diode lasers is in optical communication systems. However, there are many other potential applications. In particular, semiconductor lasers are being considered for high-speed optical recording,⁶ high-speed printing,⁷ single- and multimode database distribution systems,⁸ long-distance transmission,⁹ submarine cable transmission,¹⁰ free-space communications,¹¹ local area networks,¹² Doppler optical radar,¹³ optical signal processing,¹⁴ high-speed optical microwave sources,¹⁵ pump sources for other solid-state lasers,¹⁶ fiber amplifiers,¹⁷ and medical applications.¹⁸

For very high-speed optical recording systems (>100 MB/s), laser diodes operating at relatively short wavelengths ($\lambda < 0.75 \mu\text{m}$) are required. In the past few years, much progress has been made in developing short-wavelength semiconductor lasers, although the output powers are not yet as high as those of more standard semiconductor lasers.

One of the major applications for lasers with higher power and wide temperature of operation is in local area networks. Such networks will be widely used in high-speed computer networks, avionic systems, satellite networks, and high-definition TV. These systems have a large number of couplers, switches, and other lossy interfaces that determine the total system loss. In order to maximize the number of terminals, a higher-power laser diode will be required.

Wide temperature operation and high reliability are required for aerospace applications in flight control and avionics. One such application involves the use of fiber optics to directly link the flight control computer to the flight control surfaces, and is referred to as fly-by-light (FBL). A second application involves the use of a fiber-optic data network for distributing sensor and video information.

Finally, with the advent of efficient high-power laser diodes, it has become practical to replace flash lamps for the pumping of solid-state lasers such as Nd:YAG. Such an approach has the advantages

of compactness and high efficiency. In addition, the use of strained quantum well lasers operating at $0.98\ \mu\text{m}$ has opened significant applications for high-gain fiber amplifiers for communications operating in the $1.55\text{-}\mu\text{m}$ wavelength region.

19.4 BASIC OPERATION

Lasing in a semiconductor laser, as in all lasers, is made possible by the existence of a gain mechanism plus a resonant cavity. In a semiconductor laser the gain mechanism is provided by light generation from the recombination of holes and electrons (see Fig. 1). The wavelength of the light is determined by the energy bandgap of the lasing material. The recombining holes and electrons are injected, respectively, from the p and n sides of a p - n junction. The recombing carriers can be generated by optical pumping or, more commonly, by electrical pumping, i.e., forward-biasing the p - n junction. In order for the light generation to be efficient enough to result in lasing, the active region of a semiconductor laser, where the carrier recombination occurs, must be a direct bandgap semiconductor. The surrounding carrier injection layers, which are called *cladding layers*, can be indirect bandgap semiconductors. For a discussion of semiconductor band levels see any solid-state physics textbook such as that by Kittel.¹⁹ For a more detailed discussion of carrier recombination see Chap. 17 in this *Handbook*.

For a practical laser, the cladding layers have a wider bandgap and a lower index of refraction than the active region. This type of semiconductor laser is called a *double heterostructure* (DH) laser, since both cladding layers are made of a different material than the active region (see Fig. 2). The first semiconductor lasers were homojunction lasers,^{21–24} which did not operate at room temperature; it is much easier to achieve lasing in semiconductors at low temperatures. Today, however, all semiconductor lasers contain heterojunctions. The narrower bandgap of the active region confines carrier recombination to a narrow optical gain region. The sandwich of the larger refractive index active region surrounded by cladding layers forms a waveguide, which concentrates the optical modes generated by lasing in the active region. For efficient carrier recombination the active region must be fairly thin, typically on the order of $1000\ \text{\AA}$, so a significant fraction of the optical mode spreads into the cladding layers. In order to completely confine the optical mode in the semiconductor structure, the cladding layers must be fairly thick, usually about $1\ \mu\text{m}$.

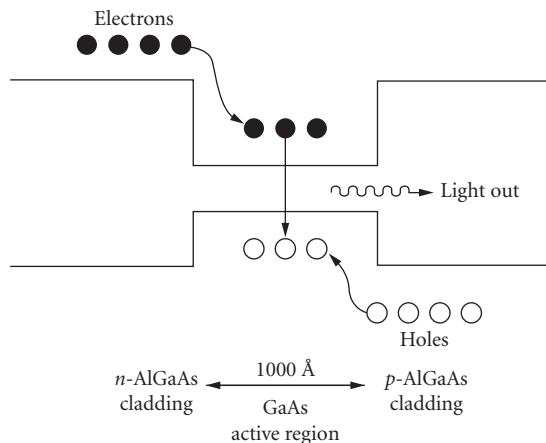


FIGURE 1 Schematic diagram of the recombination of electrons and holes in a semiconductor laser.

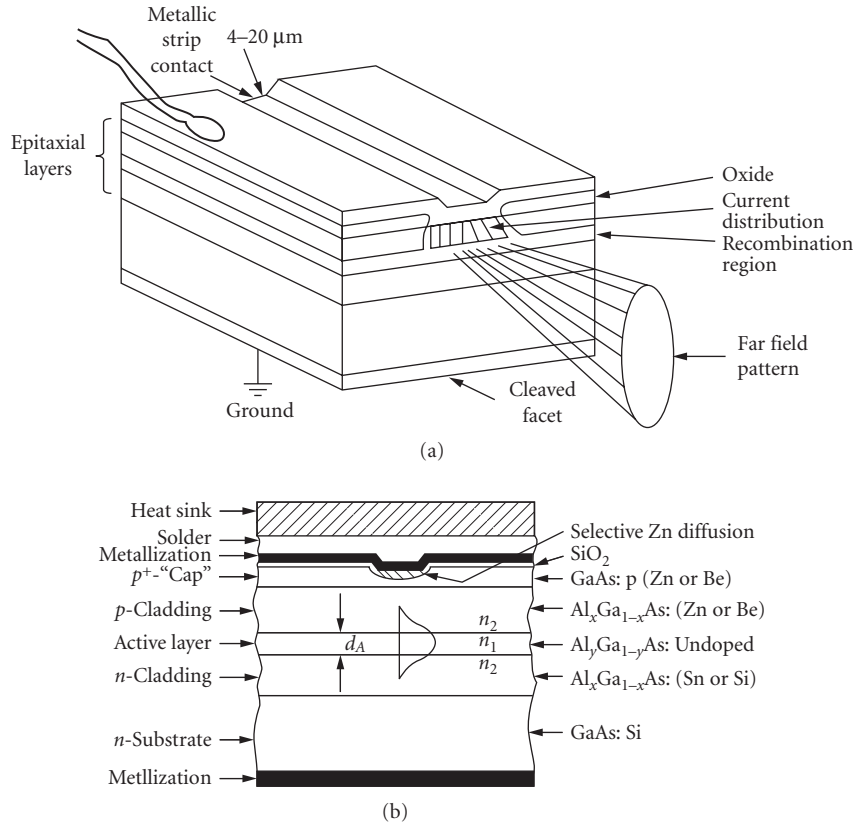


FIGURE 2 (a) Schematic diagram of a simple double heterostructure laser and (b) cross-sectional view showing the various epitaxial layers. (After Ref. 20.)

The resonant cavity of a simple semiconductor laser is formed by cleaving the ends of the structure. Lasers are fabricated with their lasing cavity oriented perpendicular to a natural cleavage plane. For typical semiconductor materials this results in mirror reflectivities of about 30 percent. If necessary, the reflectivities of the end facets can be modified by applying dielectric coatings to them.²⁵ For applications where it is not possible to cleave the laser facets, it is also possible to etch them,²⁶⁻²⁸ although this is much more difficult and usually does not work as well. Laser cavity lengths can be anywhere from about 50 to 2000 μm , although commercially available lasers are typically 200 to 1000 μm long.

Unpumped semiconductor material absorbs light of energy greater than or equal to its bandgap. When the semiconductor material is pumped optically or electrically, it reaches a point at which it stops being absorbing. This point is called *transparency*. If it is pumped beyond this point, it will have optical gain, which is the opposite of absorption. A semiconductor laser is subject to both internal and external losses. For lasing to begin, i.e., to reach threshold, the gain must be equal to these optical losses. The threshold gain per unit length is given by

$$g_{\text{th}} = \alpha_i + \frac{1}{2L} \ln \left(\frac{1}{R_F R_R} \right) \quad (1)$$

where α_i is the internal loss per unit length, L is the laser cavity length, and R_F and R_R are the front and rear facet reflectivities. (For semiconductor lasers, gain is normally quoted as gain per unit

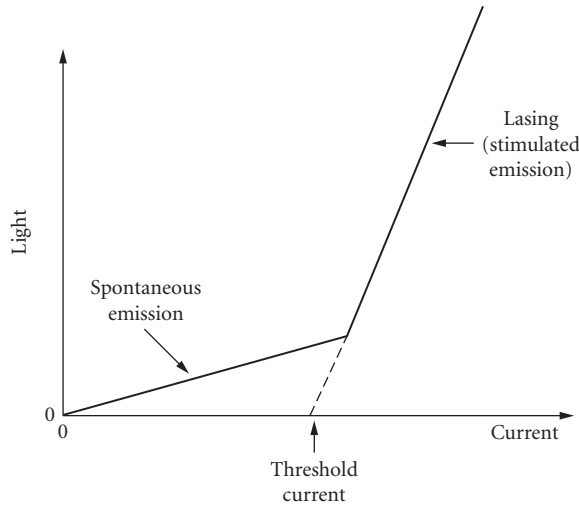


FIGURE 3 An example of the light-versus-current relationship of a semiconductor laser illustrating the definition of threshold current.

length in cm^{-1} . This turns out to be very convenient, but unfortunately is confusing for people in other fields, who are used to gain being unitless.)

The internal loss is a material parameter determined by the quality of the semiconductor layers. Mechanisms such as free-carrier absorption and scattering contribute to α_i .¹ The second term in Eq. (1) is the end loss. A long laser cavity will have reduced end loss, since the laser light reaches the cavity ends less frequently. Similarly, high facet reflectivities also decrease the end loss, since less light leaves the laser through them.

For biases below threshold, a semiconductor laser emits a small amount of incoherent light spontaneously (see Fig. 3). This is the same type of light emitted by an LED (see Chap. 17). Above threshold, stimulated emission results in lasing. The relationship between lasing emission and the bias current of a healthy semiconductor laser is linear. To find the threshold current of a laser this line is extrapolated to the point at which the stimulated emission is zero (see Fig. 3).

For further discussion of optical gain in semiconductors, see under “Quantum Well Lasers” later in this chapter. For more detail, see one of the books that has been written about semiconductor lasers.¹⁻⁵

19.5 FABRICATION AND CONFIGURATIONS

In order to fabricate a heterostructure laser, thin semiconductor layers of varying composition must be grown on a semiconductor substrate (normally GaAs or InP). There are three primary epitaxial methods for growing these layers: liquid phase epitaxy (LPE), molecular beam epitaxy (MBE), and metalorganic chemical vapor deposition (MOCVD), which is also called organometallic vapor phase deposition (OMVPE).

Most of the laser diode structures which have been developed were first grown by LPE.²⁹ For a description of LPE see Chap. 17. Most commercially available lasers are grown by LPE; however, it is not well suited for growing thin structures such as quantum well lasers, because of lack of control and uniformity, especially over large substrates.³⁰ MBE and MOCVD are better suited for growths of thin, uniform structures.

MOCVD^{31,32} is basically a specialized form of chemical vapor deposition. In MOCVD, gases reacting over the surface of a substrate form epitaxial layers; some of the gases are metalorganics. MOCVD is well suited for production environments, since epitaxial layers can be grown simultaneously on multiple large-area substrates and quickly, compared to MBE. It is expected that more commercial laser diodes will be grown by MOCVD in the future.

In the simplest terms, MBE^{30,33–35} is a form of vacuum evaporation. In MBE growth occurs through the thermal reaction of thermal beams of atoms and molecules with the substrate, which is held at an appropriate temperature in an ultrahigh vacuum. MBE is different from simple vacuum evaporation for several reasons: the growth is single crystalline; the growth is much more controlled; and the vacuum system, evaporation materials, and substrate are cleaner.

With MOCVD the sources are gases, while with MBE they are solids. There are advantages and disadvantages to both types of sources. With gaseous sources the operator must work with arsine and/or phosphine, which are extremely hazardous gases. Solid-source phosphorus, however, is very flammable. Also, with MBE balancing the ratios of arsenic and phosphorus is extremely difficult; therefore, MOCVD is the preferred method for growth of GaInAsP and InP. MBE is a slower growth process (on the order of 1 to 2 $\mu\text{m}/\text{h}$) than MOCVD. MBE therefore has the control necessary to grow very thin structures (10 \AA), but MOCVD is more efficient for production. MBE has a cleaner background environment, which tends to make it better suited for growths at which background impurities must be kept at a minimum. Newer growth techniques,^{36,37} which combine some of the advantages of both MBE and MOCVD, are gas source MBE, metalorganic MBE (MOMBE), and chemical beam epitaxy (CBE). In these growth techniques the background environment is that of MBE, but some or all of the sources are gases, which makes them more practical for growth of phosphorus-based materials.

Double heterostructure (DH) semiconductor lasers can be fabricated from a variety of lattice-matched semiconductor materials. The two material systems most frequently used are GaAs/Al_xGa_{1-x}As and In_{1-x}Ga_xAsP_{1-y}/InP. All of these semiconductors are III-V alloys. The GaAs/Al_xGa_{1-x}As material system has the advantage that all compositions of Al_xGa_{1-x}As are closely lattice-matched to GaAs, which is the substrate. For GaAs-based lasers, the active region is usually GaAs or low-Al-concentration Al_xGa_{1-x}As ($x < 0.15$), which results in lasing wavelengths of 0.78 to 0.87 μm . Al_xGa_{1-x}As quantum well lasers with wavelengths as low as 0.68 μm have been fabricated,³⁸ but the performance is reduced.

In the In_{1-x}Ga_xAsP_{1-y}/InP material system, the active region is In_{1-x}Ga_xAsF_{1-y} and the cladding layers and substrate are InP. Not all compositions of In_{1-x}Ga_xAsP_{1-y} are lattice-matched to InP; x and y must be chosen appropriately to achieve both lattice match and the desired lasing wavelength.⁴ The lasing wavelength range of InP-based lasers, 1.1 to 1.65 μm , includes the wavelengths at which optical fibers have the lowest loss (1.55 μm) and material dispersion (1.3 μm). This match with fiber characteristics makes In_{1-x}Ga_xAsP_{1-y}/InP lasers the preferred laser for long-distance communication applications. InP-based lasers can also include lattice-matched In_{1-x-y}Al_xGa_yAs layers,^{39–42} but the performance is reduced.

There is a great deal of interest in developing true visible lasers for optical data storage applications. (Al_xGa_{1-x})_{0.5}In_{0.5}P lasers^{43–46} lattice-matched to GaAs have proven superior to very short wavelength GaAs/Al_xGa_{1-x}As lasers. Higher Al concentration layers are cladding layers and a low Al concentration layer or Ga_{0.5}In_{0.5}P is the active region. In this material, system lasers with a lasing wavelength as low as 0.63 μm which operate continuously at room temperature have been fabricated.⁴⁶

In order to fabricate a blue semiconductor laser, other material systems will be required. Recently, lasing at 0.49 μm at a temperature of 77 K was demonstrated in a ZnSe- (II-VI semiconductor) based laser.⁴⁷

Very long-wavelength (>2 μm) semiconductor lasers are of interest for optical communication and molecular spectroscopy. The most promising results so far have been achieved with GaInAsSb/AlGaAsSb lattice-matched to a GaSb substrate. These lasers have been demonstrated to operate continuously at 30°C with a wavelength of 2.2 μm .⁴⁸

Lead salt lasers (Pb_{1-x}Eu_xSeTe_{1-y}, Pb_{1-x}Sn_xSe, PbS_{1-x}Se, PbS_{1-x}Sn_xTe, Pb_{1-x}Sr_xS) can be fabricated for operation at even longer wavelengths,^{4,49–52} but they have not been demonstrated at room temperature. Progress has been made, however, increasing the operating temperature; currently Pb_{1-x}Eu_xSeTe_{1-y}/PbTe lasers operating continuously at 203 K with a lasing wavelength of 4.2 μm have been

demonstrated.⁵³ Other very long-wavelength lasers are possible; recently, HgCdTe lasers with pulsed operation at 90 K and a lasing wavelength of 3.4 μm have been fabricated.⁵⁴

Laser Stripe Structures

We have discussed the optical and electrical confinement provided by a double heterostructure parallel to the direction of epitaxial growth; practical laser structures also require a confinement structure in the direction parallel to the substrate.

The simplest semiconductor laser stripe structure is called an *oxide stripe laser* (see Fig. 4a). The metallic contact on the *n*-doped side of a semiconductor laser is normally applied with no definition for current confinement; current confinement is introduced on the *p* side of the device. For a wide-stripe laser, a dielectric coating (usually SiO_2 or Si_3N_4) is evaporated on the *p* side of the laser. Contact openings in the dielectric are made through photolithography combined with etching of the dielectric. The *p* metallic contact is then applied across the whole device, but makes electrical contact only at the dielectric openings. A contact stripe works very well for wide stripes, but in narrow stripes current spreading is a very significant drawback, because there is no mechanism to prevent current spreading after the current is injected. Also, since the active region extends outside of the stripe, there is no mechanism to prevent optical leakage in a contact-stripe laser. Lasers like this, which provide electrical confinement, but no optical confinement, are called *gain-guided lasers*.

Another type of gain-guide laser is an ion bombardment stripe (see Fig. 4b). The material outside the stripe is made highly resistive by ion bombardment or implantation, which produces lattice defects.⁵⁵ Implantation causes optical damage,⁵⁶ so implantation should not be heavy enough to reach the active region.

A more complicated stripe structure with electrical and optical confinement is required for an efficient narrow-stripe laser. A number of structures which accomplish the necessary confinement have been developed. These structures are called *index-guided lasers*, since optical confinement is achieved through a change in refractive index.

The buried heterostructure laser (BH) was first developed by Tsukada.⁵⁷ To form a BH stripe, a planar laser structure is first grown. Stripe mesas of the laser structure are formed by photolithography combined with etching. For a GaAs-based BH laser, AlGaAs is then regrown around the lasing stripe. Figure 4c is a schematic diagram of a buried heterostructure. Since the active region is completely surrounded by AlGaAs, a BH has tight optical confinement. If the regrown layers are doped to produce a reverse-biased junction or are semi-insulating, a BH laser can also provide good current confinement. There are many variations on the BH structure. In some cases the active region is grown in the second growth step (see Fig. 4d). The tight optical confinement of BH lasers allows practical fabrication of very narrow stripes, on the order of 1 to 2 μm .

There are many other stripe structures that provide weaker optical confinement than a buried heterostructure. One of the simplest and most widely used of these is the ridge waveguide laser (RWG) (Fig. 4e). After epitaxial growth, most of the *p*-cladding layer is etched away, leaving a mesa where the lasing stripe will be. Only this mesa is contacted, which provides electrical confinement. The change in surrounding refractive index produces an effective change in refractive index in the active region beyond the mesa and provides optical confinement. Other stripe structures are described later in this chapter under "High-Power Semiconductor Lasers."

Another type of laser stripe is one in which confinement is provided by the *p-n* junction. The best-known laser of this type is the transverse junction stripe⁵⁸⁻⁶⁰ (see Fig. 4f). In order to fabricate a TJS laser, both cladding layers are grown as *n*-AlGaAs. Zn diffusion is then used to create a *p-n* junction and contacts are applied on either side of the junction. In this laser the current flows parallel to the substrate rather than perpendicular to it. In a TJS laser the active region is limited to the small region of GaAs in which the Zn diffusion front ends.

The examples of laser stripe structures described here are GaAs/AlGaAs lasers. Long-wavelength laser structures (InP-based) are very similar,⁴ but the active region is InGaAsP and the cladding layers are InP. With an *n*-InP substrate the substrate can be used as the *n*-type cladding, which allows greater flexibility in designing structures such as that illustrated in Fig. 4d. For a more detailed discussion of GaAs-based laser stripe structures see Casey and Panish² or Thompson.³

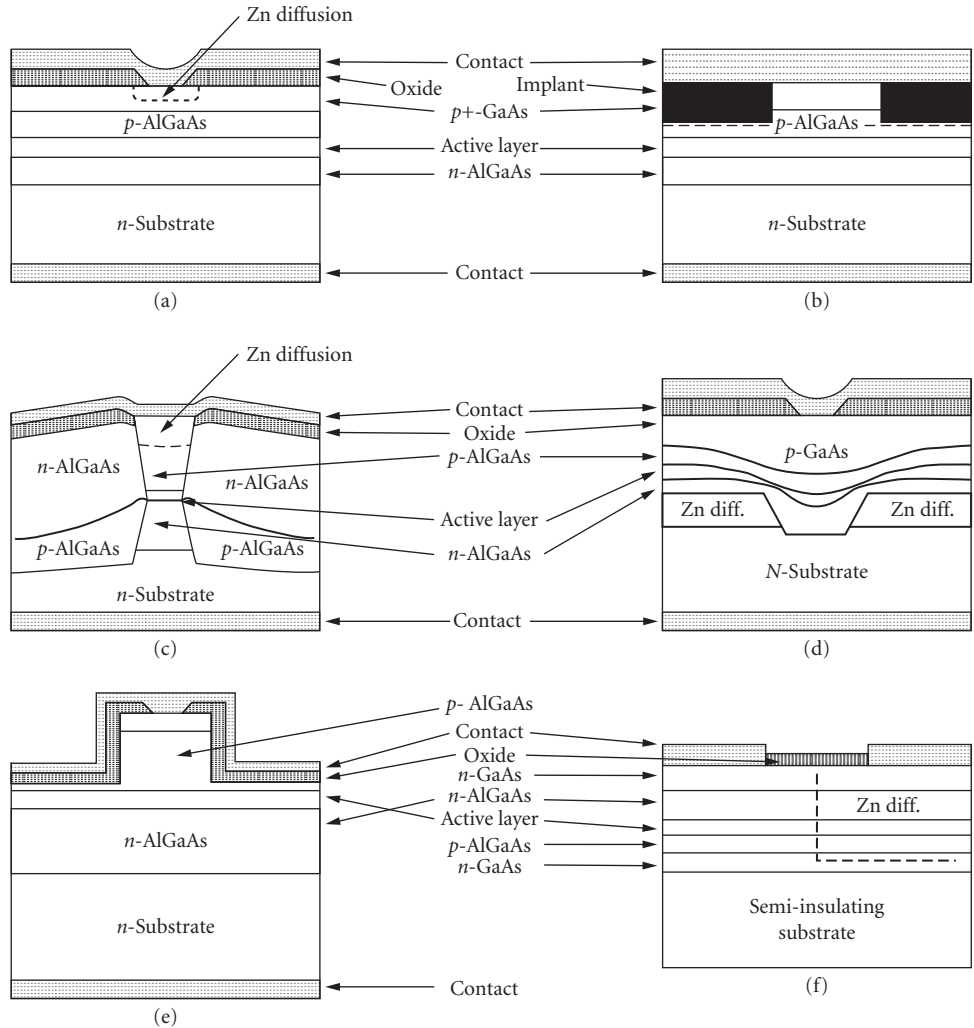


FIGURE 4 Schematic diagrams of GaAs/AlGaAs stripe laser structures. (a) Oxide stripe laser; (b) ion bombardment laser; (c) buried heterostructure (BH) laser; (d) variation on buried heterostructure laser; (e) ridge waveguide (RWG) laser; and (f) transverse junction stripe (TJS) laser.

19.6 QUANTUM WELL LASERS

The active region in a conventional DH semiconductor laser is wide enough ($\sim 1000 \text{ \AA}$) that it acts as bulk material and no quantum effects are apparent. In such a laser the conduction band and valence band are continuous (Fig. 5a). In bulk material the density of states, $D(E)$, for a transition energy E per unit volume per unit energy is²⁶

$$D(E) = \sum_{i=l,h} \frac{m_r^i}{\pi^2 \hbar^3} \sqrt{2m_r^i(E - E_g)} \quad E > E_g \quad (2)$$

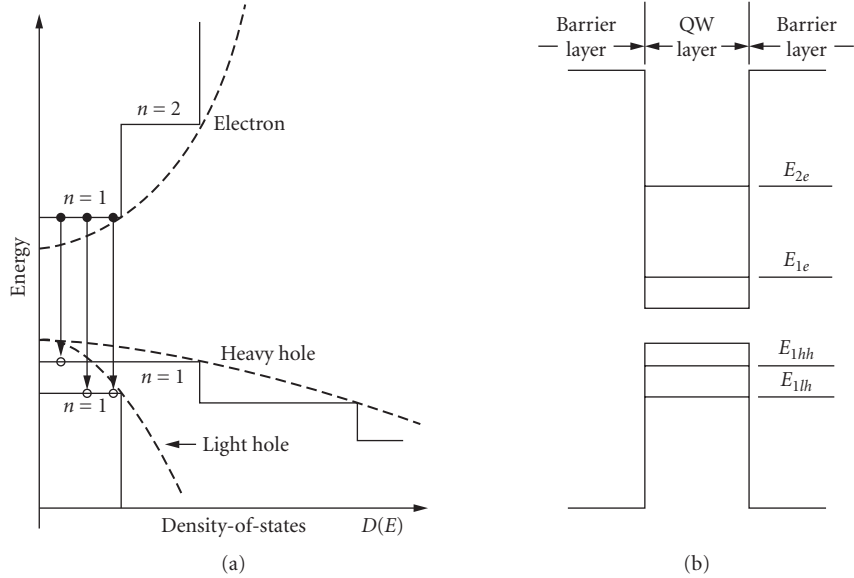


FIGURE 5 Schematic diagrams of (a) the density of states for a quantum well (solid line) and for a bulk DH (dotted line) and (b) quantized energy levels in a quantum well for $n = 1$ and 2 for the conduction band and for the light and heavy hole bands. (After Ref. 63.)

where E_g is the bandgap energy, \hbar is Planck's constant divided by 2π , l and h refer to light and heavy holes, and m_r is the effective mass of the transition which is defined as

$$\frac{1}{m_r} = \frac{1}{m_c} + \frac{1}{m_v} \quad (3)$$

where m_c is the conduction band mass and m_v is the valence band mass. (The split-off hole band and the indirect conduction bands are neglected here and have a negligible effect on most semiconductor laser calculations.)

If the active region of a semiconductor laser is very thin (on the order of the DeBroglie wavelength of an electron) quantum effects become important. When the active region is this narrow (less than ~ 200 Å) the structure is called a *quantum well* (QW). (For a review of QWs see Dingle,⁶¹ Holonyak et al.,⁶² Okamoto,⁶³ or the book edited by Zory.⁶⁴) Since the quantum effects in a QW are occurring in only one dimension they can be described by the elementary quantum mechanical problem of a particle in a one-dimensional quantum box.⁶⁵ In such a well, solution of Schrödinger's equation shows that a series of discrete energy levels (Fig. 5b) are formed instead of the continuous energy bands of the bulk material. With the approximation that the well is infinitely deep, the allowed energy levels are given by

$$E_n = \frac{(n\pi\hbar)^2}{2mL_z^2} \quad (4)$$

where $n = 1, 2, 3, \dots, m$ is the effective mass of the particle in the well, and L_z is the quantum well thickness. Setting the energy at the top of the valence band equal to zero, the allowed energies for an electron in the conduction band of a semiconductor QW become $E = E_g + E_n^c$, where E_n^c is E_n with m

equal to m_c . The allowed energies for a hole in the valence band are then $E = -E_n^v$, where E_n^v is E_n with m equal to m_v . The allowed transition energies E are limited to

$$E = E_g + E_n^c + E_n^v + \frac{\hbar^2 \mathbf{k}^2}{2m_r} \quad (5)$$

where \mathbf{k} is the wavevector, rigorous \mathbf{k} -selection is assumed, and transitions are limited to those with $\Delta n = 0$.

This quantization of energy levels will, of course, change the density of states. For a QW the density of states is given by

$$D(E) = \sum_{i=l,h} \sum_{n=1}^{\infty} \frac{m_r^i}{L_z \pi \hbar^2} H(E - E_g - E_n^c - E_n^v) \quad (6)$$

where H is the heavyside function. The difference in the density of states directly affects the modal optical gain generated by the injection of carriers. The modal gain is proportional to the stimulated emission rate:^{1,66,67}

$$g(E, N) = \alpha \frac{\Gamma D(E) |M|^2 (f_c(E, N) - f_v(E, N))}{E} \quad (7)$$

where I is the optical confinement factor, $|M|^2$ is the matrix element for the transition, N is the carrier density of either electrons or holes (the active region is undoped so they have equal densities), and $f_c(E, N)$ and $f_v(E, N)$ are the Fermi occupation factors for the conduction and valence bands. (For a detailed review of gain in semiconductor lasers see Ref. 67.)

The optical confinement factor Γ is defined as the ratio of the light intensity of the lasing mode within the active region to the total intensity over all space. Since a QW is very thin, Γ_{QW} will be much smaller than Γ_{DH} . Γ_{DH} is typically around 0.5 whereas for a single QW, Γ_{QW} is around 0.03.

The Fermi occupation functions describe the probability that the carriers necessary for stimulated emission have been excited to the states required. They are given by^{1,19}

$$f_c(E_c, N) = \frac{1}{1 + \exp((E_c + E_c - F_c)/kT)} \quad (8)$$

and

$$f_v(E_v, N) = \frac{1}{1 + \exp(-(E_v - F_v)/kT)} \quad (9)$$

where k is the Boltzmann constant, T is temperature, E_c is the energy level of the electron in the conduction band relative to the bottom of the band (including both the quantized energy level and kinetic energy), E_v is the absolute value of the energy level of the hole in a valence band, and F_c and F_v are the quasi-Fermi levels in the conduction and valence bands. Note that E_c and E_v are dependent on E , so f_c and f_v are functions of E . f_c and f_v are also functions of N through F_c and F_v . F_c and F_v are obtained by evaluating the expressions for the electron and hole densities:

$$N = \int D_c(E_c) f_c(E_c) dE_c \quad (10)$$

and

$$N = \int D_v(E_v) f_v(E_v) dE_v \quad (11)$$

where $D_c(E_c)$ and $D_v(E_v)$ are the densities of states for the conduction and valence bands and follow the same form as $D(E)$ for a transition.

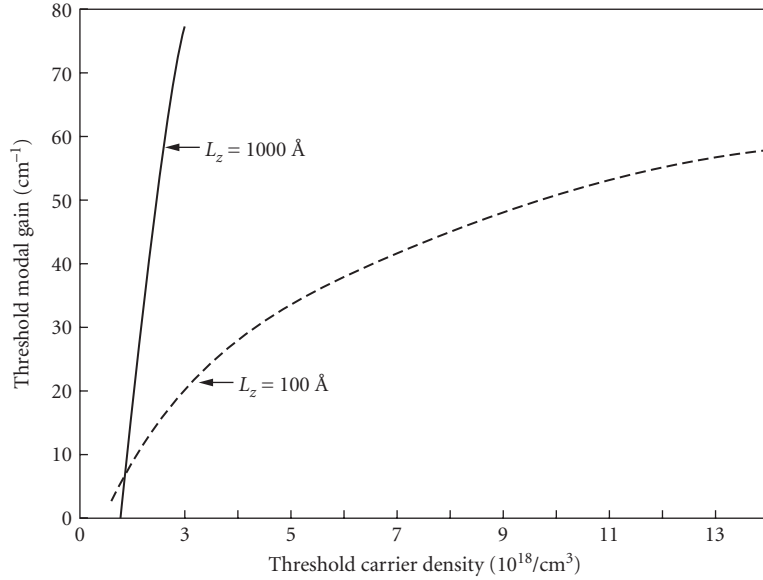


FIGURE 6 Threshold modal gain as a function of threshold carrier density for a conventional (Al, Ga)As double heterostructure with an active region thickness of 1000 Å and for a 100-Å single quantum well. (From Ref. 68.)

In Fig. 6, the results of a detailed calculation⁶⁸ based on Eq. (7) for the threshold modal gain as a function of threshold carrier density are plotted for a 100-Å single QW. The corresponding curve for a DH laser with an active region thickness of 1000 Å is also shown. The gain curves for the QWs are very nonlinear because of saturation of the first quantized state as the carrier density increases. The transparency carrier density N_0 is the carrier density at which the gain is zero. From Fig. 6 it is clear that the transparency carrier densities for QW and DH lasers are very similar and are on the order of $2 \times 10^{18} \text{ cm}^{-3}$.

The advantage of a QW over a DH laser is not immediately apparent. Consider, however, the transparency current density J_o . At transparency²⁸

$$J_o = \frac{N_0 L_z e}{\tau_s} \quad (12)$$

where L_z is the active layer thickness, e is the charge of an electron, and τ_s is the carrier lifetime near transparency. τ_s is approximately 2 to 4 ns for either a QW or a DH laser. Since N_0 is about the same for either structure, any difference in J_o will be directly proportional to L_z . But L_z is approximately 10 times smaller for a QW; therefore, J_o will be approximately 10 times lower for a QW than for a conventional DH laser. (A lower J_o will result in a lower threshold current density since the threshold current density is equal to J_o plus a term proportional to the threshold gain.) Note that this result is not determined by the quantization of energy levels; it occurs because fewer carriers are needed to reach the same carrier density in a QW as in a DH laser. In other words, this result is achieved because the QW is thin!

In this discussion we are considering current density instead of current. The threshold current density (current divided by the length and width of the stripe) is a more meaningful measure of the relative quality of the lasing material than is current. Current depends very strongly on the geometry and stripe fabrication of the device. In order to eliminate geometry-induced variations, current density is normally measured on broad-area (50 to 150 μm wide) oxide stripe lasers (see earlier under

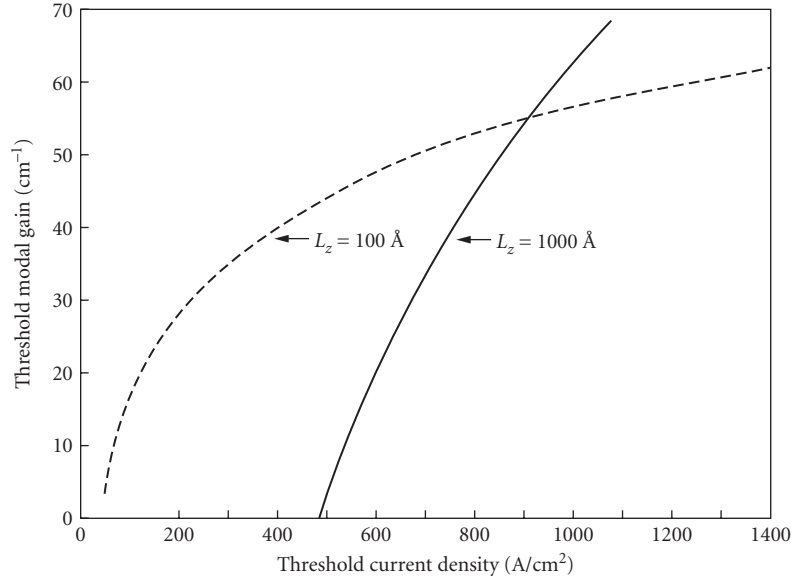


FIGURE 7 Threshold modal gain as a function of threshold current density for a conventional (Al, Ga)As double heterostructure with an active region thickness of 1000 Å and for a 100-Å single quantum well. (From Ref. 68.)

“Fabrication and Configurations”). With a narrow stripe the current spreads beyond the intended stripe width, so it is difficult to accurately measure the current density.

Figure 7 shows the results of a detailed calculation⁶⁸ of the threshold modal gain versus the threshold current density for a DH laser with an active region thickness of 1000 Å and for a 100-Å single QW. The potential for lower threshold current densities for QW lasers is clear for threshold gains less than that where the curve for the DH laser intercepts those of the QWs. With low losses, the threshold current of a QW laser will be substantially lower than that of a DH laser, since the threshold gain will be below the interception point.

To get an appreciation for how the threshold current density of a single QW will compare to that of a DH laser, consider that near transparency, the modal gain is approximately linearly dependent on the current density:

$$g(J) = A(J - J_o) \quad (13)$$

where A is a constant which should have a similar value for either a QW or a DH laser (this can be seen visually on Fig. 7). Taking Eq. (13) at threshold we can equate it to Eq. (1) and solve for J_{th} (the threshold current density):

$$J_{th} = J_o + \frac{\alpha_i}{A} + \frac{1}{2AL} \ln \left(\frac{1}{R_F R_R} \right) \quad (14)$$

α_i is related primarily to losses occurring through the interaction of the optical mode with the active region. In a QW, the optical confinement is lower, which means that the optical mode interacts less with the active region and α_i tends to be smaller. Let's substitute in the numbers in order to get an idea for the difference between a QW and a DH laser. Reasonable values are⁶⁹ $A_{QW} \sim 0.7 \text{ A}^{-1} \text{ cm}$, $A_{DH} \sim 0.4 \text{ A}^{-1} \text{ cm}$, $J_o^{QW} \sim 50 \text{ A/cm}^2$, $J_o^{DH} \sim 500 \text{ A/cm}^2$, $\alpha_i^{QW} \sim 2 \text{ cm}^{-1}$, $\alpha_i^{DH} \sim 15 \text{ cm}^{-1}$, $L \sim 400 \text{ } \mu\text{m}$, and for uncoated facets $R_F = R_R = 0.32$. Substituting in we get $J_{th}^{QW} \sim 95 \text{ A/cm}^2$ and $J_{th}^{DH} \sim 610 \text{ A/cm}^2$.

It is clear that changes in the losses will have a more noticeable effect on threshold current for a QW than for a DH laser since losses are responsible for a more significant portion of the threshold current of a QW laser. The gain curve of a QW laser saturates due to the filling of the first quantized energy level, so operating with low losses is even more important for a QW than is illustrated by the above calculation. When the gain saturates, the simple approximation of Eq. (13) is invalid. Operating with low-end losses is also important for a QW, since they are a large fraction of the total losses. This explains why threshold current density results for QW lasers are typically quoted for long laser cavity lengths (greater than 400 μm), while DH lasers are normally cleaved to lengths on the order of 250 μm . High-quality broad-area single QW lasers (without strain) have threshold current densities lower than 200 A/cm^2 (threshold current densities as low as 93 A/cm^2 have been achieved^{69–71}), while the very best DH lasers have threshold current densities around 600 A/cm^2 .⁷² The end loss can also be reduced by the use of high-reflectivity coatings.²⁵ The combination of a single QW active region with a narrow stripe and high-reflectivity coatings has allowed the realization of submilliampere threshold current semiconductor lasers^{68,69,73} and high-temperature operation.^{74,75,76}

A disadvantage of QW lasers compared to DH lasers is the loss of optical confinement. One of the advantages of a DH laser is that the active region acts as a waveguide, but in a QW the active region is too thin to make a reasonable waveguide. Guiding layers are needed between the QW and the (Al, Ga)As cladding layers. As the bandgap diagram of Fig. 8 illustrates, a graded layer of intermediate aluminum content can be inserted between the QW and each cladding layer. The advantage of this structure, which is called a *graded-index separate-confinement heterostructure* (GRIN SCH),^{77,78} is separate optical and electrical confinement. The carriers are confined in the QW, but the optical mode is confined in the surrounding layers. The grading can be either parabolic (as illustrated in Fig. 8) or linear. Experimentally it has been found that the optimum AlAs mole fraction x for layers around a GaAs QW is approximately 0.2.⁷⁹ Typically, each additional layer is on the order of 2000 \AA thick. In order to confine the optical mode, the cladding layers need a low index of refraction compared to that of an $x = 0.2$ layer. In a simple DH laser, the cladding layers typically have x between 0.3 and 0.4, but for good confinement in an $x = 0.2$ layer, more aluminum should be incorporated into the cladding layers; x should be between 0.5 and 0.7.

In the discussion so far we have considered only single QWs. Structures in which several quantum wells are separated by thin AlGaAs barriers are called *multi-quantum wells* (MQWs) and also have useful properties. For a given carrier density, an MQW with n QWs of equal thickness, L_z has gain which is approximately n times the gain for a single QW of the same thickness L_z , but the current

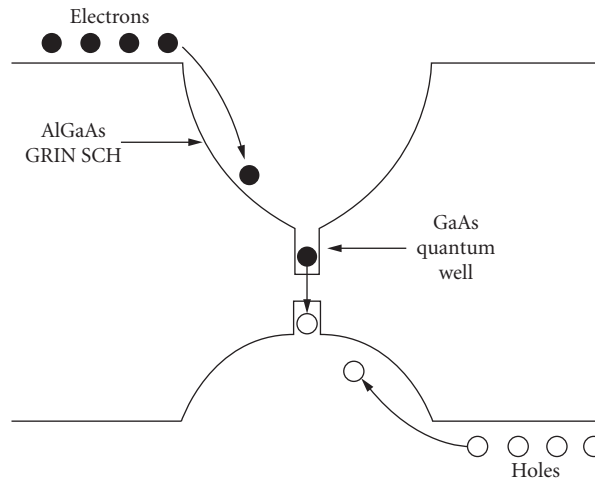


FIGURE 8 Schematic energy-band diagram for a GRIN SCH single QW.

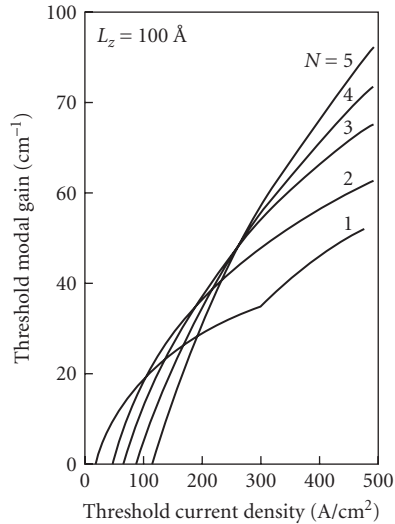


FIGURE 9 Threshold modal gain as a function of threshold current density for (Al, Ga) As single QW and MQWs with 2, 3, 4, and 5 QWs. Each QW has a thickness of 100 Å. (From Ref. 80.)

density is also approximately increased by a factor of n .⁸⁰ The transparency current density will be larger for the MQW than for the single QW since the total active region thickness is larger. Figure 9 shows that, as a function of current density, the gain in the single QW will start out higher than that in the MQW because of the lower transparency current, but the gain in the MQW increases more quickly so the MQW gain curve crosses that of the single QW at some point.⁸¹

Which QW structure has a lower threshold current will depend on how large the losses are for a particular device structure and on where the gain curves cross. The best structure for low threshold current in a GaAs-based laser is normally a single QW, but for some applications involving very large losses and requiring high gain an MQW is superior. For applications in which high output power is more important than low current an MQW is appropriate. An MQW is also preferred for high-modulation bandwidth (see “Spectral Properties” later in this chapter).

An advantage of a QW structure over a bulk DH laser is that the lasing wavelength, which is determined by the bulk bandgap plus the first quantized energy levels, can be changed by changing the quantum well thickness [see Eq. (4)]. A bulk GaAs laser has a lasing wavelength of about 0.87 μm , while a GaAs QW laser of normal thickness (60 to 120 Å) has a lasing wavelength of 0.83 to 0.86 μm . Further bandgap engineering can be introduced with a strained QW.^{67,81,82}

Strained Quantum Well Lasers

Normally, if a semiconductor layer of significantly different lattice constant is grown in an epitaxial structure, it will maintain its own lattice constant and generate misfit dislocations. If this layer is very thin, below a certain critical thickness,^{81,83–85} it will be distorted to match the lattice constant (perpendicular to the substrate) of the surrounding layers and will not generate misfit locations. A layer with thickness above the critical thickness is called “relaxed,” one below is called “strained.”

Straining a semiconductor layer changes the valence-band structure. Figure 10a shows the band structure of an unstrained III–V semiconductor, while Fig. 10b and c show the band structure under biaxial tension and compression, respectively.^{81,82} For the unstrained semiconductor, the light and

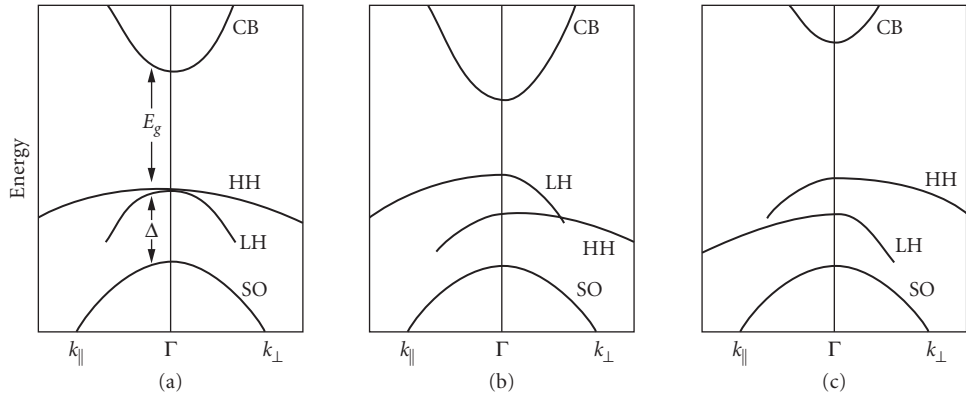


FIGURE 10 Schematic diagrams of the band structure of a III-V semiconductor: (a) unstrained; (b) under tensile strain; and (c) compressively strained. (After Ref. 81.)

heavy hole bands are degenerate at $\mathbf{k}=0(\Gamma)$. Strain lifts this degeneracy and changes the effective masses of the light and heavy holes. In the direction parallel to the substrate the heavy hole band becomes light and the light hole band becomes heavy. Under biaxial tension (Fig. 10b) the bandgap decreases and the “heavy” hole band lies below the “light” hole band. Under biaxial compression the bandgap increases and the “heavy” hole band lies above the “light” hole band.

Figure 10 is a simplification of the true band structure.^{81,82} The bands are not really exactly parabolic, especially the hole bands, and strain increases the nonparabolicity of the hole bands. The details of the band structure can be derived using $\mathbf{k}\cdot\mathbf{p}$ theory.^{86–90}

For a GaAs-based QW, strain can be introduced by adding In to the QW. Since InAs has a larger lattice constant than GaAs, this is a compressively strained QW. In the direction of quantum confinement the highest quantized hole level is the first heavy hole level. This hole level will, therefore, have the largest influence on the density of states and on the optical gain. The effective mass of holes in this level confined in the QW is that parallel to the substrate and is reduced by strain. The reduction of the hole mass within the QW results in a reduced density of hole states [see Eq. (6)].

The reduction in density of hole states is a significant improvement. In order for a semiconductor laser to have optical gain (and lase), $f_c(E, N) - f_v(E, N)$ must be greater than zero [see Eq. (7)]. In an unstrained semiconductor the electrons have a much lighter mass than the heavy holes; the holes therefore have a higher density of states than the electrons. F_c and $f_c(E, N)$ change much more quickly with the injection of carriers than F_v and $f_v(E, N)$. Since approximately equal numbers of holes and electrons are injected into the undoped active layer, reducing the mass of the holes by introducing compressive strain reduces the carrier density required to reach transparency and therefore reduces the threshold current of a semiconductor laser.⁹¹

This theoretical prediction is well supported by experimental results. Strained InGaAs single-QW lasers with record-low threshold current densities of 45 to 65 A/cm² have been demonstrated.^{92–95} These very high-quality strained QW lasers typically have lasing wavelengths from 0.98 to 1.02 μm , QW widths of 60 to 70 \AA , and In concentrations of 20 to 25 percent. InGaAs QWs with wavelengths as long as 1.1 μm have been successfully fabricated,^{96,97} but staying below the critical thickness of the InGaAs layers becomes a problem since the wavelength is increased by increasing the In concentration. (With higher In concentration the amount of strain is increased and the critical thickness is reduced.)

Strained InGaAs QWs have another advantage over GaAs QWs. Strained QW lasers are more reliable than GaAs lasers, i.e., they have longer lifetimes. Even at high temperatures (70 to 100°C), they are very reliable.^{75,76} The reasons for this are not well understood, but it has been suggested that the strain inhibits the growth of defects in the active region.^{98–100} Improving the reliability of GaAs-based lasers is of great practical significance since GaAs lasers are generally less reliable than InP-based lasers.^{4,101,102}

Up to this point our discussion of QW lasers has been limited to GaAs-based QW lasers. QW lasers can also be fabricated in other material systems. GaInP/AlGaInP visible lasers have been improved significantly with the use of a single strained QW active region.^{103–105} These are also compressively strained QWs formed by adding excess In to the active region. This is a much less developed material system than GaAs, so recent results such as 215 A/cm² for a single strained Ga_{0.43}In_{0.57}P QW¹⁰³ are very impressive.

Long-Wavelength (1.3 and 1.55 μm) Quantum Well Lasers

Long-wavelength (InGaAsP/InP) QWs generally do not perform as well as GaAs-based QWs; however, with the advent of strained QW lasers significant progress has been made. Narrow bandgap lasers are believed to be significantly affected by nonradiative recombination processes such as Auger recombination^{4,106–108} and intervalence band absorption.¹⁰⁹ In Auger recombination (illustrated in Fig. 11) the energy from the recombination of an electron and a hole is transferred to another carrier (either an electron or a hole). This newly created carrier relaxes by emitting a phonon; therefore, no photons are created. In intervalence band absorption (IVBA) a photon is emitted, but is reabsorbed to excite a hole from the split-off band to the heavy hole band. These processes reduce the performance of long-wavelength QW lasers enough to make an MQW a lower threshold device than a single QW. As illustrated by Fig. 9, this means that the threshold gain is above the point where the gain versus current density curve of a single QW crosses that of an MQW. Good threshold current density results for lasers operating at 1.5 μm are 750 A/cm² for a single QW and 450 A/cm² for an MQW.¹¹⁰

Long-wavelength QW lasers can be improved by use of a strained QW. For these narrow bandgap lasers strain has the additional benefits of suppressing Auger recombination^{111,112} and intervalence band absorption.¹¹¹ Several groups have demonstrated excellent results with compressively strained InGaAsP/InP QW lasers.^{113–115} Compressively strained single QW lasers operating at 1.5 μm have been demonstrated with threshold current densities as low as 160 A/cm².¹¹⁵

Surprisingly, tensile strained InGaAsP/InP QW lasers also show improved characteristics.^{115–117} Tensile strained QW lasers operating at 1.5 μm have been fabricated with threshold current densities as low as 197 A/cm².¹⁶ These results had not been expected (although some benefit could be expected through suppression of Auger recombination), but have since been explained in terms of TM-mode lasing^{118,119} (normally, semiconductor lasers lase in the TE mode) and suppression of spontaneous emission.¹¹⁸

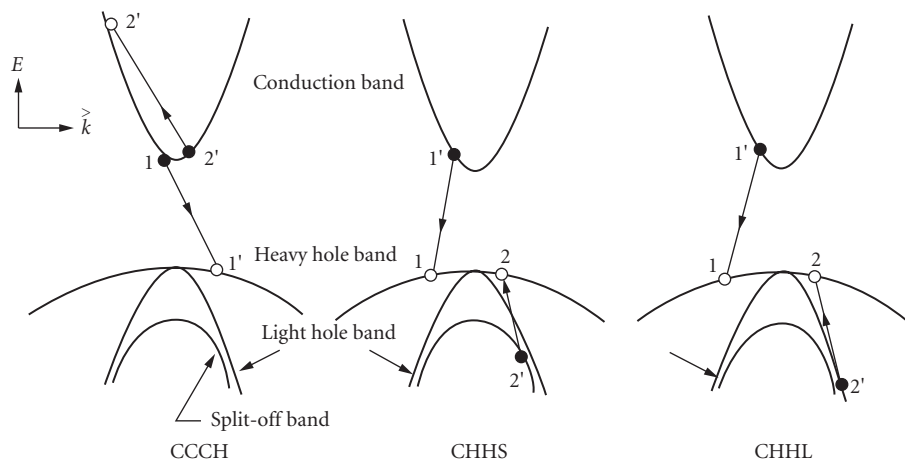


FIGURE 11 Schematic diagrams of band-to-band Auger recombination processes. (After Ref. 4.)

Long-wavelength semiconductor lasers are more sensitive to increases in operating temperature than GaAs-based semiconductor lasers. This temperature sensitivity has been attributed to the strong temperature dependence of Auger recombination^{106–108} and intervalence-band absorption.¹⁰⁹ The use of a strained QW should therefore improve the high-temperature operation of these lasers. This is in fact the case the best results are reported for tensile strained QW lasers with continuous operation at 140°C.^{115,117}

In summary, the use of QW active regions has significantly improved the performance of semiconductor lasers. In this section we have emphasized the dramatic reductions in threshold current density. Improvements have also been realized in quantum efficiency,¹¹⁹ high-temperature operation,^{74–76,115,117} modulation speed (discussed later in this chapter), and spectral linewidth (discussed later). We have limited our discussion to quantum wells. It is also possible to provide quantum confinement in two directions, which is called a *quantum wire*, or three directions, which is called a *quantum dot* or *quantum box*. It is much more difficult to fabricate a quantum wire than a QW, but quantum wire lasers have been successfully demonstrated.¹²¹ For a review of these novel structures we refer the reader to Ref. 122.

19.7 HIGH-POWER SEMICONDUCTOR LASERS

There are several useful methods for stabilizing the lateral modes of an injection laser.^{123–129} In this section, we will discuss techniques for the achievement of high-power operation in a single spatial and spectral mode. There are several physical mechanisms that limit the output power of the injection laser: spatial hole-burning effects lead to multispatial mode operation and are intimately related to multispectral mode operation, temperature increases in the active layer will eventually cause the output power to reach a maximum, and catastrophic facet damage will limit the ultimate power of the laser diode (GaAlAs/GaAs). Thus, the high-power laser designer must optimize these three physical mechanisms to achieve maximum power. In this section, we discuss the design criteria for optimizing the laser power.

High-Power Mode-Stabilized Lasers with Reduced Facet Intensity

One of the most significant concerns for achieving high-power operation and high reliability is to reduce the facet intensity while, at the same time, providing a method for stabilizing the laser lateral mode. Over the years, researchers have developed four approaches for performing this task: (1) increasing the lasing spot size, both perpendicular to and in the plane of the junction, and at the same time introducing a mechanism for providing lateral mode-dependent absorption loss to discriminate against higher-order modes; (2) modifying the facet reflectivities by providing a combination of high-reflectivity and low-reflectivity dielectric coatings; (3) eliminating or reducing the facet absorption by using structures with nonabsorbing mirrors (NAM); (4) using laser arrays and unstable resonator configurations to increase the mode volume. Techniques 1 and 2 are the commonly used techniques and will be further discussed in this section. Techniques 3 and 4 (laser arrays) will be discussed shortly.

Given the proper heat sinking, in order to increase the output power of a semiconductor GaAlAs/GaAs laser, we must increase the size of the beam and thus reduce the power density at the facets for a given power level. The first step in increasing the spot size involves the transverse direction (perpendicular to the junction). There are two approaches for accomplishing this, with the constraint of keeping threshold current low: (1) thinning the active layer in a conventional double-heterostructure (DH) laser (Fig. 12a) below 1000 Å and (2) creating a large optical cavity structure (Fig. 12b).

Thinning the active layer from a conventional value of 0.2 to 0.03 μm causes the transverse-mode spot size to triple for a constant index of refraction step, Δn_r .¹³⁰ The catastrophic power level is proportional to the effective beam width in the transverse direction, d_{eff} , the asymmetric large optical

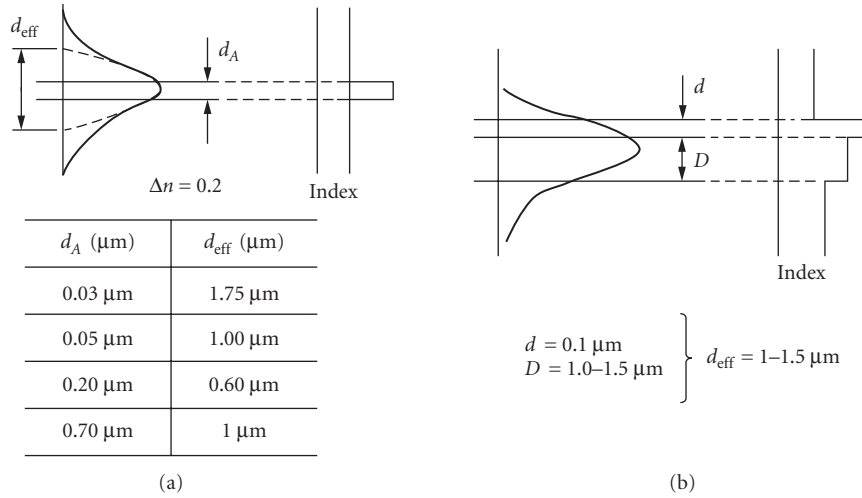


FIGURE 12 Schematic diagrams of the two most commonly used heterostructure configurations for fabricating high-power laser diodes: (a) DH structure and (b) layer large-optical cavity structure. The d_{eff} calculations are after Botez.¹³⁰

cavity (A-LOC) concept^{3,131} involves the epitaxial growth of an additional cladding layer referred to as the *guide layer* (d_G), with index of refraction intermediate between the n -AlGaAs cladding layer and the active layer. By using a relatively small index of refraction step ($\Delta n_r = 0.1$) versus 0.20 to 0.30 for DH lasers, it is possible to force the optical mode to propagate with most of its energy in the guide layer. The effective beam width for the A-LOC can be approximately expressed as

$$d_{\text{eff}} = d_A + d_G \quad (15)$$

where d_A is the active layer thickness. Mode spot sizes in the transverse direction of approximately 1.5 μm can be achieved.

Important Commercial High-Power Diode Lasers

In the last few years, several important high-power laser geometries have either become commercially available or have demonstrated impressive laboratory results. Table 1 summarizes the characteristics of the more important structures. It is evident that the structures that emit the highest cw (continuous wave) power (>100 mW), QW ridge waveguide (QWR), twin-ridge structure (TRS), buried TRS (BTRS), current-confined constricted double-heterostructure large optical cavity (CC-CDH-LOC), and buried V-groove-substrate inner stripe (BVSIS)) have several common features: (1) large spot size (CDH-LOC, TRS, BTRS, QW ridge), (2) low threshold current and high quantum efficiency, and (3) a combination of low- and high-reflectivity coatings. All the lasers with the highest powers, except for the CDH-LOC, use the thin active laser design. A recent trend has been the use of quantum well-active layers.

Figure 13 contains schematic diagrams for five of the more common DH laser designs for high cw power operation, and Fig. 14 shows plots of output power versus current for various important geometries listed in Table 1.

The CC-CDH-LOC device with improved current confinement¹³⁶ (Fig. 13a) is fabricated by one-step liquid-phase epitaxy (LPE) above a mesa separating two dovetail channels. Current confinement is provided by a deep zinc diffusion and a narrow oxide stripe geometry. The final cross-sectional

TABLE 1 Summary of Mode-Stabilized High-Power Laser Characteristics (GaAlAs/GaAs)*†

Manufacturer [Reference]	Geometry	Construction	Rated Power (mW)	Max. cw Power (mW)	Spectral Qual (cw)	Spatial Qual (cw)	I _{th} (mA)	Slope EEF (mW/mA)	Far Field
General Optonics [132]	CNS-LOC	Two-step LPE	—	60	SLM (50)	SSM (50)	50	0.67	12° × 26°
Hitachi [127]	CSP	One-step LPE (TA)	30	100	SLM (40)	SSM (40)	75	0.5	(10–12)° × 27°
MATS. [133]	TRS	One-step LPE	25	115	SLM (50)	SSM (80)	90	0.43	6° × 20°
MATS. [134]	BTRS	Two-step LPE	40	200	SLM (50)	SSM (100)	50	0.8	6° × 16°
NEC [135]	BCM	Two-step LPE	—	80	SLM (80)	SSM (80)	40	0.78	7° × 20°
RCA [136]	CC-CDH	One-step LPE	—	165	SLM (50)	SSM (50)	50	0.77	6° × 30°
RCA [137]	CSP	One-step LPE	—	190	SLM (70)	SSM (70)	50	—	6.5° × 30°
Sharp [138]	V/SIS	Two-step LPE	30	100	SLM (50)	SSM (50)	50	0.74	12° × 25°
Sharp [139]	BV/SIS	Two-step LPE	—	100	—	SSM (70)	50	0.80	12° × 25°
HP [140]	TCSM	One-step MOCVD	—	65	SLM (65)	SSM (40)	60	0.4	—
TRW [141]	ICSP	Two-step MOCVD (AH/HR)	—	100	SLM (30)	150 (50% duty cycle)	75	0.86	(8–11)° × 35°
Ortel [142]	BH/LOC (NAM)	Two-step LPE (AH/HR)	30	90	—	90	30–50	0.85	—
Spectra Diode [143]	QWR	MOCVD	—	500	SLM(100)	SSM(180)	16	1.3	8° × 22°
BN(STC) [144]	QWR	MOCVD	—	300	SLM(150)	SSM(175)	—	0.8	—

*BH Buried heterostructure

†BTRS Buried TRS

BV/SIS Buried VSIS

CC-CDH Current-confined constricted double heterostructure

CNS Channeled narrow planar

CSP Channeled substrate planar

ICSP Inverted CSP

LOC Large optical cavity

NAM Nonabsorbing mirror

QWR Quantum well ridge

SLM Single longitudinal mode

SSM Single spatial mode

TCSM Twin-channel-substrate mesa

TRS Twin-ridge substrate

V/SIS V-groove-substrate inner stripe

*Approaches for achieving high-power GaAlAs lasers:

- Thin active (TA) or A-LOC layer to decrease facet power density
- Tight current confinement to produce high current utilization
- Combination of low-/high-reflectivity facet coatings (AR/HR) to produced high differential efficiency and lower facet intensity
- Quantum well design with long cavity

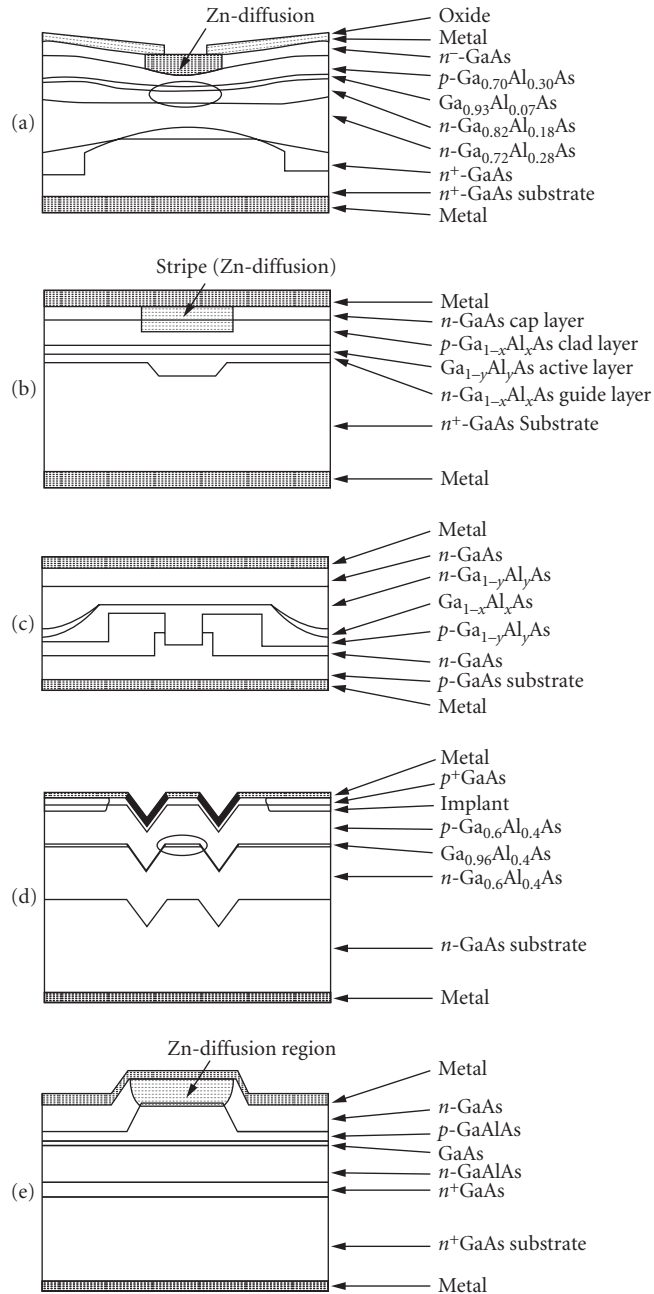


FIGURE 13 Geometries for several important high-power diode lasers. (a) Constricted double-heterostructure large-optical cavity laser (CDH-LOC);¹³⁶ (b) channel substrate planar laser (CSP);¹²⁷ (c) broad-area twin-ridge structure (BTRS);¹³⁴ (d) twin-channel substrate mesa (TCSM); and (e) inverted CSP.

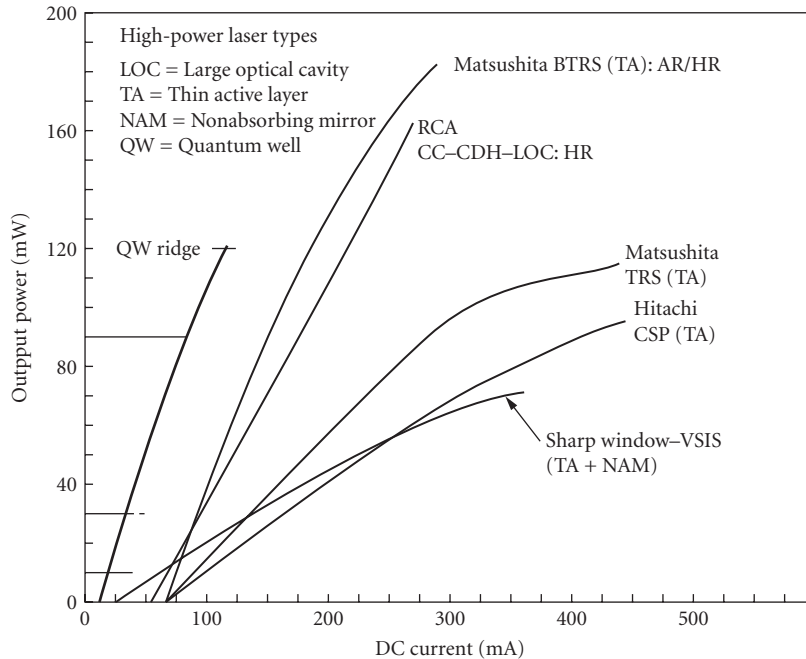


FIGURE 14 Plots showing output power versus cw current for the major high-power laser diodes. The maximum power in a single spatial mode is in the range of 100 to 150 mW and total cw power can approach 200 mW.

geometry of the device is very dependent on the exact substrate orientation and LPE growth conditions.¹⁴⁵ By properly choosing these conditions, it is possible to grow a convex-lens-shaped active layer ($\text{Al}_{0.07}\text{Ga}_{0.93}\text{As}$) on top of a concave-lens-shaped guide layer ($\text{Al}_{0.21}\text{Ga}_{0.79}\text{As}$). The combination of the two leads to a structure with antiwaveguiding properties and a large spot size. Discrimination against higher-order modes is provided by a leaky-mode waveguide. The cw threshold current is in the range of 50 to 70 mA. Single-mode operation has been obtained to 60 mW under 50 percent duty cycle, and the maximum cw power from the device is 165 mW. The power conversion efficiency at this power level is 35 percent considering only the front facet.

The channeled substrate planar (CSP) laser¹²⁷ (Fig. 13b) is fabricated by one-step LPE above a substrate channel. The current stripe is purposely made larger than the channel to ensure uniform current flow across the channel. However, this leads to some waste of current and thus a lower differential efficiency than other similar high-power laser structures (BVSIS, BTRS). Lateral mode control is very effectively obtained by the large difference in the absorption coefficient α between the center and edges of the channel and by changes in the index of refraction that result from changes in the geometry. By proper control of the active and n -cladding layer thicknesses, it is possible to obtain $\Delta\alpha \cong 1000 \text{ cm}^{-1}$ (see Ref. 146) and $\Delta n_s \cong 10^{-2}$. Threshold currents are in the range of 55 to 70 mA. The transverse far field is relatively narrow due to the very thin active layer. Researchers from RCA have obtained power levels in excess of 150 mW (cw) with a CSP-type laser.¹³⁷ A detailed study of the CSP laser has been presented by Lee et al., in a recent publication.¹⁴⁷

Matsushita¹³³ has also developed a CSP-like structure called the twin-ridge structure (TRS) that uses a 400-Å active layer thickness (Fig. 13c). The structure has demonstrated fundamental-mode cw power to 200 mW and single-longitudinal-mode cw power to 100 mW. The maximum available power for the TRS laser is 115 mW, and threshold currents are in the range of 80 to 120 mA. It appears that even though their geometry is similar to that of the CSP, lasers with ultrathin and

planar-active layers have been fabricated. It should be further pointed out that one of the keys to achieving ultrahigh power from CSP-like structures is the achievement of ultrathin ($<1000 \text{ \AA}$) active layers that are highly uniform in thickness. Small nonuniformities in the active layer thickness lead to a larger Δn_r difference, and thus a smaller lateral spot size, which will lead to lower power levels and reduced lateral mode stability.

Metalorganic chemical vapor deposition (MOCVD, discussed earlier) has been used to fabricate lasers with higher layer uniformity, which leads to a reduced spectral width and more uniform threshold characteristics. Several MOCVD laser structures with demonstrated high-power capability are schematically shown in Fig. 13, and their characteristics are summarized in Table 1. Figure 13*d* shows the twin-channel substrate mesa (TCSM) laser.¹⁴⁰ The fabrication consists of growing a DH laser structure over a chemically etched twin-channel structure using MOCVD. Optical guiding is provided by the curvature of the active layer. The TCSM laser has achieved cw powers of 40 mW in a single spatial mode and 65 mW in a single longitudinal mode. The inverted channel substrate planar (ICSP) laser¹³⁵ is schematically shown in Fig. 13*e*. This structure is one MOCVD version of the very successful CSP structure (Fig. 13*b*).¹⁴¹ The ICSP laser has achieved powers in excess of 150 mW (50 percent duty cycle) in a single spatial mode and a 100-mW (cw) catastrophic power level.

More recently, quantum well lasers using the separate carrier and optical confinement (see previous section, “Quantum Well Lasers”) and ridge waveguide geometries have been used for producing power levels in excess of 150 mW (cw) in a single spatial mode.^{143,144} The QW ridge resembles a standard RWG (Fig. 4*e*), but with a QW active region. Such laser structures have low threshold current density and low internal absorption losses, thus permitting higher-power operation.

Future Directions for High-Power Lasers

Nonabsorbing Mirror Technology The catastrophic facet damage is the ultimate limit to the power from a semiconductor laser. In order to prevent catastrophic damage, one has to create a region of higher-energy bandgap and low surface recombination at the laser facets. Thus, the concept of a laser with a nonabsorbing mirror (NAM) was developed. The first NAM structure was demonstrated by Yonezu et al.¹⁴⁸ by selectively diffusing zinc along the length of the stripe, except near the facets. This created a bandgap difference between the facet and bulk regions and permitted a three- to fourfold increase in the cw facet damage threshold and a four- to fivefold increase in pulse power operation.¹⁴⁹ More recent structures have involved several steps of liquid-phase epitaxy.^{150,151}

The incorporation of the NAM structure is strongly device dependent. For example, in the diffused device structures, such as deep-diffused stripe (DDS)¹⁴⁸ and transverse junction stripe (TJS) lasers, NAM structures have been formed by selective diffusion of zinc in the cavity direction.¹⁴⁹ The *n*-type region will have a wider bandgap than the diffused region, and thus there will be little absorption near the facets. However, most index-guided structures require an additional growth step for forming the NAM region.^{150,151} The NAM structures in the past have suffered from several problems: (1) Due to their complex fabrication, they tend to have low yields. Furthermore, cw operation has been difficult to obtain. (2) Cleaving must be carefully controlled for NAM structures having no lateral confinement, in order to avoid excessive radiation losses in the NAM region. The NAM length is a function of the spot size. (3) The effect of the NAM structure on lateral mode control has not been documented, but could lead to excessive scattering and a rough far-field pattern.

It is now becoming more clear that the use of a NAM structure will be required for the reliable operation of high-power GaAlAs/GaAs laser diodes. Experimental results¹⁵² appear to indicate that laser structures without a NAM region show a decrease in the catastrophic power level as the device degrades. However, most of the approaches currently being implemented require elaborate processing steps. A potentially more fundamental approach would involve the deposition of a coating that would reduce the surface recombination velocity and thus enhance the catastrophic intensity level.^{153,154} Such coatings have been recently used by researchers from Sharp and the University of Florida to increase the uncoated facet catastrophic power level by a factor of 2.^{155,156}

Recently, the use of NAM technology has been appearing in commercial products. The crank transverse junction stripe (TJS) laser (a TJS laser with NAM) can operate reliably at an output

power of 15 mW (cw), while the TJS laser without the NAM can operate only at 3 mW (cw).¹⁴⁷ The Ortel Corporation has developed a buried heterostructure (BH) laser with significantly improved output power characteristics compared to conventional BH lasers.¹⁴² The NAM BH laser is rated at 30 mW (cw)¹⁴² compared to 3 to 10 mW for the conventional BH/LOC device.

Last, the use of alloy disordering, whereby the bandgap of a quantum well laser can be increased by diffusion of various types of impurities (for example, Zn and Si),¹⁵⁷ can lead to a very effective technique for the fabrication of a NAM structure. Such structures have produced an enhancement of the maximum pulsed power by a factor of 3 to 4.

High-Power 1.3/1.48/1.55- μm Lasers Previous sections have discussed high cw power operation from (GaAl)As/GaAs laser devices. In the past several years there have been reports of the increasing power levels achieved with GaInAsP/InP lasers operating at $\lambda = 1.3 \mu\text{m}$. The physical mechanisms limiting high-power operation in this material system are quite different than those for GaAlAs/GaAs lasers. The surface recombination at the laser facets is significantly lower than in GaAlAs/GaAs, and thus catastrophic damage has not been observed. Maximum output power is limited by either heating or carrier leakage effects. With the advent of structures having low threshold current density and high quantum efficiency, it was just a matter of time before high-power results would become available. Furthermore, since facet damage is not a problem, the only real need for facet coatings is for improving the output power from one facet and sealing the device for improved reliability.

In Fig. 15, we schematically show the two most common long-wavelength laser structures that have demonstrated high cw power operation. In Fig. 15a, the double-channel planar buried heterostructure (DC-PBH) is systematically shown.^{158,159} The structure requires a two-step LPE growth process. The first step is the growth of the first and top cladding layers in addition to the active layer. This is followed by the etching of the structure, which is followed by a regrowth to form the blocking and contact layers. LPE growth of this material system is such that if the mesa region is narrow enough, no growth occurs on top of it during the deposition of the blocking layer, and this occurs for mesa widths of less than $\sim 5 \mu\text{m}$. Low threshold current is achieved due to the narrow mesa geometry and the good carrier and current confinement.

The DC-PBH has proved to be a laser structure with excellent output characteristics and high reliability. NEC has been able to obtain thresholds as low as 10 mA with 70 percent quantum efficiency. Degradation rates of the order of $10^{-6}/\text{h}$ for an output power of 5 mW at a temperature of 70°C have also been obtained. More recently, NEC has obtained 140-mW power in a single spatial mode.¹⁵⁸ Lasers at 50 mW and 25°C have been placed on lifetest and show relatively low degradation rates after several hundred hours. Degradation rates at 20 and 30 mW (50°C) are $1.3 \times 10^{-5}/\text{h}$ and $2.22 \times 10^{-5}/\text{h}$, respectively. TRW has also worked with DC-PBH/PBC-(planar buried crescent) type laser diodes and has obtained 100 mW(cw).¹⁵⁹ A summary of the various high-power $\lambda = 1.3\text{-}\mu\text{m}$ laser diode structures and characteristics is given in Table 2.

The other structure that has demonstrated high cw power is the buried crescent laser first investigated by Mitsubishi (Fig. 15b).¹⁶⁰ The structure is grown using a two-step LPE process and a p substrate. The final structure resembles a channel laser with an active layer that tapers to zero near the edges of the channel. The tapering provides good carrier and optical confinement. Researchers from Oki with a structure similar to the Mitsubishi structure have demonstrated maximum power levels of 200 and 140 mW in a single spatial mode.¹⁶² Lifetests¹⁶³ on these lasers demonstrated a mean time to failure of $\sim 7 \times 10^5/\text{h}$ (at 20°C) at 75 percent of the maximum cw output power (maximum = 25 to 85 mW). These results appear to indicate that $1.3\text{-}\mu\text{m}$ lasers are reliable for high-power applications.

A more recent development has been the use of multiquantum well (MQW) high-power lasers in the 1.5 to $1.55\text{-}\mu\text{m}$ wavelength band. The use of an MQW ridge waveguide structure has produced power levels of ~ 170 mW (cw).¹⁶⁴ The MQW structure consists of five wells of InGaAs, 60 Å thick, separated by four GaInAsP barriers of 100-Å thickness. The two thicker, outermost barriers of GaInAsP provide a separate confinement heterostructure (SCH) waveguide. In addition, buried heterostructure (BH) lasers¹⁶⁵ with power levels in excess of 200 mW have been achieved by incorporating strain into MQW structures. A review article by Henshall describes the state of the art in more detail.¹⁶⁶

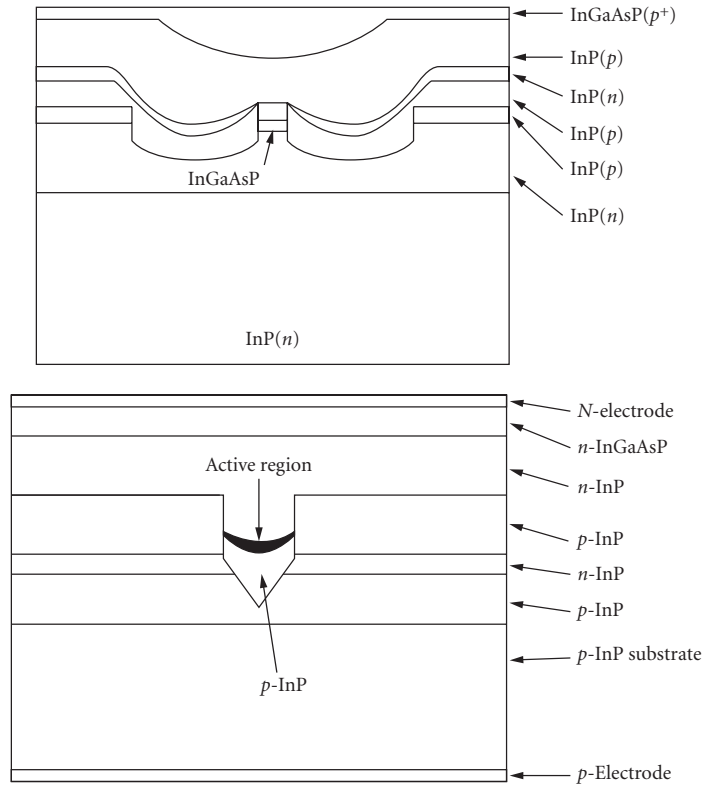


FIGURE 15 Schematic diagrams of the most prominent LPE-based 1.3- μm laser structures: (a) Double-channel planar-buried heterostructure (DC-PBH)^{158,159} and (b) buried crescent.^{160,161}

TABLE 2 Summary of Mode-Stabilized High-Power Laser Characteristics (GaInAsP/InP)^{**†}

Manufacturer [Reference]	Geometry	Construction	Max. Power (mW)	Spatial Quality	I_{th} (mA)
Mitsubishi (160)	PBC	Two-step LPE (p-subst.)	140	SSM (70)	10–30
NEC [157, 158]	DC-PBH	Two-step LPE	140	SSM (140)	10–30
OKI [161, 162]	VIPS	Two-step LPE (p-subst.)	200	SSM (200)	10–30
TRW/EORC [159]	DC-PBH type	Two-step LPE	100	SSM (70)	10–30
TRW/EORC [159]	PBC	Two-step LPE (p-subst.)	107	SSM (78)	10–30
STC [164]	MQW	MOCVD	170	—	—
ATT [165]	MQW BH	MBE	200	—	—

^{*}DC-PBH Double-channel planar buried heterostructure

MQW Multiquantum well

PBC Planar buried crescent

SSM Single spatial mode

[†]Approaches for high-power GaInAsP/InP lasers:

- Tight current confinement to reduce the threshold current
- Facet coatings (reflector/low reflecting front facet)
- Diamond heat sinking
- Long cavity length

TABLE 3 High-Power GaInAs Strained Layer Quantum Well Lasers

Laser Group [Reference]	Ridge Width (μm)	Wavelength (μm)	Threshold Current (mA)	Max. Power in Single Spatial Mode (mW)	Max. cw Power (mW)
JPL [168]	6	0.984–0.989	13	—	24
JPL [168]	3	0.978	8	116	400
NTT [169]	3	0.973–0.983	9	115	500
Spectra Diode [170]	4	0.9–0.91	~20	180	350
Boeing [76]	4	0.98	10–15	150	440

High-Power Strained Quantum Well Lasers Over the last several years, there has been extensive research in the area of strained layer quantum well high-power lasers. As with GaAs QW high-power lasers, the geometry is typically a QW ridge. Table 3 summarizes some of the latest single-spatial-mode high-power results.

Thermal Properties An important parameter in the operation of high-power laser diodes is the optimization of thermal properties of the device. In particular, optimizing the laser geometry for achieving high-power operation is an important design criterion. Arvind et al.¹⁶⁷ used a simple one-dimensional thermal model for estimating the maximum output power as a function of laser geometry (cavity length, active layer thickness substrate type, etc.). The results obtained for GaInAsP/InP narrow stripe PH lasers were as follows:

- Maximum output power is achieved for an optimum active layer thickness in the 0.15- μm region. This result applies only to nonquantum well lasers.
- Significantly higher output powers (25 to 60 percent) are obtained for lasers fabricated on p substrates compared to those on n substrates. The result is based on the lower electrical resistance of the top epitaxial layers in the p substrate compared to n substrate.
- Significantly higher output powers (~60 percent) are obtained for lasers mounted on diamond rather than silicon heat sinks as a result of the higher thermal conductivity of diamond compared to silicon, 22 versus 1.3 W/($^{\circ}\text{C}\cdot\text{cm}$).
- Significantly higher output powers (~100 percent) are obtained for lasers having a length of 700 μm compared to the conventional 300 μm . The higher power results from the reduced threshold current density and thermal resistance for the longer laser devices.

A plot of the calculation and experimental data from Oki¹⁶² is given in Fig. 16. Note that there are no adjustable parameters in the calculation.

An important conclusion from the thermal modeling is that longer cavity semiconductor lasers (700 to 1000 μm) will be able to operate at higher heat sink temperatures when the power level is nominal (~5 mW) compared to shorter cavity devices (~100 to 300 μm). In addition, the reliability of the longer cavity devices is also expected to be better. More recent calculations and experimental results using strained layer lasers have verified this.⁹⁵

Semiconductor Laser Arrays

One of the most common methods used for increasing the power from a semiconductor laser is to increase the width of the emitting region. However, as the width is increased, the occurrence of multilateral modes, filaments, and lateral-mode instabilities becomes more significant. A far-field pattern is produced that is not diffraction-limited and has reduced brightness. The most practical method to overcome this problem is to use a monolithic array of phase-locked semiconductor lasers. Such lasers have been used to generate powers in excess of 10 W (cw)¹⁷⁰ and over 200 W¹⁷¹ from a single laser bar.

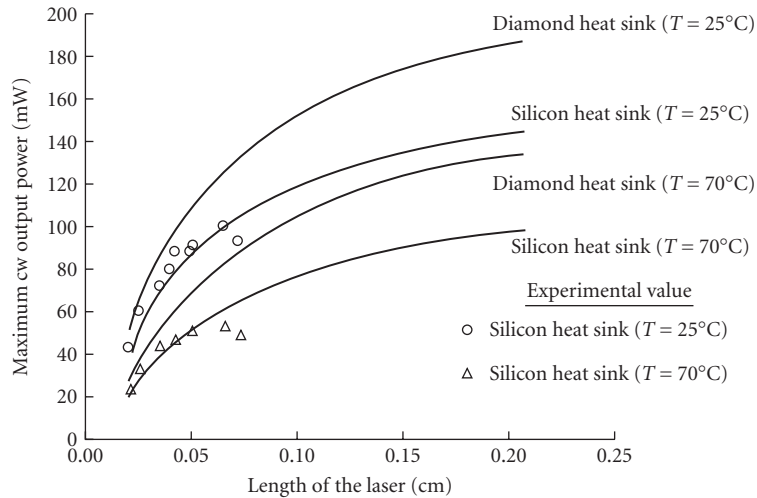


FIGURE 16 Calculation of maximum output power per facet as a function of device length for n - p -substrate-type lasers and different heat sinks.¹⁶⁷ Note the increased power level achieved for longer lasers and p -type substrates.

It was not until 1978 that Scifres and coworkers¹⁷² first reported on the phase-locked operation of a monolithic array consisting of five closely coupled proton-bombarded lasers. The original coupling scheme involved branched waveguides, but this was quickly abandoned in favor of evanescent field coupling by placing the individual elements of the array in close proximity (Fig. 17a). More recently, arrays of index-guided lasers¹⁷³ have been fabricated; one example is shown in Fig. 17b.

Recent emphasis has been on achieving higher cw power and controlling the output far-field distribution. Some of the more significant events in the development of practical semiconductor laser arrays are summarized in Table 4.

The subject of array-mode stability has become of great interest. In a series of significant papers, Butler et al.¹⁹⁰ and Kapon et al.¹⁹¹ recognized that to a first approximation, an array can be modeled as a system of n weakly coupled waveguides. The results indicate that the general solution for the field amplitudes will consist of a superposition of these array modes. The analytic results permitted, for the first time, a simple explanation for the observed far-field patterns and provided a means for designing device structures that would operate in the fundamental array mode (i.e., all elements in phase). This particular mode will provide the greatest brightness.

Many techniques have been used for improving array-mode selection^{179,188,192,197} and thus achieving a well-controlled spatial mode. Two of the more successful earlier techniques involve (1) incorporation of optical gain in the interelement regions (gain coupling of the laser array^{179-184,193}) and (2) use of interferometric techniques that involves Y -coupled junctions.^{186,188,192}

The gain-coupled arrays achieve mode selectivity by introducing optical gain in the interelement regions and thus increasing the gain of the fundamental array mode since this mode has a significant portion of its energy in the interelement regions. The first demonstration of this approach was the twin-channel laser (TCL) developed by researchers from TRW;^{180,181} since then, there have been other demonstrations.^{182,184,193}

The theoretical foundations of the Y -coupled junction were first described in a paper by Chen and Wang.¹⁹² Mode-selectivity is accomplished because the in-phase mode adds coherently at each Y junction, while the out-of-phase mode has destructive interference, since the single waveguides after the Y junction can support only the fundamental mode. Similar interferometric and mode-selective techniques have been used in the development of optical modulators.¹⁹⁸

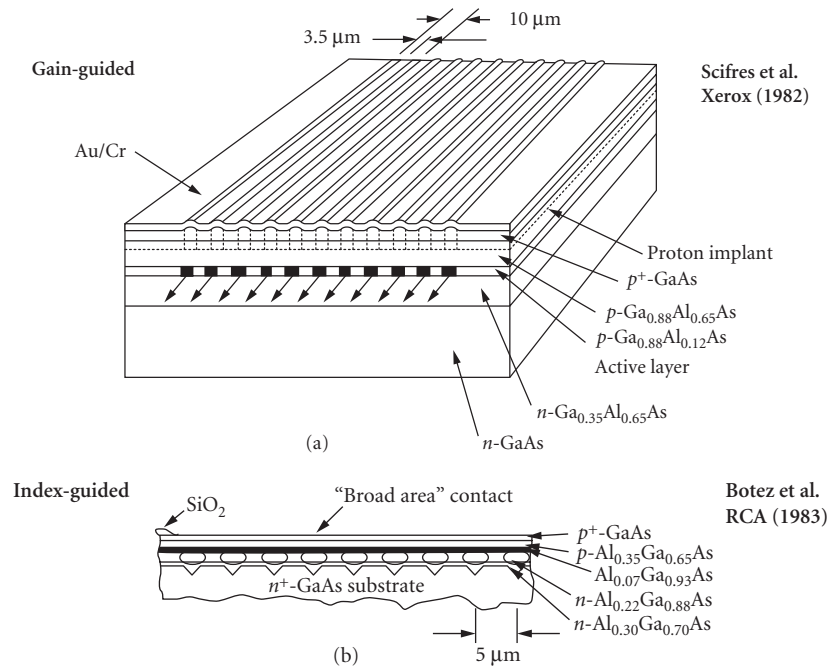


FIGURE 17 Schematic diagrams showing two types of laser array structures: (a) Gain-guided phased array using quantum well-active layers and grown by MOCVD¹⁷⁵ and (b) index-guided phased array using CSP-LOC structures grown by LPE.¹⁷³

TABLE 4 Summary of High-Power Phase-Locked Laser Arrays*

Laser Group [Reference]	No. of Elements	Material System	Type of Array*	Max. Power (mW)	Max. Power (mW)	Far Field
Xerox, 1978 [172]	5	GaAlAs-GaAs	GG	60 (P)	130 (P)	SL (2°)
HP, 1981 [174]	10	GaAlAs-GaAs	IG	1W(P)	1400 (P)	DL
Xerox, 1982 [175]	10	GaAlAs-GaAs	GG	200 (P)	270 (cw)	SL (1°)
Xerox, 1983 [176]	40	GaAlAs-GaAs	GG	800 (P)	2600 (cw)	DL
RCA, 1983 [173]	10	GaAlAs-GaAs	IG	400 (P)	1000 (P)	DL
Siemens, 1984 [177]	40	GaAlAs-GaAs	GG	—	1600 (cw)	DL
Bell Labs, 1984 [178]	10	GaAlAs-GaAs	IG	—	—	DL
TRW, 1984 [179–181]	2	GaAlAs-GaAs	IG	75 (cw)	115 (cw)	SL (4–6°)
UC Berkeley, 1982 [182]	10	GaAlAs-GaAs	IG	—	200	SL (2–7°)
Cal-Tech, 1984 [183]	5	GaAlAs-GaAs	IG	—	—	SL (3°)
Xerox/Spectra Diode, 1985 [184]	10	GaAlAs-GaAs	Offset stripe GG	575 (P)	—	SL (1.9°)
Bell Labs, 1985 [185]	10	InGaAsP-InP	GG	100 (cw)	600	SL (4°)
Sharp, 1985 [186]	2	GaAlAs-GaAs	IG; Y-C	65 (cw)	90 (cw)	SL (4.22°)
Mitsubishi, 1985 [187]	3	GaAlAs-GaAs	IG	100 (cw)	150 (cw)	SL (3.6°)
Xerox/Spectra Diode, 1986 [188]	10	GaAlAs-GaAs	IG(Y-C) stripe GG	200 (cw)	575 (P)	SL (3°)
TRW [189]	10	GaAlAs-GaAs	ROW	380 (cw) 1500 (P)	—	SL (0.7°)

*GG Gain guided
 IG Index guided
 Y-C Y-coupled
 ROW Resonant optical waveguide
 P Pulsed
 DL Double lobe
 SL Single lobe

Finally, the most recent mode control mechanism for laser arrays involves the resonant phase-locking of leaky-mode elements.¹⁸⁹ With a properly optimized geometry, the fundamental array mode has significant mode discrimination analogous to the high discrimination found in single-element leaky-mode devices. Recent results indicate power levels in excess of 360 mW (cw) in the fundamental array mode, and the beam broadens to ~2 times diffraction limited for output powers of ~500 mW (cw).

At the present time, it is not clear which technique will be most useful for achieving stable, fundamental array mode operation. The gain-coupling concept works well for two or three elements.¹⁹⁹ However, array-mode selection described by the difference in gain between the first and second array modes rapidly decreases as the number of array elements increases beyond two or three.¹⁹⁹ The Y junction and leaky-mode approaches do not appear to have the same limitations. The resonant leaky-mode arrays appear to have the most promising performance at high-power levels; however, the structures are complex and thus yield and reliability need to be more fully addressed.

Two-Dimensional, High-Power Laser Arrays

There has been a significant amount of research activity in the past few years in the area of very high-power diode lasers.^{200–205} The activity has been driven by the significant reductions in threshold current density of GaAlAs/GaAs lasers that can be achieved with metalorganic chemical vapor deposition (MOCVD), utilizing a quantum well design. Threshold current densities as low as 200 to 300 A/cm² with external efficiencies exceeding 80 percent have been achieved using GRIN-SCH quantum well lasers.^{200,201} CW powers of ~6 to 9 W have been achieved from single-laser bars. Such power levels correspond to a maximum of 11 W/cm from a single-laser bar. Table 5 lists some of the more recent results on very high-power diode laser arrays.

In order to increase the output power from laser array structures, researchers have investigated the use of a two-dimensional laser array. One particular configuration, referred to as the “rack-stack approach,” is schematically shown in Fig. 18. In essence, the approach involves stacking a linear array of edge emitters into a two-dimensional array. The two-dimensional arrays are fabricated²⁰³ by (1) cleaving linear arrays of laser diodes from a processed wafer, (2) mounting the bars on heat sinks, and (3) stacking the heat sinks into a two-dimensional array.

As shown in Table 5, the main players in this business are McDonnell-Douglas and Spectra Diode Labs. The largest stacked²⁰⁴ two-dimensional array has been manufactured by McDonnell-Douglas and has an active area with five laser bars 8 mm in length. The array was operated with 150- μ s pulses to the limit of the driver²⁰⁴ at pulse repetition rates of 20 to 666 Hz. Approximately 2.5 kW/cm² was obtained at 20 Hz (average power ~300 W) and 0.9 kW/cm² at 666 Hz (average power ~92 W). Higher output powers will be obtained as a result of achieving the ultimate limits in

TABLE 5 Summary of High-Power Laser Array Results*

Laboratory [Reference]	Array Type	Maximum Output Power	Power Efficiency (%)	Slope Efficiency (W/A)	Power Density (W/cm ²)
General Electric [201]	ID array	80 W (200- μ s pulse; 10–100 Hz)	20	0.9	80
McDonnell-Douglas [203, 204]	Broad stripe ($L = 1200 \mu\text{m}$)	6 W (cw) ($W = 300 \mu\text{m}$)	38	0.91	200
	2D: 4 bars, 8 mm	15 W (cw)	15	—	50
	2D: 5 bars, 8 mm	320 W (0.3% DF [†])	—	—	2560
Spectra Diode [170]	ID array	8 W (cw)	—	—	—
	ID array	134 W (150- μ s pulse)	49	1.26	134

*All high-power laser structures are fabricated using MOCVD and quantum well design.

[†]DF = Duty factor.

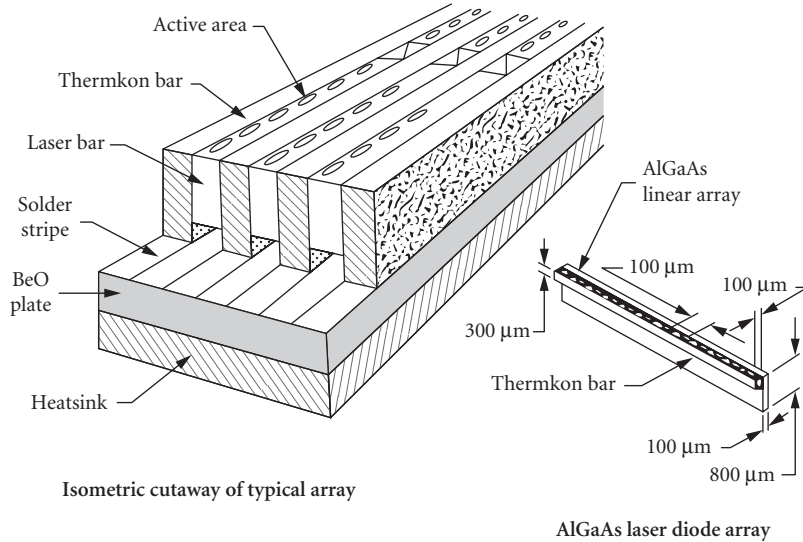


FIGURE 18 Schematic diagram of the two-dimensional rack-stack laser array architecture. (Courtesy of R. Solarz, Lawrence Livermore.)

threshold current density and optical losses in the individual with advances in nonabsorbing mirror and active cooling techniques. Some recent progress has been seen in the latter with the use of etched silicon grooves for fluid flow, which function as radiator elements to remove the heat.

19.8 HIGH-SPEED MODULATION

In many applications semiconductor lasers are modulated in order to carry information. Semiconductor laser dynamics are usually described by the rate equations for the photon and carrier densities:^{3-5,205,206}

$$\frac{dN}{dt} = \frac{I}{edLw} - \frac{cg}{n_r \Gamma} P - \frac{N}{\tau_s} \quad (16)$$

$$\frac{dP}{dt} = \frac{c}{n_r} g P - \frac{P}{\tau_p} + \Gamma \beta \frac{N}{\tau_s} \quad (17)$$

where N is the carrier density, P is the photon density, I is the current, e is the charge of an electron, d is the active layer thickness, L is the laser cavity length, w is the laser stripe width, c is the speed of light, n_r is the refractive index of the active region, g is the threshold modal gain, Γ is the optical confinement factor, τ_s is the carrier lifetime, β is the spontaneous emission factor, and τ_p is the photon lifetime of the cavity. [When these equations are written in terms of total gain instead of modal gain, Γ does not appear in Eq. (16), but multiplies the gain in Eq. (17).]

$$\tau_p = \frac{n_r}{c(\alpha_i + (1/2L)\ln(1/R_F R_R))} \quad (18)$$

where α_i is the internal loss and R_F and R_R are the front- and rear-facet reflectivities, β is the ratio of spontaneous emission power into the lasing mode to the total spontaneous emission rate.²⁰⁷ (Do not confuse β with the other spontaneous emission factor, which is used in linewidth theory and is defined as the ratio of the spontaneous emission power into the lasing mode to the stimulated emission power of the mode.)

When a semiconductor laser is modulated there is some delay before it reaches a steady state. Because it takes time for a carrier population to build up, there will be a time delay τ_d before the final photon density P_{on} is reached (see Fig. 19). Once P_{on} is reached, additional time is required for the carrier and photon populations to come into equilibrium. The output power therefore goes through relaxation oscillations before finally reaching a steady state. This type of oscillation has many parallels in other second-order systems,²⁰⁸ such as the vibration of a damped spring or an RLC circuit.

The frequency of these relaxation oscillations, f_r is called the relaxation, resonance, or corner frequency. By considering small deviations from the steady state where $N = N_{\text{th}} + \Delta N$ and $P = P_{\text{on}} + \Delta P$, we can solve Eqs. (16) and (17) for f_r with the result:^{5,205}

$$f_r \cong \frac{1}{2\pi} \sqrt{\frac{c}{n_r \Gamma} \frac{dg}{dN} \frac{P_{\text{on}}}{\tau_p}} = \frac{1}{2\pi} \sqrt{\frac{c}{n_r} \frac{dg}{dN} \frac{(I - I_{\text{th}})}{edLw}} \quad (19)$$

where dg/dN is the differential modal gain. For bulk double-heterostructure lasers the gain is linearly dependent on the carrier density and $(c/n_r)dg/dN$ is replaced by A , where A is a constant. As discussed earlier in the chapter, the gain-versus-carrier-density relationship of a quantum well is nonlinear, so A is not a constant for a QW laser.

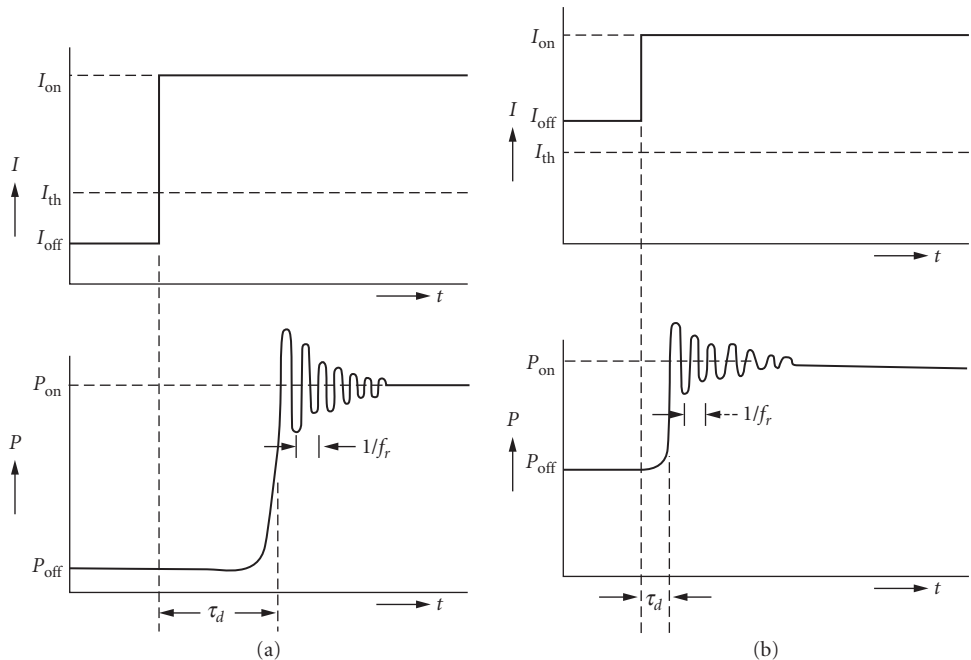


FIGURE 19 Schematic diagrams of the turn on delay and relaxation oscillations for a semiconductor laser (a) prebiased below threshold and (b) prebiased above threshold. (After Ref. 206.)

Let us return our attention to the turn on delay τ_d between a current increase and the beginning of relaxation oscillations as illustrated in Fig. 19. If the initial current I_{off} is below I_{th} , the initial photon density P_{off} can be neglected in Eq. (16). Assuming an exponential increase in carrier density we can derive²⁰⁶

$$\tau_d = \tau_s \ln \left(\frac{I_{\text{on}} - I_{\text{off}}}{I_{\text{on}} - I_{\text{th}}} \right), \quad I_{\text{off}} < I_{\text{th}} < I_{\text{on}} \quad (20)$$

Since τ_s is on the order of several nanoseconds, τ_d is usually very large for semiconductor lasers with I_{off} below I_{th} . For example, with $\tau_s = 4$ ns, $I_{\text{on}} = 20$ mA, $I_{\text{th}} = 10$ mA, and $I_{\text{off}} = 5$ mA, τ_d will be 1.6 ns. If a short current pulse is applied to a laser biased below threshold, τ_d may be so long that the laser barely responds (see Fig. 20). When a current pulse ends, it takes time for the carrier population to decay. If another identical pulse is applied before the carrier population decays fully, it will produce a larger light pulse than the first current pulse did. This phenomenon, which is illustrated in Fig. 20, is called the *pattern effect*.²⁰⁹ The pattern effect is clearly undesirable as it will distort information carried by the laser modulation. The pattern effect can be eliminated by prebiasing the laser at a current sufficient to maintain a carrier population; for most semiconductor lasers this will mean prebiasing at or above threshold. In order to modulate with a prebias below threshold, I_{on} must be much greater than I_{th} . For most lasers this will require an unpractically large I_{on} , but if I_{th} is very low, it may be possible.^{211,212}

Even a semiconductor biased above threshold will have a nonzero τ_d before it reaches its final photon density. Equations (16) and (17) may be solved for τ_d above threshold:²⁰⁶

$$\tau_d = \frac{1}{2\pi f_r} \sqrt{2 \ln \left(\frac{P_{\text{on}}}{P_{\text{off}}} \right)} = \frac{1}{2\pi f_r} \sqrt{2 \ln \left(\frac{I_{\text{on}} - I_{\text{th}}}{I_{\text{off}} - I_{\text{th}}} \right)} \quad I_{\text{th}} < I_{\text{off}} < I_{\text{on}} \quad (21)$$

τ_d will be much shorter for a prebias above threshold. For example if $f_r = 5$ GHz, $I_{\text{on}} = 40$ mA, $I_{\text{off}} = 15$ mA, and $I_{\text{th}} = 10$ mA, $\tau_d = 60$ ps. τ_d will be shortest for large P_{on} and P_{off} , so the shortest time

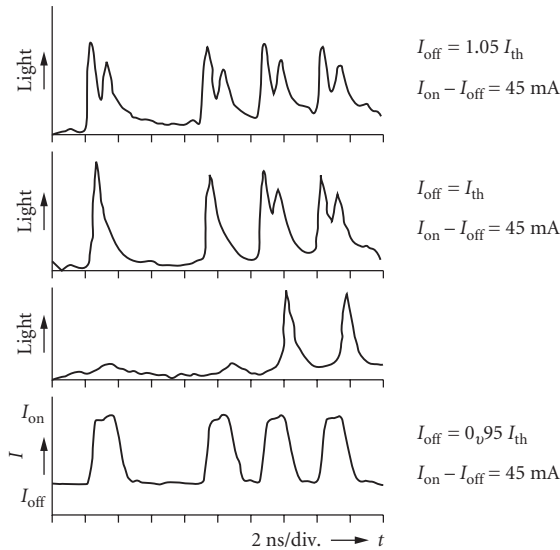


FIGURE 20 An illustration of the pattern effect for a AlGaAs laser diode with a 280 Mbit/s 10111 return to zero pattern. (From Ref. 210.)

delays will be achieved for small-scale modulation at high photon density. For digital applications in which fairly large-scale modulation is required, the maximum modulation speed of a semiconductor laser is to a large extent determined by τ_d .

For very high-speed microwave applications, lasers are prebiased at a current greater than the threshold current and modulated at high frequencies through small amplitudes about the continuous current prebias. The frequency response of a semiconductor laser has the typical shape expected from a second-order system. (For a discussion of the frequency response of a second-order system see Ref. 208.) The laser amplitude response is fairly uniform at frequencies less than the relaxation oscillation frequency. At f_r the response goes through a resonance and then drops off sharply. The relaxation oscillation frequency is therefore the primary intrinsic parameter determining the modulation bandwidth. The actual useful bandwidth is generally considered to be the frequency at which the response of the laser drops by 3 dB. Figure 21 is an example of the frequency response of a semiconductor laser under amplitude modulation.^{213,214} In this example, the maximum 3-dB bandwidth is 16 GHz. If the 3-dB bandwidth is measured in electrical dB, as in our example, it is located at approximately $1.55f_r$. Sometimes the 3-dB frequency is quoted as that at which the optical power is reduced by a factor of 2; this actually corresponds to 6 dB in electrical power and occurs at approximately $1.73f_r$. The 0-dB frequency occurs at approximately $1.41f_r$.²¹³⁻²¹⁵

The description of the relaxation oscillation frequency given here is rather simplistic since it is based on rate equations, which consider only one type of carrier, neglect the spatial dependences of the carrier and photon distributions, and neglect the effects of carrier diffusion and nonlinear gain. In addition, the spontaneous emission term of Eq. (17) was neglected in the derivation. The neglected effects are particularly important when considering damping of the

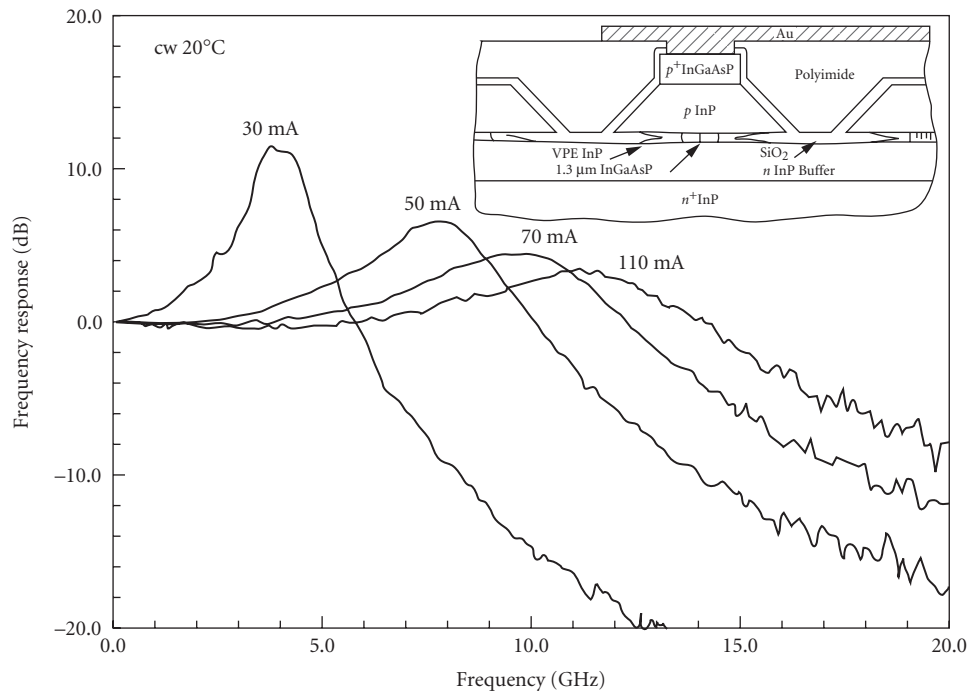


FIGURE 21 The small-signal modulation response of a 1.3- μm InGaAsP-constricted mesa laser for different bias levels. The cavity length is 170 μm and the stripe width is 1 μm . Inset: Schematic diagram of a 1.3- μm InGaAsP-constricted mesa laser. (From Refs. 213 and 214.)

relaxation oscillations.^{204–206} With significant damping²¹⁵ the measured peak frequency f_p will be more accurately determined by

$$f_p^2 = f_r^2 - \frac{f_d^2}{4} \quad (22)$$

where f_d is the damping frequency.

We have also neglected, however, the electrical parasitics of the laser and its operating circuit (bonding wires, etc.). Figure 22 is a simple equivalent circuit, which describes the parasitic elements influencing a semiconductor laser. Here L is the inductance of bond wire, R is the laser resistance including contact resistance, and C is capacitance primarily due to bonding-pad capacitance and capacitance of the current-confining structure of the laser stripe.^{213,214} The 50- Ω resistance is included to represent a 50- Ω drive. The resonant frequency of this circuit is

$$f_0 = \frac{1}{2\pi} \sqrt{\frac{R+50}{LRC}} \quad (23)$$

The circuit is strongly damped so no resonance peak occurs; instead, the response simply drops off.^{205,218} The amplitude response of a modulated semiconductor laser will begin to drop off at frequencies at which the response of the parasitics drops off even if the intrinsic peak frequency of the laser is higher. Therefore, the maximum practical modulation bandwidth may be determined by f_0 instead of f_r . Figure 23 shows the modulation response of a semiconductor laser strongly affected by parasitics.²¹⁹

The most significant parasitic limiting the performance of a semiconductor laser is normally the capacitance.²⁰⁵ In order to achieve high speeds, the laser stripe must be very narrow [see Eq. (18)]; practical narrow stripe lasers are often some form of buried heterostructure (see Fig. 4c). The substrate doping and confinement-layer doping of buried heterostructure form a parallel plate capacitor. In order to reduce the capacitance the laser can be fabricated on a semi-insulating substrate,^{205,218} the confinement layers can be semi-insulating,¹²⁴ the active area of the device can be isolated from the confinement layers by etching trenches on either side of it,^{124,125} or the confinement layers can be replaced by a thick dielectric layer such as polyimide.^{213,214} The inset on Fig. 21 is a schematic diagram of the high-speed laser stripe whose frequency response is shown in Fig. 21.

Assuming that the parasitics have been minimized, consider how a semiconductor laser can be optimized for high-speed operation. As already mentioned, minimizing the stripe width is desirable.

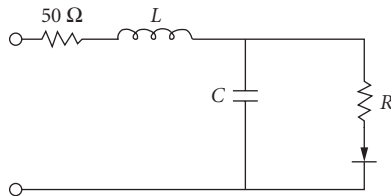


FIGURE 22 Simple equivalent circuit of the parasitics affecting the modulation of a semiconductor laser.

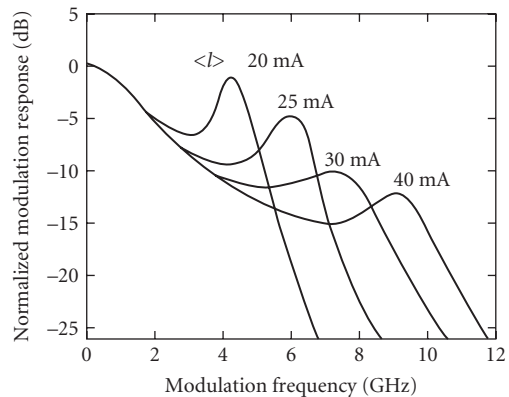


FIGURE 23 The modulation response of a 1.3- μm double-channel, planar-buried heterostructure laser with a cavity length of 80 μm and a threshold current of 18 mA. The effect of parasitics is apparent. (From Ref. 219.)

Figure 21 illustrates that increasing the photon density (or equivalently the current) increases the speed. Of course, there will be a limit as to how much the photon density can be increased; when the photon density of a semiconductor laser is increased, eventually a maximum power is reached at which the laser fails due to catastrophic facet damage. InGaAsP/InP lasers have a higher threshold for catastrophic facet damage than AlGaAs/GaAs lasers; InGaAsP/InP lasers, therefore, tend to have higher bandwidths.²¹⁴

Consideration of Eq. (19) shows that decreasing the cavity length will also increase the speed.^{205,213,214,218,219} (Decreasing the length reduces τ_p .) Increasing dg/dN will also increase the speed. With a bulk active region dg/dN is approximately constant, so use of a short cavity length will not affect it. As illustrated in Fig. 6, however, dg/dN of a QW laser will decrease with increasing threshold gain, and therefore with decreasing cavity length. A single QW laser will, therefore, make a relatively poor high-speed laser. An MQW laser, however, will have higher dg/dN than a single QW. In the past, the highest modulation bandwidths were achieved with InGaAsP/InP bulk active region InGaAsP/InP lasers^{213,214,222,223} with the best results on the order of 24 GHz²²³ for room-temperature CW measurements. The advent of strained MQW lasers has, however, recently resulted in higher bandwidths because strain increases dg/dN .^{224–226} Strained GaAs-based In_{0.3}Ga_{0.7}AsMQW lasers with bandwidths as high as 28 GHz have been demonstrated.²²⁶

So far our discussion has dealt only with the amplitude response to modulation. The phase and lasing wavelength (or optical frequency) are also affected by modulation. Figure 24 is an example of the phase response which accompanies the amplitude response of a semiconductor laser under modulation.

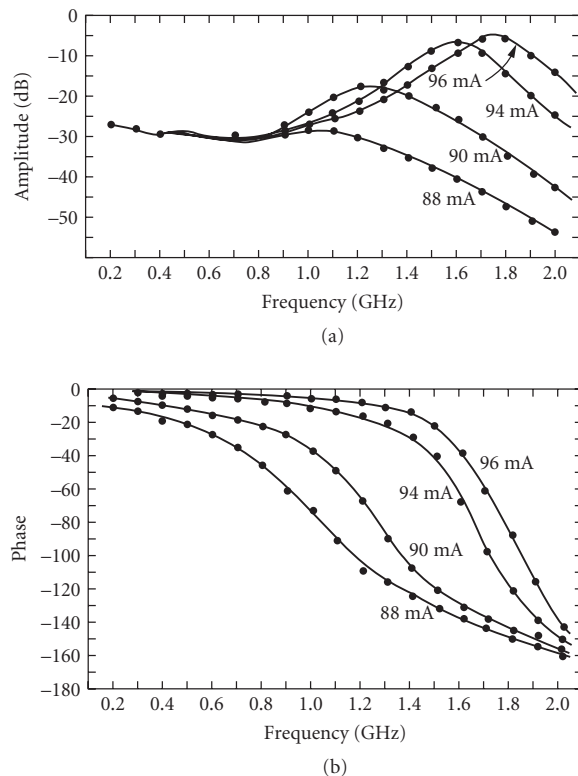


FIGURE 24 The modulation response of a proton stripe laser at various bias currents: (a) amplitude response and (b) phase. The threshold current is approximately 80 mA. (From Ref. 205.)

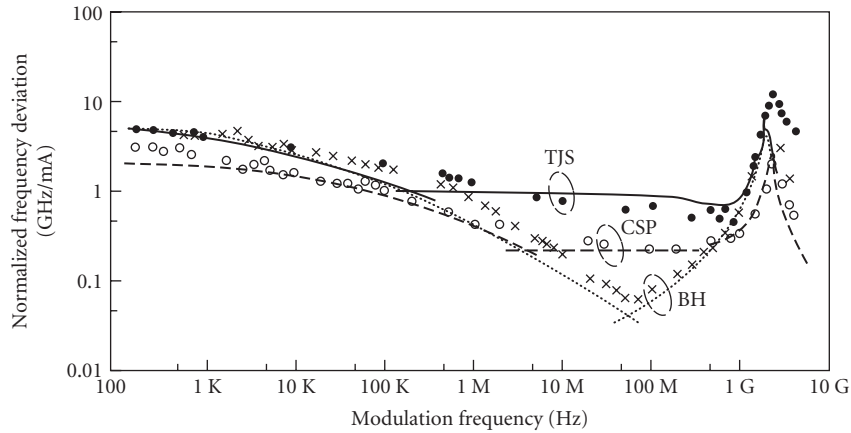


FIGURE 25 The FM response of TJS, BH, and channeled substrate planar (CSP) lasers. All the lasers are biased at 1.2 times threshold. (From Ref. 229.)

Assuming that a semiconductor laser lases in a single longitudinal mode under static conditions, high-frequency modulation can cause it to shift to another mode or to become multimode.²²⁷ The tendency to become multimode increases with the depth of modulation. For many applications single-mode operation under modulation is required. In this case a laser with built-in frequency selectivity such as a distributed feedback laser²²⁸ (DFB) (see “Spectral Properties” following) can be used to maintain single-mode operation. Even with single-mode operation under modulation the linewidth of the lasing mode will be broadened.^{206,227} This broadening which is often called *chirp*, is discussed in more detail in the next section.

While the frequency changes associated with small-scale modulation are generally undesirable in amplitude modulation (AM), they can be utilized for frequency modulation^{206,229,230} (FM). In digital systems, FM is often called frequency shift keying (FSK). FM requires only a very small amplitude modulation, so it normally refers to the effect of modulation on a single mode. In FM modulation, very fine shifts in optical frequency are detected; frequency stabilized lasers such as DFB lasers²³⁰ as well as standard semiconductor lasers will show an FM response. Typically, FM response shows a low frequency decay below f_r and a resonance at f_r ^{206,229} (see Fig. 25).

If the reader requires more in-depth information on high-speed modulation of semiconductor lasers, the book by Petermann²⁰⁶ or the review by Lau and Yariv²⁰⁵ will be particularly helpful. For a recent review of the state of the art see the tutorial by Bowers.²¹⁵ References 3, 4, and 5 also contain chapters on high-speed modulation.

19.9 SPECTRAL PROPERTIES

One of the most important features of a semiconductor laser is its high degree of spectral coherence. There are several aspects to laser coherence. First, the laser must have spatial coherence in the various transverse directions. This is usually accomplished by controlling both the geometry and the lateral-mode geometry using a structure with a built-in index as discussed earlier under “Fabrication and Configurations.” In order to achieve high spectral coherence, the semiconductor laser must operate in a single longitudinal mode. There are four technical approaches for accomplishing this:^{231,232} (1) coupled cavity, (2) frequency selective feedback, (3) injection locking, and (4) geometry control. The various techniques for achieving spectral control are described in Fig. 26.

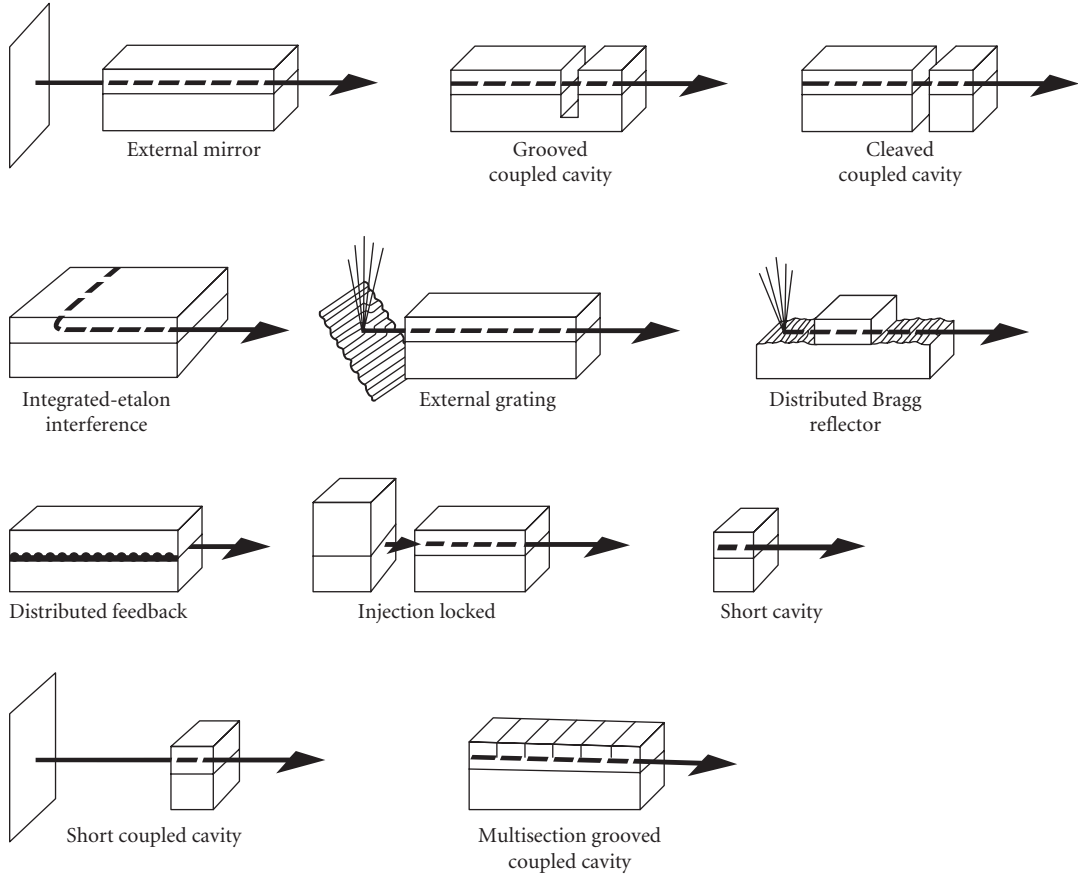


FIGURE 26 Eleven major designs for single-frequency lasers. The three in the top row and the first in the second row are coupled cavity lasers; the next three are frequency-selective-feedback lasers; the next is an injection-locked laser; the last one in the third row is a geometry-controlled laser. The left two are hybrid designs. (From Ref. 232.)

The temporal coherence of the laser is related to the spectral width of the stimulated emission spectrum by

$$L_c = K \frac{\lambda_0^2}{\Delta\lambda_{1/2}} \quad (24)$$

where K is a constant dependent on the distribution of spectral output function, L_c is coherence length, λ_0 is the wavelength of the stimulated emission peak, and $\Delta\lambda_{1/2}$ is the halfwidth of the spectral emission. $K = 1$ for rectangular, $K = 0.32$ for a lorentzian, and $K = 0.66$ for a gaussian.

The spectral linewidth, $\Delta f_{1/2}$, for a single longitudinal mode can be expressed as²³³

$$\Delta f_{1/2} = \frac{n_{sp}}{4\pi\tau_p} \left(\frac{J}{J_{th}} - 1 \right)^{-1} (1 + \alpha^2) \quad (25)$$

where n_{sp} is the spontaneous emission factor, defined as the ratio of spontaneous to stimulated emission in the lasing mode, τ_p is the cavity lifetime, and α is the linewidth enhancement factor.

Typical linewidths for a solitary single-longitudinal-mode laser are in the range of 5 to 20 MHz. Narrower linewidths can be achieved by using some of the techniques described in Fig. 26. More recently, the use of QW lasers, as described earlier, has led to a significant reduction in the linewidth enhancement factor and the corresponding laser linewidths.²³⁴ Typical linewidths in the range of 0.9 to 1.3 MHz have been achieved.

Coupled cavity lasers make up a family of devices, whereby spectral control is achieved by reinforcing certain wavelengths which resonate in several cavities.²³⁵ A typical configuration is shown in Fig. 26, whereby the long cavity is cleaved into two smaller cavities. By properly controlling the length ratios and the gap width, good longitudinal mode discrimination (better than 20-dB side-mode suppression) can be obtained.

Another important technique is the use of an external resonant optical cavity (frequency-selective feedback) as shown in Fig. 27. This technique has been used by researchers at Boeing to achieve extremely narrow linewidth ($\Delta f_{1/2} \sim 1-2$ KHz) single-longitudinal-mode operation.²³⁶

Frequency-selective feedback can also be achieved by using either a distributed feedback (DFB) or distributed Bragg reflector (DBR) laser. As shown in Fig. 26, it differs from other types of lasers in that the feedback is provided by a grating internal to the diode laser. By using a DFB/DBR in combination with a long external cavity, it is possible to achieve linewidths below 1 MHz in a monolithic diode.²³⁷

Injection-locked lasers have also been under investigation at several research labs.²³⁸ In this technique, a low-power, single-frequency laser, which does not have to be a semiconductor laser, is coupled to a single-mode semiconductor laser by injecting the continuous wave emission of a single wavelength of radiation into the laser's cavity.

The last technique for achieving single-longitudinal-mode operation involves the geometry-controlled cavity. Basically, this involves a short cavity 50 μm or less in length, since the longitudinal-mode spacing $\Delta\lambda_L$ in a semiconductor laser is given by¹²⁸

$$\Delta\lambda_L = \frac{\lambda_0^2}{2n_{\text{eff}}L} \tag{26}$$

where n_{eff} is the effective index of refraction and L is the cavity length. Then if $L < 50$ mm, $\Delta\lambda_L > 20 \text{ \AA}$ and the gain available to modes away from the gain maximum falls rapidly, the laser operates in a

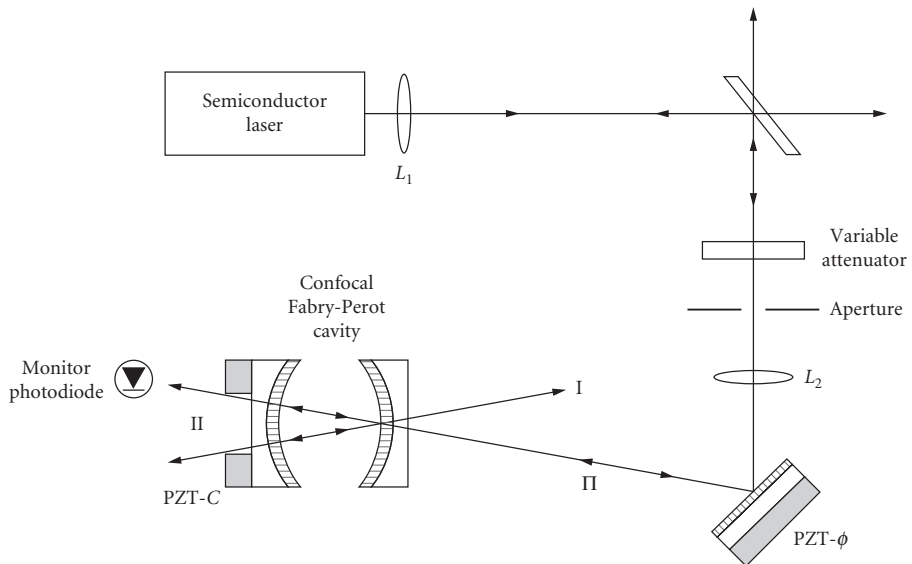


FIGURE 27 Semiconductor laser using resonant optical feedback.

single longitudinal mode. However, the width of the spectral model can still be rather large as dictated by Eq. (25), unless special precautions are made (e.g., ultrahigh mirror reflectivities).

19.10 SURFACE-EMITTING LASERS

Monolithic two-dimensional (2D) laser arrays are key to many important applications such as massive parallel data transfer, interconnect, processing, computing, and high-power, diode-pumped, solid-state lasers. Conventional lasers, as described in previous sections, require a pair of parallel crystalline facets (by cleaving) for delineating the laser cavity, thus limiting laser emission parallel to the junction plane. In this section, we describe laser structures and fabrication techniques which allow light to emit perpendicular to the junction plane, namely, surface-emitting lasers (SEL). SEL structures are compatible with monolithic 2D laser array integration and requirements.

There are three designs for the fabrication of surface-emitting lasers and arrays: (1) in-plane laser with a 45° mirror, (2) in-plane laser with a distributed grating coupler, and (3) vertical cavity laser. The main body of the first two structures is very similar to the conventional cavity design with the axis parallel to the junction plane (in-plane). Light is coupled out from the surface via an integrated mirror or grating coupler. The third structure is an ultrashort cavity (10 μm) “microlaser” requiring no cleaving and compatible with photodiode and integrated-circuit processing techniques. High-density, surface-emitting laser arrays of this type have been demonstrated jointly by AT&T and Bellcore.²³⁹ The following subsections will summarize each of the three structures.

Integrated Laser with a 45° Mirror

The development of this SEL structure requires the wafer processing of two 90° laser mirrors as well as a 45° mirror for deflecting the laser output from the junction plane as shown in Fig. 28. Dry etching techniques such as reactive ion (beam) etching (RIE), chemical-assisted ion beam etching (CAIBE), and ion beam milling are usually used for the fabrication. In combination with a mass-transport process,²⁴⁰ a smooth parabolic sidewall has been demonstrated for the 45° mirror of InGaAsP/InP lasers.

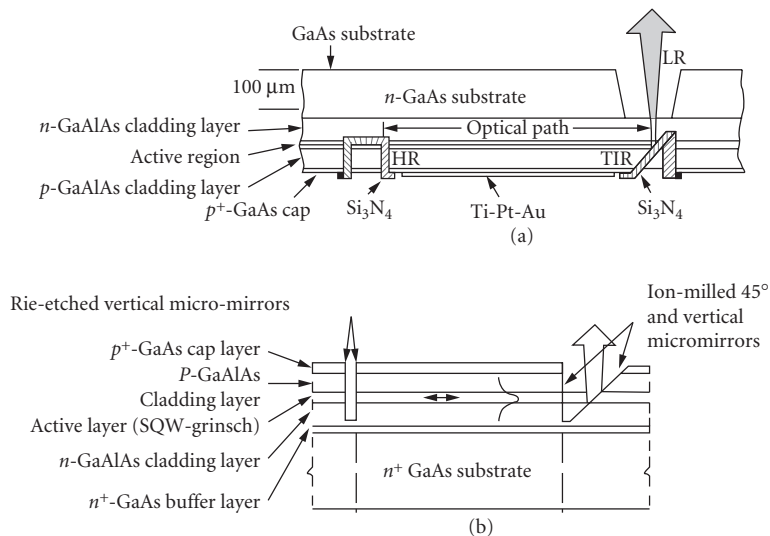


FIGURE 28 Monolithic in-plane-cavity surface-emitting lasers with 45° mirrors: (a) Junction-down and (b) junction-up configurations. (From Ref. 241.)

This approach for an SEL takes advantage of well-established layer structure growth for typical lasers. The laser performance relies on the optical quality and reliability of the facet mirrors formed. The etched-mirror lasers have been improved over the past decade to the stage very comparable with the cleaved lasers. Two-dimensional, high-power (over 1 W) laser arrays have been demonstrated by both TRW²⁴¹ and MIT Lincoln Laboratory.²⁴⁰ These structures would require injection-locking or other external optical techniques in order to achieve coherent phased array operation as mentioned in the high-power laser section.

Distributed Grating Surface-Emitting Lasers

Distributed feedback (DFB) (see under “Spectral Properties,” discussed earlier) and distributed Bragg reflector (DBR) lasers were proposed and demonstrated in the early 1970s. It is well known that for the second-order gratings fabricated in the laser, the first-order diffraction will be normal to the grating surface, as shown in Fig. 29. Since early demonstrations, it has taken over 10 years for both the applications and processing techniques to become mature. Low-threshold, high-reliability DFB lasers with true single-mode characteristics are readily fabricated. The critical issue involved in the fabrication of the laser structure is the fabrication of the gratings with a period on the order of 2000 Å. Holographic interference techniques with an Ar⁺ or He-Cd laser are generally used in many laboratories. The fabrication of large-area gratings with good throughput can be easily achieved by this technique. Another technique involves direct electron-beam writing, which is effective for design iterations.

The development of DBR structure with second-order gratings for surface-emitting lasers did not occur until it was funded by the U.S. Air Force pilot program. This type of laser does not require discrete mirrors for the laser action, so that one could link an array of the lasers with residual in-plane light injection (leaking) across neighboring lasers for coherent operation. A near-diffraction-limited array operation has been demonstrated with this type of SEL. The concept was recently used

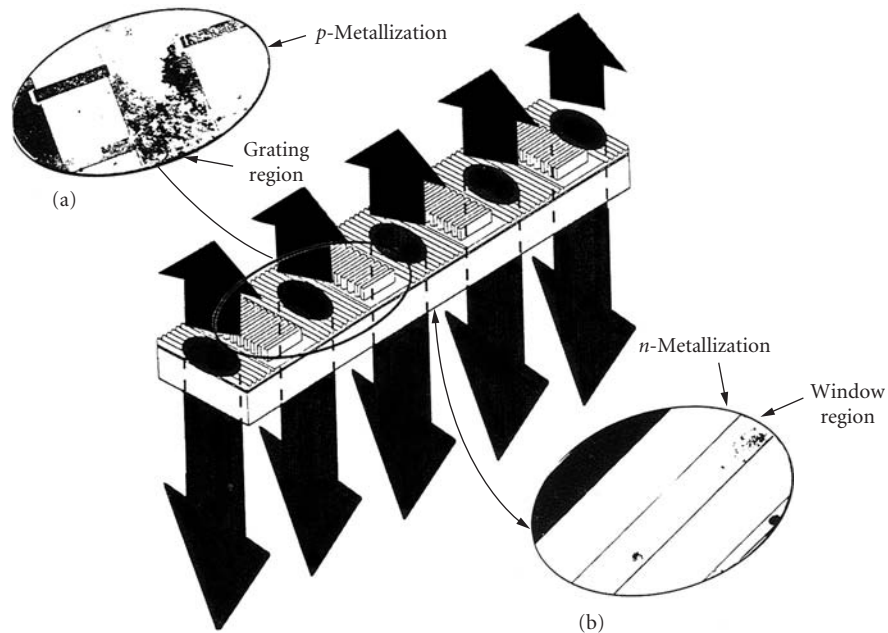


FIGURE 29 Grating surface-emitting laser array. Groups of 10 parallel ridge-guide lasers are laterally coupled in each gain section. (From Ref. 242.)

for a high-power MOPA (master oscillator power amplifier) laser amplifier demonstration with the slave lasers' grating slightly off the resonant second-order diffraction condition. High-power laser arrays of this type have been demonstrated by SRI DSRC²⁴² and Spectra Diode Laboratories.²⁴³ Coherent output powers of 3 to 5 W with an incremental quantum efficiency in excess of 30 percent have been obtained with the array.

Vertical Cavity Lasers

The term "vertical" refers to the laser axis (or cavity) perpendicular to the wafer surface when it is fabricated. Conventional lasers have a relatively long cavity, on the order of 250 μm . It is not practical to grow such a thick layer for the laser. From the analysis, if we reduce the cavity length down to 10 μm , one needs to have a pair of very high reflectivity mirrors to make it lase at room temperature. To satisfy these conditions, researchers at Tokyo Institute of Technology²⁴⁴ have used a metal thin-film or a quarter-wavelength stack of dielectric layers (Bragg reflectors) of high- and low-index material for the mirror post to the growth of laser layers. The advances in epitaxial growth techniques allow an accurate control of semiconductor layer compositions and thicknesses such that Bragg reflectors with 99.9 percent reflectivity can be attained. Therefore, a complete vertical cavity laser structure as shown in Fig. 30 consisting of a gain medium- and high-reflectivity (more than 10 periods of alternate layers due to incremental index difference) mirrors can be grown successfully in one step by MBE or MOCVD techniques.

It is important to optimize the structure for optical gain. To maximize the modal gain, one can locate the standing wave field peak at the thin quantum well-active layer(s) (quantum well lasers were discussed earlier in this chapter) to form a resonant periodic gain structure.²⁴⁵ The issue associated with the semiconductor superlattice Bragg reflectors is the built-in carrier resistance across the abrupt heterojunction discontinuity. Without modifying the structure, the series resistance is on the order of several hundred to a thousand ohms. There have been two techniques applied to lower the resistance, namely, the use of graded junctions²⁴⁶ and peripheral Zn diffusion²⁴⁷ for conducting current to the active region. Both have demonstrated improvement over the original design.

The laser size is defined by etching into a circular column that can be mode-matched to a single-mode fiber for high coupling efficiency. It is desirable that the lasers can be planarized. Proton-bombardment-defined lasers²⁴⁸ with good performance and high yield have been obtained. Meanwhile, the small size of the laser has resulted in low threshold currents close to 1 mA.²⁴⁹ The differential quantum efficiency has been improved from a few percent to more than 30 percent; the output power level, modulation frequency, and maximum operating temperature have also increased over the past several years. As mentioned previously, the advantage of this SEL structure is the potential of high packing density. Bellcore researchers²⁵⁰ have demonstrated a novel WDM (wavelength division multiplexing) laser source with a good histogram of wavelength distribution. The grading of layer thickness across the wafer during a portion of growth translates into different lasing wavelengths. Two-dimensional, individually addressable lasers in a matrix form have also been demonstrated.²⁵¹ In the future, 2D laser arrays operating at a visible wavelength will be very useful for display and optical recording/reading applications. The performance characteristics of vertical cavity SELs reported are shown in Table 6.

19.11 CONCLUSION

In this chapter we have introduced the basic properties of semiconductor lasers and reviewed some areas of the field, including high-power operation, high-speed operation, coherence, and surface-emitting lasers. We have particularly emphasized the advantages of quantum well lasers and strained quantum well lasers. Up until very recently, all the major laser diodes were fabricated using GaAs/GaAlAs and GaInAsP/InP heterostructures. However, there have been such significant advances in the use of strained quantum wells that these lasers have performance levels which exceed, in many

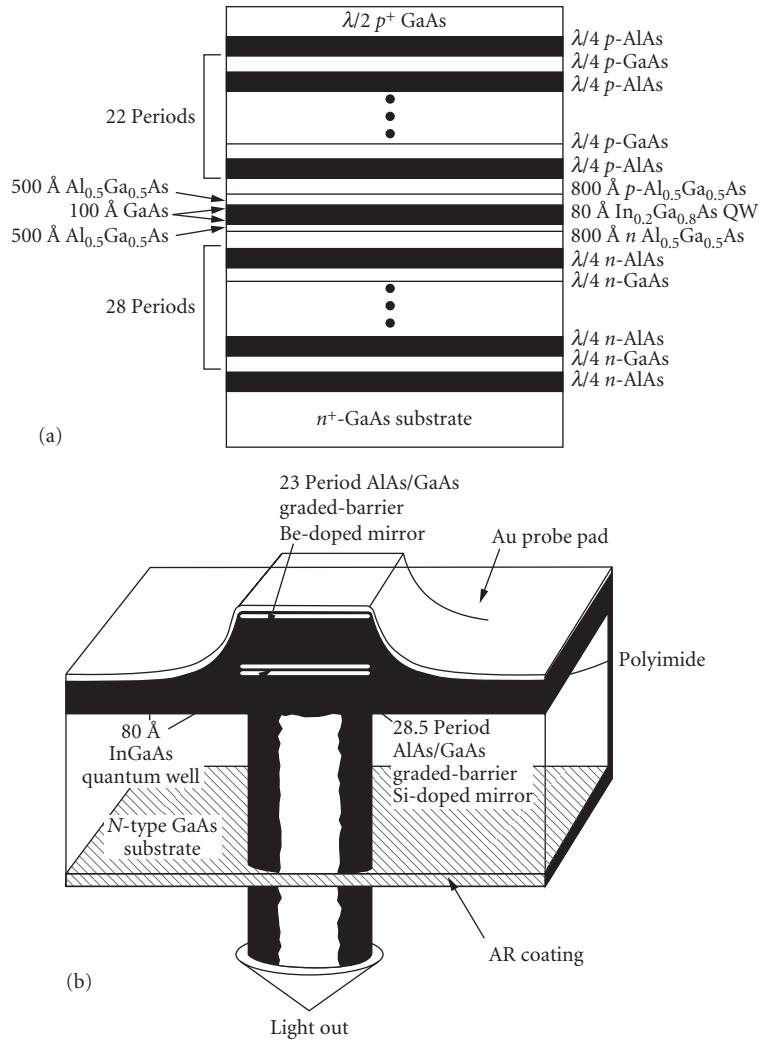


FIGURE 30 Vertical cavity surface-emitting laser with (a) layer structure and (b) device geometry. (From Ref. 249.)

cases, that found for the unstrained lasers. Coupled with the measured excellent reliability results and better output power/temperature performance, these types of lasers will experience a high demand in the future. High-power strained InGaAs/GaAs quantum well lasers are also of interest because their lasing wavelength range includes $0.98 \mu\text{m}$, which makes them useful for pumping erbium-doped fiber amplifiers.¹⁷

Another trend in the future will be the extension of commercial semiconductor lasers to a wider variety of lasing wavelengths. Just as with the standard lasers, strained quantum wells result in significant performance improvements in these more novel laser systems. Shorter wavelength ($\lambda < 0.75 \mu\text{m}$) lasers with high output powers are of interest for high-density optical recording and printing. In the last few years, much progress has been made in developing short-wavelength semiconductor lasers. In the near future, practical visible wavelength lasers will be in the red-to-yellow range, but progress has begun on even shorter wavelengths in the blue. Recent work on very

TABLE 6 Performance Characteristics of Vertical Cavity Surface-Emitting Lasers Developed in Various Research Laboratories*

J_{th} (kA/cm ²)	I_{th} (mA)	V_{th} (V)	T_0 (K)	P_{max} (mW)	η (%)	Size (μ m)	Structure	Reference
22.6	40		115	1.2	1.2	15 ϕ	0.5 mm, DH-DBR	Ref. 246, Tai, AT&T
6.6	2.5			0.3	3.9	7 ϕ	3 \times 8 nm MQW strained columnar μ laser	
6.0	1.5					5 \times 5	10 nm SQW strained columnar μ laser	Ref. 240 Jewell, AT&T
4.1	0.8					5 ϕ	Passivated 3 \times 8 nm MQS strained columnar μ laser	
3.6	3.6	3.7		0.7	4.7	10 \times 10	3 \times 8 nm MQW strained ion-implanted μ laser	Ref. 252 Lee 1990, AT&T
2.8	2.2	7.5		0.6	7.4	10 ϕ	4 \times 10 nm MQW proton- implanted μ laser	
1.4	0.7	4.0				7 \times 7	8 nm SQW strained columnar μ laser	Ref. 249 Geels, UCSB
1.2	4.8			3.0 [†]	12 [†]	20 \times 20	Strained columnar μ laser	Ref. 253 Clausen 1990 Bellcore
1.1	7.5	4.0		3.2	8.3	30 ϕ	4 \times 10 nm MQW proton- implanted μ laser	Ref. 254 Tell 1990, AT&T
0.8	1.1	4.0				12 \times 12	8 nm SQW strained columnar μ laser	Ref. 249 Geels, UCSB
				1.5	14.5	10 ϕ	4 \times 10 nm GRIN-SCH proton-implanted μ laser	Ref. 254 Tell 1990, AT&T

* J_{th} = threshold current density; I_{th} = threshold current; V_{th} = threshold voltage; P_{max} = maximum power; η = overall efficiency; T_0 = characteristic temperature;

[†]Pulsed.

long-wavelength ($\lambda > 2.0 \mu\text{m}$) GaInAsSb/AlGaAsSb lasers is also very promising. Long wavelengths ($> 1.55 \mu\text{m}$) are of interest for eye-safe laser radar, metrology, and medical instrumentation.

Currently, commercial lasers are of the edge-emitting variety, but two-dimensional surface-emitting laser arrays have advanced considerably in the past few years. When they reach maturity, they will be used for pixel interconnect and display applications.

In a limited chapter such as this it is impossible to cover all areas of the field of semiconductor lasers in depth. One of the most important areas neglected is that of tunable lasers.²⁵⁵⁻²⁵⁷ Single-mode tunable DFB and DBR lasers are of great interest for future coherent optical transmission systems. These lasers lase in a single longitudinal mode, but that mode can be tuned to a range of frequencies.

Another area not discussed in detail here is amplifiers. Amplifiers are of great interest for long-haul communication systems, for example, submarine cable systems. Amplifiers can be laser pumped fiber amplifiers¹⁷ or laser amplifiers.²⁵⁸⁻²⁶² A laser amplifier has a structure similar to that of a semiconductor laser and has some optical gain, but only enough to amplify an existing signal, not enough to lase on its own.

It is hoped that the information presented in this chapter will satisfy readers who are interested in the basics of the subject and will give readers interested in greater depth the understanding necessary to probe further in order to satisfy their specific requirements.

19.12 REFERENCES

1. H. C. Casey, Jr. and M. B. Panish, *Heterostructure Lasers, Part A: Fundamental Principles*, Academic Press, Orlando, 1978.
2. H. C. Casey, Jr. and M. B. Panish, *Heterostructure Lasers, Part B: Materials and Operating Characteristics*, Academic Press, Orlando, 1978.

3. G. H. B. Thompson, *Physics of Semiconductor Laser Devices*, John Wiley & Sons, New York, 1980.
4. G. P. Agrawal and N. K. Dutta, *Long-Wavelength Semiconductor Lasers*, Van Nostrand Reinhold, New York, 1986.
5. H. Kressel and J. K. Butler, *Semiconductor Lasers and Heterojunction LEDs*, Academic Press, New York, 1977.
6. R. A. Bartolini, A. E. Bell, and F. W. Spang, *IEEE J. Quantum Electron.* **QE-17**:69 (1981).
7. R. N. Bhargava, *J. Cryst. Growth* **117**:894 (1992).
8. C. Lin (ed.), *Optoelectronic Technology and Lightwave Communication Systems*, Van Nostrand, Reinhold, New York, 1989.
9. S. E. Miller and I. Kaminow (eds.), *Optical Fiber Telecommunications*, Academic Press, Orlando, 1988.
10. M. Katzman (ed.), *Laser Satellite Communications*, Prentice-Hall, Englewood Cliffs, N.J., 1987.
11. M. Ross, *Proc. SPIE* **885**:2 (1988).
12. J. D. McClure, *Proc. SPIE* **1219**:446 (1990).
13. G. Abbas, W. R. Babbitt, M. de La Chappelle, M. L. Fleshner, and J. D. McClure, *Proc. SPIE* **1219**:468 (1990).
14. J. W. Goodman, *International Trends in Optics*, Academic Press, Orlando, 1991.
15. R. Olshansky, V. Lanzisera, and P. Hill, *J. Lightwave Technology* **7**:1329 (1989).
16. R. L. Byer, *Proc. of the CLEO/IQEC Conf.*, Plenary Session, Baltimore, Md., 1987.
17. K. Nakagawa, S. Nishi, K. Aida, and E. Yonoda, *J. Lightwave Technology* **9**:198 (1991).
18. T. F. Deutch, J. Boll, C. A. Poliafito, K. To, *Proc. of the CLEO Conf.*, San Francisco, Calif., 1986.
19. C. Kittel, *Introduction to Solid State Physics*, John Wiley & Sons, New York, 1976.
20. L. Figueroa, "Semiconductor Lasers," *Handbook of Microwave and Optical Components*, K. Chang (ed.), J. Wiley & Sons, New York, 1990.
21. R. N. Hall, G. E. Fenner, J. D. Kingsley, T. J. Soltys, and R. O. Carlson, *Phys. Rev. Lett.* **9**:366 (1962).
22. M. I. Nathan, W. P. Dumke, G. Burns, F. H. Dill, Jr., and G. Lasher, *Appl. Phys. Lett.* **1**:62 (1962).
23. N. Holonyak, Jr. and S. F. Bevacqua, *Appl. Phys. Lett.* **1**:82 (1962).
24. T. M. Quist, R. H. Rediker, R. J. Keyes, W. E. Krag, B. Lax, A. L. McWhorter, and H. J. Zeigler, *Appl. Phys. Lett.* **1**:91 (1962).
25. T. R. Chen, Y. Zhuang, Y. J. Xu, P. Derry, N. Bar-Chaim, A. Yariv, B. Yu, Q. Z. Wang, and Y. Q. Zhou, *Optics & Laser Tech.* **22**:245 (1990).
26. G. L. Bona, P. Buchmann, R. Clauberg, H. Jaeckel, P. Vettiger, O. Voegeli, and D. J. Webb, *IEEE Photon. Tech. Lett.* **3**:412 (1991).
27. A. Behfar-Rad, S. S. Wong, J. M. Ballantyne, B. A. Stolz, and C. M. Harding, *Appl. Phys. Lett.* **54**:493 (1989).
28. N. Bouadma, J. F. Hogrel, J. Charil, and M. Carre, *IEEE J. Quantum Electron.* **QE-23**:909 (1987).
29. M. B. Panish, J. Sumski, and I. Hayashi, *Met. Trans.* **2**:795 (1971).
30. W. T. Tsang (ed.), *Semiconductors and Semimetals*, vol. 22, part A, Academic Press, New York, 1971, pp. 95–207.
31. R. D. Dupuis and P. D. Dapkus, *IEEE J. Quantum Electron.* **QE-15**:128 (1979).
32. R. D. Burnham, W. Streifer, T. L. Paoli, and N. Holonyak, Jr., *J. Cryst. Growth* **68**:370 (1984).
33. A. Y. Cho, *Thin Solid Films* **100**:291 (1983).
34. K. Ploog, *Crystal Growth, Properties and Applications*, vol. 3, H. C. Freyhardt (ed.), Springer-Verlag, Berlin, 1980, pp. 73–162.
35. B. A. Joyce, *Rep. Prog. Phys.* **48**:1637 (1985).
36. W. T. Tsang, *J. Cryst. Growth* **105**:1 (1990).
37. W. T. Tsang, *J. Cryst. Growth* **95**:121 (1989).
38. T. Hayakawa, T. Suyama, K. Takahashi, M. Kondo, S. Yamamoto, and T. Hijikata, *Appl. Phys. Lett.* **51**:707 (1987).
39. A. Kasukawa, R. Bhat, C. E. Zah, S. A. Schwarz, D. M. Hwang, M. A. Koza, and T. P. Lee, *Electron. Lett.* **27**:1063 (1991).
40. J. I. Davies, A. C. Marchall, P. J. Williams, M. D. Scott, and A. C. Carter, *Electron. Lett.* **24**:732 (1988).

41. W. T. Tsang and N. A. Olsson, *Appl. Phys. Lett.* **42**:922 (1983).
42. H. Temkin, K. Alavi, W. R. Wagner, T. P. Pearsall, and A. Y. Cho, *Appl. Phys. Lett.* **42**:845 (1983).
43. D. P. Bour, *Proc. SPIE* **1078**:60 (1989).
44. M. Ishikawa, K. Itaya, M. Okajima, and G. Hatakoshi, *Proc. SPIE* **1418**:344 (1991).
45. G. Hatakoshi, K. Itaya, M. Ishikawa, M. Okajima, and Y. Uematsu, *IEEE J. Quantum Electron.* **QE-27**:1476 (1991).
46. H. Hamada, M. Shono, S. Honda, R. Hiroyama, K. Yodoshi, and T. Yamaguchi, *IEEE J. Quantum Electron.* **QE-27**:1483 (1991).
47. M. A. Haase, J. Qiu, J. M. DePuydt, and H. Cheng, *Appl. Phys. Lett.* **59**:1272 (1991).
48. H. K. Choi and S. J. Eglash, *Appl. Phys. Lett.* **59**:1165 (1991).
49. D. L. Partin, *IEEE J. Quantum Electron.* **QE-24**:1716 (1988).
50. Z. Feit, D. Kostyk, R. J. Woods, and P. Mak, *J. Vac. Sci. Technol.* **B8**:200 (1990).
51. Y. Nishijima, *J. Appl. Phys.* **65**:935 (1989).
52. A. Ishida, K. Muramatsu, H. Takashiba, and H. Fujiasu, *Appl. Phys. Lett.* **55**:430 (1989).
53. Z. Feit, D. Kostyk, R. J. Woods, and P. Mak, *Appl. Phys. Lett.* **58**:343 (1991).
54. R. Zucca, M. Zandian, J. M. Arias, and R. V. Gill, *Proc. SPIE* **1634**:161 (1992).
55. K. Wohlleben and W. Beck, *Z. Naturforsch.* **A21**:1057 (1966).
56. J. C. Dymont, J. C. North, and L. A. D'Asaro, *J. Appl. Phys.* **44**:207 (1973).
57. T. Tsukada, *J. Appl. Phys.* **45**:4899 (1974).
58. H. Namizaki, H. Kan, M. Ishii, and A. Ito, *J. Appl. Phys.* **45**:2785 (1974).
59. H. Namizaki, *IEEE J. Quantum Electron.* **QE-11**:427 (1975).
60. C. P. Lee, S. Margalit, I. Ury, and A. Yariv, *Appl. Phys. Lett.* **32**:410 (1978).
61. R. Dingle, *Festkörper Probleme XV (Advances in Solid State Physics)*, H. Queisser (ed.), Pergamon, New York, 1975, pp. 21–48.
62. N. Holonyak, Jr., R. M. Kolbas, R. D. Dupuis, and P. D. Dapkus, *IEEE J. Quantum Electron.* **QE-16**:170 (1980).
63. N. Okamoto, *Jpn. J. Appl. Phys.* **26**:315 (1987).
64. P. Zory (ed.), *Quantum Well Lasers*, Academic Press, Orlando, 1993.
65. C. Cohen-Tannoudji, B. Diu, and F. Lalöe, *Quantum Mechanics*, vol. 1, John Wiley & Sons, New York, 1977.
66. G. Lasher and F. Stern, *Phys. Rev.* **133**:A553 (1964).
67. S. W. Corzine, R. H. Yan, and L. A. Coldren, *Quantum Well Lasers*, P. Zory (ed.), Academic Press, Orlando, 1993.
68. P. L. Derry, *Properties of Buried Heterostructure Single Quantum Well (Al, Ga)As Lasers*, thesis, Calif. Inst. of Tech., Pasadena, Calif., 1989.
69. P. L. Derry, A. Yariv, K. Y. Lau, N. Bar-Chaim, K. Lee, and J. Rosenberg, *Appl. Phys. Lett.* **50**:1773 (1987).
70. H. Z. Chen, A. Ghaffari, H. Morkoç, and A. Yariv, *Appl. Phys. Lett.* **51**:2094 (1987).
71. H. Chen, A. Ghaffari, H. Morkoç, and A. Yariv, *Electron. Lett.* **23**:1334 (1987).
72. R. Fischer, J. Klem, T. J. Drummond, W. Kopp, H. Morkoç, E. Anderson, and M. Pion, *Appl. Phys. Lett.* **44**:1 (1984).
73. P. L. Derry, T. R. Chen, Y. Zhuang, J. Paslaski, M. Middlestein, K. Vahala, A. Yariv, K. Y. Lau, and N. Bar-Chaim, *Optoelectronics—Dev. and Tech.* **3**:117 (1988).
74. P. L. Derry, R. J. Fu, C. S. Hong, E. Y. Chan, K. Chiu, H. E. Hager, and L. Figueroa, *Proc. SPIE* **1634**:374 (1992).
75. R. J. Fu, C. S. Hong, E. Y. Chan, D. J. Booher, and L. Figueroa, *IEEE Photon. Tech. Lett.* **3**:308 (1991).
76. R. J. Fu, C. S. Hong, E. Y. Chan, D. J. Booher, and L. Figueroa, *Proc. SPIE* **1418**:108 (1991).
77. W. T. Tsang, *Appl. Phys. Lett.* **40**:217 (1982).
78. W. T. Tsang, *Appl. Phys. Lett.* **39**:134 (1981).
79. W. T. Tsang, *Appl. Phys. Lett.* **39**:786 (1981).

80. Y. Arakawa and A. Yariv, *IEEE J. Quantum Electron.* **QE-21**:1666 (1985).
81. E. P. O'Reilly, *Semicond. Sci. Technol.* **4**:121 (1989).
82. E. P. O'Reilly and A. Ghiti, *Quantum Well Lasers*, P. Zory (ed.), Academic Press, Orlando, 1993.
83. J. W. Matthews and A. E. Blakeslee, *J. Cryst. Growth* **27**:118 (1974).
84. I. J. Fritz, S. T. Picraux, L. R. Dawson, T. J. Drummon, W. D. Laidig, and N. G. Anderson, *Appl. Phys. Lett.* **46**:967 (1985).
85. T. G. Andersson, Z. G. Chen, V. D. Kulakovskii, A. Uddin, and J. T. Vallin, *Appl. Phys. Lett.* **51**:752 (1987).
86. M. Altarelli, *Heterojunctions and Semiconductor Superlattices*, G. Allan, G. Bastard, N. Boccarda, M. Lannoo, and M. Voos (eds.), Springer-Verlag, Berlin, 1985, p. 12.
87. J. M. Luttinger and W. Kohn, *Phys. Rev.* **97**:869 (1955).
88. P. Lawaetz, *Phys. Rev.* **B4**:3640 (1971).
89. D. Ahn and S. L. Chuang, *IEEE J. Quantum Electron.* **QE-24**:2400 (1988).
90. S. L. Chuang, *Phys. Rev.* **B43**:9649 (1991).
91. E. Yablonovitch and E.O. Kane, *J. Lightwave Tech.* **LT-4**:504 (1986).
92. H. K. Choi and C. A. Wang, *Appl. Phys. Lett.* **57**:321 (1990).
93. N. Chand, E. E. Becker, J. P. van der Ziel, S. N. G. Chu, and N. K. Dutta, *Appl. Phys. Lett.* **58**:1704 (1991).
94. C. A. Wang and H. K. Choi, *IEEE J. Quantum Electron.* **QE-27**:681 (1991).
95. R. L. Williams, M. Dion, F. Chatenoud, and K. Dzurko, *Appl. Phys. Lett.* **58**:1816 (1991).
96. S. L. Yellen, R. G. Waters, P. K. York, K. J. Beerink, and J. J. Coleman, *Electron Lett.* **27**:552 (1991).
97. P. K. York, K. J. Beerink, J. Kim, J. J. Alwan, J. J. Coleman, and C. M. Wayman, *J. Cryst. Growth* **107**:741 (1991).
98. R. G. Waters, D. P. Bour, S. L. Yellen, and N. F. Ruggieri, *IEEE Photon. Tech. Lett.* **2**:531 (1990).
99. K. Fukagai, S. Ishikawa, K. Endo, and T. Yuasa, *Japan J. Appl. Phys.* **30**:L371 (1991).
100. J. J. Coleman, R. G. Waters, and D. P. Bour, *Proc. SPIE* **1418**:318 (1991).
101. S. Tsuji, K. Mizuishi, H. Hirao, and M. Nakamura, *Links for the Future: Science, Systems & Services for Communications*, P. Dewilde and C. A. May (eds.), IEEE/Elsevier Science, North Holland, 1984, p. 1123.
102. H. D. Wolf and K. Mettler, *Proc. SPIE* **717**:46 (1986).
103. J. Hashimoto, T. Katsyama, J. Shinkai, I. Yoshida, and H. Hayashi, *Appl. Phys. Lett.* **58**:879 (1991).
104. H. B. Serreze, Y. C. Chen, and R. G. Waters, *Appl. Phys. Lett.* **58**:2464 (1991).
105. D. F. Welch and D. R. Scifres, *Electron. Lett.* **27**:1915 (1991).
106. J. I. Pankove, *Optical Processes in Semiconductors*, Dover Publ. Inc., New York, 1971.
107. N. K. Dutta, *J. Appl. Phys.* **54**:1236 (1983).
108. N. K. Dutta and R. J. Nelson, *J. Appl. Phys.* **53**:74 (1982).
109. A. R. Adams, M. Asada, Y. Suematsu, and S. Arai, *Jpn. J. Appl. Phys.* **19**:L621 (1980).
110. T. Tanbun-Ek, R. A. Logan, H. Temkin, K. Berthold, A. F. J. Levi, and S. N. G. Chu, *Appl. Phys. Lett.* **55**:2283 (1989).
111. A. R. Adams, *Electron. Lett.* **22**:249 (1986).
112. Y. Jiang, M. C. Teich, and W. I. Wang, *Appl. Phys. Lett.* **57**:2922 (1990).
113. H. Temkin, R. A. Logan, and T. Tanbun-Ek, *Proc. SPIE* **1418**:88 (1991).
114. C. E. Zah, R. Bhat, R. J. Favire, Jr., S. G. Menocal, N. C. Andreadakis, K. W. Cheung, D. M. Hwang, M. A. Koza, and T. P. Lee, *IEEE J. Quantum Electron.* **27**:1440 (1991).
115. P. J. A. Thijs, L. F. Tiemeijer, P. I. Kuindersma, J. J. M. Binsma, and T. Van Dongen, *IEEE J. Quantum Electron.* **27**:1426 (1991).
116. C. E. Zah, R. Bhat, B. Pathak, C. Caneau, F. J. Favire, Jr., N. C. Andreadakis, D. M. Hwang, M. A. Koza, C. Y. Chen, and T. P. Lee, *Electron. Lett.* **27**:1414 (1991).
117. P. J. A. Thijs, J. J. M. Binsma, E. W. A. Young, and W. M. E. Van Gils, *Electron. Lett.* **27**:791 (1991).
118. E. P. O'Reilly, G. Jones, A. Ghiti, and A. R. Adams, *Electron. Lett.* **27**:1417 (1991).
119. S. W. Corzine and L. A. Coldren, *Appl. Phys. Lett.* **59**:588 (1991).

120. A. Larsson, M. Mittelstein, Y. Arakawa, and A. Yariv, *Electron. Lett.* **22**:79 (1986).
121. S. Simhony, E. Kapon, E. Colas, R. Bhat, N. G. Stoffel, and D. M. Hwang, *IEEE Photon. Tech. Lett.* **2**:305 (1990).
122. K. J. Vahala, J. A. Lebens, C. S. Tsai, T. F. Kuech, P. C. Sercel, M. E. Hoenk, and H. Zarem, *Proc. SPIE* **1216**:120 (1990).
123. N. Chinone, *J. Appl. Phys.* **48**:3237 (1978).
124. P. A. Kirby, A. R. Goodwin, G. H. B. Thompson, D. F. Lovelace, and S. E. Turley, *IEEE J. Quantum Electron.* **QE-13**:720 (1977).
125. R. Lang, *IEEE J. Quantum Electron.* **QE-15**:718 (1979).
126. S. Wang, C. Y. Chen, A. S. Liao, and L. Figueroa, *IEEE J. Quantum Electron.* **QE-17**:453 (1981).
127. K. Aiki, N. Nakamura, T. Kurada, and J. Umeda, *Appl. Phys. Lett.* **30**:649 (1977).
128. M. Nakamura, *IEEE Trans. Circuits Syst.* **26**:1055 (1979).
129. D. Botez, *IEEE Spectrum* **22**:43 (1985).
130. D. Botez, *RCA Rev.* **39**:577 (1978).]
131. H. C. Casey, M. B. Panish, W. O. Schlosser, and T. L. Paoli, *J. Appl. Phys.* **45**:322 (1974).
132. R. J. Fu, C. J. Hwang, C. S. Wang, and B. Lolevic, *Appl. Phys. Lett.* **45**:716 (1984).
133. M. Wada, K. Hamada, H. Himuza, T. Sugino, F. Tujiri, K. Itoh, G. Kano, and I. Teramoto, *Appl. Phys. Lett.* **42**:853 (1983).
134. K. Hamada, M. Wada, H. Shimuzu, M. Kume, A. Yoshikawa, F. Tajiri, K. Itoh, and G. Kano, *Proc. IEEE Int. Semicond. Lasers Conf.*, Rio de Janeiro, Brazil, 1984, p. 34.
135. K. Endo, H. Kawamo, M. Ueno, N. Nido, Y. Kuwamura, T. Furese, and I. Sukuma, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 38.
136. D. Botez, J. C. Connolly, M. Ettenberg, and D. B. Gilbert, *Electron. Lett.* **19**:882 (1983).
137. B. Goldstein, J. K. Butler, and M. Ettenberg, *Proc. CLEO Conf.*, Baltimore, Md., 1985, p. 180.
138. Y. Yamamoto, N. Miyauchi, S. Maci, T. Morimoto, O. Yamamoto, S. Yomo, and T. Hijikata, *Appl. Phys. Lett.* **46**:319 (1985).
139. S. Yamamoto, H. Hayashi, T. Hayashi, T. Hayakawa, N. Miyauchi, S. Yomo, and T. Hijikata, *Appl. Phys. Lett.* **42**:406 (1983).
140. D. Ackley, *Electron. Lett.* **20**:509 (1984).
141. J. Yang, C. S. Hong, L. Zinkiewicz, and L. Figueroa, *Electron. Lett.* **21**:751 (1985).
142. J. Ungar, N. Bar-Chaim, and I. Ury, *Electron. Lett.* **22**:280 (1986).
143. D. F. Welch, W. Streifer, D. R. Scifres, *Proc. SPIE* **1043**:54 (1989).
144. D. R. Daniel, D. Buckley, B. Garrett, *Proc. SPIE* **1043**:61 (1989).
145. D. Botez, *IEEE J. Quantum Electron.* **QE-17**:2290 (1981).
146. T. Kuroda, M. Nakamura, K. Aiki, and J. Umeda, *Appl. Opt.* **17**:3264 (1978).
147. S. J. Lee, L. Figueroa, and R. Rammaswamy, *IEEE J. Quant. Electron.* **25**:1632 (1989).
148. H. Yonezu, M. Ueno, T. Kamejima, and I. Hayashi, *IEEE J. Quantum Electron.* **15**:775 (1979).
149. H. Kumabe, T. Tumuka, S. Nita, Y. Seiwa, T. Sugo, and S. Takamija, *Jpn. J. Appl. Phys.* **21**:347 (1982).
150. H. Blauvelt, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **40**:1029 (1982).
151. D. Botez and J. C. Connolly, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 36.
152. H. Matsubara, K. Ishiki, H. Kumabe, H. Namazaki, and W. Susaki, *Proc. CLEO.*, Baltimore, Md., 1985, p. 180.
153. F. Capasso, and G. F. Williams, *J. Electrochem. Soc.* **129**:821 (1982).
154. H. H. Lee and L. Figueroa, *J. Electrochem. Soc.* **135**:496 (1988).
155. H. Kawanishi, H. Ohno, T. Morimoto, S. Kaneiwa, N. Miyauchi, H. Hayashi, Y. Akagi, Y. Nakajima, *Proc. SPIE* **1219**:309 (1990).
156. J. Yoo, H. Lee, and P. Zory, *IEEE Photonics Lett.* **3**:594 (1991).
157. Y. Suzuki, Y. Horikoshi, M. Kobayashi, and H. Okamoto, *Electron. Lett.* **20**:384 (1984).

158. M. Yamaguchi, H. Nishimoto, M. Kitumara, S. Yamazaki, I. Moto, and K. Kobayashi, *Proc. CLEO*, Baltimore, 1988, p. 180.
159. C. B. Morrison, D. Botez, L. M. Zinkiewicz, D. Tran, E. A. Rezek, and E. R. Anderson, *Proc. SPIE* **893**:84 (1988).
160. Y. Sakakibara, E. Oomura, H. Higuchi, H. Namazaki, K. Ikeda, and W. Susaki, *Electron. Lett.* **20**:762 (1984).
161. K. Imanaka, H. Horikawa, A. Matoba, Y. Kawai, and M. Sakuta, *Appl. Phys. Lett.* **45**:282 (1984).
162. M. Kawahara, S. Shiba, A. Matoba, Y. Kawai, and Y. Tamara, *Proc. Opt. Fiber Commun.* (OFC 1987), paper ME1, 1987.
163. S. Oshiba, A. Matoba, H. Horikawa, Y. Kawai, and M. Sakuta, *Electron. Lett.* **22**:429 (1986).
164. B. S. Bhumbra, R. W. Glew, P. D. Greene, G. D. Henshall, C. M. Lowney, and J. E. A. Whiteaway, *Electron. Lett.* **26**:1755 (1990).
165. T. Tanbun-Ek, R. A. Logan, N. A. Olsson, H. Temkin, A. M. Sergent, and K. W. Wecht, *International Semiconductor Laser Conference*, paper D-3, Davos, 1990.
166. G. D. Henshall, A. Hadjifotiou, R. A. Baker, and K. J. Warwick, *Proc. SPIE* **1418**:286 (1991).
167. M. Arvind, H. Hsing, and L. Figueroa, *J. Appl. Phys.* **63**:1009 (1988).
168. A. Larsson, S. Forouher, J. Cody, and R. J. Lang, *Proc. SPIE* **1418**:292 (1991).
169. M. Okayasu, M. Fukuda, T. Takeshita, O. Kogure, T. Hirone, and S. Uehara, *Proc. of Optical Fiber Communications Conf.*, 29, 1990.
170. D. F. Welch, C. F. Schaus, S. Sun, M. Cardinal, W. Streifer, and D. R. Scifres, *Proc. SPIE* **1219**:186 (1990).
171. G. L. Harnagel, J. M. Haden, G. S. Browder, Jr., M. Cardinal, J. G. Endriz, and D. R. Scifres, *Proc. SPIE* **1219**:186 (1990).
172. D. R. Scifres, R. D. Burnham, and W. Steifer, *Appl. Phys. Lett.* **33**:1015 (1978).
173. D. Botez and J. C. Connally, *Appl. Phys. Lett.* **43**:1096 (1983).
174. D. E. Ackley and R. G. Engelmann, *Appl. Phys. Lett.* **39**:27 (1981).
175. D. R. Scifres, R. D. Burnham, W. Streifer, and M. Bernstein, *Appl. Phys. Lett.* **41**:614 (1982).
176. D. R. Scifres, C. Lindstrom, R. D. Burnham, W. Streifer, and T. L. Paoli, *Appl. Phys. Lett.* **19**:160 (1983).
177. F. Kappeler, H. Westmeier, R. Gessner, M. Druminski, and K. H. Zschauer, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 90.
178. J. P. Van der Ziel, H. Temkin, and R. D. Dupuis, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 92.
179. L. Figueroa, C. Morrison, H. D. Law, and F. Goodwin, *Proc. Int. Electron Devices Meeting*, 1983, p. 760.
180. L. Figueroa, C. Morrison, H. D. Law, and F. Goodwin, *J. Appl. Phys.* **56**:3357 (1984).
181. C. Morrison, L. Zinkiewicz, A. Burghard, and L. Figueroa, *Electron. Lett.* **21**:337 (1985).
182. Y. Twu, A. Dienes, S. Wang, and J. R. Whinnery, *Appl. Phys. Lett.* **45**:709 (1984).
183. S. Mukai, C. Lindsey, J. Katz, E. Kapon, Z. Rav-Noy, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **45**:834 (1984).
184. D. F. Welch, D. Scifres, P. Cross, H. Kung, W. Streifer, R. D. Burnham, and J. Yaeli, *Electron. Lett.* **21**:603 (1985).
185. N. Dutta, L. A. Kozzi, S. G. Napholtz, and B. P. Seger, *Proc. Conf. Lasers Electro-Optics (CLEO)*, Baltimore, Md., 1985, p. 44.
186. M. Taneya, M. Matsumoto, S. Matsui, Y. Yano, and T. Hijikata, *Appl. Phys. Lett.* **47**:341 (1985).
187. J. Ohsawa, S. Himota, T. Aoyagi, T. Kadowaki, N. Kaneno, K. Ikeda, and W. Susaki, *Electron. Lett.* **21**:779 (1985).
188. D. F. Welch, P. S. Cross, D. R. Scifres, W. Streifer, and R. D. Burnham, *Proc. CLEO*, San Francisco, Calif., 1986, p. 66.
189. L. Mawst, D. Botez, E. R. Anderson, M. Jansen S. Ou, M. Sargent, G. L. Peterson, and T. J. Roth, *Proc. SPIE* **1418**:353 (1991).
190. J. K. Butler, D. E. Ackley, and D. Botez, *Appl. Phys. Lett.* **44**:293 (1984).
191. E. Kapon, J. Katz, and A. Yariv, *Opt. Lett.* **10**:125 (1984).
192. K. L. Chen and S. Wang, *Electron. Lett.* **21**:347 (1985).

193. W. Streifer, A. Hardy, R. D. Burnham, and D. R. Scifres, *Electron. Lett.* **21**:118 (1985).
194. S. Chinn and R. J. Spier, *IEEE J. Quantum Electron.* **20**:358 (1985).
195. J. Katz, E. Kapon, C. Lindsey, S. Margalit, U. Shreter, and A. Yariv, *Appl. Phys. Lett.* **42**:521 (1983).
196. E. Kapon, C. P. Lindsey, J. S. Smith, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **45**:1257 (1984).
197. D. Ackley, *Electron. Lett.* **20**:695 (1984).
198. T. R. Ranganath and S. Wang, *IEEE J. Quantum Electron.* **13**:290 (1977).
199. L. Figueroa, T. Holcomb, K. Burghard, D. Bullock, C. Morrison, L. Zinkiewicz, and G. Evans, *IEEE J. Quantum Electron.* **22**:241 (1986).
200. L. J. Mawst, M. E. Givens, C. A. Zmudzinski, M. A. Emanuel, and J. J. Coleman, *IEEE J. Quantum Electron.* **QE-23**:696 (1987).
201. P. S. Zory, A. R. Reisinger, R. G. Walters, L. J. Mawst, C. A. Zmudzinski, M. A. Emanuel, M. E. Givens, and J. J. Coleman, *Appl. Phys. Lett.* **49**:16 (1986).
202. R. G. Walters, P. L. Tihanyi, D. S. Hill, and B. A. Soltz, *Proc. SPIE* **893**:103 (1988).
203. M. S. Zediker, D. J. Krebs, J. L. Levy, R. R. Rice, G. M. Bender, and D. L. Begley, *Proc. SPIE* **893**:21 (1988).
204. C. Krebs and B. Vivian, *Proc. SPIE* **893**:38 (1988).
205. K. Y. Lau and A. Yariv, *Semiconductors and Semimetals Volume 22: Lightwave Communications Technology*, W. T. Tsang (ed.), Academic Press, New York, 1985, pp. 69–151.
206. K. Petermann, *Laser Diode Modulation and Noise*, Kluwer Academic Publ., Dordrecht, The Netherlands, 1988.
207. K. Petermann, *IEEE J. Quantum Electron.* **QE-15**:566 (1979).
208. G. F. Franklin, J. D. Powell, and A. Emami-Naeini, *Feedback Control of Dynamic Systems*, Addison-Wesley Publishing Company, Reading, 1986.
209. T. P. Lee and R. M. Derosier, *Proc. IEEE* **62**:1176 (1974).
210. G. Arnold, P. Russer, and K. Petermann, *Topics in Applied Physics, vol. 39: Semiconductor Devices for Optical Communication*, H. Kressel (ed.), Springer-Verlag, Berlin, 1982, p. 213.
211. K. Y. Lau, N. Bar-Chaim, P. L. Derry, and A. Yariv, *Appl. Phys. Lett.* **51**:69 (1987).
212. K. Y. Lau, P. L. Derry, and A. Yariv, *Appl. Phys. Lett.* **52**:88 (1988).
213. J. E. Bowers, B. R. Hemenway, A. H. Gnauck, and D. P. Wilt, *IEEE J. Quantum Electron.*, **QE-22**:833 (1986).
214. J. E. Bowers, *Solid-State Electron.* **30**:1 (1987).
215. J. Bowers, *Conference on Optical Fiber Communication: Tutorial Sessions*, San Jose, Calif. 1992, p. 233.
216. K. Furuya, Y. Suematsu, and T. Hong, *Appl. Optics* **17**:1949 (1978).
217. D. Wilt, K. Y. Lau, and A. Yariv, *J. Appl. Phys.* **52**:4970 (1981).
218. K. Y. Lau and A. Yariv, *IEEE J. Quantum Electron.* **QE-21**:121 (1985).
219. R. S. Tucker, C. Lin, C. A. Burrus, P. Besomi, and R. J. Nelson, *Electron. Lett.* **20**:393 (1984).
220. W. H. Cheng, A. Appelbaum, R. T. Huang, D. Renner, and K. R. Cioffi, *Proc. SPIE* **1418**:279 (1991).
221. C. B. Su and V. A. Lanzisera, *IEEE J. Quantum Electron.* **QE-22**:1568 (1986).
222. R. Olshansky, W. Powazink, P. Hill, V. Lanzisera, and R. B. Lauer, *Electron. Lett.* **23**:839 (1987).
223. E. Meland, R. Holmstrom, J. Schlafer, R. B. Lauer, and W. Powazink, *Electron. Lett.* **26**:1827 (1990).
224. S. D. Offsey, W. J. Schaff, L. F. Lester, L. F. Eastman, and S. K. McKernan, *IEEE J. Quantum Electron.* **27**:1455 (1991).
225. R. Nagarajan, T. Fukushima, J. E. Bowers, R. S. Geels, and L. A. Coldren, *Appl. Phys. Lett.* **58**:2326 (1991).
226. L. F. Lester, W. J. Schaff, X. J. Song, and L. F. Eastman, *Proc. SPIE* **1634**:127 (1992).
227. K. Y. Lau, C. Harder, and A. Yariv, *IEEE J. Quantum Electron.* **QE-20**:71 (1984).
228. Y. Sakakibara, K. Furuya, K. Utaka, and Y. Suematsu, *Electron. Lett.* **16**:456 (1980).
229. S. Kobayashi, Y. Yamamoto, M. Ito, and T. Kimura, *IEEE J. Quantum Electron.* **QE-18**:582 (1982).
230. P. Vankwikelberge, F. Buytaert, A. Franchois, R. Baets, P. I. Kuindersma, and C. W. Fredrksz, *IEEE J. Quantum Electron.* **25**:2239 (1989).

231. T. Ikegami, "Longitudinal Mode Control in Laser Diodes," *Opto-Electronics Technology and Lightwave Communication Systems*, Van Nostrand Reinhold, New York, 1989, p. 264.
232. T. E. Bell, *IEEE Spectrum* **20**(12):38 (December 1983).
233. M. Osinsky and J. Boos, *IEEE J. Quantum Electron.* **QE-23**:9 (1987).
234. S. Takano, T. Sasaki, H. Yamada, M. Kitomura, and I. Mito, *Electron Lett.* **25**:356 (1989).
235. W. T. Tsang, N. A. Olson, R. A. Linke, and R. A. Logan, *Electron Lett.* **19**:415 (1983).
236. R. Beausoleil, J. A. McGarvey, R. L. Hagman, and C. S. Hong, *Proc. of the CLEO Conference*, Baltimore, Md., 1989.
237. S. Murata, S. Yamazaki, I. Mito, and K. Koboyashi, *Electron Lett.* **22**:1197 (1986).
238. L. Goldberg and J. F. Weller, *Electron Lett.* **22**:858 (1986).
239. J. L. Jewell, A. Scherer, S. L. McCall, Y. H. Lee, S. J. Walker, J. P. Harbison, and L. T. Florez, *Electron. Lett.* **25**:1123 (1989).
240. Z. L. Liao and J. N. Walpole, *Appl. Phys. Lett.* **50**:528 (1987).
241. M. Jansen, J. J. Yang, S. S. Ou, M. Sergeant, L. Mawst, T. J. Roth, D. Botez, and J. Wilcox, *Proc. SPIE* **1582**:94 (1991).
242. G. A. Evans, D. P. Bour, N. W. Carlson, et al., *IEEE J. Quantum Electron.* **27**:1594 (1991).
243. D. Mehuys, D. Welch, R. Parke, R. Waarts, A. Hardy, and D. Scifres, *Proc. SPIE* **1418**:57 (1991).
244. K. Iga, F. Koyama, and S. Kinoshita, *IEEE J. Quantum Electron.* **24**:1845 (1988).
245. M. Y. A. Raja, S. R. J. Brueck, M. Osinski, C. F. Schaus, J. G. McInerney, T. M. Brennan, and B. E. Hammons, *IEEE J. Quantum Electron.* **25**:1500 (1989).
246. K. Tai, L. Yang, Y. H. Wang, J. D. Wynn, and A. Y. Cho, *Appl. Phys. Lett.* **56**:2496 (1990).
247. Y. J. Yang, T. G. Dziura, R. Fernandez, S. C. Wang, G. Du, and S. Wang, *Appl. Phys. Lett.* **58**:1780 (1991).
248. M. Orenstein, A. C. Von Lehmen, C. Chang-Hasnain, N. G. Stoffel, J. P. Harbison, L. T. Florez, E. Clausen, and J. E. Jewell, *Appl. Phys. Lett.* **56**:2384 (1990).
249. R. S. Geels and L. A. Coldren, *Appl. Phys. Lett.* **57**:1605 (1990).
250. C. J. Chang-Hasnain, J. R. Wullert, J. P. Harbison, L. T. Florez, N. G. Stoffel, and M. W. Maeda, *Appl. Phys. Lett.* **58**:31 (1991).
251. A. Von Lehmen, M. Orenstein, W. Chan, C. Chang-Hasnain, J. Wullert, L. Florez, J. Harbison, and N. Stoffel, *Proc. of the CLEO Conference*, Baltimore, Md., 1991, p. 46.
252. Y. H. Lee, J. L. Jewell, B. Tell, K. F. Brown-Goebeler, A. Scherer, J. P. Harbison, and L. T. Florez, *Electron. Lett.* **26**:225 (1990).
253. E. M. Clausen, Jr., A. Von Lehmen, C. Chang-Hasnain, J. P. Harbison, and L. T. Florez, *Techn. Digest Postdeadline Papers, OSA 1990 Annual Meeting*, Boston, Mass., 1990, p. 52.
254. B. Tell, Y. H. Lee, K. F. Brown-Goebeler, J. L. Jewell, R. E. Leibenguth, M. T. Asom, G. Livescu, L. Luther, and V. D. Matterna, *Appl. Phys. Lett.* **57**:1855 (1990).
255. T. P. Lee, *IEEE Proceedings* **79**:253 (1991).
256. M.-C. Amann and W. Thulke, *IEEE J. Selected Areas Comm.* **8**:1169 (1990).
257. K. Kobayashi and I. Mito, *J. Lightwave Tech.* **6**:1623 (1988).
258. T. Saitoh and T. Mukai, *IEEE Global Telecommunications Conference and Exhibition*, San Diego, Calif., 1990, p. 1274.
259. N. A. Olsson, *J. Lightwave Tech.* **7**:1071 (1989).
260. A. F. Mitchell and W. A. Stallard, *IEEE Int. Conference on Communications*, Boston, Mass., 1989, p. 1546.
261. T. Saitoh and T. Mukai, *J. Lightwave Tech.* **6**:1656 (1988).
262. M. J. O'Mahony, *J. Lightwave Tech.* **6**:531 (1988).

ULTRASHORT OPTICAL SOURCES AND APPLICATIONS

Jean-Claude Diels

*Departments of Physics and Electrical Engineering
University of New Mexico
Albuquerque, New Mexico*

Ladan Arissian

*Texas A&M University
College Station Texas, and
National Research Council of Canada
Ottawa, Ontario, Canada*

20.1 INTRODUCTION

It is considered an easy task to control waveforms down to a few cycles with electronic circuits, at frequencies in the megahertz range. Ultrafast optics has seen the development of the same capability at optical frequencies, i.e., in the peta Hertz range. Laser pulses of a few optical cycles (pulse duration of a few femtoseconds) are routinely generated, with a suboptical cycle accuracy. The high power of these ultrashort bursts of electromagnetic radiation have led to new type of high field interactions. Electrons ejected from an atom/molecule by tunnel or multiphoton ionization can be recaptured by the next half optical cycle of opposite sign. The interaction of the returning electron with the atom/molecule is rich of new physics, including high harmonic generation, generation of single attosecond pulses of attosecond pulse trains, scattering of returning electrons by the atom/molecule, etc. Generation, amplification, control, and manipulation of optical pulses is an important starting point for these high field studies.

As compared to fast electronics, ultrafast optical pulses have reached a considerable higher level of accuracy. Pulse trains can be generated, of which the spacing between pulses (of the order of nanoseconds) is a measurable number of optical cycles (one optical cycle being approximately 2 fs in the visible). The frequency spectrum of these pulse trains is a frequency comb, of which each tooth can be an absolute standard with a subhertz accuracy. These frequency combs have numerous applications in metrology and physics—for instance, determining the eventual drift of physical constants, or in astronomy, a considerable improvement in the determination of Doppler shifts of various sources. In addition to the high level of accuracy and control in time and frequency, the femtosecond sources have a remarkable amplitude stability. This stability is the result of nonlinear intracavity losses being minimum for a particular intensity.

This chapter starts with a detailed description of an optical pulse and an optical pulse train. Nonlinear mechanisms are described that can be exploited to control pulse duration, chirp, intensity

of the mode-locked lasers. In particular, a mode-locked laser with two intracavity pulses will be discussed, and its analogy with a quantum mechanical two-level system.

20.2 DESCRIPTION OF OPTICAL PULSES AND PULSE TRAINS

Single Optical Pulse

In this first section we will summarize the essential notations and definitions used throughout the chapter. Ideally, a mode-locked laser emits a continuous train of identical ultrashort pulses. To this infinite series of identical pulses corresponds, in the frequency domain, a finite (but large) number of equally spaced modes, generally referred to as a *frequency comb*. Inside the laser typically, only one pulse circulates. The shape of an intracavity pulse results from a steady-state equilibrium between various mechanisms of pulse stretching (saturable gain, dispersion), compression (saturable absorption, combination of self-phase modulation, and dispersion), amplification, and losses.

The pulse is characterized by measurable quantities which can be directly related to the electric field. A complex representation of the field amplitude is particularly convenient in dealing with propagation problems of electromagnetic pulses.

The complex spectrum of the pulse $\tilde{E}(\Omega)^*$ is defined by taking the complex Fourier transform \mathcal{F} of the real electric field $E(t) = \varepsilon(t)\cos[\omega t + \varphi(t)]$:

$$\tilde{E}(\Omega) = \mathcal{F}\{E(t)\} = \int_{-\infty}^{\infty} E(t)e^{-i\Omega t} dt = \left| \tilde{E}(\Omega) \right| e^{i\Phi(\Omega)} = \tilde{\varepsilon}(\Omega - \omega) = \tilde{\varepsilon}(\Delta\Omega) \quad (1)$$

In the definition (1), $|\tilde{E}(\Omega)|$ denotes the spectral amplitude and $\Phi(\Omega)$ is the spectral phase. Since $E(t)$ is a real function, its Fourier transform is symmetric, and its negative frequency part can be considered as redundant information. We will therefore choose to represent the light pulse by either the positive frequency function $\tilde{E}(\Omega) = E(\Omega)e^{i\Phi(\Omega)}$ (defined as being equal to zero for $\Omega < 0$) or its complex inverse Fourier transform in the time domain

$$\tilde{E}(t) = \frac{1}{2\pi} \int_0^{\infty} \tilde{E}(\Omega)e^{i\Omega t} d\Omega = \frac{1}{2} \tilde{\varepsilon}(t)e^{i\omega t} = \frac{1}{2} \varepsilon(t)e^{i[\omega t + \varphi(t)]} \quad (2)$$

The relation with the real physical measurable field $E(t)$ is

$$E(t) = \tilde{E}(t) + \text{c.c.} = \varepsilon(t)\cos[\omega t + \varphi(t)] \quad (3)$$

The latter part of Eq. (2) defines a pulse envelope function $\varepsilon(t)$, a carrier frequency ω and a phase $\varphi(t)$. The decomposition is somewhat arbitrary, since the instantaneous frequency is given by

$$\omega(t) = \omega + \frac{d}{dt}\varphi(t) \quad (4)$$

In general, the carrier frequency ω will be chosen such that the average contribution from the phase factor $\varphi(t)$ is zero:

$$\langle \varphi(t) \rangle = \frac{\int_{-\infty}^{\infty} \varepsilon(t)^2 \dot{\varphi}(t) dt}{\int_{-\infty}^{\infty} \varepsilon(t)^2 dt} = 0 \quad (5)$$

*Complex quantities related to the field will be represented with a tilde.

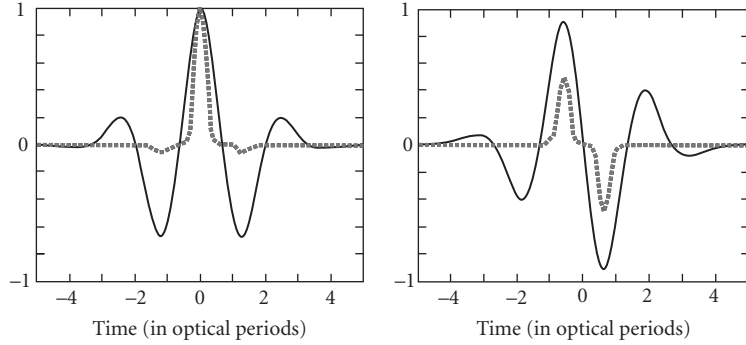


FIGURE 1 Comparison of a two-cycle pulse with $\varphi_e = 0$ (left) and $\varphi_e = \pi/2$ (right). The solid line traces the instantaneous electric field normalized to the peak value of its envelope, as a function of time in units of the optical period. The dotted lines correspond to the seventh power of the electric field, which would be driving a seven photon process.

For pulses of a few optical cycles, the variation of the phase factor can often be neglected across the pulse, and $\varphi(t) = \varphi_e$ is constant. Even for a single pulse, the phase factor φ_e is of practical significance, when a nonlinear phenomena traces the electric field under the envelope of a pulse of only a few cycle duration. If the phase φ_e is zero, the time dependence of the electric field is symmetric, with a peak in the center at $t = 0$, larger than the two opposite maxima at $t = \pm T/2$. If the phase $\varphi_e = \pi/2$, the time dependence of the electric field is antisymmetric, with equal opposite extrema at $t = \pm\pi/4$. These two pulses can give a different response in highly nonlinear phenomena. Let us consider for instance the shortest pulse that can be generated at 800 nm, which has a full width at half maximum (FWHM) of the intensity of 2.5 fs. Its complex electric field envelope can be written as $\tilde{\mathcal{E}}(t) = (\varepsilon_0/2)\exp[-(t/2T)^2 + i\varphi_e]$, which corresponds to the real electric field $E(t) = \varepsilon_0 \exp[-(t/2T)^2] \cos[2\pi(t/2T) + \varphi_e]$ which is plotted as a solid line in Fig. 1 for $\varphi_e = \pi/2$ (left) and $\varphi_e = \pi/2$ (right). If we consider that this pulse is used to excite a seven-photon process (for instance a seven-photon ionization), the driving function for that process is the seventh power of the field, which is plotted as a dotted line in Fig. 1. One can see that the different values of φ_e make a significant difference on how the process is driven. For $\varphi_e = 0$, the excitation is a single spike, as close approximation as practical to a δ -function. In the case of $\varphi_e = \pi/2$ (right), the excitation consists in a succession of positive and negative spikes.

Train of Pulses

The “ideal” mode-locked laser emits a train of identical pulses, at equal time interval. The period of the pulse train is τ_{RT} , defined as the separation between two successive envelopes. In the particular case that the pulse separation is an integer number of optical cycles $\tau_{RT} = NT = N/\nu$ (T being the light period and $\nu = \omega/(2\pi)$ the optical frequency) the successive pulses are identical. This will generally not be the case, and there will be a phase shift $\varphi_p = \omega\tau_{RT} \neq 2N\pi$ between successive pulses. The complex electric field of the total pulse train \tilde{E}_{pt} is

$$\tilde{E}_{pt}(t) = e^{i\omega t} [\tilde{\mathcal{E}}(t) + \tilde{\mathcal{E}}(t - \tau_{RT})e^{i\varphi_p} + \tilde{\mathcal{E}}(t - 2\tau_{RT})e^{2i\varphi_p} + \dots] \quad (6)$$

where $\tilde{\epsilon}(t) = \epsilon(t)e^{i\varphi_c}$ is the electric field of one particular pulse. The n th pulse has the phase factor $\exp[i(\varphi_c + n\varphi_p)]$, different from the previous and next pulse. To the change in phase between successive pulses φ_p , corresponds a frequency:

$$f_0 = \frac{1}{2\pi} \frac{\varphi_p}{\tau_{RT}} \quad (7)$$

This frequency is called the *carrier to envelope offset*. The *carrier to envelope offset* is an important parameter of pulse trains, where the change in phase from pulse to pulse is a measurable quantity, independent of the duration of the individual pulse in the train.

One can “idealize” to the extreme the concept of a pulse train, by considering an infinite train of δ -functions, equally spaced by the period of the train τ_{RT} , as shown in Fig. 2a. The Fourier transform of this ideal pulse train shown in Fig. 2b is an identical picture in the frequency domain: a comb of infinite extent (because the pulses were δ -function in time), with δ -function teeth (because of the infinite extent of the train).

Since the comb extends to infinity, there is no particular tooth that can be called an average frequency. Each mode ν_m of index m carries the same weight, and corresponds in the time domain to an infinite sine wave, which is a particular term of a Fourier series representation of δ -function. The first tooth at frequency $\nu_0 = f_0$ represents the carrier to envelope offset defined above. The corresponding carrier to envelope phase φ_c defined previously can be identified in the time domain, even with a train of δ -functions. The harmonic wave corresponding to the mode ν_2 is shown in Fig. 2a,* and the phase φ_e is identified as the phase at which each δ -function crosses the harmonic field. In the sketch of Fig. 2a, $\varphi_e = 0$ for the first pulse, and φ_p is then the carrier to envelope phase φ_e of the second pulse as indicated in the figure.

A somewhat more mundane train of pulses of finite duration τ^\dagger is sketched in Fig. 3a. In the frequency domain (Fig. 3b), the infinite pulse train is represented by a finite frequency comb. The

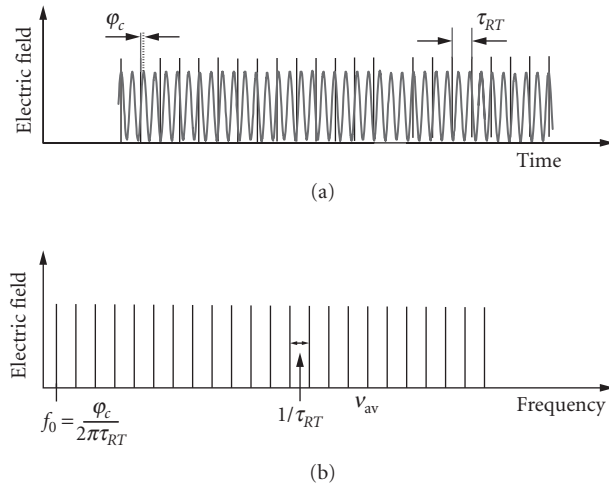


FIGURE 2 Idealized infinite train of δ -function pulses (a), and its Fourier transform (b). In (a), the carrier to envelope phase φ_c of the first pulse is assumed to be zero.

*This harmonic wave sketched is associated with ν_2 because there are two periods between pulses. In the Fourier spectrum of a train of δ -functions, any mode ν_n can be chosen as being the “average frequency”.

†When not otherwise specified, the pulse duration will be the full width at half maximum (FWHM) of the pulse intensity profile.

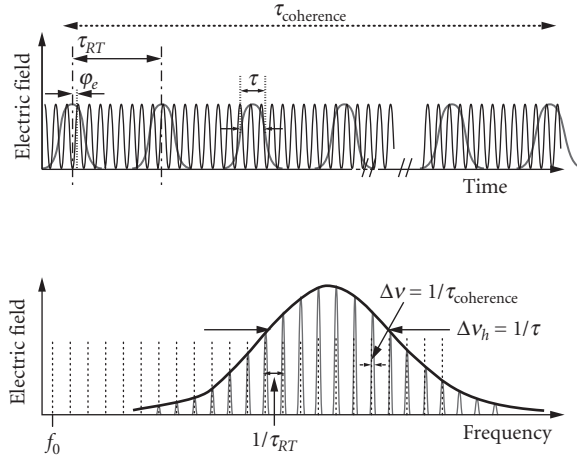


FIGURE 3 (a) Train of pulses of finite duration τ . The successive pulse envelope repeats every τ_{RT} . Within the coherence time of the train $T_{coherence}$, it is the same carrier at the optical frequency that is modulated by the successive envelopes. (b) The Fourier transform of the pulse train shown in (a).

envelope of the comb is the Fourier transform of the envelope of a single pulse of the train, thus of extension $\approx 1/\tau$. The teeth of the frequency comb are no longer δ -functions, but sharp peaks of width $1/\tau_c$, where τ_c is the coherence time of the pulse train. The carrier to envelope phase ϕ_p is indicated for a pulse of the time sequence in Fig. 3a. Note that this phase is changing from one pulse to the next. The rate of change $\phi_p/(2\pi\tau_{RT})$ is the frequency f_0 , which is indicated in the frequency picture by the lowest frequency tooth of the extension (dashed lines in Fig. 3b).

The angular frequency ω_m of the m th mode of the comb is given by

$$\omega_m = 2\pi f_0 + m \frac{2\pi}{\tau_{RT}} \quad (8)$$

In the case of a train of pulses of finite duration, the frequency f_0 is no longer a real tooth of the comb, but the first mode of an extension of the frequency comb beyond the pulse bandwidth as shown by the dotted line in (Fig. 3b).

It can be seen from this definition that f_0 is indeed the change of phase per round-trip between the envelope and the carrier. By definition of τ_{RT} , the pulse envelope peaks exactly at the same locations after one round-trip. With respect to this envelope, the shift of phase of the mode m is obtained by multiplying Eq. (8) by τ_{RT} :

$$\omega_m \tau_{RT} = 2\pi f_0 \tau_{RT} + 2m\pi \quad (9)$$

which, after substitution of the definition of f_0 Eq. (7), is indeed the phase ϕ_p defined earlier.

Soliton Solution and Steady-State Pulse Train

As mentioned in the introduction, if a laser is to generate a pulse of well defined duration and shape, there has to be compression and broadening mechanisms that balance each other, and lead to

a stable pulse. The mechanisms that lead to an emergence of a pulse out of noise in a laser cavity are usually dissipative, i.e., the pulse that emerges dissipates a minimum amount of energy by nonlinear loss mechanisms and extracts the maximum amount of gain from the active laser medium. We will consider in this section only nondissipative interaction that plays a dominant role for the stable formation of the shortest pulses. First we consider the evolution of a single pulse as it propagates through a cavity. Next we will study the formation of a pulse train with a similar nondissipative model.

Evolution of a Single Pulse in an “Ideal” Cavity When a laser is in continuous operation, the cavity gain and losses are in equilibrium. In the case of femtosecond mode-locked lasers, the major pulse-shaping mechanism is a combination of self-phase modulation and dispersion at each round-trip. The self-phase modulation results from a nonlinear index of refraction $n_2 I$ (I being the intensity in W/cm^2 , n_2 being the nonlinear index in cm^2/W of a nonlinear element of length ℓ in the cavity). The dispersion k'' results from the frequency dependence of the average index of refraction n_{av} , defined previously, and is characterized by the second derivative of a cavity averaged k vector with respect to frequency. In what follows, for simplicity, we will neglect higher-order terms in Kerr effect and in dispersion. The evolution of a pulse in the mode-locked laser cavity can be considered as a propagation (of a nondiffracting beam) through an infinite lossless medium, with a positive Kerr nonlinearity ($n_2 > 0$) and a negative dispersion (as can be introduced with intracavity prisms¹ or chirped mirrors² in the cavity). The pulse evolution generally converges toward a steady-state solution, designated as “solitons,” which can be explained as follows. The nonlinearity is responsible for spectral broadening and up-chirp. Because of the anomalous dispersion, $k'' < 0$, the high-frequency components produced in the trailing part of the pulse, travel faster than the low-frequency components of the pulse leading edge. Therefore, the tendency of pulse broadening owing to the exclusive action of group velocity dispersion can be counterbalanced. To determine the approximate parameters of that solution, let us assume a Gaussian pulse $\mathcal{E}(t) = \mathcal{E}_0 \exp[-(t/\tau_{G0})^2]$, and let us state that the chirp produced in the pulse center by the nonlinearity and the dispersion are of equal magnitude (but of opposite sign). Under this equilibrium condition the pulse circulates in the cavity without developing a frequency modulation and spectral broadening. The effect of group velocity dispersion is to create a pulse broadening and a down-chirp (in a medium of negative dispersion). The change (per round-trip) of the second derivative of the phase versus time, at the center of the pulse, is given, in a first-order approximation, by Ref. 3:

$$\Delta \left(\frac{\partial^2 \varphi(t)}{\partial t^2} \Big|_{(t=0)} \right) = \frac{4k''_{\text{av}}}{\tau_{G0}^4} P \quad (10)$$

where $k''_{\text{av}} = d^2 k_{\text{av}}/d\Omega^2$ is the second-order dispersion averaged over the cavity of perimeter P , which is < 0 for an optical element with negative dispersion. Assuming that the cavity contains an element with a nonlinear index $n = n_0 + n_2 I$ of length ℓ_{Kerr} , the phase induced by self-phase modulation, near the center of the Gaussian pulse, is

$$\Delta \varphi(t) = -k_{\text{NL}} \cdot \ell_{\text{Kerr}} \Big|_{(t=0)} = \frac{2\pi n_2}{\lambda} \ell_{\text{Kerr}} I \approx \frac{4\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda} \frac{t}{\tau_{G0}^2} \quad (11)$$

where we have used a quadratic approximation for the Gaussian near $t = 0$. Taking the second derivative yields the chirp induced by phase modulation at the pulse center:

$$\Delta \left(\frac{\partial^2 \varphi(t)}{\partial t^2} \right) \approx \frac{8\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda \tau_{G0}^2} \quad (12)$$

The peak intensity of the pulse (at $t = 0$) is $I_0 = \epsilon_0^2 / 2\eta$; $\eta = \sqrt{\mu_0 / \epsilon}$ being the characteristic impedance of the medium. Expressing that the chirps induced by phase modulation [Eq. (12)] and dispersion [Eq. (10)] should cancel each other leads to

$$I_0 \tau_{G0}^2 = -\frac{\lambda k''_{av}}{2\pi n_2} \frac{P}{\ell_{Kerr}} \quad (13)$$

This expression leads to a value for the pulse duration given by

$$\tau_{G0}^2 = \frac{\lambda |k''_{av}|}{2\pi n_2 I_0} \frac{P}{\ell_{Kerr}} \quad (14)$$

In a Ti:sapphire laser, $n_2 = 10.5 \cdot 10^{16}$ cm²/W; the crystal length is typically $\ell = 4$ mm, and a well-designed laser can produce a train of 10 fs pulses with an intracavity average power of the order of 10 W. These numbers can be used in Eq. (14) to determine the negative cavity dispersion k''_{av} required for stable laser operation.

How Unequally Spaced Modes Lead to a Perfect Frequency Comb The physical reason for a non-zero carrier to envelope offset f_0 is the dispersion of the laser cavity in which a pulse is circulating. The components of the laser cavity impose different group velocities on the pulse. We can define an average group velocity $v_g = P / \tau_{RT}$ of the pulse envelope, where P represents twice the length of a linear cavity, or the perimeter of a ring cavity. That group velocity is different from the phase velocity c/n_{av} (n_{av} being the linear index of refraction averaged over the laser components). The two quantities are related by

$$\frac{1}{v_g} = \frac{n_{av}}{c} + \frac{\omega}{c} \frac{dn_{av}}{d\Omega} \bigg|_{\omega} \quad (15)$$

Note that these quantities n_{av} and v_g are function of the spectral frequency of the pulse. The “ideal mode-locked laser” considered in this section already poses a conceptual dilemma. Mode-locking is generally described as putting the modes of a laser cavity in phase. If the cavity has dispersion, we have seen that the mode-comb issued from the laser does not start at zero frequency but with a frequency offset f_0 . Keeping in mind that a cavity with dispersion has unequally spaced modes, is contradictory to the fact that the frequency comb has rigorously equally spaced teeth.* To resolve this apparent contradiction, we will look at the pulse train formation, and discuss how an initially irregular set of modes can lead to a perfect frequency comb.

As shown above, a minimum negative cavity dispersion k''_{av} is required for stable mode-locked operation. Such a cavity dispersion implies that the index of refraction n_{av} is frequency (wavelength) dependent, hence the spacing of the cavity modes $c/[n_{av}(\Omega)P]$ varies across the pulse spectrum.

The laser is modeled by a circulating pulse, which enters a Kerr medium of thickness ℓ , resulting in phase modulation at each passage, and a medium that represents the linear dispersive properties of the cavity. We will assume that the balance of gain and losses maintains a constant Gaussian shape for the envelope of the circulating pulse. At each passage through the cavity, the phase of the pulse is modified in the time domain through the Kerr effect, and in the frequency domain through dispersion. We consider first the modulation in the time domain:

$$\varphi(t) = -k_{NL} \ell_{Kerr} = -\frac{2\pi n_2 \ell_{Kerr}}{\lambda} I_0 e^{-2(t/\tau_G)^2} \quad (16)$$

*A fact that has been verified experimentally with millihertz accuracy.⁴

where τ_G is the $1/e$ half-width of the pulse electric field envelope (the FWHM of the intensity is $\tau_p = \sqrt{2 \ln 2} \tau_G$). Ignoring at this point the influence of dispersion (which will be introduced after Fourier transformation into the frequency domain), the pulse *train* issued from the laser can be represented by

$$\sum_{q=0}^{\infty} \mathcal{E}(t - \tau_q) e^{iq\varphi(t - \tau_q)} e^{i\omega t} \quad (17)$$

where τ_q is the time of arrival of the center of gravity of the successive pulses. At this point τ_q is not set to any value. It is assumed here that at $t = 0$, the first pulse is unmodulated. Using a parabolic approximation for the Gaussian intensity profile, the time-dependent phase is

$$\varphi(t - \tau_q) \approx \frac{4\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda} \left(\frac{t - \tau_q}{\tau_G} \right)^2 = a \left(\frac{t - \tau_q}{\tau_G} \right)^2 \quad (18)$$

The Fourier transform of the pulse train given by Eq. (17) is

$$\mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega\tau_q} e^{iq\Delta\Omega^2\tau_k^2} \right] \quad (19)$$

where

$$\begin{aligned} \Delta\Omega &= \Omega - \omega \\ \mathcal{E}(\Delta\Omega) &= \frac{\mathcal{E}_0 \sqrt{\pi} \tau_G}{\sqrt[4]{1+a^2}} \exp\left\{ -\frac{\Delta\Omega^2 \tau_G^2}{4(1+a^2)} \right\} \\ \tau_k^2 &= \frac{a\tau_G^2}{4(1+a^2)} \end{aligned} \quad (20)$$

The width of the Gaussian pulse spectrum, broadened by the Kerr effect, is the inverse of the characteristic time τ_k . Let us now take dispersion into account. The operation representing the dispersion of the cavity is a product of the spectral field by $\exp[-ik_{\text{av}}(\Delta\Omega)P]$, where $-k_{\text{av}}(\Delta\Omega)P^*$ is the phase change per round-trip. The combined Kerr effect and dispersion, in the frequency domain, leads to the output spectral field:

$$\mathcal{E}_{\text{out}}(\Delta\Omega) = \mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega\tau_q} e^{iq\Delta\Omega^2\tau_k^2} e^{-iqk_{\text{av}}(\Delta\Omega)P} \right] \quad (21)$$

Expanding the wave vector $k_{\text{av}}(\Delta\Omega)$ in series, to second order:

$$\begin{aligned} k_{\text{av}}(\Delta\Omega)P &= k_{\text{av}}(\Delta\Omega=0)P + \Delta\Omega k'_{\text{av}}P + \frac{\Delta\Omega^2}{2} k''_{\text{av}}P \\ &= k_{\text{av}}(\Delta\Omega=0)P + \Delta\Omega \tau_{RT} + \frac{k''_{\text{av}}P}{2} \Delta\Omega^2 \end{aligned} \quad (22)$$

*In the argument of k_{av} , the light frequency ω is taken as origin ($\Delta\Omega = 0$) of the frequency scale.

where the derivatives k'_{av} and k''_{av} are calculated at the light frequency $\omega(\Delta\Omega=0)$. Note that $k'_{av}=1/v_g=\tau_{RT}/P$ [cf. Eq. (15)] are material properties independent of the index q , as is the cavity perimeter P . The modes of the cavity are not equally spaced. The parameter k'' characterizes the departure from equal spacing. Substituting (22) in Eq. (21),

$$\mathcal{E}_{out}(\Delta\Omega) = \mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega(\tau_q - q\tau_{RT})} e^{iq\Delta\Omega^2(\tau_k^2 - k''_{av}P/2)} \right] \quad (23)$$

The conditions

$$\tau_k^2 = -\frac{k''_{av}P}{2} \quad (24)$$

$$\tau_q = (q+1)\tau_{RT} \quad (25)$$

leads to modes that are exactly equally spaced. The inverse Fourier transform of the frequency comb becomes then

$$\tilde{\mathcal{E}}_{out}(t) = \mathcal{E}(t) + \mathcal{E}(t - \tau_{RT})e^{-ik_{av0}P} + \mathcal{E}(t - 2\tau_{RT})e^{-2k_{av0}P} + \dots \quad (26)$$

This last equation corresponds indeed to the description of the ideal frequency comb, with equally spaced pulses in time and frequency, and a carrier to envelope phase shift of $\varphi_p = -k_{av0}P$. In the case of small Kerr modulation, $a \ll 1$, it can easily be verified that the condition in Eq. (24) is identical to the soliton Eq. (14). Indeed, substituting

$$\tau_k^2 = \frac{a\tau_G^2}{4(1+a^2)} \approx \frac{4\pi n_2 I_0 \ell_{Kerr} \tau_G^2}{4\lambda} = \frac{k''_{av}P}{2} \quad (27)$$

which is indeed equivalent to Eq. (14). One can thus conclude that the mechanism that leads to an equal spacing for the teeth of the frequency comb emitted by the laser is the same Kerr effect responsible for creating maximum intracavity pulse compression.

20.3 PULSE EVOLUTION TOWARD STEADY STATE

A Simple Model

In the previous section we have considered the dispersive mechanisms that ultimately give the final shape in amplitude and phase to the steady-state pulse. This mechanism dominates in the sub-picosecond regime, where dissipative mechanisms have reached equilibrium. Other elements play a decisive role in initiating the mode-locking, which are usually referred to as the passive mode-locking elements. The latter can most often be represented by intensity-dependent intracavity loss. Larger losses at low intensity imply that the laser has less gain—and may be below threshold—for low-intensity continuous wave (cw) radiation than for pulses with higher peak intensity. This leads to the emergence of a pulse out of the amplified spontaneous emission noise of the laser. Rather than concentrating on the primary process of formation of a precursor of a pulse from random noise, let us follow the evolution of the pulse from its birth from noise until it has blossomed into a fully shaped stable laser pulse. In this intermediate stage of the evolution toward steady state, the main shaping elements are dissipative, as opposed to the purely dispersive interaction considered in the previous section. We will look for simple evolution equations for the pulse energy $W = \int_{-\infty}^{\infty} I(t) dt$, with $I(t)$ being the pulse intensity. The element responsible for saturable losses (gain) should have

typically a *linear* loss (gain) factor at low energies, and a *constant* loss (gain) at higher energies. We thus have, at low energies: $dW/dz = \mp \alpha W$. At large energies $W \gg W_s$: $dW/dz = \mp W_s$ where W_s is the saturation energy for the chosen geometry. The simplest differential equation to combine these two limits is

$$\frac{dW}{dz} = \alpha_g W_{sg} \left[1 - e^{W/W_{sg}} \right] \quad (28)$$

Equation (28) is written for a medium with a linear gain α_g and a saturation energy W_{sg} . It can be integrated to yield the energy W_2 at the end of the amplifier of thickness d_g , as a function of the input energy W_1 :

$$W_2 = G(W_2, W_1)W_1 = W_{sg} \ln \left\{ 1 - e^{\alpha_g d_g} (1 - e^{W_1/W_{sg}}) \right\} \quad (29)$$

A similar equation applies to the saturable absorber, with a negative absorption coefficient $-\alpha_a$ and a smaller saturation energy W_{sa} :

$$W_2 = A(W_2, W_1)W_1 = W_{sa} \ln \left\{ 1 - e^{\alpha_a d_a} (1 - e^{W_1/W_{sa}}) \right\} \quad (30)$$

The dominant linear loss element is the output coupler, with (intensity) reflectivity r . The transfer function for that element is simply

$$W_2 = L(W_2, W_1)W_1 = rW_1 \quad (31)$$

and the energy of the output pulse is $(1-r)W_1$. The evolution of the pulse energy in a single round-trip can be simply calculated from the product of all three transfer functions given by Eqs. (29), (30), and (31). For instance, if we consider a ring laser with the sequence: mirror, gain, and absorber, the pulse energy W_4 after the absorber is given by the product $A(W_4, W_3)G(W_3, W_2)L(W_2, W_1)$. One can also express the relation between the energy W_4 and the pulse energy W_1 before the output mirror by the algebraic relation:

$$1 + a[e^{W_4/W_{sa}} - 1] = \left\{ 1 + g[e^{rW_1/W_{sg}} - 1] \right\}^{W_{sg}/W_{sa}} \quad (32)$$

where $a = \exp\{-\alpha_a d_a\}$ is the *linear* small-signal attenuation of the passive element and $g = \exp\{\alpha_g d_g\}$ is the *linear* small-signal amplification.

High-Gain Oscillators

Unlike laser amplifiers, where it is desirable to use a gain medium with as high a saturation energy density as possible, mode-locked oscillators will often use high-gain laser media. These are opposite requirements: the larger the amplification cross section σ_g , the larger the gain $\alpha_g = \Delta N \sigma_g$, and the smaller the saturation energy density $W_s = \hbar \omega / (2\sigma_g)$. Both numbers a and g can be large, and the reflectivity of the output coupler r can be even lower than 50 percent. Examples are dye lasers, semiconductor lasers with tapered amplifiers, and to a smaller extent the Ti:sapphire laser. As a result, the order of the elements matters in the design of the laser, and in its performances. To illustrate this point, let us assume that the passive element is totally saturated in normal operation. In full saturation, the input energy $W(0)$ is related to the output energy $W(d)$ by

$$W(d) = W(0) - \alpha_a d_a W_{sa} \quad (33)$$

Given an initial energy W_1 , the energy W_4 for a single passage through three different sequences of the same elements is given below. For the sequence mirror-absorber-gain

$$e^{W_4/W_{sg}} = 1 - g \left[1 - e^{(rW_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \quad (34)$$

For the sequence absorber-mirror-gain

$$e^{W_4/W_{sg}} = 1 - g \left[1 - e^{r(W_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \quad (35)$$

For the sequence absorber-gain-mirror

$$e^{W_4/W_{sg}} = \left\{ 1 - g \left[1 - e^{(W_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \right\}^{1/r} \quad (36)$$

Let us take a numerical example for a high-gain system such as the flash-lamp pumped Nd:glass laser, with $W_{sa}/W_{sg} = 0.1$, $\alpha_a d_a = 1$, $\alpha_g d_g = 1.5$ and output coupling of $r = 0.5$. For an initial energy $W_1/W_{sg} = 0.5$, we find for the sequence absorber-mirror-gain $W_4/W_{sg} = 0.689$, and for the sequence absorber-gain-mirror $W_4/W_{sg} = 0.582$. The order of the elements, the relative saturation of the gain to that of the passive element, as well as the output coupling influence the stability and output power of the laser, as shown in Ref. 5. The evolution of the pulse in the cavity can be calculated by repeated applications of products of operations such as $A(W_4, W_3)G(W_3, W_2)L(W_2, W_1)$ for the sequence (passive element, gain, output coupler), starting from a minimum value of W_1 above threshold for pulsed operation, and recycling at each step the value of W_4 as the new input energy W_1 . Figure 4 shows the growth of intracavity pulse energy as a function of the round-trip index j , for different orders of the elements. The initial pulse energy is 1 percent of the saturation energy W_{sg} in the gain medium. The saturation energy and optical thickness of the absorbing medium are, respectively, $W_{sa} = 0.8W_{sg}$ and $\alpha_a d_a = 1.2$. The linear gain is $\alpha_g d_g = 1.5$ and the output coupling $r = 0.8$.

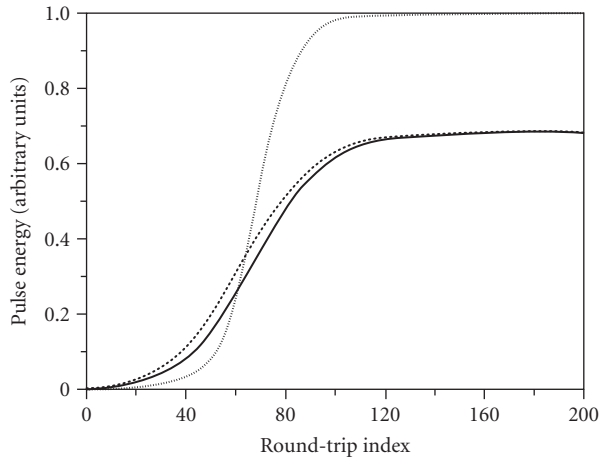


FIGURE 4 Intracavity pulse energy versus round-trip index j . The solid, dashed, and dotted lines correspond to the sequence *m*(irror)-*a*(bsorber)-*g*(ain), *a*-*g*-*m* and *a*-*m*-*g*, respectively.

The main point of this exercise is that the order of the elements is important, a fact confirmed by measurement on moderate-gain lasers such as $\text{Ti:Al}_2\text{O}_3$. Such dependence on the order of the elements indicates that analytical theories based on the approximation of infinitesimal change per element and per cavity round-trip are not quite adequate. Numerical codes have been developed that attempt to include all physical phenomena affecting pulse shape and duration. Unfortunately, the sheer number of these mechanisms makes it difficult to reach a physical understanding of the pulse-generation process, or even identify the essential parameters. Therefore, the most popular approach is to construct a simplified analytical model on a selected mechanism.

20.4 COUPLING CIRCULATING PULSES INSIDE A CAVITY

Pulse Train Interferometry

A pulse train combines the temporal resolution of a single pulse, and the spectral resolution of a cw beam. It is therefore not surprising that new interferometric technique can be developed exploiting these properties.

Because the ratio of pulse duration to the period of the train is generally less than $1:10^5$, interferences of pulse trains of different repetition rates will not be considered here. We will instead focus on a situation where two pulse trains of identical repetition rate, but different carrier to envelope phase, are made to interfere. It will be shown in the next section how pulse trains of identical repetition rate are generated. The experimental arrangement for pulse train interferometry is depicted in Fig. 5. The two pulse trains are combined by a beam splitter, and their relative delay adjusted in an optical delay line, in order to have superposition of the pulse envelopes on the detector. If the carrier to envelope phase is identical for both pulse trains, the detector will simply register a constant signal, with an amplitude dependent on the relative phase of the two trains at the detector.

If instead the two pulse trains have a different carrier to envelope phase, successive pulses will interfere differently. The envelope of the interfering pulse trains, as seen by the detector, will be modulated at the frequency $f_{01} - f_{02}$, where f_{01} and f_{02} are the carrier to envelope offsets of either pulse trains.

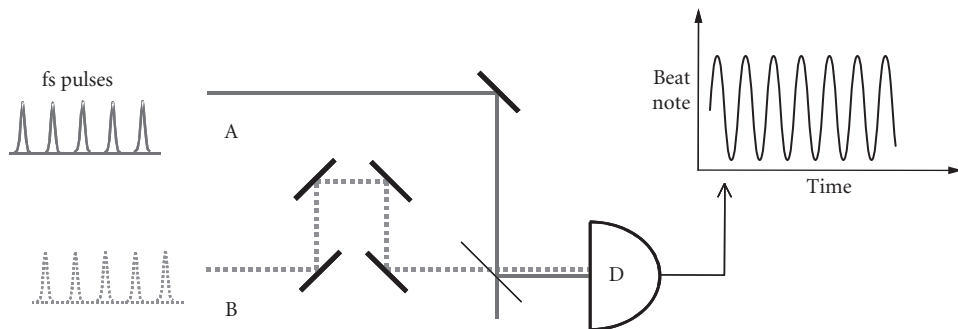


FIGURE 5 Interference of two pulse trains of the same repetition rate. An optical delay line is required to ensure temporal overlap of the pulses in either train.

Interwoven Pulse Trains Generated by Two Intracavity Pulses

We have seen the interference of pulse trains. These can be generated by a laser cavity. Inside the laser, there are then two or more pulses circulating in the cavity. The interactions of these two pulses will determine the relative properties of the pulse trains. The common picture of the mode locked laser is that of a resonator cavity with gain in which a single pulse circulates. It is interesting to consider, both from the point of view of fundamental understanding of the laser operation and applications to sensors, the situation where several pulses circulate in the cavity. Mode-locking with multiple pulses/round-trip is sometimes referred to as “harmonic mode-locking” 6–10. Techniques of harmonic mode-locking have been developed for telecommunications, where a pulse rate of over a GHz is desirable. The fundamental clock remains the round-trip time, which may typically be of the order of 100 MHz. If there are m pulses in the cavity, there will be m values of interpulse delay, which, in the frequency domain, will mean some splitting of the modes, which, is small, may be seen as a slight broadening of the tooth of the comb.

In the following we will consider only the case of two pulses circulating in the cavity, separated exactly by half the cavity round-trip time. Such a situation is encountered in bidirectional ring lasers, but also in linear cavities. Of particular interest is to determine the type of coupling that may or may not exist between the two pulses, and the resulting correlation between the two pulse trains. As will be the case in most sensor applications,¹¹ we will assume in the following that the two pulses experience a relative shift in phase δ at each round-trip.

Locking Two Pulse Trains by Backscattering

Whether in a ring or linear laser, the two circulating pulses will meet at two points of the cavity. Unless the meeting point is in vacuum, there will be some coupling introduced by the medium in which the pulses meet. The most common case is that of a medium with random scattering. Using the plane wave description of Eq. (2), the backscattering component $\tilde{r}_{ij} = r \exp(\theta_{ij})$ of scattering will couple the pulse with field envelope \tilde{e}_j into the pulse with field envelope \tilde{e}_i :

$$\begin{aligned}\frac{\partial \tilde{e}_1}{\partial t} &= i \frac{\delta}{2\tau_{RT}} \tilde{e}_1 + \frac{1}{\tau_{RT}} \tilde{r}_{12} \tilde{e}_2 \\ \frac{\partial \tilde{e}_2}{\partial t} &= \frac{1}{\tau_{RT}} \tilde{r}_{21} \tilde{e}_1 - i \frac{\delta}{2\tau_{RT}} \tilde{e}_2\end{aligned}\quad (37)$$

Of particular interest here is the impact of the coupling on the phase of the two fields, since the balance of saturable gain and losses will in general restore a steady-state value of the pulses energy. Expressing the fields in terms of amplitude and phase as in Eq. (2) in the system of Eqs. (37), and taking the difference of the imaginary parts, yields:

$$\frac{\partial(\varphi_2 - \varphi_1)}{\partial t} = \frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} + \frac{r}{\varepsilon_1 \varepsilon_2 \tau_{RT}} [\varepsilon_1^2 \sin(\theta_{21} - \psi) - \varepsilon_2^2 \sin(\theta_{12} + \psi)] \quad (38)$$

In the absence of coupling ($r=0$), the carrier frequency of the two pulses would differ by $\dot{\psi} = \delta/\tau_{RT}$, a frequency difference that can easily be detected by beating the two output pulse trains of the laser (corresponding to either pulse in the cavity) against each other on a detector.

Distributed Backscattering In presence of a sufficient coupling $r \neq 0$, there is generally a solution $\partial\psi/\partial t$ to Eq. (38), such that the two circulating pulses are identical, only differing by a phase factor. If we take for instance the particular case of $\theta_{12} = \theta_{21} = 0$, the constant phase difference is ψ_0 given by

$$\sin \psi_0 = \frac{2\varepsilon_1\varepsilon_2}{\varepsilon_1^2 + \varepsilon_2^2} \frac{\delta}{2r} \quad (39)$$

Any backscattering such that

$$r \geq \frac{\delta\varepsilon_1\varepsilon_2}{\varepsilon_1^2 + \varepsilon_2^2} \quad (40)$$

will lock the carrier frequency of the two waves to each other. This implies that the mode frequencies, the repetition rates, and the CEO of the two pulse trains are identical.

Interface Coupling A reciprocal backscattering, where $\theta_{12} = \theta_{21}$ is the norm when dealing with distributed scattering of a solid, liquid, or gaseous medium.^{12,13} The situation is different however in a short pulse laser, where the meeting points of the two pulses are localized rather than being distributed over the whole length of the laser resonator. In the case of the mode-locked laser, the backscattering can be due to an interface, in which case $\tilde{r}_{21} = -\tilde{r}_{12}^* = -\tilde{r}^*$; and $\theta_{21} = \theta_{12} + \pi = \theta + \pi$. This type of coupling does not prevent lock-in, since Eq. (38) becomes

$$\frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} - \frac{r}{\varepsilon_1\varepsilon_2\tau_{RT}} [\varepsilon_1^2 + \varepsilon_2^2] \sin(\theta + \psi) \quad (41)$$

which, for sufficient large r , still has a lock-in solution ψ_0 for which $\partial\psi/\partial t = 0$.

Phase Conjugated Coupling Not all couplings lead to identical mode frequencies of the two output pulse trains of the laser. In a phase conjugated coupling, a fraction r_c of the complex conjugate of one field is coupled into the other field. Such a phase conjugated coupling¹⁴ does preserve the phase identity of each intracavity pulse. The coupled equations for the two pulses are then

$$\begin{aligned} \frac{\partial\tilde{\varepsilon}_1}{\partial t} &= i\frac{\delta}{2\tau_{RT}}\tilde{\varepsilon}_1 + \frac{r_c}{\tau_{RT}}\tilde{\varepsilon}_2^* \\ \frac{\partial\tilde{\varepsilon}_2}{\partial t} &= \frac{r_c}{\tau_{RT}}\tilde{\varepsilon}_1^* - i\frac{\delta}{2\tau_{RT}}\tilde{\varepsilon}_2 \end{aligned} \quad (42)$$

Subtracting the imaginary parts of this equation:

$$\frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} \quad (43)$$

and there is no “lock-in” possible with this type of coupling.

Repetition Rate Coupling It has been observed that the repetition rate in both directions can be locked by a saturable absorber. The mechanism by which the average group velocity of the two pulses are locked to each other is described below. This mechanism leaves the carrier frequencies of the two pulses uncoupled.

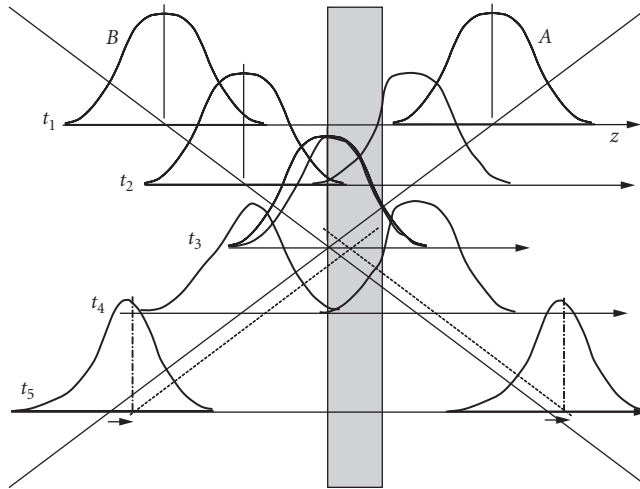


FIGURE 6 Representation of the intracavity pulses entering a saturable absorber. The pulses are plotted as a function of space (z) at successive times. The saturable absorber is initially at the left of the pulse crossing point. Because of mutual saturation, there is only significant absorption when only one of the pulses is present in the absorber. Therefore, the leading edge of pulse A is attenuated more, resulting in an apparent slowing down of the pulse. Similarly, the trailing edge of pulse B is absorbed more, resulting in an apparent acceleration of that pulse. The effect of the absorption combined with mutual saturation is to “pull” the pulse crossing point toward the center of the absorber.

To appreciate this envelope coupling, let us consider a saturable absorber of smaller longitudinal dimensions than the optical pulse, as sketched in Fig. 6. In the figure, the meeting point of the two pulses is on the left of the saturable absorber. Therefore, the pulse A entering from the left enters first the absorber, and its leading edge is attenuated. The absorption is saturated when the two pulses meet in the absorber. The pulse B coming from the left is still partly in the absorber when pulse A has left the absorber. Therefore, its tail will be more absorbed. The net effect is a shift of the center of gravity of both pulses, such that at the next round-trip they will meet closer to the middle of the absorber.

20.5 DESIGNS OF CAVITIES WITH TWO CIRCULATING PULSES

The properties of interwoven intracavity pulses have been discussed in the preceding sections. A few examples of laser cavity designs where two pulses are circulating independently are presented next.

Ring Dye Laser

In a bidirectional ring laser, two pulses circulate in opposite direction. The first realization was a dye laser,^{15,16} in which the gain is provided by a jet of Rhodamine 6G in ethylene glycol (pumped by an argon ion laser). A saturable absorber jet of DODCI* has three functions: (i) to mode-lock the laser,

*Di-oxa-di-carbo-cyanide-iodide.

(ii) to ensure bidirectional operation of the laser, and (iii) to define and maintain the pulse crossing point of the envelopes. Bidirectional operation is favored over unidirectional operation because, for the same pulse energy and duration, a single pulse will suffer more intracavity losses than two pulses creating a standing wave by crossing in the absorber.³ The pulse crossing point is set and maintained by the mechanism just described in Section “Locking Two Pulse Trains by Backscattering”. The saturable absorber has to be located approximately $1/4$ cavity perimeter ($P/4$) away from the gain jet, in order for each counter-circulating pulse to enter the gain medium at equal time interval ($P/2c$). The phase difference per round-trip between the clockwise (CW) and counter-clockwise (CCW) pulses is measured by combining the CW and CCW pulses on a detector.

Ti:Sapphire Ring Laser with Saturable Absorber

A similar cavity configuration has been used with a Ti:sapphire laser as a gain medium, and a saturable absorber jet of HITCI* for repetition rate synchronization between the two counter-circulating pulses. A sketch of such a cavity is shown in Fig. 7.

The nonlinear index of the gain medium results in a lensing effect, which can be approximated by a positive lens collocated with the gain medium. In mode-locked operation, Kerr lensing induces a positive lens at the location of the gain medium, which modifies the beam size distribution. For the empty cavity, the beam size versus position is represented by the solid line in the graph of Fig. 8. The beams size distribution modified by the self-lensing in the gain rod is indicated by the dotted line in the figure. An aperture located at the position A_1 will favor mode-locked operation, since the losses will decrease with intensity. An aperture located at A_2 will create increasing losses with intensity. This negative feedback stabilizes the mode-locked laser operation. Kerr lensing, which could be considered as an instantaneous saturable absorption,³ is the technique commonly used to generate the shortest pulses with Ti:sapphire lasers. It is however not a preferred technique for achieving stable bidirectional operation, when, as is typically the case, the active element for Kerr lensing is the gain medium. There is a competition in the gain medium between mutual Kerr lensing, favoring bidirectional operation (with the pulses crossing in the gain medium) and mutual gain saturation favoring unidirectional operation, with the latter generally dominating. An experimental study of a Kerr-lens mode-locked Ti:sapphire laser¹⁸ showed unidirectional operation, switching direction periodically (approximately every 0.1 second). The operation became bidirectional after insertion of a dilute saturable absorber jet inside the cavity.¹⁸

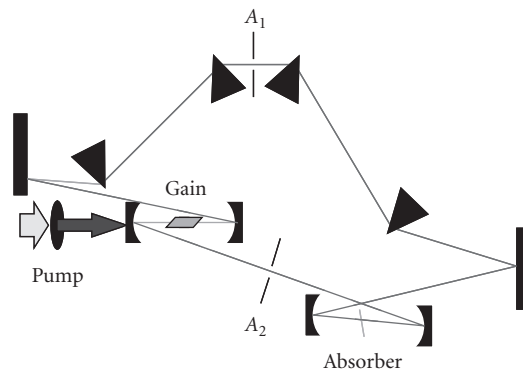


FIGURE 7 Ti:sapphire mode-locked ring laser mode-locked with a saturable absorber jet. Four prisms are used for the control of cavity dispersion.¹⁷

*Hexa-methyl-indo-carbo-cyanide-iodide.

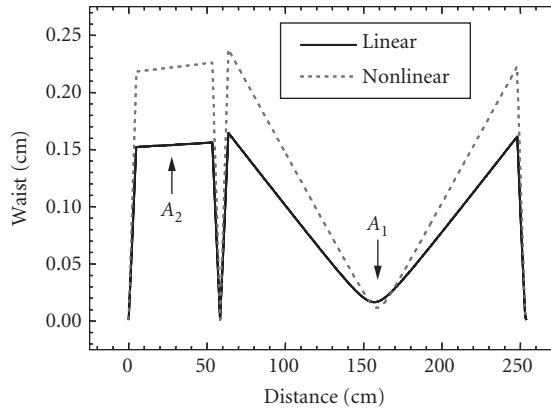


FIGURE 8 Calculation of the cavity mode as a function of position along the cavity, for the empty cavity (solid line) and the cavity modified by nonlinear lensing at the position of the gain medium (dotted line).

The use of a liquid saturable absorber flowing at high velocity (several meters per second) through a narrow nozzle (typically $100\ \mu\text{m}$ thick and 5 mm wide) is essential to ensure the absence of phase coupling between the two pulses. If the saturable absorber were a nonmoving solid, the coupling by scattering would result in a mutual locking of the carrier frequency of the two pulses, as discussed in section “Locking Two Pulse Trains by Backscattering.” In the case of the moving fluid of the absorbing dye jet, θ_{12} and θ_{21} are both random functions of time, varying much faster than ψ . Over the time scale that the variation of ψ is negligible, the last terms of Eq. (38) average to zero. Therefore, the dead band has been eliminated, as has been verified experimentally.¹⁹

Ring Laser with Additional Kerr Crystal

Another technique to achieve bidirectional operation for a Kerr-lens mode-locked laser^{20,21} is to insert a nonlinear crystal (for which the nonlinear phase shift is larger than that produced in the gain medium) $1/4$ cavity perimeter away from the gain medium. The laser cavity is sketched in Fig. 9.

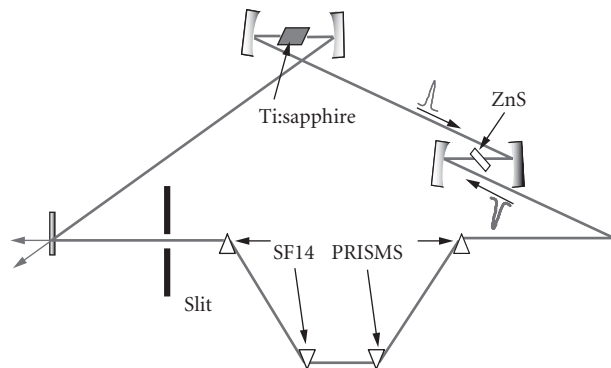


FIGURE 9 Ring laser mode-locked with a nonlinear crystal (ZnS).

The modification to the nonlinear index of refraction due to two counter-propagating fields of direction i and j is

$$n_2^i = n_2(I_i + 2I_j) \quad (44)$$

The factor 2 on the right hand side of Eq. (44) reflects the well-known act that cross-phase modulation is twice as effective as self-phase modulation for the same intensity.²² Assuming equal intensity in the counter-propagating fields and similar waists for single pulse operation versus bidirectional operation yields the following approximate relationship for the nonlinear index of refraction in the ZnS crystal for single pulse (unidirectional) versus double pulse (bidirectional) operation:

$$\Delta n_{\text{bidirectional}} \sim \frac{3}{2} \Delta n_{\text{single}} \left[\frac{\text{interaction length}}{\text{crystal length}} \right] \quad (45)$$

Because the overlap region of the pulse envelopes is approximately the pulsewidth/ c , the mutual Kerr-lens will dominate only if the length of the Kerr medium is not much longer than the pulsewidth. Depending on the pulsewidth and crystal thickness, there may be enough nonlinearity to distinguish between single-pulse operation and bidirectional operation. To enhance the mutual Kerr effect, a crystal of ZnS was chosen, because its nonlinear index is 50 times larger than that of Ti:sapphire.²¹ Pulses as short as 60 fs were generated, meeting in the ZnS crystal used as a nonlinear element. Because the two countercirculating pulses have different intensities, a differential phase shift between the two directions results in a difference of cavity modes of 60 kHz. In addition to the relative shift of the modes, this laser system has the additional complexity that the center of spectral envelope of either countercirculating pulse is shifted by 2 nm.

Imperfection in the ZnS crystal used as nonlinear element resulted in coupling of the beams by scattering. The amount of coupling can be measured through the spectrum of the beat note of the two pulse trains. Such a spectrum shows for instance a fundamental at 60 kHz and harmonics at 120 kHz and 180 kHz. Such a spectrum can lead to the lock-in frequency between the two beams,²³ which in the case of the experiment cited is 1.8 kHz.²¹

Linear Lasers

Considering the ring laser sketched on the top of Fig. 7, it is possible to visualize stretching the cavity by the two mirrors at the extreme left and right, while keeping the perimeter constant. The limit of the stretched out ring is a linear cavity in which two pulses circulate.

As an example, let us consider a linear cavity used to measure with high accuracy the electro-optic coefficient.²⁴ The laser cavity is similar to the typical linear cavity mode-locked Ti:sapphire laser, but with a saturable absorber (a jet of HITCI dye dissolved in ethylene glycol) placed in the center of the cavity, as sketched in Fig. 10. As in the case of the ring laser discussed above, Kerr-lens mode-locking does not appear to be possible with double pulse operation. Instead, a saturable absorber is positioned in the middle of the cavity by translating one of the end mirrors. The distance that the end mirror can be translated while maintaining double pulses is about 2 cm, in excess of the pulse length of approximately 0.6 mm (2-ps pulses). The 2-cm distance is a 120-ps delay and corresponds to the lifetime of the dye. The dye concentration is not a critical parameter and can be varied over a broad range without affecting the performance of the laser. The 140-MHz output from one end of the laser, detected on a fast photodiode, is filtered, its frequency divided by 2 in an ECL logic. The resulting 70-MHz signal is which yielded a 70-MHz sinusoidal signal. Finally the signal is amplified again and applied to synchronize the measurement to be performed. In the case of measurement of an electro-optic coefficient, the 70-MHz signal is applied directly to electrode on the crystal to be measured, in parallel with a 50-ohm terminator.²⁴

Whether to use a ring or linear cavity depends on the quantity to be measured. The ring laser is sensitive to rotation, and to fresnel drag in one of its arms.²⁵ Without any rotation or modulation a

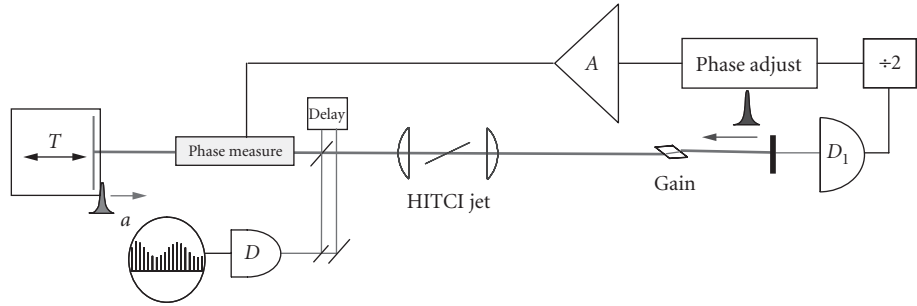


FIGURE 10 Linear laser mode-locked with a saturable absorber to produce two pulses/cavity round-trip. The end mirror on the left is on a translation stage T , in order to set the position of the saturable absorber in the middle of the cavity. The saturable absorber is a jet of HITCI dissolved in ethylene glycol, between two focusing elements. The output pulse train, recorded on a photodiode, reduced to half frequency, is used to synchronize a phase measurement. Each of the two intracavity pulses experience a different phase shift per round-trip, hence a different carrier frequency. The two intracavity pulses are extracted from the cavity with a beam splitter, and recombined after an appropriate optical delay. The two interfering pulse trains show a beat signal on the detector D .

mode-locked ring laser normally has a beat frequency offset of at least 100 Hz and often as high as 100 kHz.²⁶ This is a result of the asymmetry in the CW and CCW pulse. Because of the nonlinear intracavity elements, the order in which the pulse encounters the optical elements will affect the pulsewidth and pulse amplitude.^{3,27,28} Any variation in pulse amplitude or pulsewidth will be seen as a beat signal. Since the pulses in a linear cavity travel through the same optical elements in the same order, there is no asymmetry. Therefore, one advantage of a linear cavity versus a ring geometry is the improvement in the frequency offset.

As the electronic delay of the signal applied to the sample is varied, the beat note shows a sinusoidal dependence, as shown in Fig. 11, which is a plot of the beat frequency versus the delay. The optimum timing occurs when one pulse sees a voltage on the sample of $+V_0$ and the second pulse sees a voltage of $-V_0$ at the sample. The line plotted in Fig. 11 is not a fit, but a plot of $V_0 \left| \sin\left(\frac{2\pi c}{2L} \tau - \phi_0\right) \right|$, where the fixed phase ϕ_0 was the only free parameter.

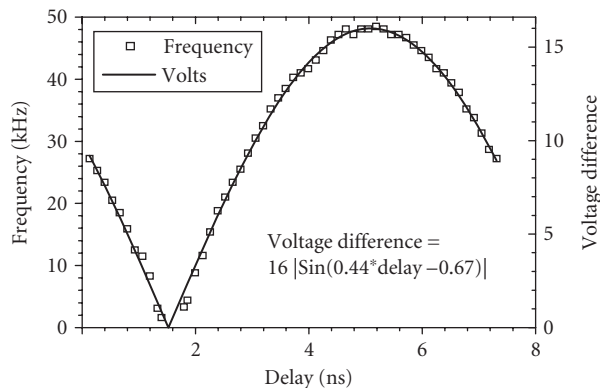


FIGURE 11 Delay dependence of the beat note frequency.

Optical Parametric Oscillators

The coupling between oppositely circulating beams in a ring laser is eliminated if a laser operates with ultrashort pulses* circulating in opposite direction and crossing in a nonscattering medium such as vacuum or air. It is however quite a challenge to find means to couple the pulse envelopes, without introducing any phase coupling. One solution described in the previous sections is to couple the pulses in amplitude in a medium that moves transversally to the beam, such as a jet of liquid saturable absorber. Another solution is the synchronously pumped optical parametric oscillator (OPO), which offers the possibility to decouple relative phase and repetition rates of the oscillating signals, without the need for any moving element.

A simple configuration is that of an OPO-pumped extracavity by a Ti:sapphire mode-locked laser, as sketched in Fig. 12. The position of the crossing point of the two circulating pulses in the OPO is simply determined by the timing of the pump pulses, rather than by a saturable absorber¹⁹ or a nonlinear crystal.²¹ The mode frequencies are still set by the cavity. Another advantage over other systems is the tunability, which is important for applications such as detecting ultra low magnetic fields, where the laser radiation has to be tuned to a narrow atomic transition.

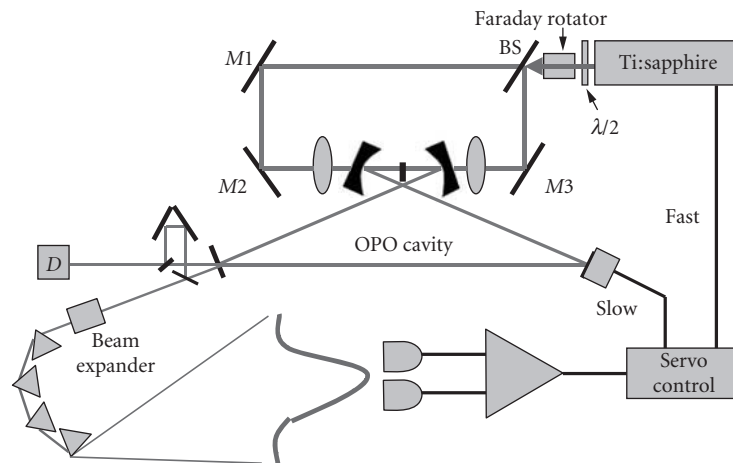


FIGURE 12 Illustration of the OPO cavity pumped by the Ti:sapphire laser. The reflected and transmitted parts of the beam splitter BS are focused into the periodically poled lithium niobate crystal via the two branches of an antiresonant ring. The beams should destructively interfere at the antiresonant ring output, which is monitored with a CCD. Since all the radiation is, in exact alignment, reflected back into the laser, a two-stage optical isolator is required to prevent disruption of the mode-locked operation. The nonlinear crystal is a 0.8-mm-long periodically poled lithium niobate (PPLN) crystal with a period of $19.75 \mu\text{m}$, temperature stabilized at 353 K to prevent photorefractive damage and achieve quasi-phase-matching condition for generation of a signal of $1.35 \mu\text{m}$ with an average power of 30 mW per direction. The difference between the two optical paths from the beam splitter to the crystal determines the crossing point of the signal pulses in the OPO cavity. The two output pulses are made to interfere on a detector *D* after an optical delay line brings them in coincidence. A four-prism sequence sends the pulse spectrum onto a pair of detectors. The difference between the two detected signals is amplified and applied to a piezo to stabilize the cavity length.

*Pulse duration τ in the range from femtoseconds to a few picoseconds, pulse length $C\tau$ much shorter than any linear dimension in the cavity.

A periodically poled LiNbO_3 crystal (PPLN) is excited by a “pump pulse” at an optical frequency ω_p to provide gain at a “signal” frequency ω_s through the process $\omega_p = \omega_s + \omega_i$ where ω_i is the “idler.” The OPO is an oscillator which uses the gated gain at ω_s . Therefore, as opposed to a conventional gain medium, in an OPO, the timing, direction, and position of the gain are determined by the time of arrival, k vector and focal spot location of the pump pulse in the parametric crystal. A Ti:sapphire laser provides a train of gating pump pulses of 200-fs duration at 789 nm, 143-MHz repetition rate, and 400-mW average power. The cavity of the pump laser is in a ring configuration, and operates unidirectionally. Such a configuration is less sensitive to feedback than a linear cavity. A double stage Faraday isolator (providing -60 -db isolation) is still required to prevent the feedback from the OPO antiresonant ring pumping arrangement (Fig. 12)^{29,30} from destroying the mode-locked operation. The sensitivity of the OPO wavelength to cavity mismatch can be exploited to stabilize the synchronously pumped OPO. Indeed, since the repetition rate of the OPO is fixed by the pump laser, the signal wavelength will adjust to a value for which the round-trip rate matches the pump rate.^{31,32} As a result, any fluctuation of OPO cavity length relative to that of the pump cavity will be translated into a change in wavelength of the OPO laser. In the arrangement sketched in Fig. 12, motions of the spectrum are detected by spectrally dispersing (with 4 prisms) an expanded output from the counter-clockwise OPO beam. The signal spectrum, centered at $1.35 \mu\text{m}$, is split into two parts and collected by a pair of lenses into two infrared photodiodes (Fig. 12). The difference signal of the two detectors monitoring two spectral components on either side of the pulse spectrum is sent through a high-gain amplifier ($\omega_{3db} = 1 \text{ kHz}$) to drive piezoelectric transducers (PZT) translating an OPO (slow servo loop) and a Ti:sapphire mirror (fast servo loop).

The beat note observed with the configuration of Fig. 12 has a bandwidth of tens of kilohertz,³³ because of the extreme (nanometer) sensitivity of the OPO to the pump spot position, as demonstrated experimentally in reference.³⁴ The beat note bandwidth is thus fundamentally due to fluctuations in the gain spot position for either circulating pulse, due to the beam pointing instability of the pump laser. The basic remedy is to make the two pump spots part of the same spatial mode of a cavity. One solution that has been implemented³⁴ is to insert the OPO crystal inside the cavity of the pump laser. Implementation of an OPO pumped intracavity by a linear Ti:sapphire laser is shown in Fig. 13. Four LaKL21 prisms are incorporated in the pump cavity to compensate the group velocity dispersion (GVD) from the Ti:sapphire crystal, the PPLN crystal, and other intracavity elements such as lenses and mirrors. This four-prisms configuration was necessitated by the desire to have large GVD compensation (needed because of the large positive GVD of LiNbO_3) and a reasonably short cavity length ($1/2$ of the perimeter of the OPO cavity). Two quantum wells (MQW) of AlGaAs on top of a mirror structure are used in the cavity as a saturable absorber to mode-lock the laser. The

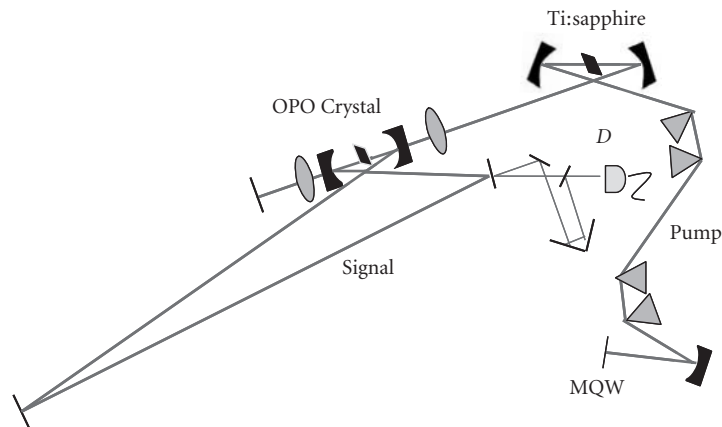


FIGURE 13 Illustration of the intracavity OPO pumped by the Ti:sapphire laser. The main control of the GVD compensation is the prism spacing L_2 .

Ti:sapphire laser radiation consists of 200-fs pulses centered at 785 nm, repetition rate of 95 MHz. The OPO crystal is a 3-mm-long Brewster cut PPLN crystal (HC Photonics, Taiwan) with a period of 19.4 μm (quasi-phase matching for signal near 1.36 μm), which is temperature stabilized at 408 K to prevent photorefractive damage. Attempts to use OPO crystals cut near normal incidence and antireflection coated failed, because the Ti:sapphire mode-locked operation was prevented by the smallest feedback from the antireflection coating. This feedback problem was completely eliminated through the use of a Brewster angle cut, but at the price of a considerably more difficult alignment procedure. Because of the Brewster angle cut of LiNbO_3 , the idler, the pump radiation and its second harmonic all exited the crystal at a different angle, and could not be used for alignment of the OPO cavity (the OPO cavity mirrors were reflecting at the signal wavelength of 1.4 μm and at the second harmonic of the pump).

20.6 ANALOGY OF A TWO-LEVEL SYSTEM

In the previous section, some methods for mode-locking a laser with two intracavity pulses were described. Inside the cavity, these two pulses can couple to each other in diverse ways which were analyzed in Sec. 20.4. The purpose of circulating two pulses in a laser cavity is that intracavity perturbations created by the quantity to be measured will alter the carrier frequency of a pulse, a shift in frequency that can easily be detected by interfering the output pulse trains. The laser itself is used as an interferometer, with the remarkable properties that the phase shift occurring in the cavity is transformed in a frequency shift. There are numerous factors influencing the accuracy as well as the sensitivity of a measurement performed intracavity. A better understanding of the two pulse per cavity laser can be reached by noting the complete analogy with a quantum mechanical two-level system. In the quantum mechanical situation an atomic or molecular system can be in one or two quantum states $|k\rangle$, with $k = 1$ and 2, of energy $\pm\omega_0/2$. Each of these states correspond to pulse $|2\rangle$ and pulse $|1\rangle$ in the laser cavity. For instance, in a ring laser, one of the states would correspond to a counterclockwise circulating pulse, the other to the clockwise circulating pulse. The interaction of a two-level system with a near resonant field is the most thoroughly studied problem in atomic and molecular physics. Techniques developed to achieve sublinewidth resolution in atomic physics may be transposed to the laser situation, and, thanks to the analogy, lead to methods to enhance the resolution of intracavity laser sensors.

Review of Coherent Interaction of Two-Level Systems

Considering the case of two level with a dipole allowed transition, in presence of a near resonant electromagnetic field $E = 1/2\tilde{E}(t)\exp(i\omega t) + \text{c.c.}$ In presence of this electric field, the state of the atomic/molecular system is described by the wavefunction ψ , a solution of the time-dependent Schrödinger equation:

$$H\psi = i\hbar \frac{\partial\psi}{\partial t} \quad (46)$$

with the total Hamiltonian given by

$$H = H_0 + H' = H_0 - p \cdot E(t) \quad (47)$$

where p is the dipole moment. In the standard technique for solving time dependent problems, the wave function ψ is written as a linear combination of the basis functions $|k\rangle$:

$$\psi(t) = \sum_k a_k(t) |k\rangle \quad (48)$$

This expression for ψ is inserted in the time dependent Schrödinger Eq. (46). Taking into account the normalization conditions for the basis functions ψ_k , one finds the coefficients a_k have to satisfy the following set of differential equations:

$$\frac{da_k}{dt} = -i\omega_k a_k + \sum_j \frac{i}{2\hbar} p_{k,j} [\tilde{\mathcal{E}} e^{i\omega_j t} + \tilde{\mathcal{E}}^* e^{-i\omega_j t}] a_j \quad (49)$$

where $p_{k,j}$ are the components of the dipole coupling matrix for the transition $k \rightarrow j$, and a_k are the amplitudes of the eigenstates. Phase and amplitude relaxation have been neglected so far and will be introduced later. It should be noted that Eq. (49) is of a quite general nature, is ideally suited to numerical integration, and is not limited to two level systems. Similarly, the laser analogy can be extended to lasers with more than two intracavity pulses. We will consider here only two levels, with a very small detuning $\Delta\omega = \omega_0 - \omega \ll \omega_0$:

$$\Delta\omega = \omega_0 - \omega \quad (50)$$

Consistent with the approximation of small detuning, we replace the set of coefficients a_k , which have temporal variations at optical frequencies, by the “slowly varying” set of coefficients c_k , using the transformation:

$$a_k = e^{-ik\omega_0 t} c_k \quad (51)$$

Inserting in the pair of Eqs. (49), leads to the pair of differential equations for the two coefficients c_k :

$$\frac{d}{dt} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} i\frac{\Delta\omega}{2} & i\frac{1}{2\hbar} p \tilde{\mathcal{E}} \\ -i\frac{1}{2\hbar} p \tilde{\mathcal{E}}^* & -i\frac{\Delta\omega}{2} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \quad (52)$$

where $\kappa|\tilde{\mathcal{E}}| = (p/\hbar)|\tilde{\mathcal{E}}|$ (p being the dipole moment of the single photon transition) is the Rabi frequency.

The Laser as a Two-Level System

The two-level system considered above is isolated, hence the total population is conserved. The analogue of the electromagnetic field coupling states $|1\rangle$ and $|2\rangle$ is a conservative intracavity coupling $\tilde{r}_{12} = -\tilde{r}_{21}^*$. Such a type of coupling can only be considered in the case of mode-locked lasers, where the localization of the radiation in the cavity enables one to select a truly conservative coupling. The coupling, localized at the crossing point of the two circulating pulses, can be produced by the back-scattering at a dielectric interface[†] between two media a and b , for which $\tilde{r}_{ab} = \tilde{r}$ and $\tilde{r}_{ba} = -\tilde{r}^*$. It can easily be verified that the total intensity change introduced by this coupling is zero, as expected for a conservative coupling. In fact, the phase relation between the two reflections at either sides of the interface is a consequence of energy conservation.

In the analogy of the laser, the coefficients $c_i(t)$ correspond to the complex field amplitudes $\tilde{\mathcal{E}}_i$ (the tilde indicating a complex quantity) of each pulse circulating in the ring cavity (round-trip

[†]In the case of a linear laser with two pulses/round-trip, the conservative coupling is only possible when there are two crossing points in the cavity, and that the interface is located at one of the crossing points.

time τ_{RT}). The state of the system is also defined by $\psi(t) = \tilde{\epsilon}_1(t)|1\rangle + \tilde{\epsilon}_2(t)|2\rangle$. The evolution equation of these fields are

$$\frac{d}{dt} \begin{pmatrix} \tilde{\epsilon}_1 \\ \tilde{\epsilon}_2 \end{pmatrix} = \frac{1}{\tau_{RT}} \begin{pmatrix} \tilde{r}_{11} & \tilde{r}_{12} \\ \tilde{r}_{21} & \tilde{r}_{22} \end{pmatrix} \begin{pmatrix} \tilde{\epsilon}_1 \\ \tilde{\epsilon}_2 \end{pmatrix} = \frac{1}{\tau_{RT}} ||R|| \cdot ||E|| \quad (53)$$

In order to have an equivalence between Eqs. (52) and (53), the matrix $||R||$ should be *Anti-Hermitian*, which, in addition to the condition $\tilde{r}_{21} = -\tilde{r}_{12}^*$, imposes and that \tilde{r}_{kk} be purely imaginary. It can also easily be verified that this is the only form of interaction matrix for which energy is conserved $d/dt(|\tilde{\epsilon}_1|^2 + |\tilde{\epsilon}_2|^2) = 0$. The real parts of the diagonal elements of the matrix $||R||$ represent gain and loss in the cavity. In steady state, the gain and loss are in equilibrium, and the real parts of \tilde{r}_{kk} are zero. A gain (or absorber) with a recovery (relaxation) time longer than $\tau_{RT}/2$ will cause transients in population. The laser equivalent to the detuning $\Delta\omega$ is a differential phase shift for the pulses $|1\rangle$ and $|2\rangle$ in the cavity. Such a differential phase shift is introduced either by rotation in a ring laser,^{15,16} or with an electro-optic modulator in a linear laser.²⁴ In the latter case, the electro-optic phase modulator imposes an opposite phase shift ($\Delta\phi/2$ and $-\Delta\phi/2$) for either pulse, thereby modifying the resonance of the cavity for the pulse $\tilde{\epsilon}_1$ by $\Delta\omega/2 = \Delta\phi/(2\tau_{RT})$, and for pulse $\tilde{\epsilon}_2$ by $-\Delta\omega/2 = -\Delta\phi/(2\tau_{RT})$. These detuning terms contribute to the diagonal terms of the matrix $||R||$: $\tilde{r}_{11} = -\tilde{r}_{22} = i\Delta\phi/2$.

The analogy between the laser and a two-level system applies also to the set of density matrix equations. These can be obtained by rewriting Eq. (53) in terms of the intensities in either sense of rotation $\rho_{22} = \tilde{\epsilon}_2 \tilde{\epsilon}_2^*$ and $\rho_{11} = \tilde{\epsilon}_1 \tilde{\epsilon}_1^*$, and the quantities $\rho_{12} = \tilde{\epsilon}_1 \tilde{\epsilon}_2^*$ and $\rho_{21} = \tilde{\epsilon}_2 \tilde{\epsilon}_1^*$:

$$\frac{d(\rho_{22} - \rho_{11})}{dt / \tau_{RT}} = -4\text{Re}(\tilde{r}_{12}\rho_{21}) - \frac{(\rho_{22} - \rho_{11})}{T_1} \quad (54)$$

$$\frac{d\rho_{21}}{dt / \tau_{RT}} = -i\Delta\omega\tau_{RT}\rho_{21} + \tilde{r}_{12}^*(\rho_{22} - \rho_{11}) - \frac{\rho_{21}}{T_2} \quad (55)$$

where, as in the case of the two-level system interacting with a near resonant field, phenomenological relaxation times T_1 and T_2 have been introduced. One recognizes here Bloch's equation for a two-level system driven off-resonance by a step function Rabi frequency of amplitude \tilde{r}/τ_{RT} .³⁵ The difference in intensities $(\rho_{22} - \rho_{11})$ is the direct analogue of the population difference between the two levels. The off-diagonal matrix element ρ_{21} is the interference signal obtained by beating the two outputs of the laser on a detector. As in the case of the two-level system, one can introduce phenomenological relaxation times T_1 for the energy relaxation (diagonal matrix element) and T_2 for the coherence relaxation (off-diagonal matrix elements). As for the quantum mechanical two-level system, $1/T_2$ is the homogeneous component of the linewidth of the beat note between the two pulse trains. There is also an "inhomogeneous" component to that linewidth, which has as physical origin the mechanical vibration of the laser components, causing random fluctuations of the beat note. Because of mechanical vibrations, each pulse sees random differences in the cavity length caused by mirror motion over a time of $\tau_{RT}/2$.

Table 1 summarizes the main points of the analogy between a laser with two pulses/cavity and the coherent interaction of a two-level system with a near resonant electromagnetic field.

Experimental Demonstration of the Analogy

The most typical manifestation of a two-level system interacting with a step function electromagnetic pulse is Rabi cycling, which is a periodic transfer of population from one state to the other. To observe such a periodic transfer, the system should be in one of the two states at $t = 0$. One method to prepare the ring laser with one state (direction) dominating, is to feedback one direction into the other outside of the cavity (Fig. 14). The output pulse from one direction is extracted, and fed back (<1 percent) with a mirror, after appropriate optical delay, into the opposite direction. By using a fast switch (turn-off time of less than the cavity round-trip time of 10 ns) at the Pockel's cell, the coupling can be turned off, to let the counter-circulating fields evolve in the cavity.

TABLE 1 Summary of the Analogy between a Two-Level System and the Laser with Two Circulating Intracavity Pulses

	Two-Level System	Laser
Basic states	$ k\rangle; k=1, 2$	$ k\rangle; k=1, 2$
corresponding to	energy level $\pm \frac{\omega_0}{2}$	Intracavity pulses 1, 2 selected by geometry
Coupling through	Near-resonant E-field at ω	Backscattering at interface
Detuning	$\Delta\omega = \omega_0 - \omega$	$\Delta\omega = \Delta\phi \tau_{RT}$
Slowly varying	$\Delta\omega \ll \omega$	$\Delta\omega \ll 1/\tau_{RT}$
Wave function	$\psi(t) = a_1(t) 1\rangle + a_2(t) 2\rangle$	$\psi(t) = \varepsilon_1(t) 1\rangle + \varepsilon_2(t) 2\rangle$
Density matrix	$\rho_{kk} = a_k a_k^*$ Populations	$\rho_{kk} = \tilde{\varepsilon}_k \varepsilon_k^*$ (Intensities)
elements	$\rho_{ij} = -a_i a_j^*$	$\rho_{ij} = \tilde{\varepsilon}_i \varepsilon_j^*$ (beat signal)

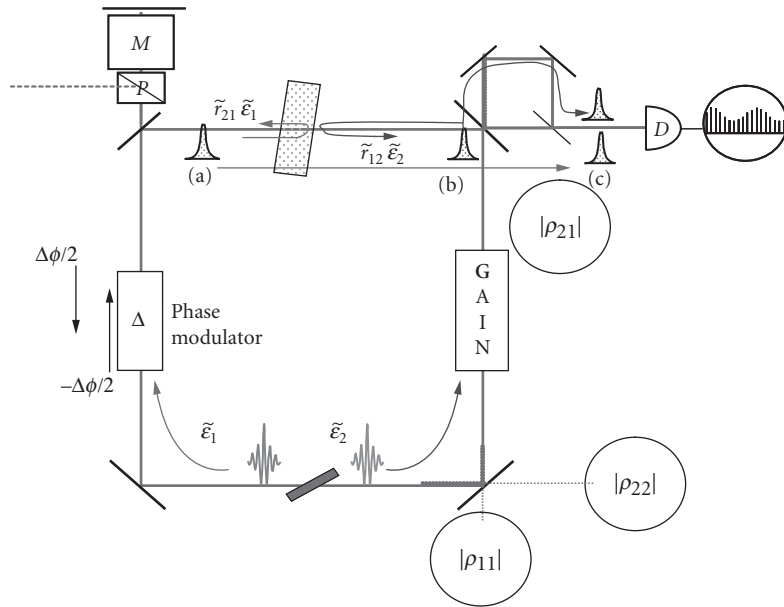


FIGURE 14 Sketch of the ring laser used to demonstrate the analogy with a two-level system. In the bidirectional mode-locked ring laser two circulating pulses meet in a saturable absorber jet. Three successive positions: (a), (b), and (c) of the two pulses are shown. An interface, positioned at or near the opposite crossing point of the two pulses, controls the amplitude of the coupling parameter \tilde{r}_{ij} . The circulating intensities in the laser, measured for each direction by quadratic detectors, are the diagonal elements (populations) of the density matrix of the equivalent two-level system. The absence of phase modulation corresponds to the two levels being on resonance, driven at the Rabi frequency $(p/\hbar)\varepsilon$ by a resonant field (the Rabi frequency $(p/\hbar)\varepsilon = \kappa\varepsilon$ corresponds to the frequency r_{12}/τ_{RT} in the ring laser analogy). The backscattering at the interface provides thus coherent coupling (Rabi cycling) between the two states, while other noncoherent decays tend to equalize the population in the two directions, and washes out the phase information. The detuning $\Delta\omega$ corresponds to the phase difference per round-trip $\Delta\phi/\tau_{RT}$, imposed by an electro-optic phase modulator driven exactly at the cavity round-trip time. A beat-note detector measuring the interference between the two fields, records the off-diagonal matrix element. A combination of a Pockel's cell M and polarizer P controls a feedback of the clockwise pulse into the counterclockwise one.

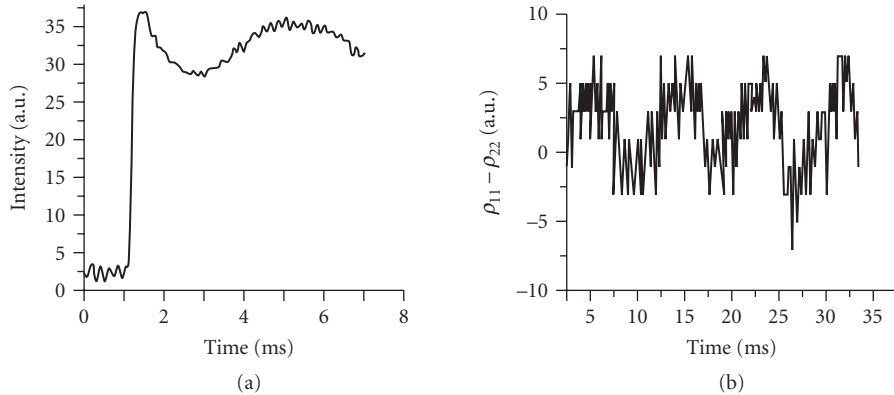


FIGURE 15 The evolution of the intensities are shown after switching the Pockels' cell. (a) The counterclockwise direction is shown—the intensity at clockwise direction is 180° out of phase with this graph, with population dropping from the maximum initial value. The fast initial transient reflects the gain and cavity dynamics associated with the sudden change in cavity loss at the switching time. Thereafter, a slow oscillation due to population transfer or Rabi oscillation between two directions is observed. (b) Population difference showing the Rabi cycling.

Rabi Cycling on Resonance In the measurements that follow, the system is “at resonance”; (i.e., $\Delta\omega=0$). An example of “Rabi cycling” is shown in Fig. 15. The counterclockwise intensity (ρ_{22}) is plotted as a function of time [(Fig. 15a)]. The clockwise intensity ρ_{11} (not shown) is complementary. The system is prepared so that the ρ_{11} is initially populated ($\rho_{11}=0.8, \rho_{22}=0.2$). As the feedback that creates the initial state is switched off at $t=1$ ms, there is a fast (approximately $10\ \mu\text{s}$) transient. This risetime reflects combined dynamics of the gain and cavity, as the laser adapts to the different (now symmetrical) cavity losses. This risetime corresponds roughly to the fluorescence lifetime of the upper state of Ti:sapphire. The “Rabi cycling” of the “population difference” $\rho_{22}-\rho_{11}$ is plotted in Fig. 15b. One can also record the beat note frequency (off-diagonal element $|\rho_{12}|$) as sketched in Fig. 14. As can easily be seen from the Bloch vector model of Feynman et al.,³⁵ the oscillation of the diagonal elements and the off-diagonal elements are 90° out of phase. This property can indeed be seen in Fig. 16a. The Rabi frequency $|\tilde{r}|/\tau_{RT}$ can be varied by changing the position of the scattering surface, as shown in Fig. 16b. The maximum value measured³⁶ for this interface corresponds to a

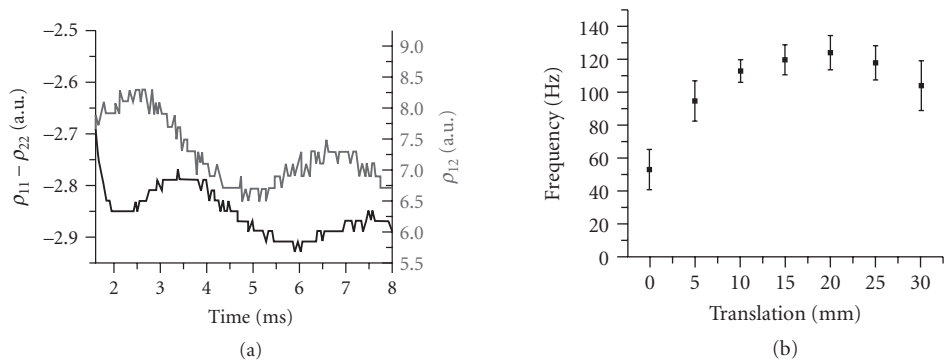


FIGURE 16 (a) Comparison of the oscillation of the population difference $\rho_{22}-\rho_{11}$ and the off-diagonal element (beat note) ρ_{12} . (b) Rabi frequency as a function of position of the glass at the meeting point of the two directions. Translation of the glass-air interface along the beam results in different values of coupling \tilde{r} .

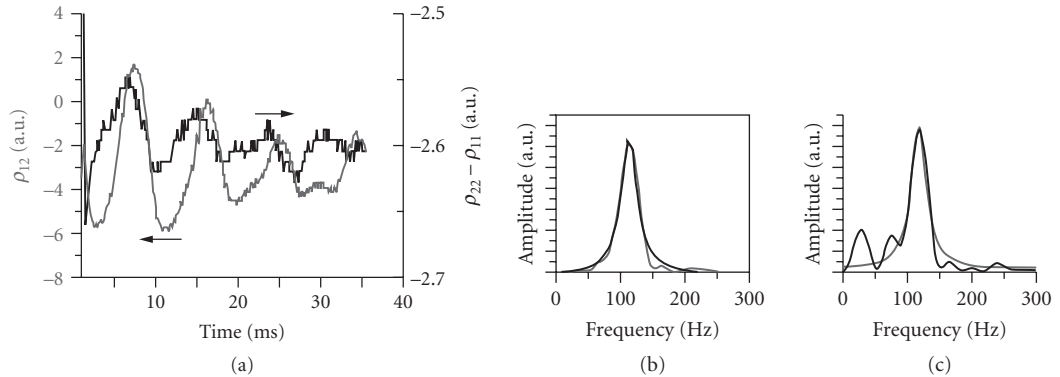


FIGURE 17 (a) Measurement of the decay of the Rabi oscillations $\rho_{22} - \rho_{11}$ and ρ_{21} . (b and c) The Fourier transforms of the relative measurements are shown on the right.

backscattering coefficient of $|\tilde{r}| \approx 1 \cdot 10^{-6}$. Note that the Rabi frequency provides a direct measurement of very minute backscattering coefficients, without the need to trace a complete gyroscopic response as in Refs. 36 and 37.

In the case of a two-level system, the phenomenological “longitudinal” and “transverse” relaxation times have been identified as energy relaxation time (fluorescence decay) and phase relaxation time (due for instance to atomic collisions). Figure 17 shows a measurement of the decay of the Rabi oscillation for the diagonal and off-diagonal elements. The decay is measured by fitting the Fourier transform of the measurement to a Lorentzian, and measuring its FWHM. The values are 27 and 30 Hz. As noted previously, there are at least two origins to the decay of the off-diagonal element: vibration of mirrors, and coupling through absorption (gain). The latter affects equally the diagonal and off-diagonal elements. The former can be seen as a type of “inhomogeneous broadening,” since it has its origin in random cavity length fluctuation, expressed as randomness in the value of $\Delta\omega$. The approximately 30-Hz bandwidth of both decays is consistent with 0.3- μm amplitude vibrations at 100 Hz of cavity components, causing differential cavity fluctuations of 0.3 pm/round-trip.

The role of the gain coupling is particularly important in the present Ti:sapphire laser because of the long-gain recovery time. A better candidate for this analogy would be a dye or semiconductor laser for which the gain lifetime is shorter than the cavity round-trip time. An OPO provides an even better situation, since the gain exists only at the time of pumping.

Rabi Cycling Off-Resonance

If the radiation of amplitude \mathcal{E} (Rabi frequency $\kappa\mathcal{E}$) is off-resonance with a two-level system by an amount $\Delta\omega$, the Rabi frequency becomes $\sqrt{\kappa^2\mathcal{E}^2 + \Delta\omega^2}$. In the case of the ring laser, we can control the off-resonance amount $\Delta\omega$ with a Pockel’s cell (Fig. 14), the initial condition is set favorable to the counter-clockwise direction as shown in Fig. 15a. The Rabi cycling is measured indeed to correspond to $\sqrt{|\tilde{r}|^2/\tau^2 + \Delta\omega^2}$. In resonance case ($\Delta\omega = 0$), measurement of ρ_{12} leads to $r/\tau_{RT} = 138 \text{ Hz} \pm 15 \text{ Hz}$. With $\Delta\omega$ of $171 \text{ Hz} \pm 12 \text{ Hz}$, the off-resonant measurement is $r/\tau_{RT} = 237 \text{ Hz} \pm 21 \text{ Hz}$, which follows the behavior of an off-resonance two-level system.

Impact of the Analogy

The analogy between two-level and laser systems may be more than just a scientific curiosity. A thorough understanding of two-level systems led to powerful spectroscopic techniques, using sophisticated

pulse sequences. Rather than using optical pulses, sublinewidth molecular spectroscopy has been successfully realized by pulsing the detuning by the Stark shifts.^{38–40} In the case of the ring laser, similar pulse sequences can be applied to the detuning. The information sought in spectroscopy is contained in the measurement of $|\rho_{12}|$, as a function of the driving field (measurement of the Rabi frequency $\kappa\epsilon$ leading to the determination of the dipole moment) or detuning $\Delta\omega$. In the case of the ring laser, the measurement of $|\rho_{12}|$ is linked to the properties of some sample inserted in the cavity.¹¹ Any resolution enhancing technique that has been devised in spectroscopy, such as the Ramsey fringes,^{41,42} could be transposed to a laser phase sensor with two intracavity pulses.

20.7 CONCLUSION

This chapter started with a mathematical description of short optical pulses and optical pulse train. Basic physics of short pulse generation in a mode-locked laser are discussed. It is shown in particular that the mechanism by which a steady-state pulse is generated inside the laser, is also responsible for creating equally spaced modes in the frequency domain. It is shown that pulse train interferometry combines the properties of temporal and spatial resolution. The laser can be used as a most sensitive interferometer, when the reference and sample pulses are circulating in the same cavity. Measurements of extreme sensitivity can be performed by interfering the two pulse trains emitted by such a laser. The exquisite sensitivity to phase results from the fact that a phase shift is transposed into a frequency shift inside the active cavity. Exploitation of such lasers as sensors requires a thorough understanding of the coupling between the two intracavity pulses. A new modeling of the laser with two intracavity pulses is introduced by making an analogy with a quantum mechanical two-level system. Beyond its physical elegance, this analogy inspires new sensitivity enhancement techniques for the use of the two-pulse per cavity laser sensor.

20.8 REFERENCES

1. Ladan Arissian and Jean-Claude Diels. "Carrier to Envelope and Dispersion Control in a Cavity with Prism Pairs." *Phys. Rev. A* **75**:013814–013824, 2007.
2. F. X. Kärtner, N. Matuschek, T. Schibli, U. Keller, H. A. Haus, C. Heine, R. Morf, V. Scheuer, M. Tilsch, and T. Tschudi. "Design and Fabrication of Double Chirped Mirrors." *Opt. Lett.* **22**:831–833, 1997.
3. J.-C. Diels and Wolfgang Rudolph. *Ultrashort Laser Pulse Phenomena*, 2d ed. Elsevier, Boston, 2006, ISBN 0-12-215492-4.
4. Th. Udem, J. Reichert, R. Holzwarth, and T.W. Hänsch. "Accurate Measurement of Large Optical Frequency Differences with a Mode-Locked Laser." *Opt. Lett.* **24**:881–883, 1999.
5. J.-C. Diels. Femtosecond Dye Lasers. In F. Duarte and L. Hillman, (eds.) *Dye Laser Principles: With Applications*, Academic Press, Boston, 1990. ISBN 0-12-215492-4, chapter 3, pages 41–132.
6. C. M. Depriest, T. Yilmaz, P. J. Delfyett, S. Etemad, A. Braun, and J. H. Abeles. "Ultralow Noise and Supermode Suppression in an Actively Mode-Locked External-Cavity Semiconductor Diode Ring Laser." *Opt. Lett.* **27**:719–721, 2002.
7. B. Resan and P. J. Delfyett. "Dispersion-Managed Breathing-Mode Semiconductor Mode-Locked Ring Laser: Experimental Characterization and Numerical Simulations." *IEEE J. of Quantum Elect.* **40**:214–220, 2004.
8. T. Yilmaz, C. M. Depriest, and P. J. Delfyett. "Complete Noise Characterisation of External Cavity Semiconductor Laser Hybridly Modelocked at 10 GHz." *Elect. Lett.* **22**:1338–1339, 2003.
9. T. Yilmaz, C. M. Depriest, A. Braun, and J. H. Abeles and P. J. Delfyett. "Residual Phase Noise and Longitudinal Mode Linewidth Measurements of Hybridly Modelocked External Linear Cavity Semiconductor Laser." *Opt. Lett.* **27**:872–874, 2002.
10. T. Yilmaz, C. M. Depriest, A. Braun, J. H. Abeles, and P. J. Delfyett. "Noise in Fundamental and Harmonic Mode-Locked Semiconductor Lasers: Experiments and Simulations." *IEEE J. of Quantum Elect.* **39**:838–849, 2003.

11. J.-C. Diels, Jason Jones, and Ladan Arissian. "Applications to Sensors of Extreme Sensitivity." In Jun Ye and Stephen Cundiff, (eds.) *Femtosecond Optical Frequency Comb: Principle, Operation and Applications*, Springer, New York, 2005, chapter 13, 333–354.
12. F. Aronowitz and R. J. Collins. "Mode Coupling due to Backscattering in a He-Ne Traveling-Wave Ring Laser." *Appl Phys. Lett.* **9**:55–58, 1966.
13. F. Aronowitz. "The Laser Gyro." In Ross, (ed.) *Laser Applications*, Academic Press, New York, 1971, 133–200.
14. J.-C. Diels, I. C. McMichael, J. J. Fontaine, and C. Y. Wang. "Subpicosecond Pulse Shape Measurement and Modeling of Passively Mode locked Dye Lasers Including Saturation and Spatial Hole Burning." In K. B. Eisenthal, R. M. Hochstrasser, W. Kaiser, and A. Laubereau, (eds.) *Picosecond Phenomena III*, Springer-Verlag, New York, 1982, 116–119.
15. M. L. Dennis, J.-C. Diels, and M. Lai. "The Femtosecond Ring Dye Laser: A Potential New Laser Gyro." *Opt. Lett.* **16**:529–531, 1991.
16. Ming Lai, Jean-Claude Diels, and Michael Dennis. "Nonreciprocal Measurements in fs Ring Lasers." *Opt. Lett.* **17**:1535–1537, 1992.
17. Ladan Arissian and Jean-Claude Diels. "Repetition Rate Spectroscopy of the Dark Line Resonance in Rubidium." *Opt. Comm.* **264**:169–173, 2006.
18. Matthew J. Bohn and Jean-Claude Diels. "Bidirectional Kerr-Lens Mode-Locked Femtosecond Ring Laser." *Opt. Comm.* **141**:53–58, 1997.
19. Scott Diddams, Briggs Atherton, and Jean-Claude Diels. "Frequency Locking and Unlocking in a Femtosecond Ring Laser with the Application to Intracavity Phase Measurements." *Appl. Phys. B* **63**:473–480, 1996.
20. Czeslaw Radzewicz, Gary W. Pearson, and Jerzy S. Krasinski. "Use of ZnS as an Additional Highly Nonlinear Intracavity Self-Focusing Element in a Ti:sapphire Self-Modelocked Laser." *Opt. Comm.* **102**:464–468, 1993.
21. M. J. Bohn, R. J. Jones, and J.-C. Diels. "Mutual Kerr-Lens Mode-Locking." *Opt. Comm.* **170**:85–92, 1999.
22. G. P. Agrawal. *Nonlinear Fiber Optics*. Academic Press, Boston, 1995, ISBN 0-12-045142-5.
23. G. E. Stedman, Z. Li, C. H. Rowe, A. D. McGregor, and H. R. Bilger. "Harmonic Analysis in a Large Ring Laser with Backscatter-Induced Pulling." *Physical Review A* **51**(6), June 1995.
24. Matthew J. Bohn, Jean-Claude Diels, and R. K. Jain. "Measuring Intracavity Phase Changes Using Double Pulses in a Linear Cavity." *Opt. Lett.* **22**:642–644, 1997.
25. Ming Lai and Jean-Claude Diels. "Wave-Particle Duality of a Photon in Emission." *J. of the Opt Soc. Am. B* **9**:2290–2294, 1992.
26. D. Gnass, N. P. Ernsting, and F. P. Schaefer. "Sagnac Effect in the Colliding-Pulse-Mode-Locked Dye Ring Laser." *Appl. Phys. B* **53**:119–120, 1991.
27. F. Krausz, Ch. Spielman, T. Brabec, E. Wintner, and A. J. Schmidt. "Generation of 33-fs Optical Pulses from a Solid-State Laser." *Opt. Lett.* **17**:204, 1992.
28. C. Spielmann, P. F. Curley, T. Brabec, and F. Krausz. "Ultrabroadband Femtosecond Lasers." *IEEE J. Quant. Elec.* **QE-30**:1100–1114, 1994.
29. A. E. Siegman. "An Antiresonant Ring Interferometer for Coupled Laser Cavities, Laser Output Coupling, Mode-Locking, and Cavity Dumping." *IEEE J. Quantum Electron.* **QE-9**:247–250, 1973.
30. N. Jamasbi, J.-C. Diels, and L. Sarger. "Study of a Linear Femtosecond Laser in Passive and Hybrid Operation." *J. of Modern Optics* **35**:1891–1906, 1988.
31. D. T. Reid, M. Padgett, C. McGowan, W. E. Sleat, and W. Sibbett. "Light-Emitting Diodes as Measurement Devices for Femtosecond Laser Pulses." *Opt. Lett.* **22**:233–235, 1997.
32. E. S. Wachman, D. C. Edelstein, and C. L. Tang. "Continuous-Wave Mode-Locked and Dispersion Compensated fs Optical Parametric Oscillator." *Opt. Lett.* **15**:136–139, 1990.
33. Xianmei Meng, Jean-Claude Diels, Dietrich Kuehlke, Robert Batchko, and Robert Byer. "Bidirectional, Synchronously Pumped, Ring Optical Parametric Oscillator." *Opt. Lett.* **26**:265–267, 2001.
34. Xianmei Meng, Raphael Quintero, and Jean-Claude Diels. "Intracavity Pumped Optical Parametric Bidirectional Ring Laser as a Differential Interferometer." *Opt. Comm.* **233**:167–172, 2004.
35. R. P. Feynman, F. L. Vernon, and R. W. Hellwarth. "Geometrical Representation of the Schroedinger Equation for Solving Maser Problems." *J. Appl. Phys.* **28**:49–52, 1957.
36. Rafael Quintero-Torres, Mark Ackerman, Martha Navarro, and Jean-Claude Diels. "Scatterometer Using a Bidirectional Ring Laser." *Opt. Comm.* **241**:179–183, 2004.

37. M. Navarro, O. Chalus, and Jean-Claude Diels. "Mode-Locked Ring Lasers for Backscattering Measurement of Mirror." *Opt. Lett.* **31**:2864–2866, 2006.
38. R. G. Brewer and R. L. Shoemaker. "Photon Echoes and Optical Nutation in Molecules." *Phys. Rev. Lett.* **27**:631–634, 1971.
39. R. G. Brewer and R. L. Shoemaker. "Optical Free Induction Decay." *Phys. Rev. A* **6**:2001–2007, 1972.
40. P. R. Berman, J. M. Levy, and R. G. Brewer. "Coherent Optical Transient Study of Molecular Collisions: Theory and Observations." *Phys. Rev.* **11**:1668–1688, 1975.
41. N. F. Ramsey. "A Molecular Beam Resonance Method with Separated Oscillating Fields." *Phys. Rev.* **78**:695–699, 1950.
42. M. M. Salour and C. Cohen-Tannoudji. "Observation of Ramsey's Interference Fringes in the Profile of Doppler-Free Two-Photon Resonances." *Phys. Rev. Lett.* **38**:757–760, 1977.

Zenghu Chang

*Department of Physics
Kansas State University
Cardwell Hall
Manhattan, Kansas*

21.1 GLOSSARY

A	electric field envelope of a laser pulse
E	kinetic energy of an electron in a laser field
ϵ_L	electric field strength of a laser pulse at a given time
E_{\max}	maximum kinetic energy of an electron
ϵ_x	electric field strength of an attosecond XUV pulse at a given time
f_0	carrier-envelope offset frequency of a frequency comb
f_{rep}	repetition frequency of a pulse train
G	temporal gate function
φ_{CE}	carrier-envelope phase (also called absolute phase) of a laser pulse
h	Planck constant
\hbar	Planck constant divided by 2π
I	laser
I_p	ionization potential of an atom
λ_0	center wavelength of a laser pulse
S	trace of the frequency-resolved optical gating
τ	time delay between a laser pulse and an attosecond XUV pulse
U_p	ponderomotive potential of an electron in a laser field
ν_c	frequency of the cutoff harmonic order
ω_0	carrier angular frequency of a laser pulse

21.2 INTRODUCTION

Since the invention of the laser in 1960, the duration of coherent optical pulses has decreased from hundreds of microseconds^{1,2} to 6 femtoseconds in the first 27 years.³ Such tremendous progress was driven by the desire to generate high peak power, study dynamics in matter, increase the speed of telecommunications, and many other applications. However, by the year 1987, the optical pulse length was approaching the limit, i.e., one optical cycle of visible light, which is a few femtoseconds. The bandwidth required to support such few-cycle pulses is generated by perturbative nonlinear interactions such as self-phase modulation.

The characteristic time scale of electron motion in atoms is one atomic unit of time, which is 24.2 attoseconds. One attosecond is 10^{-18} seconds. In the Bohr's model of the atom, the electron orbital time around the hydrogen nucleus is 152 attoseconds. The study of electron dynamics in atoms and molecules called for optical pulses with attosecond duration.⁴⁻⁶ In the frequency domain, a transform-limited Gaussian pulse with 24 attosecond full width at half maximum (FWHM) corresponds to a 73 eV FWHM power spectrum, which is much broader than the entire visible light range. In other words, attosecond pulses are inherently XUV light or x rays. The duration of such extremely short pulses was first measured in 2001.^{7,8} The required ultrabroad spectrum was obtained by using a nonperturbative nonlinear optics process called high-order harmonic generation, discovered in 1987–1988.^{9,10}

High Harmonic Generation

When a linearly polarized, short-pulse laser beam with an intensity on the order of 10^{14} W/cm² interacts with noble gases, odd harmonics of the fundamental frequency—up to tens or even hundreds in order—emerge in the output beam,^{11,12} as depicted in Fig. 1. The intensity of the first few order harmonics decreases quickly as the order increases, then the intensity remains almost unchanged over many harmonic orders, forming a plateau. Finally, the signal cuts off abruptly at the highest order. The broad width of the plateau provided the required spectral bandwidth to support attosecond pulses. The appearance of the intensity plateau is the signature of this nonperturbative laser-atom interaction, which can be described by a semiclassical model.

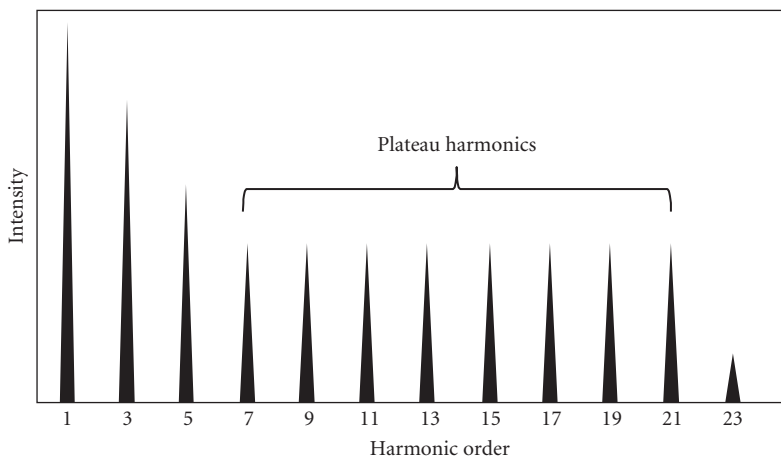


FIGURE 1 High-order harmonic spectrum.

Semiclassical Model

It is also called three-step or two-step model, rescattering model. The electric field acting on an atom changes sinusoidally within one laser cycle. As the laser intensity reaches the level of 10^{14} W/cm², the field near the peak of the oscillation is comparable to the atomic Coulomb field. The superposition of the laser field and the Coulomb field transforms the potential well that binds the electron into a potential barrier. As a result, the electron in the ground state tunnels through the barrier (the first step). The freed electron moves in the laser field like a classical particle and its trajectory can be calculated using Newton's second law. In one laser cycle, the electron first moves away from the nucleus, then is driven back when the force changes direction. During the returning journey, the electron can acquire kinetic energy up to hundreds of electron volts (the second step). Finally, the electron recombines with the parent ion with the emission of a photon (the third step).^{13,14}

When all electrons released near one peak of a laser cycle are considered, the emitted photons form an attosecond pulse. Since there are two field maxima in one laser cycle, two attosecond pulses are generated. For a laser pulse that contains many cycles, an attosecond pulse train is produced. The pulse train corresponds to discrete harmonic peaks in the frequency domain. In other words, high harmonic generation and the attosecond pulse train are two manifestations of the same nonperturbative interaction.

The photon energy of the cutoff harmonic order, $h\nu_c$, is determined by the maximum kinetic energy of the electron gained in the laser field, E_{\max} . It can be shown that $h\nu_c = I_p + E_{\max} \approx I_p + 3U_p$. Here I_p is the ionization potential of the atom and U_p is the ponderomotive potential of the electron in the laser field. Apparently, the width of the plateau and therefore the minimum attosecond pulse duration is limited by $h\nu_c$. The cutoff order is also affected by the depletion of the ground state population due to the ionization.¹⁵

Ponderomotive Potential

The ponderomotive potential is the cycle-averaged kinetic energy of an electron in a laser field, $U_p [eV] = 9.33 \times 10^{-14} I \cdot \lambda_0^2$, where I is the laser intensity in W/cm² and λ_0 is the center wavelength of the laser in micrometer. It is clear that the cutoff photon energy of the high harmonic spectrum can be extended by using longer wavelength driving lasers.¹⁶

Strong Field Approximation

A fully quantum three-step model was developed in 1994.¹⁷ It is valid when the ponderomotive potential is much larger than the ionization potential. It assumes that the harmonic emission is the result of the dipole transition between the ground state and the continuum states only, with the excitation states playing no role. Solving the Schrödinger equation results in an analytical solution of the dipole moment, from which one can obtain both the phase and the intensity of each harmonic order. The model reveals that there are two quantum trajectories that contribute to each plateau harmonic. One is called the long trajectory and the other is the short trajectory. The phase of each harmonic depends on the laser intensity. The intensity dependence of the dipole phase (also called intrinsic phase) is different for the two trajectories.

Quantum Trajectories

By solving the equation of motion, it can be shown that an electron released right at the peak of the laser field will return to the starting point one cycle later, with zero kinetic energy. As the releasing time from the field peak increases, the returning energy increases first, reaches the maximum value ($3U_p$), then decreases to zero. Therefore electrons releasing at two different moments can come back to the parent ion with the same kinetic energy, which corresponds to the same harmonic order.^{13,14}

The electron that starts the journey earlier returns later. Its path is called the long trajectory. The other one is the short trajectory. Quantum mechanically, there are many trajectories contribute to each harmonics (Feynman's path-integral), but the dominating contributions are from the two trajectories corresponding to the classical ones.¹⁷

Phase-Matching

The semiclassical model and the strong field approximation describe the single atom response. To generate a high-intensity high harmonic beam, many atoms must contribute to the output constructively.¹⁸ Ionization of the atom is unavoidable in high harmonic generation because it is the first step of the process. In highly ionized gas targets, the phase velocity of the laser field (and thus the polarization) is greater than that of the harmonic field. The resulting phase mismatch can be compensated for by several approaches. One of them utilizes the intensity dependent phase.¹⁹ In most cases, only the short trajectory is phase matched. Nevertheless, low laser to harmonic conversion efficiency is still a major problem that needs to be solved. The spatial coherence of the high harmonic/attosecond train beam is excellent when the phase matching conditions are fulfilled. The divergence angle of the XUV beam is smaller than the driving laser beam.¹²

Single Isolated Pulses

The attosecond pulse train corresponding to high order harmonics is useful for some applications. In general, however, single isolated attosecond pulses are required for performing pump-probe experiments with arbitrary delay between the pump and the probe pulses. Such pulses can be generated by suppressing all the pulses in the train except one, which can be accomplished by using single-cycle driving lasers²⁰ or pulse extraction switches with a subcycle opening time.²¹ Also the pulses from the gas target are positively chirped.²² Dispersion compensation over a broad XUV spectral range is a major challenge. By 2008, the shortest single isolated pulses, which were generated from neon gas by using 3.3-fs driving lasers centered at 720 nm, were 80 attoseconds and contained ~0.5 nJ of energy.²³ Their spectrum was centered at 80 eV.

21.3 THE DRIVING LASER

There are several basic requirements on the driving lasers for the generation of single isolated attosecond pulses. First, the intensity at the focus must be high enough, on the order of 10^{14} to 10^{15} W/cm², which is a fraction of an atomic unit of intensity (3.55×10^{16} W/cm²). The corresponding pulse energy is 100 μ J or higher. The spectral bandwidth of the attosecond pulses is proportional to the driving laser intensity. Second, the laser pulse duration must be short enough. The ionization of the target atoms by the laser field before the cycle where the attosecond pulse is generated must not deplete the ground state population completely. Depending on the generation scheme, acceptable laser pulses range from 3 to 30 fs. Third, the carrier-envelope phase needs to be stabilized. Since the single attosecond pulses are generated in a fraction of the laser cycle, a shift in the carrier-envelope phase results in shot-to-shot variations of the attosecond pulses. Finally, the repetition rate of the laser should be high, on the order of kilohertz. Many attosecond characterization and application schemes rely on photoelectron measurements. There is an upper limit on the number of electrons per shot to avoid the space charge effect. Thus the signal count rate is primarily determined by the repetition rate. The energy stability of high-repetition-rate lasers is also better than those with low repetition rates.

High power laser pulses with duration around 30 fs can be generated with chirped pulse amplification. Pulses down to ~4 fs with submillijoule energy can be obtained by spectral broadening in hollow-core fibers filled with gases, followed by dispersion compensation using chirped mirrors or phase modulators,^{24–27} as illustrated by the block diagram in Fig. 2.

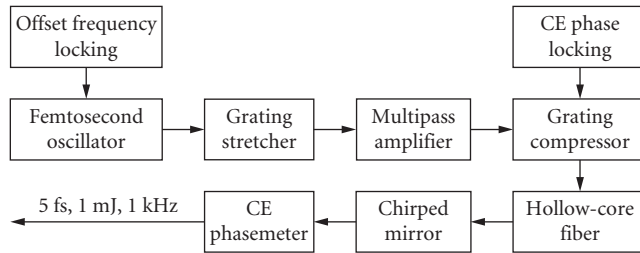


FIGURE 2 A carrier-envelope phase stabilized, few-cycle laser system.

Chirped Pulse Amplification (CPA)

Ti:Sapphire is a commonly used gain medium for femtosecond lasers primarily because of its broad gain bandwidth.²⁸ Its center wavelength is ~ 800 nm, which corresponds to a ~ 2.6 -fs optical period. Femtosecond oscillators that use Ti:Sapphire as the gain medium can generate pulses with nanojoule level energy. Direct amplification of the pulse to millijoule level may cause damage to the laser crystal. In a chirped amplifier, the pulses from the oscillator are stretched to hundreds of picoseconds to lower the peak power. Then the pulses are amplified in multipass or regenerative amplifiers. The high energy pulses are finally compressed to femtosecond duration.^{29,30} Most high energy (>5 mJ) lasers use grating pairs to stretch and compress pulses.

Carrier-Envelope Phase

The effects of the magnetic field of the laser on high harmonic generation can be ignored because the electron velocity during the three-step journey is much slower than the speed of light. The electric field of a linearly polarized laser pulse can be expressed as $\epsilon_L(t) = A(t)\cos(\omega_0 t + \varphi_{CE})$. Here $A(t)$ is the pulse envelope, ω_0 is the carrier frequency, φ_{CE} is the carrier-envelope (CE) phase that specifies the offset between the peak of the pulse envelope and the closest maximum of the oscillating field. To stabilize the carrier-envelope phase of the laser pulses from chirped pulse amplifiers followed by the hollow-core fiber compressors, the carrier-envelope offset frequency of the oscillator must first be stabilized. Furthermore, the CE phase drifts in the chirped pulse amplifier and in the hollow-core fiber compressor must be compensated.^{31,32}

Carrier-Envelope Offset Frequency

The technique for stabilizing carrier-envelope offset frequency was originally developed for frequency metrology in 2000.^{33,34} Most femtosecond oscillators used for seeding amplifiers work at a repetition rate $f_{\text{rep}} \sim 80$ MHz. The oscillator output is a femtosecond pulse train that corresponds to a frequency comb. The frequency of the n th tooth of the comb is $f_0 + nf_{\text{rep}}$, where f_0 is the frequency of the zeroth tooth. The rate of carrier-envelope phase change is determined by the offset frequency f_0 , which can be stabilized by using f -to- $2f$ technology. For this to work, the laser spectrum must cover an octave. The f_0 is measured by beating the $2n$ th tooth with the frequency doubled n th tooth, that is $f_0 = 2[f_0 + nf_{\text{rep}}] - (f_0 + 2nf_{\text{rep}})$. It was found that f_0 can be stabilized by controlling the pump power to the gain medium of chirped mirror based oscillators. When f_0 is fixed to f_{rep}/m , the CE phase of every m th pulse from the oscillator is the same. It is common practice to choose $m = 4$.

Carrier-Envelope Phase of Chirped Pulse Amplifiers

The oscillator pulses with the same CE phase are switched out by a Pockels cell and sent to Ti:Sapphire amplifiers that operate at kilohertz repetition rates. When grating pairs are used to stretch and compress laser pulses for the chirped pulse amplification, a submicrometer change of separation between gratings can lead to a 2π CE phase shift.³⁵ This effect has been used to correct the slow CE phase variation introduced by the amplifier components. It was accomplished by measuring the CE phase variation after the amplifier and using the measured signal for feedback control of the grating separation.³⁶ The CE phase error of CPA systems can be controlled to <200 mrad over hours. The relative CE phase variation can be measured by a single shot f -to- $2f$ interferometer, whereas the absolute phase value can be determined by a phasemeter (discussed below), which measures electrons from the above-threshold ionization of atoms by the laser pulses.

Single Shot f -to- $2f$ Interferometer

The laser pulse from the hollow-core fiber compressor is a white-light continuum that can cover an octave spectral range. One can select a narrow range near 1000 nm and frequency double it to interfere with the light around 500 nm. The interferogram in the frequency domain is a sinusoidal fringe. The period of the fringe pattern is inversely proportional to the delay between the two interfering pulses. A CE phase shift will cause the fringes to shift. Thus by measuring the interferogram with a spectrometer, the CE phase variation can be measured.³⁷

Carrier-Envelope Phasemeter

In the three-step semiclassical model, the attosecond photon pulses are generated by the recombination of the returning electrons. A returning electron can also scatter away from the parent ion. The kinetic energy distribution of the rescattered electrons after the laser field vanishes can extend to $10U_p$. There is also a plateau in the electron spectrum similar to the high harmonic spectrum. This electron emission process is called above-threshold ionization. The angular distribution of the electrons is concentrated along the field polarization direction. When the laser pulse is only a few cycles long, the number of plateau electrons flying to one direction can be different from those to the opposite direction. The asymmetry depends on the carrier-envelope phase of the laser. Thus, by simultaneously measuring electrons in two directions, the absolute CE phase value can be determined.³⁸

21.4 ATTOSECOND PULSE GENERATION

A typical attosecond pulse generation setup consists of a kilohertz femtosecond Ti:Sapphire laser system, a vacuum chamber where the gas target is located, and an XUV spectrometer/attosecond streak camera that characterizes the pulses in the spectral domain and the time domain, as shown in Fig. 3. The attosecond pulses are XUV or soft x-ray light that cannot propagate in air because of high absorption. The gas density in the laser interaction region is on the order of 10^{17} to 10^{18} atoms/cm³. The interaction length is typically a few millimeters for gas cells or gas jets. It should be smaller than the Rayleigh range of the focusing laser beam, so that the carrier-envelope phase does not change significantly due to the Gouy phase shift inside the target. The target is located after the focal point to achieve good phase-matching.

Attosecond Pulse Train

Such pulses are generated with linearly polarized laser pulses that contain many optical cycles. When only the fundamental frequency is used, the spacing between two neighboring harmonic

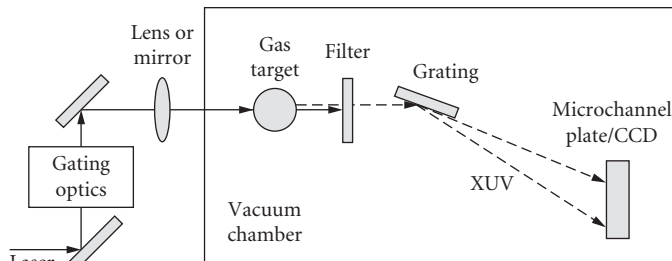


FIGURE 3 Setup for generating attosecond pulses and measuring their spectrum.

peaks is two-photon energies. In the time domain, the spacing between adjacent pulses is one-half of an optical cycle; 1.3 fs for Ti:Sapphire lasers.^{7,39} The amplitude changes from pulse to pulse. Since many-cycle lasers (>20 fs) can be generated directly from high power (terawatt) chirped pulses amplifiers, the attosecond pulse energy can be high enough to perform nonlinear physics experiments.⁴⁰

Two-Color Gating

When the many-cycle driving laser is a combination of the fundamental frequency and its second harmonic, the breaking of symmetry of the laser field leads to the generation of both odd and even high harmonics and thus the spacing between two neighboring harmonic peaks is one photon-energy. In the time domain, the spacing between adjacent pulses becomes a full optical cycle.⁴¹ Such pulse trains are useful for performing experiments using the powerful attosecond streaking technique.

Amplitude Gating

As the driving laser approaches a single optical cycle, the cycle-to-cycle field amplitude variation becomes significant. When the carrier-envelope phase of the pump laser is set to zero, the spectrum of the attosecond pulses generated near the peak of the laser pulse envelope extends to a shorter XUV wavelength range as compared to the adjacent attosecond pulses emitted when the laser field is weaker. As a result, the high-order harmonic spectrum becomes a continuum in the cutoff region.²⁰ Discrete harmonics remain in other portions of the spectrum. The shorter the driving laser is, the broader the XUV continuum becomes. A single isolated attosecond pulse as short as 80 attoseconds was obtained by selecting the continuum region of the XUV spectrum with a high-pass filter, using <4-fs pump lasers.²¹ The scaling of the attosecond pulse energy is limited by the maximum energy of the driving laser from the hollow-core fiber compressor. Combining this type of gating with the two-color gating can relax the requirement of the laser pulse duration.

Polarization Gating

In the plateau region, four attosecond pulses are produced in one laser cycle taking into account both the long and short trajectory's contributions. However, only the two pulses from the short trajectory can be phase matched on axis. Consequently, the spacing between pulses is still half of a laser cycle. Single isolated attosecond pulses can be extracted by a scheme called polarization gating.²¹ It uses a laser field with a rapid change of ellipticity. Since XUV attosecond pulses can only be efficiently generated with linearly polarized driving fields, a single attosecond pulse is emitted if the laser field

is linearly polarized in only a short time range and elliptically polarized in the other portion of the driving pulse. The time range over which the attosecond pulse is generated is called the polarization gate. So far, single isolated XUV pulses as short as 130 attoseconds were generated with this method using 5 fs pump lasers.⁴² For the same driving laser pulse duration, polarization gating has the potential to generate shorter attosecond pulses because it can create a broader continuum in the plateau region.⁴³

Double Optical Gating

The few-cycle laser pulses used in amplitude gating and polarization gating are difficult to generate daily. A method called double optical gating was proposed to allow the generation of single isolated attosecond pulses with longer pump lasers.⁴⁴ It is a combination of the two-color gating and the polarization gating. A second harmonic field is added to the fundamental field in order to break the symmetry of the field and increases the spacing between the adjacent attosecond pulses to one optical cycle. When the polarization gating is applied, the polarization gate width equals one optical cycle to select one isolated XUV pulse. The depletion of the ground state population by the leading edge of the laser pulses can be significantly reduced with this scheme; as a result, multicycle lasers can be used. This scheme has been demonstrated with laser pulses as long as 20 fs. Since such lasers do not necessarily need hollow-core fiber compressors, they are much easier to operate. The laser pulse energy can also be much higher, which is important for the energy scaling of the attosecond pulses.

21.5 ATTOSECOND PULSE CHARACTERIZATION

Measurement of the optical pulse duration requires a temporal gate. For femtosecond lasers, nonlinear optics phenomena such as second harmonic generation can serve as the gating, which is the foundation of widely implemented autocorrelation and the frequency-resolved optical gating (FROG) techniques.⁴⁵ The intensity of the attosecond pulses is not high enough to generate second harmonic light yet. Most of the methods for determining the width of the attosecond pulses require the measurements of photoelectrons or ions. The XUV beam is focused to a gas target to generate the photoelectrons/ions. The charged particles are detected by a time-of-flight spectrometer. A second beam, either an XUV or an intense laser beam is also focused to the same target, overlapping spatially and temporally with the first beam. The interaction of the two pulses in the gas serves as the temporal gate. A typical setup is shown in Fig. 4, where the attosecond XUV pulses are generated in the first gas target and are measured in the second gas target. Similar apparatus have been used for studying electron dynamics in atoms.

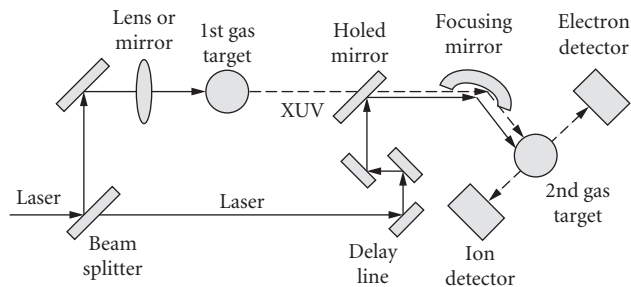


FIGURE 4 Setup for measuring the attosecond pulse duration.

Second-Order Autocorrelator or FROG

This technique resembles the second harmonic autocorrelation in femtosecond optics. Ionization of atoms (such as helium) or Coulomb explosion of molecules (such as N_2) by nonresonant two-photon absorption can serve as the nonlinearity. The ion signal as a function of the time delay between the two attosecond pulses is measured to yield the second order autocorrelation function.⁴⁰ By assuming a certain pulse shape, the pulse duration can be obtained by fitting the autocorrelation trace. The interferometric second-order autocorrelations have been used to characterize attosecond pulse trains generated with low-repetition rate femtosecond lasers with tens of millijoule pulse energy.

When the photoelectron kinetic energy spectrum is measured as a function of the delay, a two-dimensional frequency-resolved optical gating pattern is obtained. Both the phase and pulse profile of subfemtosecond pulses can be reconstructed using this method.⁴⁶

RABITT (Reconstruction of Attosecond Beating by Interference of Two-Photon Transition)

When the attosecond pulse train is generated with the fundamental wave of kilohertz lasers, the intensity of the XUV light may not be strong enough to cause measurable nonlinear effects. A cross-correlation method based on the two-color above-threshold ionization was developed to determine the duration of the pulses in the train.⁷ A high harmonic beam interacting with atomic gases alone will generate photoelectron peaks separated by two laser photon-energies. Adding a dressing laser with intensity of 10^{11} W/cm² generates an electron sideband located in the middle of two peaks. By measuring sideband intensity as a function of the delay between the XUV pulse and the dressing laser, relative phase between adjacent harmonics can be determined. Combining this with the high harmonic power spectrum, one can deduce the attosecond pulse duration. This approach assumes that the width of all the pulses in the train is the same.

Attosecond Streak Camera

The photoelectron replicas generated by attosecond XUV pulses have durations shorter than the optical cycle of the driving lasers. When the photoelectrons are released in the presence of a laser field, their momentum after the laser pulse is gone will be different from the initial value. The momentum shift is determined by the vector potential of the laser field at the time the electron is released. Thus, the leading edge of the electron pulse will gain an additional momentum that is different from the electron in the trailing edge. By measuring the momentum distribution of photoelectrons, the width of the photoelectron pulse (and thus the XUV pulse) can be determined.^{47,48} The required laser intensity is on the order of 10^{12} to 10^{13} W/cm². This approach is similar to the picosecond optical streak camera. It has been used to measure single isolated attosecond pulses and the pulse trains generated from the two-color gating. It is, however, difficult to measure pulses with half laser cycle spacing with this method.

FROG-CRAB (Frequency-Resolved Optical Gating for Complete Reconstruction of Attosecond Bursts)

The momentum streaking of the photoelectrons in a laser field can also be described as the phase shift of an electron wave packet. The phase shift of the electron wave by the laser field can be considered as a temporal phase gate, $G(t)$. When the photoelectron spectrum is measured as a function of the delay τ between the XUV field $\epsilon_x(t)$ and the laser field, a FROG trace is obtained, given by
$$S(E, \tau) = \left| \int_{-\infty}^{\infty} dt \epsilon_x(t - \tau) G(t) e^{j(E + I_p)t/\hbar} \right|^2.$$
 Here E is the energy of the photoelectron. Such a spectrogram can be processed using a FROG retrieval algorithm to fully characterize the XUV pulse as well as the electric field of the near IR laser pulse.^{49,50} This method works well for measuring both attosecond pulse trains and single isolated pulses.

21.6 ACKNOWLEDGMENTS

This material is supported by the U.S. Army Research Office under grant number W911NF-07-1-0475.

21.7 REFERENCES

1. T. H. Maiman, "Stimulated Optical Radiation in Ruby," *Nature* **187**: 493–494 (1960).
2. A. E. Siegman, *Lasers*, University Science Books, Mill Valley, California (1986), ISBN 0-935702-11-3, p. 61.
3. R. L. Fork, C. H. B. Cruz, P. C. Becker, and C. V. Shank, "Compression of Optical Pulses to Six Femtoseconds by Using Cubic Phase Compensation," *Opt. Lett.* **12**: 483–485 (1987).
4. P. Agostini, and L. F. DiMauro, "The Physics of Attosecond Light Pulses," *Reports on Progress in Physics* **67**: 813 (2004).
5. P. B. Corkum and F. Krausz, "Attosecond Science," *Nat. Phys.* **3**: 381 (2007).
6. M. F. Kling and M. J. J. Vrakking, "Attosecond Electron Dynamics," *Annual Review of Physical Chemistry* **59**: 463 (2008).
7. P. M. Paul, E. S. Toma, P. Breger, G. Mullot, F. Auge, Ph. Balcou, H. G. Muller, and P. Agostini, "Observation of a Train of Attosecond Pulses from High Harmonic Generation," *Science* **292**: 1689 (2001).
8. M. Hentschel, R. Kienberger, Ch. Spielmann, G. A. Reider, N. Milosevic, T. Brabec, P. Corkum, U. Heinzmann, M. Drescher, and F. Krausz, "Attosecond Metrology," *Nature* **414**: 509 (2001).
9. A. McPherson, G. Gibson, H. Jara, U. Johann, T. S. Luk, I. A. McIntyre, K. Boyer, and C. K. Rhodes, "Studies of Multiphoton Production of Vacuum-Ultraviolet Radiation in the Rare Gases," *J. Opt. Soc. Am. B* **4**: 595 (1987).
10. M. Ferray, A. L'Huillier, X. F. Li, L. A. Lompré, G. Mainfray, and C. Manus, "Multiple-Harmonic Conversion of 1064 nm Radiation in Rare Gases," *J. Phys. B* **21**: L31 (1988).
11. A. L'Huillier, T. Auguste, Ph. Balcou, B. Carie, P. Monot, P. Salières, C. Altucci, et al., "High-Order Harmonics: A Coherent Source in the XUV Range," *J. Nonl. Opt. Phys. and Mat.* **4**: 647 (1995).
12. P. Salières, A. L'Huillier, P. Antoine, M. Lewenstein, "Studies of the Spatial and Temporal Coherence of High Order Harmonics," *Adv. Atom. Mol. Opt. Phys.* **41**: 83 (1999).
13. P. B. Corkum, "Plasma Perspective on Strong-Field Multiphoton Ionization," *Phys. Rev. Lett.* **71**: 1994–1997 (1993).
14. K. C. Kulander, K. J. Schafer, J. L. Krause, in *Super-Intense Laser-Atom Physics*, B. Piraux, A. L'Huillier, and K. Rzazewski (eds.) Plenum, New York (1993). NATO ASI, Ser. B, Vol. **316**: p. 95.
15. Z. Chang, A. Rundquist, H. Wang, M. M. Murnane, and H. C. Kapteyn, "Generation of Coherent Soft X Rays at 2.7 nm Using High Harmonics," *Phys. Rev. Lett.* **79**: 2967 (1997).
16. B. Shan, Z. Chang, "Dramatic Extension of the High-Order Harmonic Cutoff by Using a Long-Wavelength Pump," *Phys. Rev. A* **65**: 011804(R) (2002).
17. M. Lewenstein, Ph. Balcou, M. Yu. Ivanov, A. L'Huillier, and P. B. Corkum, "Theory of High-Harmonic Generation by Low-Frequency Laser Fields," *Phys. Rev. A* **49**: 2117–2132 (1994).
18. M. B. Gaarde, J. L. Tate, and K. J. Schafer, "Macroscopic Aspects of Attosecond Pulse Generation," *J. Phys. B: At. Mol. Opt. Phys.* **41**: 32001 (2008).
19. M. Lewenstein, P. Salières, and A. L'Huillier "Phase of the Atomic Polarization in High-Order Harmonic Generation," *Phys. Rev. A* **52**: 4747 (1995).
20. I. P. Christov, M. M. Murnane, and H. Kapteyn, "High-Harmonic Generation of Attosecond Pulses in the 'Single-Cycle' Regime," *Phys. Rev. Lett.* **78**: 1251–1254 (1997).
21. P. B. Corkum, N. H. Burnett, and M. Y. Ivanov, "Subfemtosecond Pulses," *Opt. Lett.* **19**: 1870 (1994).
22. Z. Chang, "Chirp of the Attosecond Pulses Generated by a Polarization Gating," *Phys. Rev. A* **71**: 023813 (2005).
23. E. Goulielmakis, M. Schultze, M. Hofstetter, V. S. Yakovlev, J. Gagnon, M. Uiberacker, A. L. Aquila, et al., "Single-Cycle Nonlinear Optics," *Science* **320**: 1614 (2008).
24. M. Nisoli, S. D. Silvestri, and O. Svelto, "Generation of High Energy 10 fs Pulses by a New Pulse Compression Technique," *Appl. Phys. Lett.* **68**: 2793–2975 (1996).

25. R. Szipöcs, K. Ferencz, C. Spielmann, and F. Krausz, "Chirped Multilayer Coatings for Broadband Dispersion Control in Femtosecond Lasers," *Opt. Lett.* **19**: 201–203 (1994).
26. M. Nisoli, S. D. Sverstri, O. Svelto, R. Szipöcs, K. Ferencz, Ch. Spielmann, S. Sartania, and F. Krausz, "Compression of High-Energy Laser Pulse below 5 fs," *Opt. Lett.* **22**: 522–524 (1997).
27. H. Wang, Y. Wu, C. Li, H. Mashiko, S. Gilbertson, and Z. Chang, "Generation of 0.5 mJ, Few-Cycle Laser Pulses by an Adaptive Phase Modulator," *Opt. Exp.* **16**: 14448–14455 (2008).
28. P. F. Moulton, "Spectroscopic and Laser Characteristics of $\text{Ti:Al}_2\text{O}_3$," *J. Opt. Soc. Am. B* **3**: 125 (1986).
29. D. Strickland and G. Mourou, "Compression of Amplified Chirped Optical Pulses," *Opt. Commun.* **56**: 219 (1985).
30. G. A. Mourou, T. Tajima, and S. V. Bulanov, "Optics in the Relativistic Regime," *Rev. Mod. Phys.* **78**: 309 (2006).
31. A. Baltuska, Th. Udem, M. Uiberacker, M. Hentschel, E. Goulielmakis, Ch. Gohle, R. Holzwarth, et al., "Attosecond Control of Electronic Processes by Intense Light Fields," *Nature* **421**: 611 (2003).
32. A. Baltuska, M. Uiberacker, E. Goulielmakis, R. Kienberger, V. S. Yakovlev, T. Udem, T. W. Hänsch, and F. Krausz, "Phase-Controlled Amplification of Few-Cycle Laser Pulses," *IEEE J. Sel. Topics Quantum Electron.* **9**: 972 (2003).
33. D. J. Jones, S. A. Diddams, J. K. Ranka, A. Stentz, R. S. Windeler, J. L. Hall, and S. T. Cundiff, "Carrier-Envelope Phase Control of Femtosecond Mode-Locked Lasers and Direct Optical Frequency Synthesis," *Science* **288**: 635–639 (2000).
34. A. Apolonski, A. Poppe, G. Tempea, C. Spielmann, T. Udem, R. Holzwarth, T. W. Hänsch, and F. Krausz, "Controlling the Phase Evolution of Few-Cycle Light Pulses," *Phys. Rev. Lett.* **85**: 740–743 (2000).
35. Z. Chang, "Carrier Envelope Phase Shift Caused by Grating-Based Stretchers and Compressors," *Appl. Opt.* **45**: 8350(2006).
36. C. Li, E. Moon, and Z. Chang, "Carrier-Envelope Phase Shift Caused by Variation of Grating Separation," *Opt. Lett.* **31**: 3113–3115 (2006).
37. M. Kakehata, H. Takada, Y. Kobayashi, K. Torizuka, Y. Fujihara, T. Homma, and H. Takahashi, "Single-Shot Measurement of Carrier-Envelope Phase Changes by Spectral Interferometry," *Opt. Lett.* **26**: 1436–1438 (2001).
38. G. G. Paulus, F. Grabson, H. Walther, P. Villorosti, M. Nisoli, S. Stagira, E. Priori, and S. De Silvestri, "Absolute-Phase Phenomena in Photoionization with Few-Cycle Laser Pulses," *Nature* **414**: 182–184, (2001).
39. P. Antoine, A. L'Huillier, and M. Lewenstein, "Attosecond Pulse Trains Using High-Order Harmonics," *Phys. Rev. Lett.* **77**, 1234 (1996).
40. P. Tzallas, D. Charalambidis, N. A. Papadogiannis, K. Witte, and G. D. Tsakiris, "Direct Observation of Attosecond Light Bunching," *Nature* **426**: 267–271 (2003).
41. J. Mauritsson, P. Johnsson, E. Gustafsson, A. L'Huillier, K. J. Schafer, and M. B. Gaarde, "Attosecond Pulse Trains Generated Using Two Color Laser Fields," *Phys. Rev. Lett.* **97**: 013001 (2006).
42. G. Sansone, E. Benedetti, F. Calegari, C. Vozzi, L. Avaldi, R. Flammini, L. Poletto, et al., "Isolated Single-Cycle Attosecond Pulses," *Science* **314**: 443 (2006).
43. Z. Chang, "Single Attosecond Pulse and xuv Supercontinuum in the High-Order Harmonic Plateau," *Phys. Rev. A* **70**: 043802 (2004).
44. H. Mashiko, S. Gilbertson, C. Li, S. D. Khan, M. M. Shakya, E. Moon, and Z. Chang, "Double Optical Gating of High-Order Harmonic Generation with Carrier-Envelope Phase Stabilized Lasers," *Phys. Rev. Lett.* **100**: 103906 (2008).
45. R. Trebino, D. J. Kane, "Using Phase Retrieval to Measure the Intensity and Phase of Ultrashort Pulses: Frequency-Resolved Optical Gating," *J. Opt. Soc. Am. A* **10**: 1101 (1993).
46. A. Kosuge, T. Sekikawa, X. Zhou, T. Kanai, S. Adachi, and S. Watanabe, "Frequency-Resolved Optical Gating of Isolated Attosecond Pulses in the Extreme Ultraviolet," *Phys. Rev. Lett.* **97**: 263901 (2006).
47. R. Kienberger, M. Hentschel, M. Uiberacker, Ch. Spielmann, M. Kitzler, A. Scrinzi, M. Wieland, et al., "Steering Attosecond Electron Wave Packets with Light," *Science* **297**: 1144–1148 (2002).
48. J. Itatani, F. Quéré, G. L. Yudin, M. Yu. Ivanov, F. Krausz, and P. B. Corkum, "Attosecond Streak Camera," *Phys. Rev. Lett.* **88**: 173903 (2002).
49. Y. Mairesse, and F. Quéré, "Frequency-Resolved Optical Gating for Complete Reconstruction of Attosecond Bursts," *Phys. Rev. A* **71**: 011401(R) (2005).
50. J. Gagnon, E. Goulielmakis, V. S. Yakovlev, "The Accurate FROG Characterization of Attosecond Pulses from Streaking Measurements," *App. Phys. B* **92**: 25 (2008).

This page intentionally left blank.

LASER STABILIZATION

John L. Hall, Matthew S. Taubman*, and Jun Ye

JILA

*University of Colorado and National Institute of Standards and Technology
Boulder, Colorado*

22.1 INTRODUCTION AND OVERVIEW

For laser applications in which measurement precision is a key feature, frequency-stabilized lasers are preferred, if not essential. This observation was true in the gas laser days when the 10^{-6} fractional Doppler width set the uncertainty scale. Now we have diode-pumped solid state lasers with fractional tuning range approaching 10^{-2} or more, and laser diode systems with several percent tuning. Such tuning is useful to find the exact frequency for our locking resonance, but then stabilization will be essential. Locking to cavities and atomic references can provide excellent stability, even using a widely tunable laser source. Indeed, laser frequency stability between independent systems has been demonstrated at 1×10^{-15} in 1 s averaging time, and more than a decade better at 300 seconds. This incredible performance enhancement is possible because of a feedback system, beginning from measurement of the laser's frequency error from our chosen setpoint, suitable processing of this error signal by a filter/amplifier system, and finally application of a correction signal to an actuator on the laser itself, which changes its frequency in response. While such feedback in response to performance may be the most important principle in evolution, in machines and lasers feedback enables the design of lighter, less costly systems. The accuracy is obtained, not by great bulk and stiffness, but rather by error correction, comparing the actual output against the ideal. This continuous correction will also detect and suppress the system's internal nonlinearity and noise. The performance limitation ultimately is set by imprecision of the measurement, but naturally there is a lot of care required to get into that domain: we must have a very powerful and accurate correction effort to completely hide the original sins.

This chapter is our attempt to lead the worker newly interested in frequency control of lasers on a guided tour of stabilized lasers, ideally providing enough insight for recruiting yet another colleague into this wonderful arena. As nonlinear optics becomes just part of our everyday tools, the buildup cavities which enhance the nonlinear couplings are taking on a more critical role: this is the reason that we focus on the taming of piezoelectric-based (PZT-based) systems. We then cover locking with other transducers, and present some details about their construction and use. We consider the frequency discriminator, which is a key element for these control systems. The chapter concludes with description of the design and performance of several full practical systems, including subhertz linewidth systems.

*Matthew S. Taubman is now with the Pacific Northwest National Laboratories, Richland WA.

Quantifying Frequency Stability

In thinking about the stability of our lasers, one may first wonder whether time- or frequency-domain pictures will be more powerful and instructive. Experience shows that time-domain perturbations of our lasers are usually associated with unwelcome sounds—door slamming, telephone bells, and loud voices. Eventually these time-localized troubles can be eliminated. But what remains is likely the sum of zillions of smaller perturbations: none too conspicuous, but too many in number to attack individually. This perspective leads to a frequency-domain discussion where we can add the Fourier amplitudes caused by the many little sources. Eventually we are led to idealize our case to a continuum of spectrally-described perturbations. This physical outlook is one reason we will mainly be specifying our performance measures in the frequency domain: We have already removed the few really glaring problems and now begin to see (too) many small ones.

Another important issue concerns the nifty properties of Mr. Fourier's description: in the frequency domain, cascaded elements are represented by the multiplication of their individual transfer functions. If we had chosen instead the time domain, we would need to work with convolutions, nonlocal in time. Today's result in time is the sum of all previous temporal events that have the proper delay to impact us now. So it seems clear that frequency domain is good for analysis. What about describing the results?

Frequency versus Time: Drift—the Allan Variance Method

At the other end of our laser stabilization project, describing the results, it is convenient to measure and record the frequency as a function of time. We can measure the frequency averaged over one-second gating time, for example, and stream 100 points to a file. This would be a good way to see the variations around a mean for the 1-second time intervals. This measurement could be repeated using a succession of gate times, 3 s, 10 s, 30 s, 100 s. . . . Surely it will be attractive to make this measurement just once and numerically combine the data to simulate the longer gate times. Thinking this way brings us a new freedom: we can process this data to recover more than just the mean and the standard deviation. Of course, we can expect to eventually see some drift, particularly over long times. When we look at the drift and slowly varying laser frequency, one wishes for a method to allow us to focus on the random noise effects which are still visible, even with the extended gate times. This is where the resonance physics is, while the drift is mainly due to technical problems. Dave Allan introduced the use of first differences, which has come to be called the Allan Variance method.¹ If we take the difference between adjacent samples of the measured frequency, we focus on the random processes which are averaged down to small, but not insignificant values within each gate time τ . These first differences (normalized by $1/\sqrt{2}$ to account for random noise in each entry) form a new data set which is first-order insensitive to long-term processes such as drift which dominate the directly recorded data.

Essentially the Allan Variance calculation presents us with a display of the laser's fractional frequency variation, σ_y , as a function of the time over which we are interested. At medium times, say τ of a few seconds, most laser stabilization systems will still be affected by the random measurement noise arising from shot noise and perhaps laser technical noise. At longer times the increased signal averaging implies a smaller residual fluctuation due to random processes. It is easy to show that the dependence of σ_y versus τ can be expected to be $1/\tau^{1/2}$, in the domain controlled by random (white) noise. The Allan deviation also has a great utility in compressing our statement of laser stability: we might say, for example, "the (in-)stability is 2×10^{-12} at 1 second, with the $1/\tau^{1/2}$ dependence which shows that only random noise is important out to a time of 300 s."

Allan Deviation Definition

With a counter linked to a computer, it is easy to gather a file of frequency values f_i measured in successive equal gate time intervals, t_g . Usually there is also some dead time, say t_d , while the counter-to-computer data transfers occur via the GPIB connection. This leads to a sample-to-sample time interval

of $t_s = t_g + t_d$. Allan variance is one half of the average squared difference between adjacent samples, and the usually quoted quantity, the Allan Deviation, is the square root of this averaged variance,

$$\sigma_y(\tau) = \left[\frac{1}{2(N-1)} \sum_{n=1}^{N-1} (f_{n+1} - f_n)^2 \right]^{1/2} \quad (1)$$

The dependence of σ_y upon the measuring time τ contains information essential for diagnosis of the system performance. These values for several times can be efficiently calculated from the (large) data set of frequencies observed for a fixed minimum gate time by adding together adjacent measurements to represent what would have been measured over a longer gate time. (This procedure neglects the effects of the small dead-time t_d , which are negligible for the white frequency noise $1/\sqrt{\tau}$ of usual interest but, for systems with drift and increased low-frequency noise, the dead-time effects can seriously impact the apparent results.) In any case, fewer samples will be available when the synthetic gate time becomes very long, so the uncertainty of this noise measurement increases strongly. Usually one insists on three or four examples to reduce wild variations, and so the largest synthetic gate time τ_{\max} will be taken to be the total measurement time/3. For a serious publication we might prefer 5 or 10 such synthetic measurements for the last point on the graph.

The Allan Deviation has one curiosity in the presence of a distinct sinusoidal modulation of the laser's frequency: when the gate time is 1/2 the sinusoid's period, adjacent samples will show the maximum deviation between adjacent measurements, leading to a localized peak in σ_y versus τ . Interestingly, there will be "ghosts" or aliases of this when the gate time/modulation period ratio is 1/4, 1/8, and so on. For longer gate times compared with the modulation, some fractional cycle memories can be expected also. So a clean slope of $-1/2$ for a log-log plot of σ_y versus τ shows that there is no big coherent FM process present.

Historically, Allan Variance has been valuable in locating time scales at which new physical processes must be taken into account. For example, at long times it is usual for a laser or other stable oscillator to reach a level of unchanging σ_y versus τ . We speak of this as a "flicker" floor. It arises from the interplay of two opposing trends: the first is the decreasing random noise with increasing τ (decreasing σ_y versus τ). At longer times one sees an increasing σ_y versus τ , due to drifts in the many system parameters (electronic offsets, temperature . . .), which make our lasers lock at points increasingly offset from the ideal one. If we wait long enough, ever larger changes become likely. So for several octaves of time, the combination of one decreasing and one increasing contribution leads to a flat curve. Eventually significant drift can occur even within one measurement time, and this will be mapped as a domain of rising σ_y increasing as the +1 power of τ .

It is useful to note that the frequency/time connection of the Allan Variance transformation involves very strong data compression and consequently cannot at all be inverted to recover the original data stream in the way we know from the Fourier transform pair. However in the other direction, we can obtain the Allan Deviation from the Phase Spectral Density.²

Spectral Noise Density

As noted earlier, when the number of individual contributions to the noise becomes too large to enumerate, it is convenient to move to a spectral density form of representation. To carry this idea forward, two natural quantities to use would be the frequency deviations occurring at some rate and the narrow bandwidth within which they occur. To work with a quantity that is positive definite and has additive properties, it is convenient to discuss the squared frequency deviations $\langle (f^2_N) \rangle$ which occur in a noise bandwidth B around the Fourier frequency f. This Frequency Noise Power Spectral Density, $S_f \equiv \langle (f^2_N) \rangle / B$, will have dimensions of Hz² (deviation²)/Hz (bandwidth). The summation of these deviations over some finite frequency interval can be done simply by integrating S_f between the limits of interest.

Connecting Allan Deviation and Spectral Density Sometimes one can estimate that the system has a certain spectrum of frequency variations described by $S_f(f)$, and the question arises of what Allan

Deviation this would represent. We prefer to use the Allan presentation only for experimental data. However Ref. 2 indicates the weighted transform from S_f to Allan Variance.

Connecting Linewidth and Spectral Density A small surprise is that an oscillator's linewidth generally will not be given by the summation of these frequency deviations! Why? The answer turns on the interesting properties of Frequency Modulated (FM) signals. What counts in distributing power is the Phase Modulation Index β , which is the peak modulation-induced phase shift or, equivalently, the ratio of the peak frequency excursion compared with the modulation rate. Speaking of pure tone modulation for a moment, we can write the phase-modulated field as

$$E(t) = \sin(\Omega t + \beta \sin(\omega t))$$

$$= J_0(\beta) \exp(i\Omega t) + \sum_{n=1}^{\infty} J_n(\beta) \exp(i(\Omega + n\omega)t) + \sum_{n=1}^{\infty} J_n(\beta) (-1)^n \exp(i(\Omega - n\omega)t) \quad (2)$$

where Ω is the "carrier" frequency, and $\omega = 2\pi f$ and its harmonics are the modulation frequencies. The frequency offset of one of these "sidebands," say the n th one, is n times the actual frequency of the process' frequency f . The strength of the variation at such an n th harmonic decreases rapidly for $n > \beta$ according to the Bessel function $J_n(\beta)$. We can distinguish two limiting cases.

Large excursions, slow frequency rate This is the usual laboratory regime with solid state or HeNe and other gas lasers. The dominant perturbing process is driven by laboratory vibrations that are mainly at low frequencies (5–200 Hz). The extent of the frequency modulation they produce depends on our mechanical design, basically how efficient or inefficient an "antenna" have we constructed to pick up unwanted vibrations. Clearly a very stiff, lightweight structure will have its mechanical resonances at quite high frequencies. In such case, both laser mirrors will track with nearly the same excursion, leading to small differential motion, i.e., low pickup of the vibrations in the laser's frequency. On the other hand, heavy articulated structures, particularly mirror mounts with soft springs, have resonances in the low audio band and lead to big FM noise problems. A typical laser construction might use a stiff plate, say 2 inches thick of Al or honeycomb-connected steel plates. The mirror mounts would be clamped to the plate, and provide a laser beam height of 2 inches above the plate. Neglecting air pressure variations, such a laser will have vibration-induced excursions $(\langle f_N^2 \rangle)^{1/2}$ of $\ll 100$ kHz. An older concept used low expansion rods of say 15 mm diameter Invar, with heavy Invar plates on the ends, and kinematic but heavy mirror mounts. This system may have a vibration-induced linewidth $(\langle f_N^2 \rangle)^{1/2}$ in the megahertz range. Only when the "rods" become several inches in diameter is the axial and transverse stiffness adequate to suppress the acceleration-induced forces. With such massive laser designs we have frequency excursions of tens to thousands of kilohertz, driven by low-frequency laboratory vibrations in a bandwidth $B < 1$ kHz. In this case $(\langle f_N^2 \rangle)^{1/2} \gg B$, and the resulting line shape is Gaussian. The linewidth is given by Ref. 3, $\Delta f_{\text{FWHM}} = [8 \ln(2) (\langle f_N^2 \rangle)]^{1/2} \cong 2.355 (\langle f_N^2 \rangle)^{1/2}$.

The broadband fast, small excursion limit This is the domain in which we can usually end up if we can achieve adequate servo gain to reduce the vibration-induced FM. Since the drive frequency of the perturbation is low, it is often feasible to obtain a gain above 100, particularly if we use a speedy transducer such as an acousto-optic modulator (AOM) or an electro-optic modulator (EOM). In general we will find a noise floor fixed, if by nothing else than the broadband shot noise which forms a minimum noise level in the measurement process. Here we can expect small frequency excursions at a rapid rate, $(\langle f_N^2 \rangle)^{1/2} \ll B$, leading to a small phase modulation index. If we approximate that the Spectral Noise Frequency Density $S_f = (\langle f_N^2 \rangle)/B$ is flat, with the value $S_f \text{ Hz}^2$ (deviation²)/Hz (bandwidth), then the linewidth in this domain is Lorentzian,³ with the $\Delta f_{\text{FWHM}} = \pi S_f = \pi (\langle f_N^2 \rangle)/B$.

This summary of frequency-domain measures is necessarily brief and the interested reader may find additional discussion useful.³⁻⁵ A number of powerful consequences and insights flow from reworking the above discussions in terms of a Phase Noise Power Spectral Density, $S_\phi = S_f/f^2$. The National Institute of Standards and Technology (NIST) Frequency and Time Division publishes collections

of useful tutorial and overview articles from time to time. The currently available volume² covers these topics in more detail. Vendors of rf-domain spectrum analyzers also have useful application notes.⁶

22.2 SERVO PRINCIPLES AND ISSUES^{7,8}

Bode Representation of a Servo System

We will describe our systems by transfer functions, output/input, as a function of Fourier frequency ω . We begin purely in the domain of electronics. The amplifier gain is $G(\omega)$. The electrical feedback is represented as $H(\omega)$. Both will have voltage as their physical domain, but are actually dimensionless in that they are output/input ratios. Considering that we will have to represent phase of these AC signals, both $G(\omega)$ and $H(\omega)$ will generally be complex. It will be fundamental to view these functions with their dependence on frequency, for both the amplitude and phase response.

Imagine a closed loop system with this amplifier as the forward gain $G(\omega)$ between input V_i and output V_o . Some fraction of the output is tapped off and sent back to be compared with the actual input. For more generality we will let $H(\omega)$ represent this feedback transfer ratio. The actual input, minus this sampled output will be our input to our servo amplifier $G(\omega)$. After a line of algebra we find the new gain of the closed loop—in the presence of feedback—is

$$A_{cl} = \frac{V_o}{V_i} = \frac{G(\omega)}{1 + G(\omega)H(\omega)} \quad (3)$$

A particularly instructive plot can be made for the product $G(\omega)H(\omega)$, called the “open loop gain,” which appears in the denominator. In this so-called Bode plot, the gain and phase are separately plotted. Also, from inspection of Eq. (3) we can learn one of the key advantages which feedback brings us: if the feedback factor GH were $\gg 1$, the active gain G would basically cancel out and we would be left with $A_{cl} \sim 1/H$. We imagine this feedback channel will be passive, formed from nearly ideal nondistorting components. The noise, exact value of the gain, and distortion introduced by it are seen to be nearly unimportant, according to the large magnitude of $1 + GH$. Gentle amplifier overload will lead to overtone production, but could alternatively be represented by a decrease of G with signal. Since the output doesn't depend sensitively upon G anyway, we are sure these distortion products and internally generated noise will be suppressed by the feedback. We can identify the denominator $1 + GH$ as the noise and distortion reduction factor.

What is the cost of this reduced dependence on the active components $G(\omega)$ and their defects? Basically it is that the gain is reduced and we must supply a larger input signal to obtain our desired output. For a music system one can then worry about the distortion in the preamplifier system. However, we want to make quiescent lasers, without the slightest hint of noise. So it is nice that the amplification of internal noise is reduced.

To be concrete, the circuit of Fig. 1 represents a common building block in our servo design. It also represents a simple case of feedback. We show it as a current summing input node: the subtraction at the input arises here because the sign of the gain is negative. With the nearly ideal high-gain operational amplifiers now available, $G \gg 1$ and we can closely approximate the closed loop gain by $1/H(\omega)$, yielding a flat gain above and a rising gain below some corner frequency $\omega_0 = 1/\tau_0$, with $\tau_0 = R_f C$. Remember $1/H(\omega)$ is the closed loop gain between V_o and V_i . To find the exact relationship between the signals V_s and V_o , we notice the related voltage-divider effect gives $V_i = (1 - H(\omega))V_s$, which leads to

$$\frac{V_o}{V_s} = -\frac{R_f}{R_i} \frac{(1 + j\omega/\omega_0)}{j\omega/\omega_0} = -\frac{R_f}{R_i} \frac{(1 + j\omega\tau_0)}{j\omega\tau_0} \quad (4)$$

The negative sign arises from the fact that the forward gain is negative. When the corner frequency ω_0 is chosen to be sufficiently high, we may have to consider the bandwidth issue of the OpAmp: $G(\omega)$ could

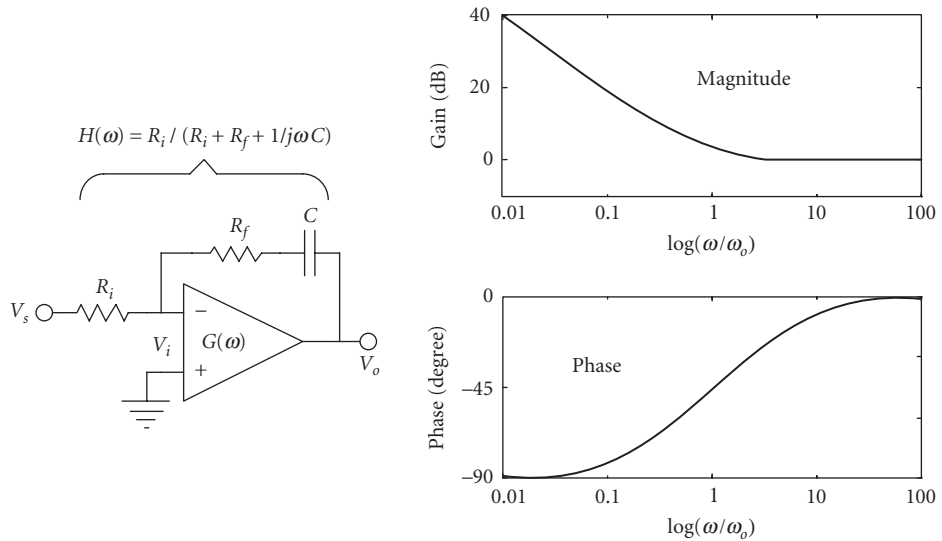


FIGURE 1 Phase and amplitude response of a Proportional-Integral (PI) amplifier circuit. The PI function is implemented using an inverting OpAmp.

start to roll off and no longer satisfy the approximation of $G \gg 1$. A more complex network is needed to compensate for the gain roll-off and that is exactly the topic of feedback we wish to cover below.

Phase and Amplitude Responses versus Frequency

We can plot⁹ the gain magnitude and phase of this elementary feedback example in Fig. 1, where we can see the flat gain at high frequencies and the rising response below ω_0 . Our laser servo designs will need to echo this shape, since the drift of the laser will be greater and greater at low frequencies, or as we wait longer. This will require larger and larger gains at low frequencies (long times) to keep the laser frequency nearby our planned lock point. The phase in Fig. 1 shows the lag approaching 90° at the lowest frequencies. (An overall minus sign is put into the subtractor unit, as our circuit shows an adder.) The time-domain behavior of this feedback system is a prompt inverted output, augmented later by the integration’s contribution.

As a first step toward modeling our realistic system, Fig. 2 shows the laser included in our control loop. The servo system’s job is to keep the laser output at the value defined by the reference or setpoint input. Some new issues will arise at the high-frequency end with the physical laser, as its piezo-electric transducer (PZT) will have time delay, finite bandwidth, and probably some resonances.

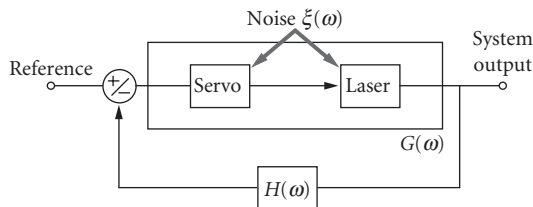


FIGURE 2 Model of laser system, including frequency noise, as part of a servo control loop.

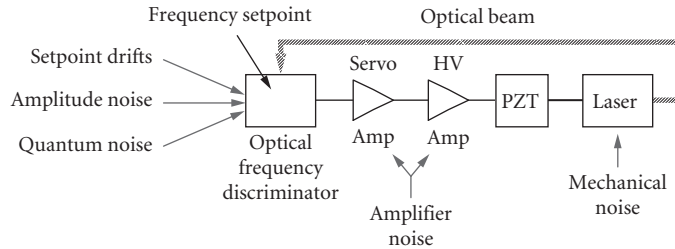


FIGURE 3 Detailed model of a frequency-controlled laser.

One way we should expand the model is to include the laser's operation as a frequency transducer, converting our control voltage to a frequency change. Probably the laser will have some unwanted frequency noises, and in Fig. 3 we can indicate their unwanted contributions arriving in the optical frequency discriminator, which functions like a summing junction. The emitted laser field encounters an optical frequency discriminator and the laser frequency is compared with the objective standard, which we will discuss below. In our diagram we show this laser frequency discriminator's role as an optical frequency-to-voltage converter element. More exactly, laser frequency *differences* from the discriminator's reference setpoint are converted to voltage outputs. Laser amplitude noises (due to the intrinsic property of the laser itself or external beam propagation) and vibration effects on the discriminator will appear as undesired additive noises also.

The first simple idea is that the feedback loop can be closed when the servo information, carried as a voltage signal in our amplifier chain, is converted to a displacement (in meters) by the PZT, then into laser frequency changes by the laser's standing-wave boundary condition. As the length changes, the "accordian" in which the waves are captive is expanded or compressed, and along with it the wavelength and frequency of the laser's light.

A second truth becomes clear as well: there is freedom in designating the division into the forward gain part and the feedback path part. Actually, we probably would like the laser to be tightly locked onto the control cavity/discriminator, and then we will tune the whole system by changing the set point, which is the discriminator's center frequency. This leads us to view the optical frequency discriminator as the summing junction, with the amplifier and PZT transducer as the forward gain part. The output is taken as an optical frequency, which would be directly compared to the setpoint frequency of the discriminator. So the feedback path $H = 1$.

We should consider some magnitudes. Let K_{PZT} represent the tuning action of the PZT transducer, expressed as displacement meter per volt. A typical value for this would be $K_{PZT} = 0.5 \text{ nm/V}$. The laser tunes a frequency interval $c/2L$ for a length change by $\lambda/2$, so the PZT tuning will be $\sim 600 \text{ V/order}$ at 633 nm.

$$K_V = K_{PZT} \frac{2}{\lambda} \frac{c}{2L} \quad (5)$$

So we obtain a tuning sensitivity $K_V \sim 800 \text{ kHz/V}$ tuning for a foot-long laser, assuming a disk-type PZT geometry. See the section below on PZT design.

Measurement Noise as a Performance Limit—It Isn't

Usually our desire for laser stability exceeds the range of the possible by many orders, and we soon wonder about the ultimate limitations. Surely the ultimate limit would be due to measurement noise. However, we rarely encounter the shot-noise-limited case, since the shot noise-limited S/N of a $100 \mu\text{W}$ locking signal is $\sim 6 \times 10^6$ in a 1 Hz bandwidth. (See section on "The Optical Cavity-Based Frequency Discriminator" later.) Rather we are dealing with the laser noise remaining because our

servo gain is inadequate to reduce the laser's intrinsic noise below the shot noise limit, the clear criterion of gain sufficiency. So our design task is to push up the gain as much as possible to reduce the noise, limited by the issue of stability of the thus-formed servo system.

Closed-Loop Performance Expectations When Transducer Resonance Limits the Usable Gain

Servo Stability: Larger Gain at Lower Frequencies, Decreasing to Unity Gain and Below Our need for high gain is most apparent in the low-frequency domain ~ 1 kHz and below. Vibrations abound in the dozens to hundreds of hertz domain. Drifts can increase almost without limit as we wait longer or consider lower Fourier frequencies. Luckily, we are allowed to have more gain at low frequencies without any costs in stability. At high frequencies, it is clear we will not help reduce our noise if our correction is applied too late and so no longer has an appropriate phase. One important way we can characterize the closed-loop behavior of our servo is by a time delay t_{delay} . Here we need to know the delay time before any servo response appears; a different (longer) time characterizes the $1/e$ response. The latter depends on the system gain, while the ultimate high-speed response possible is controlled by the delay until the first action appears. A good criterion is that the useful unity gain frequency can be as high as $f_{\tau} = 1/(2\pi t_{\text{delay}})$, corresponding to 1 rad extra phase-shift due to the delay. Below this ultimate limit we need to increase the gain—increase it a lot—to effectively suppress the laser's increased noise at low frequencies. This brings us to address the closed-loop stability issue.

Closed-Loop Stability Issues

One can usefully trace the damping of a transient input as it repetitively passes the amplifier and transducer, and is reintroduced into the loop by the feedback. Evidently stability demands that the transient is weaker on each pass. The settling dynamics will be more beautiful if the second-pass version of the perturbation is reduced in magnitude and is within say $\pm 90^\circ$ of the original phase. Ringing and long delay times result when the return phasor approaches -1 times the input signal vector, as then we are describing a sampled sinewave oscillation. These time-domain pictures are clear and intuitive, but require treatment in terms of convolutions, so here we will continue our discussion from the frequency-domain perspective that leads to more transparent algebraic forms. We can build up an arbitrary input and response from a summation of sinusoidal inputs. This leads to an output as the sum of corresponding sinusoidal outputs, each including a phase shift.

In our earlier simple laser servo example, no obvious limitation of the available closed-loop gain was visible. The trouble is we left out two fundamental laboratory parasites: time delay, as just noted, and mechanical resonances. We will usually encounter the mechanical resonance problem in any servo based on a PZT transducer. For design details, see the "Practical Issues" section. A reasonable unit could have its first longitudinal resonance at about 25 kHz, with a $Q \sim 10$. In servo terms, the actual mechanical PZT unit gives an added 2-pole roll-off above the resonance frequency and a corresponding asymptotic phase lag of 180° . Including this reality in our model adds another transfer function $R_{\text{PZT}} = \omega_0^2 / (\omega_0^2 + 2\omega\eta\omega_0 + \omega^2)$, where ω_0 is 2π times the resonance frequency, and $\eta = 1/2Q$ is the damping factor of the resonance. This response is shown in Fig. 4.

We now talk of stabilizing this system. The elements are the laser and some means to correct its frequency, a frequency discriminator to measure the difference between the actual and the setpoint frequencies, and a feedback amplifier. Here we propose to do the frequency control by means of a PZT transducer to change the laser frequency. For the present discussion, we assume the frequency discriminator has a flat response. For the feedback amplifier, the first appealing option is to try a pure integrator. The problem then is that we are limited in gain by the peakheight of the resonance which must remain entirely below unity gain to avoid instability. In Fig. 5 case (a) we see that the unity gain frequency is limited to a value of 1.5 kHz. Some margin is left to avoid excessive ringing near the resonant frequency, but it is still visible in the time domain. Techniques that help this case include a roll-off filter between the unity gain and PZT resonance frequencies.

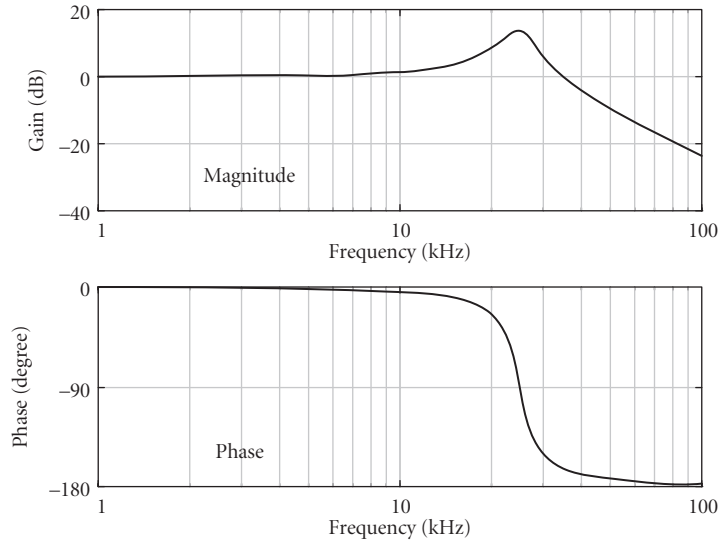


FIGURE 4 The amplitude and phase response of a tubular PZT transducer and an 8-mm-diameter by 5-mm-thick mirror. The resonance is at 25 kHz with a Q of 10.

Figure 5 shows the “open loop” gain function GH of the feedback equation, and the corresponding phase response. We already noted the dangerous response function of -1 where the denominator of Eq. (3) vanishes. In the time-domain iterative picture, the signal changes sign on successive passes and leads to instability/oscillation. We need to deal with care as we approach near this point in order to obtain maximum servo gain: it is useful to consider two stability margins. The *phase stability margin* is the phase of the open-loop function when the gain is unity. It needs to be at least 30° . The *gain*

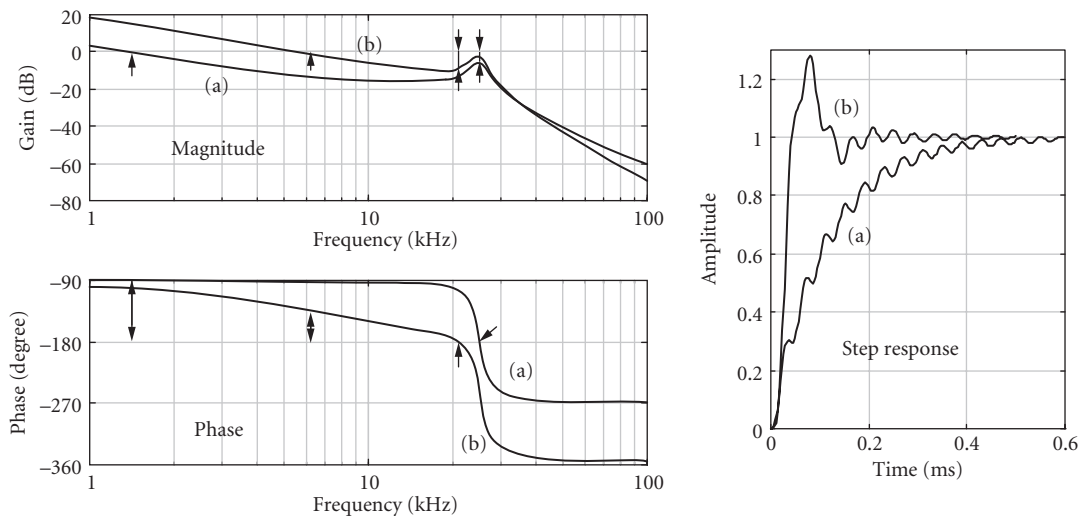


FIGURE 5 (a) Integrator gain function alone. Gain must be limited so that gain is <1 even at the resonance. (b) Single-pole low pass at 6 kHz inserted. Now unity gain can increase to 6 kHz and time response is ~ 3 -fold faster. Small arrows in the graph indicate the phase margin at the unity gain frequency (gain = 0 dB) and gain margin at a phase shift of -180° .

margin is the closed loop gain when the phase is 180° . In Fig. 5 case (a) we see that the phase is not shifted very much until we really “sense” the amplitude increase from the resonance. So this resonance may tend to fix an apparently solid barrier to further servo improvement. But as shown in Fig 5 case (b), just a low-pass to push down the PZT resonance is very helpful.

In fact, there are many ways of improving the low frequency gain of this system. They include: (1) imposing yet another high frequency roll-off (or multi-pole low pass filter) just before the resonance thus pushing its height down and allowing the open loop transfer function to come up, (2) adding lag compensators before the resonance to push the low frequency gain up while keeping the high frequency response relatively unchanged, (3) adding a lead compensator just above the resonance to advance the phase and increase the unity gain point, (4) or placing a notch at the resonant frequency to “cut it out” of the open loop transfer function. The last two options in this list are quite promising and are discussed in more detail below.

Proportional Integral Derivative (PID) Controller versus Notch Filters Like many “absolute” barriers, it is readily possible to shoot ahead and operate with a larger closed loop bandwidth than that represented by the first PZT resonance. The issue is that we must control the lagging phase that the resonance introduces. A good solution is a differentiator stage, or a phase lead compensator, which could also be called a high frequency boost/gain-step circuit. In Fig. 6 case (a) we show the Bode plot of our PZT-implemented laser frequency servo, based on a PID (Proportional Integral Differentiator) controller design. Just a few moments of design pay a huge benefit, as the unity gain frequency has now been pushed to 40 kHz, almost a factor of 2 *above* the PZT’s mechanical resonance. For this PID controller example, unity gain occurs at a 7-fold increased frequency compared with Fig. 5 case (b). Thus at the lower frequencies we would hope to have increased the servo gain by a useful factor of 7x or 17 dB. However, comparison of Fig. 5 (b) and 6 (a) shows that the low frequency gain is hardly changed, even though we greatly increased the servo bandwidth.

So, how *do* we go forward? We could in principle continue to increase the gain and unity gain frequency, but this is not really practical, however, since we will again be limited by additional structure resonances that exist beyond the first resonance. Also, the Derivatives needed to tame these resonances cost low frequency gain, and it is hard to win overall system performance. To make

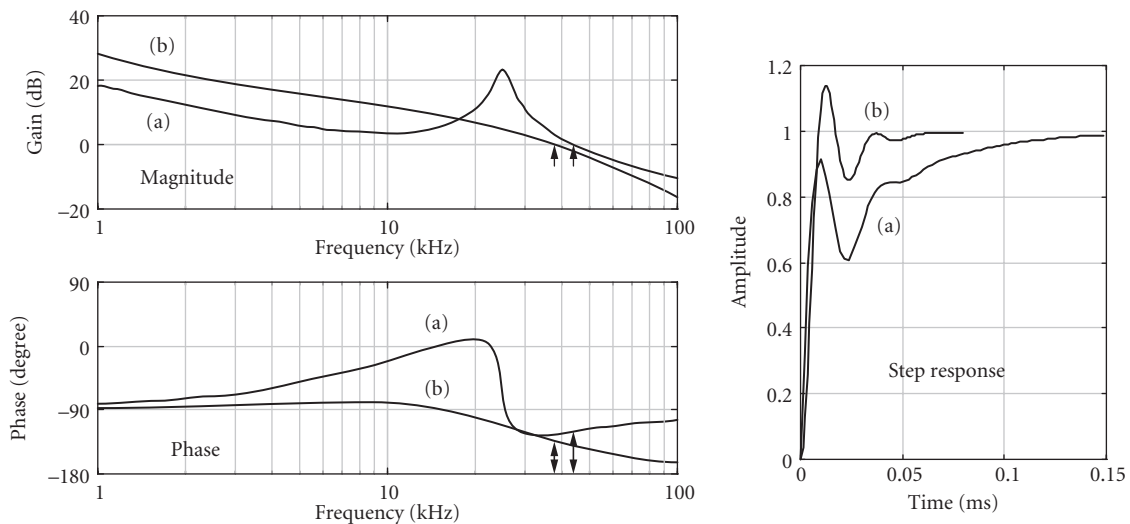


FIGURE 6 Two methods of working through and beyond a resonance. (a) PID controller where the Derivative term advances the phase near the resonance. (b) Adding a notch is a better approach, where the notch function approximates the inverse of the resonance peak. Transient response settles much more quickly. Again, we use the small arrows in the graph to indicate the phase margin at the unity gain frequency.

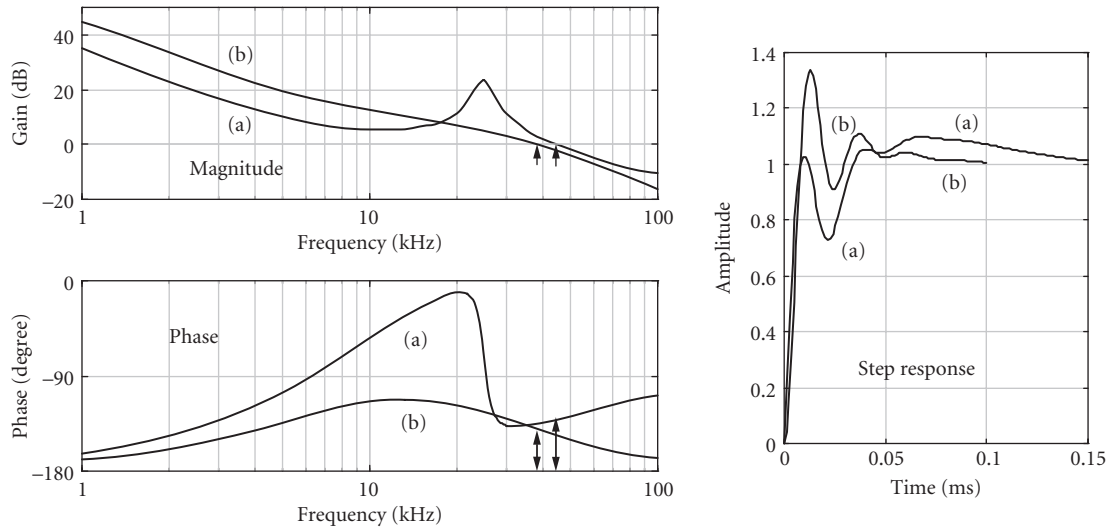


FIGURE 7 Adding an additional PI stage to (a) the PID and (b) the PI-plus-notch stabilizers of Fig. 6. Note that low frequency gain is strongly increased.

progress, we use a notch as an alternative technique to suppress the resonance. Now a D term is not needed, and we can conserve the gain at low frequencies. The notch filter, combined with a PI stage, gives unity gain at higher frequencies, and increases gain for low ones. See Fig. 6 case (b). Then Fig. 7 compares adding another PI stage to the two cases of Fig. 6, adding a PID in (a) and a Notch plus PI in (b). The time-domain approach, shown in Fig. 7, shows case (b) settles rather nicely. *And* the gain has increased more than 20 dB at frequencies of 1 kHz and below. So this is very encouraging.

While we have come to the cascaded-integrators approach cautiously in this discussion, in fact at least 2 integrators would always be used in practice. Workers with serious gain requirements, for example, the LIGO and VIRGO interferometric gravitational wave detector groups, may use the equivalent of 4 cascaded integrators! Such a design is “conditionally stable” only, meaning that the gain cannot be smoothly reduced or increased. Such aggressive stabilizer designs have their place, but not for a first design!

“Rule of Thumb” PID Design for System with a Transducer Resonance Optimizing servo performance is an elegant art, turned into science by specification of our “cost function” for the system performance shortcomings. In the case that we wish to minimize the time-integrated magnitude of the residuals following a disturbance, one comes to the case studied by Ziegler and Nichols for the PID controller used in a system with a combined roll-off and time delay.⁸ Such a case occurs also in thermal controllers. With only the P term, one first looks for the frequency f_{osc} where the system first oscillates when the gain is increased. The PD corner is then set $1.27\times$ higher than this f_{osc} , the P gain is reset at 0.6 of the oscillation gain, and the PI corner is set at 0.318 times the oscillation frequency. This “rule of thumb” design of the phase compensation produces a transient response which settles reasonably well, so as to minimize the Time Integrated Error. For phase-locking lasers, a cost function with more emphasis on long-lasting errors leads to another kind of “optimum” tuning, but with qualitatively similar results.

When a notch is used to suppress the resonance, there is no longer an anomalous gain at the resonant frequency and one is returned to the same case as in its absence. A reasonable servo approach to using two PI stages is to design with only one, achieving the desired unity gain frequency. The second PI is then added to have its corner frequency at this same point or up to 10-fold lower in frequency, depending on whether we wish the most smooth settling or need the highest feasible low frequency gain. Figures 6b and 7 show the Bode plot of such designs, along with the system’s

closed-loop transient response. An elegant strategy is to use adaptive clamping to softly turn on the extra integrator stage when the error is small enough, thus dynamically increasing the order of the controller when it will not compromise the dynamics of recovery.

22.3 PRACTICAL ISSUES

Here we offer a number of important tidbits that are useful background material for a successful application of the grand schemes discussed above.

Frequency Discriminators for Laser Locking—Overview

So far we have devoted our main effort to addressing the issues of the feedback scheme. Of equal importance is the subject of frequency reference system. After all, a good servo eliminates intrinsic noises of the plant (laser), and replaces them with the measurement noise associated with the reference system. Indeed, development of prudent strategies in high precision spectroscopy and the progress of laser stabilization have been intimately connected to each other through the years,¹⁰ with the vigorous pursuit of resolution and sensitivity resulting in amazing achievements in both fields.

To stabilize a laser, one often employs some kind of resonance information to derive a frequency/phase-dependent discrimination signal. The resonance can be of material origin, such as modes of an optical interferometer; or of natural origin, such as atomic or molecular transitions. If the desired use of a stabilized laser is to be an optical frequency standard, its long-term stability or reproducibility will be key, so the use of a natural resonance is preferred. Reproducibility is a measure of the degree to which a standard repeats itself from unit to unit and upon different occasions of operation. The ultimate reproducibility is limited to the accuracy of our knowledge of the involved transitions of free atoms or molecules. The term “free” means the resonance under study has a minimum dependence on the laboratory conditions, such as the particle moving frame (velocity), electromagnetic fields, collisions, and other perturbations. To realize these goals, modern spectroscopy has entered the realm of quantum-limited measurement sensitivities and exquisite control of internal and external degrees of freedom of atomic motions.

A careful selection of a high-quality resonance can lead to superior system performance and high working efficiency. For example, the combined product of the transition quality factor Q and the potential signal-to-noise ratio (S/N) is a major deciding factor, since this quantity controls the time scale within which a certain measurement precision (fractional frequency) can be obtained. This importance is even more obvious when one considers the waiting time for a systematic study is proportional to the inverse square of ($Q \times S/N$). A narrower transition linewidth of course also helps to reduce the susceptibility to systematic errors. The resonance line shape is another important aspect to explore. By studying the line shape we will find out whether we have come to a complete understanding of the involved transition and whether there are other unresolved small lines nearby ready to spoil our stabilization system.

Sometimes it may not be sufficient to use the natural resonance alone for stabilization work, or may not be necessary. The saturation aspect of the atomic transition limits the attainable S/N . To stabilize a noisy laser we need to use, for example, an optical resonator, which can provide a high-contrast and basically unlimited S/N of the resonance information. Careful study of the design and control of the material properties can bring the stability of material reference to a satisfactory level. See below for a more detailed discussion on this topic.

Ideally, a resonance line shape is even symmetric with respect to the center frequency of the resonance, and deviations from this ideal case will lead to frequency offsets. However, for the purpose of feedback, the resonance information needs to be converted to an odd symmetric discriminator shape: we need to know in which direction the laser is running away from the resonance. A straightforward realization of an error signal using direct absorption technique is to have the laser tuned to the side of resonance.¹¹ The slope of the line is used to convert the laser frequency noise to amplitude information

for the servo loop. This technique is essentially a DC approach and can suffer a huge loss in S/N due to the low-frequency amplitude noise of the laser. A differential measurement technique using dual beams is a requirement if one wishes to establish a somewhat stable operation. With a dual beam approach, the information about the laser noise can be measured twice and therefore it is possible to completely eliminate the technical noise and approach the fundamental limit of shot noise using clever designs of optoelectronic receivers. Conventional dual beam detection systems use delicate optical balancing schemes,¹² which are often limited by the noise and drift of beam intensities, residual interference fringes, drift in amplifiers, and spatial inhomogeneity in the detectors. Electronic auto-cancellation of the photodetector currents has provided near shot noise-limited performance.¹³ Although this process of input normalization helps to increase S/N of the resonance, the limitation on the locking dynamic range remains a problem. The servo loop simply gets lost when the laser is tuned to the tail or over the top of the resonance. Further, it is found that transient response errors basically limit the servo bandwidth to be within the cavity linewidth.¹⁴ Another effective remedy to the DC measurement of resonances is the use of zero-background detection techniques, for example, polarization spectroscopy.^{15,16} In polarization spectroscopy the resonance information is encoded in the differential phase shifts between two orthogonally polarized light beams. Heterodyne detection between the two beams can reveal an extremely small level of absorption-induced polarization changes of light, significantly improving the detection sensitivity. However, any practical polarizer has a finite extinction ratio (ϵ) which limits the attainable sensitivity. Polarization spectroscopy reduces the technical noise level by a factor of $\sqrt{\epsilon}$, with $\epsilon \sim 10^{-7}$ for a good polarizer. Polarization techniques do suffer the problem of long-term drifts associated with polarizing optics.

Modulation techniques are of course often used to extract weak signals from a noisy background. Usually noises of technical origins tend to be more prominent in the low frequency range. Small resonance information can then be encoded into a high-frequency region where both the source and the detector possess relatively small noise amplitudes. Various modulation schemes allow one to compare on-resonant and off-resonant cases in quick succession. Subsequent demodulations (lock-in detection) then simultaneously obtain and subtract these two cases, hence generating a signal channel with no output unless there is a resonance. Lorentzian signal recovery with the frequency modulation method has been well documented.¹⁷ The associated lock-in detection can provide the first, second, and third derivative type of output signals. The accuracy of the modulation waveform can be tested and various electronic filters can be employed to minimize nonlinear mixing among different harmonic channels and excellent accuracy is possible. In fact, the well-established 633-nm HeNe laser system¹⁸ is stabilized on molecular iodine transitions using this frequency dither technique and third harmonic (derivative) signal recovery. Demodulation at the third or higher order harmonics helps to reduce the influence of other broad background features.¹⁹ The shortcoming of the existence of dither on the output beam can be readily cured with an externally implemented "un-dithering" device based on an AOM.²⁰ However, in this type of modulation spectroscopy the modulation frequency is often chosen to be relatively low to avoid distortions on the spectral profile by the auxiliary resonances associated with modulation-induced spectral sidebands. An equivalent statement is that the line is distorted because it cannot reach an equilibrium steady state in the face of the rapidly tuning excitation. This low-frequency operation (either intensity chopping or derivative line shape recovery) usually is still partly contaminated by the technical noise and the achievable signal-to-noise ratio (S/N) is thereby limited. To recover the optimum signal size, large modulation amplitudes (comparable to the resonance width) are also employed, leading to a broadened spectral linewidth. Therefore the intrinsic line shape is modified by this signal recovery process and the direct experimental resolution is compromised.

A different modulation technique was later proposed and developed in the microwave magnetic resonance spectroscopy and similarly in the optical domain.²¹⁻²³ The probing field is phase-modulated at a frequency much larger than the resonance linewidth under study. When received by a square-law photodiode, the pure FM signal will generate no photocurrent at the modulation frequency unless a resonance feature is present to upset the FM balance. Subsequent heterodyne and rf phase-sensitive detection yield the desired signal. The high sensitivity associated with the FM spectroscopy is mainly due to its high modulation frequency, usually chosen to lie in a spectral region where the amplitude noise level of the laser source approaches the quantum (shot noise) limit. The redistribution of some of the carrier

power to its FM sidebands causes only a slight penalty in the recovered signal size. Another advantage of FM spectroscopy is the absence of linewidth broadening associated with low-frequency modulation processes. The wide-spread FM spectra allows each individual component to interact with the spectral features of interest and thereby preserves the ultrahigh resolution capability of contemporary narrow-linewidth lasers.

Since its invention, FM spectroscopy has established itself as one of the most powerful spectroscopic techniques available for high-sensitivity, high-resolution, and high-speed detection. The high bandwidth associated with the radio frequency (rf) modulation enables rapid signal recovery, leading to a high Nyquist sampling rate necessary for a high-bandwidth servo loop. The technique has become very popular in nonlinear laser spectroscopy,²⁴ including optical heterodyne saturation spectroscopy,²³ two-photon spectroscopy,²⁵ Raman spectroscopy,²⁶ and heterodyne four-wave mixing.²⁷ Recent developments with tunable diode lasers have made the FM technique simpler and more accessible. The field of FM-based laser diode detection of trace gas and remote sensing is rapidly growing. In terms of laser frequency stabilization, the rf sideband based Pound-Drever-Hall locking technique²⁸ has become a uniformly adopted fast stabilization scheme in the laser community. The resonance-based error signal in a high-speed operating regime is shown to correspond to the instantaneous phase fluctuations of the laser, with the atom or optical cavity serving the purpose of holding the phase reference. Therefore a properly designed servo loop avoids the response time of the optical phase/frequency storage apparatus and is limited only by the response of frequency-correcting transducers.

In practice some systematic effects exist to limit the ultimate FM sensitivity and the resulting accuracy and stability. Spurious noise sources include residual amplitude modulation (RAM), excess laser noise, and étalon fringes in the optical system.²⁹ A number of techniques have been developed to overcome these problems. In many cases FM sidebands are generated with electro-optic modulators (EOM). A careful design of EOM should minimize the stress on the crystal and the interference between the two end surfaces (using angled incidence or antireflection coatings). Temperature control of the EOM crystal is also important and has been shown to suppress the long-term variation of RAM.³⁰ The RAM can also be reduced in a faster loop using an amplitude stabilizer³¹ or a tuning filter cavity.³² The étalon fringe effect can be minimized by various optical or electronic means.³³ An additional low-frequency modulation (two-tone FM³⁴) can be used to reduce drifts and interference of the demodulated baseline.

In closing this section, we note that a laser is not always stabilized to a resonance but is sometimes referenced to another optical oscillator.³⁵ Of course the working principle does not change: one still compares the frequency/phase of the laser with that of reference. The technique for acquiring the error information is however more straightforward, often with a direct heterodyne detection of the two superposed waveforms on a fast photo detector. The meaning of the fast photo detector can be quite extensive, sometimes referring to a whole table-top system that provides THz-wide frequency gap measurement capabilities.^{36–38} Since it is the phase information that is detected and corrected, an optical phase locked loop usually provides a tight phase coherence between two laser sources. This is attractive in many measurement applications where the relative change of optical phase is monitored to achieve a high degree of precision. Other applications include phase-tracked master-slave laser systems where independent efforts can be made to optimize laser power, tunability and intrinsic noise.

The Optical Cavity-Based Frequency Discriminator

It is difficult to have both sensitive frequency discrimination and short time delay, unless one uses the reflection mode of operation: these issues have been discussed carefully elsewhere.²⁸ With ordinary commercial mirrors, we can have a cavity linewidth of 1 MHz, with a contrast C above 50 percent. We can suppose using 200 μW optical power for the rf sideband optical frequency discriminator, leading to a dc photo current i_0 of $\sim 100 \mu\text{A}$ and a signal current of $\sim 25 \mu\text{A}$. The shot noise of the dc current is $i_n = \sqrt{2ei_0}$ in a 1-Hz bandwidth, leading to an S/N of $\sim 4 \times 10^6$. The frequency noise-equivalent would then be 250 millihertz/ $\sqrt{\text{Hz}}$. If we manage to design enough useful gain in the controller to suppress

the laser's intrinsic noise below this level, the laser output frequency spectrum would be characterized by this power spectral density. Under these circumstances, according to the earlier discussion in the Introduction and Overview Section, the output spectrum would be Lorentzian, of width $\Delta\nu_{\text{FWHM}} = \pi S_f = \pi (0.25 \text{ Hz})^2/\text{Hz} \sim 0.8 \text{ Hz}$. One comes to impressive predictions in this business! But usually the results are less impressive.

What goes wrong? From measurements of the servo error, we can see that the electronic lock is very tight indeed. However, the main problem is that vibrations affect the optical reference cavity's length and hence its frequency. For example measurements show the JILA Quiet Room floor has a seismic noise spectrum which can be approximated by $4 \times 10^{-9} \text{ m rms}/\sqrt{\text{Hz}}$ from below 1 Hz to about 20 Hz, breaking there to an f^{-2} roll-off. Below 1 Hz the displacement noise climbs as f^{-3} . Horizontal and vertical vibration spectra are similar. Accelerations associated with these motions lead to forces on the reference cavity that will induce mechanical distortion and hence frequency shifts. To estimate the resulting frequency shift, simple approximate analysis leads to a dynamic fractional modulation of the cavity length l by the (colinear) acceleration a , as

$$\left. \frac{\Delta l}{l} \right|_{\text{axial}} = -\frac{\Delta f}{f} = \frac{a \rho l \varepsilon}{2Y} \quad (6)$$

where $Y \sim 70 \text{ GPa}$ is the Young's modulus and $\rho \sim 2.2 \text{ gm/cm}^3$ is the density for the ULE (or Zerodur) spacer. The factor ε ($-1 < \varepsilon < 1$) is a geometrical design factor. For example, suppose the cavity is hanging vertically, suspended from the top. Then the cavity is stretched by its weight, and $\varepsilon = 1$. Using $l = 10 \text{ cm}$ and $a = 1 \text{ g}$, we expect $\Delta l/l = -\Delta f/f \sim 1.5 \times 10^{-8} \rightarrow \sim 8.7 \text{ MHz/g}$, supposing $\lambda = 532 \text{ nm}$. (This is equivalent to 885 kHz/ms^{-2} .) If the cavity were vertical, but supported from below, it would be in compression and $\varepsilon = -1$. Evidently there is an interesting regime in which the cavity is supported near its middle height, where there will be a strong cancelation of the net vertical length change. We return below to this case where $\varepsilon \sim 0$.

First, let us suppose our reference cavity bar is uniformly supported horizontally from a flat horizontal surface. Even in this transverse case, vertically accelerating the interferometer produces length changes through the distortion coupling between the transverse compression and lengthwise extension, the effect of "extrusion of the toothpaste." So the longitudinal displacement of Eq. (6) is reduced by this Poisson ratio $\sigma = 0.17$. Also the vertical weight now comes from the cavity's height, which is now really the spacer's diameter ϕ , typically about 5-fold less than the length. So we have

$$\left. \frac{\Delta l}{l} \right|_{\text{transverse}} = \frac{a \phi \rho \sigma}{2Y} \quad (7)$$

We come to a predicted sensitivity then of $\sim 300 \text{ kHz/g}$ for vertically applied uniform force (equivalent to 30 kHz/ms^{-2}).

Some important things have been so far left out of this discussion. For one, to make a stable reference cavity the details of the mounting and cavity support can be very important, since the expansion coefficient of the metal vacuum envelope is likely three orders of magnitude greater than that of ULE near its critical temperature-stable point. To prevent the vacuum shell's dimensional expansion from causing stresses in the cavity, it makes sense to use a pendulum suspension of the cavity. With two loops around the horizontal bar, forming a dual pendulum suspension, the cavity motion is mainly restricted to the axial direction, and the horizontal acceleration forces at high frequency are filtered down. Now we have the question:

What should be the spacing B between the two suspension loops? Put them close together and the expansion of the metal outer shell has even less impact on the cavity length. But the cavity rod (or bar, or tube) now takes on a stronger static bend, which shortens the cavity and the resulting cross-term in the cosine projection leads to a first-order length response with vertical acceleration noise. Furthermore, the bending-induced misalignment of the cavity mirrors means the intracavity resonant mode will displace laterally across the mirror surface to again have the optimal standing-wave buildup. Certainly the mirrors are rather nicely polished on their surfaces, but at least one is

a curved surface. So with our greedy dream of 10^{-15} frequency stability, wiping the beam vertically across the curvature will introduce disastrous optical length changes.

What about a wider spacing of the supports? Luckily for us the “two-point suspension problem” was addressed by G. B. Airy in the nineteenth century. He established that a support-spacing-to-length ratio of $B/L = 0.577$ was an ideal design for such a suspension, as it restored the parallelism of the two end faces of the measurement bar. A series of JILA experiments explored this domain.³⁹ These showed a vertical acceleration sensitivity of the horizontally suspended bar of 2200 kHz/ms^{-2} at $B/L = 0.11$, reduced to 150 kHz/ms^{-2} at the Airy spacing $B/L = 0.577$. Our “theory” in Eq. (7) doesn’t consider static bending of the bar, but would lead to 90 kHz/ms^{-2} if scaled for the $5.7 \times 7.1 \times 27.7 \text{ cm}$ dimensions of our cavity’s spacer-bar (suspended with the 5.1-cm direction vertical). Regrettably, the sign of the vibration-induced response was not determined: cavity bending shortens the optical path, while vertical squeezing would lengthen it.

Integrating the acceleration produced by the mid-band floor vibration spectrum quoted above leads to a broadband noise of a few dozen hertz in both H and V planes. Left out however is the 1 milli-“g” vibration near 30 Hz due to ac motors in JILA (Pepsi refrigerators!). So we should have a vibration-induced linewidth of something like $1/2 \text{ kHz}$, which correlates adequately well with experience. Passive air-table antivibration measures suppress this vibration (acceleration) spectrum to $\sim 2 \times 10^{-6} \text{ ms}^{-2}/\sqrt{\text{Hz}}$, again roughly flat over 2 to 20 Hz band by filtering the floor’s vibrational noise above $\sim 2 \text{ Hz}$ Fourier frequency. The calculated Fourier frequency at which the phase modulation processes have removed $1/2$ the laser carrier power (approximate half-linewidth of the locked laser) is $\sim 1 \text{ Hz}$, but nonmodeled noise led to experimental values more like 5 Hz. Elegant passive vibration-damping suspensions at NIST have led to record-level subhertz cavity-locked laser linewidths.⁴⁰ It has been suggested that much of the remaining noise is associated with thermal mechanical displacement noise in the mirror coatings.⁴¹ Later measurements confirmed that the thermal noise was indeed the dominant source limiting the laser linewidth.⁴²

Returning to the cavity-mounting problem, we introduced the symmetry factor ε in the axial direction, because it is clear that holding the cavity in the midplane seems wise. Then the acceleration-induced net length change would tend be cancelled: one half of the length is under compression, the other half is under tension at a particular moment in the ac vibration cycle. We denote this cancellation by symmetry as ε , with $-1 \leq \varepsilon \leq 1$. Some experiments were made with short vertically mounted cavities.⁴³ The hand-assembly of the central disk limited the observed asymmetry value for our vertical mountings to a ~ 20 -fold reduction of the vibration sensitivity ($\varepsilon \geq 0.05$), to about $\sim 10 \text{ kHz/ms}^{-2}$, measured at the Nd fundamental wavelength. It was directly possible to observe subhertz laser beats! A computer design⁴⁴ for a more optimal cavity is shown as Fig. 8.

Quantum Resonance Absorption⁴⁵

Establishing a long-term stable optical frequency standard requires a natural reference of atomic or molecular origin. Historically, the use of atomic/molecular transitions was limited to those that had accidental overlap with some fixed laser wavelengths. With the advent of tunable lasers, research on quantum absorbers has flourished. A stabilized laser achieves fractional frequency stability

$$\frac{\delta\nu}{\nu} = \frac{1}{Q} \frac{1}{S/N} \frac{1}{\sqrt{\tau}}$$

where Q is the quality factor of the transition involved, S/N is the recovered signal-to-noise ratio of the resonance information, and τ is the averaging time. Clearly one wishes to explore the limits on both resolution and sensitivity of the detected signal. The nonlinear nature of a quantum absorber, while on one hand limiting the attainable S/N , permits sub-Doppler resolutions. With sensitive techniques such as FM-based signal modulation and recovery, one is able to split a MHz scale linewidth by a factor of 10^4 to 10^5 , at an averaging time of 1 s or so. Sub-Hertz long-term stability can be achieved with carefully designed optical systems where residual effects on baseline stability are minimized. However, a pressing question is: How accurate is our knowledge of the center of the resonance? Collisions, electromagnetic fringe fields, probe field wavefront curvature, and probe

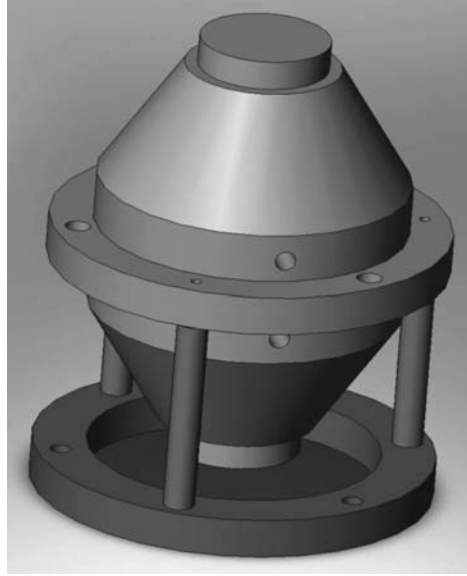


FIGURE 8. Computer model of vibration-resistant optical reference cavity.

power can all bring undesired linewidth broadening and center shifts. Distortion in the modulation process and other physical interactions can produce asymmetry in the recovered signal line shape. These issues will have to be addressed carefully before one can be comfortable talking about accuracy. A more fundamental issue related to time dilation of the reference system (second-order Doppler effect) can be solved in a controlled fashion: one simply knows the sample velocity accurately (e.g., by velocity selective Raman process), or the velocity is brought down to a negligible level using cooling and trapping techniques.

The simultaneous use of quantum absorbers and an optical cavity offers an attractive laser stabilization system. On one hand, a laser prestabilized by a cavity offers a long phase-coherence time, reducing the need of frequent interrogations of the quantum absorber. In other words, information of the atomic transition can be recovered with an enhanced S/N and the long averaging time translates into a finer examination of the true line center. On the other hand, the quantum absorber's resonance basically eliminates inevitable drifts associated with material standards. The frequency offset between the cavity mode and atomic resonance can be bridged by an AOM. In this case the cavity can be made of totally passive elements: mirrors are optically contacted to a spacer made of ultralow expansion material such as ULE or Zerodur. In case that the cavity needs to be made somewhat tunable, an intracavity Brewster plate driven by Galvo or a mirror mounted on PZT is often employed. Of course, these mechanical parts bring additional thermal and vibrational sensitivities to the cavity, along with nonlinearity and hysteresis. Temperature tuning of a resonator is potentially less noisy but slow. Other tuning techniques also exist, for example, through the use of magnetic force or pressure (change of intracavity refractive index or change of cavity dimension by external pressure). An often-used powerful technique called frequency-offset-locking brings the precision rf tuning capability to the optical world.¹⁴

Transducers

PZT Transducer Design: Disk versus Tube Designs We will usually encounter the mechanical resonance problem in any servo based on a PZT transducer: Small mirrors clearly are nice as they can have higher resonance frequencies. A mirror, say 7.75 mm $\Phi \times 4$ mm high, might be waxed onto a

PZT disk 10 mm diameter \times 0.5 mm thick. The PZT, in turn, is epoxied onto a serious backing plate. This needs to be massive and stiff, since the PZT element will produce a differential force between the mirror and the backing plate. At short times there will be a “reduced mass” kind of splitting of the motion between the mirror and the support plate. At lower frequencies, one hates to get a lot of energy coupled into the mirror mount since it will have a wealth of resonances in the sub-kHz range. For this size mirror, the backing plate might be stainless steel, 1 inch diameter by $\sim 3/4$ inch wedged thickness, and with the PZT deliberately decentered to break down high Q modes. The piston mode will be at ~ 75 kHz.

Tubular PZT Design Often it is convenient to use a tubular form of PZT, with the electric field applied radially across a thin wall of thickness t . This gives length expansion also, transverse to the field using a weaker d_{31} coefficient, but wins a big geometric factor in that the transverse field is generating a length response along the entire tube height h . The PZT tube could be $1/2$ inch diameter by $1/2$ inch length, with a wall thickness $t = 1.25$ mm. This geometry leads to a ~ 7 -fold sensitivity win, when $d_{31} \approx 0.7 d_{33}$ is included. Typical dimensions for the mirror might be 12.5 mm diameter \times 7 mm high. The PZT tube also is epoxied onto a serious backing plate. For the high voltage isolation of the PZT electrodes at the tube ends, a thin sheet (say < 0.5 mm) of stiff ceramic, alumina for example, will suffice. An alternative way to provide the electrical isolation of the ends involves removing the silver electrodes for several millimeters at the end. A new technique uses a diamond-charged tubular core drill mounted into a collet in a lathe. The active tool face projects out only 2 mm so that handheld PZT grinding leads to clean electrode removal, inside and out. This end of the PZT tube is attached to the backing mass with strong epoxy. The mirror is attached to the open PZT tube end with melted wax. This is vastly better than epoxy in that it does not warp the optic, and the small energy dissipation occurs at the best place to damp the Q of the PZT assembly. If done well, this unit will have its first longitudinal resonance at about 25 kHz, with a $Q \sim 10$. As noted above, in servo terms, the actual mechanical PZT unit gives an added 2-pole roll-off above the resonance frequency and a corresponding asymptotic phase lag of 180° . So it is useful to design for high resonant frequency and low Q.

Comparing disk and tubular designs, the disk approach can have a three-fold higher resonance frequency, while the tubular design is ~ 7 -fold more sensitive. Perhaps more important is the tube's reduced stiffness, moving the PZT/mirror resonance down into the 20-kHz domain. This brings us to the subject of spectral shaping of the amplifier gain and limitations of servo performance due to electronic issues.

Amplifier Strategies for PZT Driver We enjoy the tubular PZT for its large response per volt and its relatively high resonance frequencies. But it gives a problem in having a large capacitance, for example, of 10 nF in the above design. Even with the high sensitivity of 70 V/order, achieving a tight lock requires high frequency corrections and can lead to a problem in supplying the necessary ac current, supposing that we ask the HV amplifier alone to do the job. An apparent answer is to use a pair of amplifiers, one fast and the other HV, separately driving the two sides of the PZT. This alone doesn't solve the problem, as the big high-frequency ac current is only returned via the HV amplifier. The answer is to use a crossover network on the HV amplifier side. A capacitor to ground, of perhaps 3- or 5-fold larger value than the PZT will adequately dump the fast currents coming through the PZT's capacitance. A resistor to this PZT/shunt capacitor junction can go to the HV amplifier. Now this HV amp has indeed more capacitance to drive, but is only needed to be active below a few hundred hertz where the current demand becomes reasonable. An alternative topology sums the two inputs on one side of the PZT.

Other Useful Transducers—Slow but Powerful Commercial multiple wafer designs utilize 100 or more thin PZT sheets mechanically in series and electrically in parallel to produce huge excursions such as $10 \mu\text{m}$ for 100 V. Of course the capacitance is $\sim 0.1 \mu\text{F}$ and the stability leaves something to be desired. These are useful for applications that can tolerate some hysteresis and drift, such as grating angle tuning in a diode laser. When a large dynamic range is needed to accommodate wide tuning range or to correct for extensive laser frequency drifts at low frequencies, a galvo-driven Brewster

plate can be used inside the optical cavity. Typically a Brewster plate inflicts an insertion loss less than 0.1 percent if its angular tuning range is limited within $\pm 4^\circ$. Walk-off of the optical beam by the tuning plate can be compensated with a double-passing arrangement or using dual plates. In the JILA-designed Ti:Sapphire laser, we use the combination of PZT and Brewster plate for the long-term frequency stabilization. The correction signal applied to the laser PZT is integrated and then fed to the Brewster plate to prevent saturation of the PZT channel. At higher frequencies we use much faster transducers, such as AOM and EOM, which are discussed below.

Temperature control of course offers the most universal means to control long-term drifts. Unfortunately the time constant associated with thermal diffusion is usually slow and therefore the loop bandwidth of thermal control is mostly limited to Hertz scale. However, thoughtful designs can sometimes push this limit to a much higher value. For example, a Kapton thin-film heater tape wrapped around the HeNe plasma tube has produced a thermal control unity gain bandwidth in excess of 100 Hz.⁴⁶ The transducer response is reasonably modeled as an integrator above 0.3 Hz and excessive phase shifts associated with the thermal diffusion does not become a serious issue until ~ 200 Hz. This transfer function of the transducer can be easily compensated with an electronic PI filter to produce the desired servo loop response. Radiant heating of a glass tube by incandescent lamps has achieved a time delay < 30 ms and has also been used successfully for frequency control of HeNe lasers.⁴⁷ If a bipolar thermal control is needed, Peltier-based solid state heat pumps (thermoelectric coolers) are available and can achieve temperature differences up to 70°C , or can transfer heat at a rate of 125 W, given a proper configuration of heat sinking. Parallel use of these Peltier devices results in a greater amount of heat transfer while a cascaded configuration achieves a larger temperature difference.

Combining various servo transducers in a single feedback loop requires thorough understanding of each actuator, their gains and phase shifts, and the overall loop filter function one intends to construct. Clearly, to have an attractive servo response in the time domain, the frequency transfer functions of various gain elements need to crossover each other smoothly. A slow actuator may have some resonance features in some low-frequency domain, hence the servo action needs to be relegated to a faster transducer at frequency ranges beyond those resonances. The roll-off of the slow transducer gain at high-frequencies needs to be steep enough, so that the overall loop gain can be raised without exciting the associated resonance. On the other hand, the high-frequency channel typically does not have as large a dynamic range as the slow ones. So one has to pay attention not to overload the fast channel. Again, a steep filter slope is needed to rapidly relinquish the gain of the fast channel toward the low-frequency range. However, we stress here that the phase difference between the two channels at the crossover point needs to be maintained at less than 90° . In the end, predetermined gains and phase shifts will be assigned to each transducer so that the combined filter function resembles a smooth single channel design. Some of these issues will be addressed briefly in the section below on example designs.

Servo Design in the Face of Time Delay: Additional Transducers Are Useful As one wishes for higher servo gain, with stability, it means a higher closed-loop bandwidth must be employed. Eventually the gain is sufficiently large that the intrinsic laser noise, divided by this gain, has become less than the measurement noise involved in obtaining the servo error signal. This should be sufficient gain. However, it may not be usable in a closed-loop scenario, due to excessive time delay. If we have a time delay of t_{delay} around the loop from an injection to the first receipt of correction information, a consideration of the input and response as vectors will make it clear that no real servo noise suppression can occur unless the phase of the response at least approximates that required to subtract from the injected error input to reduce its magnitude on the next cycle through the system. A radian of phase error would correspond to a unity-gain frequency of $1/(2\pi t_{\text{delay}})$, and we find this to be basically the upper useful limit of servo bandwidth. One finds that to correct a diode laser or dye laser to leave residual phase errors of 0.1 rad, it takes about 2 MHz servo bandwidth. This means a loop delay time, at the absolute maximum, of $t_{\text{delay}} = 1/2\pi \cdot 2\text{MHz} = 80$ ns. Since several amplifier stages will be in this rf and servo-domain control amplifier chain, the individual bandwidths need to be substantially beyond the 12 MHz naively implied by the delay spec. In particular rf modulation frequencies need to be unexpectedly large, 20 MHz at least, and octave rf bandwidths need to

be utilized, considering that the modulation content can only be 1/2 the bandwidth. Suppression of even-order signals before detection is done with narrow resonant rf notches.

Of course a PZT transducer will not be rated in the nanosecond regime of time delay. Rather, one can employ an AOM driven by a fast-acting Voltage-Controlled Oscillator to provide a frequency shift. Unfortunately the acoustic time delay from the ultrasonic transducer to the optical interaction seems always to be 400 ns, and more if we are dealing with a very intense laser beam and wish to avoid damage to the delicate AO transducer. The AOM approach works well with diode-pumped solid state lasers, where the bandwidth of major perturbations might be only 20 kHz. By double-passing the AOM the intrinsic angular deflection is suppressed. Usually the AOM prefers linear polarization. To aid separation of the return beam on the input side, a spatial offset can be provided with a collimating lens and roof prism, or with a cat's-eye retroreflector. Amplitude modulation or leveling can also be provided with the AOM's dependence on rf drive, but it is difficult to produce a beam still at the shot noise level after the AOM.

The final solution is an *EOM phase modulator*. In the external beam, this device will produce a phase shift per volt, rather than a frequency shift. So we will need to integrate the control input to generate a rate of change signal to provide to the EOM, in order to have a frequency relationship with the control input.⁴⁸ Evidently this will bring the dual problems of voltage saturation when the output becomes too large, and a related problem, the difficulty of combining fast low-delay response with high-voltage capability. The standard answer to this dilemma was indicated in our PZT section, namely, one applies fast signals and high-voltage signals independently, taking advantage of the fact that the needed control effort at high frequencies tends to cover only a small range. So fast low-voltage amplifier devices are completely adequate, particularly if one multipasses the EOM crystal several times. A full discussion of the crossover issues and driver circuits will be prepared for another publication.

Representative/Example Designs

Diode-Pumped Solid State Laser Diode-pumped solid state lasers are viewed as the most promising coherent light sources in diverse applications, such as communications, remote detection, and high precision spectroscopy. The diode-pumped Nd:YAG laser is probably the most highly developed of the rapidly expanding universe of diode-pumped solid state lasers, and it has enjoyed continuous improvements in its energy efficiency, size, lifetime, and intrinsic noise levels. The laser's free-running linewidth of ~ 10 kHz makes it a straightforward task to stabilize the laser via an optical cavity or an optical phase locked loop. In our initial attempt to stabilize the laser on a high finesse ($F \sim 100,000$, linewidth ~ 3 kHz) cavity, we employ an external AOM along with the laser internal PZT which is bonded directly on the laser crystal. The frequency discrimination signal between the laser and cavity is obtained with 4-MHz FM sidebands detected in cavity reflection. The PZT corrects any slow but potentially large laser frequency noise. Using the PZT alone allows the laser to be locked on the cavity. However, the loop tends to oscillate around 15 kHz and the residual noise level is more than 100 times higher than that obtained with the help of an external AOM. The AOM is able to extend the servo bandwidth to ~ 150 kHz, limited by the propagation time delay of the acoustic wave inside the AOM crystal. The crossover frequency between the PZT and AOM is about 10 kHz. Such a system has allowed us to achieve a residual frequency noise spectral density of 20 mHz/ $\sqrt{\text{Hz}}$. The laser's linewidth relative to cavity is thus a mere 1.3 mHz,⁴⁹ even though the noise spectral density is still 100 times higher than the shot noise. This same strategy of servo loop design has also been used to achieve a microradian level phase locking between two Nd:YAG lasers.⁵⁰

It is also attractive to stabilize the laser directly on atomic/molecular transitions, given the low magnitude of the laser's intrinsic frequency noise. Of course the limited S/N of the recovered resonance information will not allow us to build speedy loops to clean off the laser's fast frequency/phase noise. Rather we will use the laser PZT alone to guide the laser for a long-term stability. An example here is the 1.064 μm radiation from the Nd:YAG, which is easily frequency doubled to 532 nm where strong absorption features of iodine molecules exist.^{51,52} The doubling is furnished with a noncritical phase-matched KNb_3O_7 crystal located inside a buildup cavity. 160 mW of green light

output is obtained from an input power of 250 mW of IR. Only mW levels of the green light are needed to probe the iodine saturated absorption signal. Low vapor pressure (~ 0.5 Pa) of the iodine cell is used to minimize the collision-induced pressure shift and to reduce the influence on baseline by the linear Doppler absorption background. The signal size decreases as the pressure is reduced. However, this effect is partly offset due to the reduced resonance linewidth (less pressure broadening) which helps to increase the slope of the frequency locking error signal. A lower pressure also helps to reduce power-related center frequency shifts since a lower power is needed for saturation. With our 1.2-m long cells, we have achieved an S/N of 120 in a 10-kHz bandwidth, using the modulation transfer spectroscopy.⁵³ (Modulation transfer is similar to FM except that we impose the frequency sideband on the saturating beam and rely on the nonlinear medium to transfer the modulation information to the probe beam which is then detected.) Normalized to 1-s averaging time, this S/N translates to the possibility of a residual frequency noise level of 10 Hz when the laser is locked on the molecular resonance, given the transition linewidth of 300 kHz. We have built two such iodine-stabilized systems and the heterodyne beat between the two lasers permits systematic studies on each system and checks the reproducibility of the locking scheme.⁵⁴ With a 1 second counter gate time, we have recorded the beat frequency between the two lasers. The standard deviation of the beat frequency noise is ~ 20 Hz, corresponding to ~ 14 -Hz rms noise per IR laser, basically a S/N limited performance. The beat record can be used to calculate the Allan standard deviation: starting at 5×10^{-14} at 1 second, decreasing with a slope of $1/\sqrt{\tau}$ up to 100 second. (τ is the averaging time.) After 100-s the deviations reach the flicker noise floor of $\sim 5 \times 10^{-15}$. At present, the accuracy of the system is limited by inadequate optical isolation in the spectrometer and the imperfect frequency modulation process (residual amplitude noise, RAM) used to recover the signal. This subject is under intense active study in our group.⁵⁵

External Cavity Diode Lasers Diode lasers are compact, reliable, and coherent light sources for many different applications.⁵⁶ The linewidth of a free-running diode laser is limited by the fundamental spontaneous emission events, enhanced by the amplitude-phase coupling inside the gain medium. With a low-noise current driver, a typical milliwatt scale AlGaAs diode laser has a linewidth of several MHz. To reduce this fast frequency noise, one typically employs an external cavity formed between one of the diode laser facets and a grating (or an external mirror that retroreflects the first-order grating diffraction).^{57–59} This optical feedback mechanism suppresses the spontaneous emission noise, replaced by much slower fluctuations of mechanical origin. The linewidth of the grating-stabilized external cavity diode laser (ECDL) is usually between 100 kHz and 1 MHz, determined by the quality factor of the optical feedback. The ECDL also offers much better tuning characteristics compared against a solitary diode. To do such tuning, the external grating (or the mirror that feeds the grating-dispersed light back to the laser) is controlled by a PZT for scanning. Synchronous tuning of the grating dispersion and the external cavity mode can be achieved with a careful selection of the grating rotation axis position. Similarly, this PZT-controlled grating can be used to stabilize the frequency of an ECDL. However, owing to the low bandwidth limited by the mechanical resonance of PZT, a tight frequency servo is possible only through fast transducers such as the laser current or intracavity phase modulators.

This hybrid electro-optic feedback system is attractive, and ECDLs have been demonstrated to show hertz level stability under a servo bandwidth of the order of 1 MHz. For a solitary diode, feedback bandwidth of tens of megahertz would have been needed in order to bring the frequency noise down to the same level. However, considering that the optical feedback has a strong impact on the laser frequency noise spectrum, one finds the frequency response of the compound laser system is clearly dependent upon the optical alignment. Therefore, for each particular ECDL system, we need to measure the frequency response function of the laser under the optimally aligned condition. We are dealing with a multichannel feedback system (e.g., PZT plus current), so that designing smooth crossovers between different transducers requires knowledge of the transfer functions of each transducer. Normally the current-induced FM of a solitary diode has a flat response up to ~ 100 kHz, and then starts to roll off in the region between 100 kHz and 1 MHz, initially with a single-pole character. This is due to the time response of the current-induced thermal change of the refractive index inside the diode. (At a faster time scale, the carrier density variation will remain and then dominate

the laser frequency response.) Design of a fast feedback loop needs to take into account this intrinsic diode response. Fortunately the time delay associated with the current response is low, typically below 10 ns.

In our example system, the frequency discrimination signal of the ECDL is obtained from a 100 kHz linewidth cavity with a sampling frequency of 25 MHz. The error signal is divided into three paths: PZT, current modulation through the driver, and direct current feedback to the diode head. The composite loop filter function is shown in Fig. 9. The crossover between the slow current channel and the PZT usually occurs around 1 kHz, in order to avoid the mechanical resonance of the PZT at a few kHz. In our system, the frequency response of the PZT/grating is 10 GHz/V. To furnish this in-loop gain of ~ 1000 at 1 kHz, we need to supply an electronic gain of 0.1, given that the error signal has a slope of 1 V/1 MHz.

Toward the lower frequency range, the PZT gain increases by 40 dB/decade (double integrators) to suppress the catastrophically rising laser frequency noise. It is obvious from Fig. 9 that the intermediate current channel tends to become unstable at a few hundred kHz, due to the excessive phase shift there. The fast current loop, bypassing the current driver to minimize additional time delay and phase shift, has a phase lead compensator to push the unity gain bandwidth to 2 MHz. With this system we can lock the ECDL robustly on the optical cavity, with a residual noise spectral density of 2 Hz/ $\sqrt{\text{Hz}}$, leading to a relative linewidth of 12 Hz. The achieved noise level is about 100 times higher than the fundamental measurement limit set by shot noise. We note in passing that when an ECDL gradually goes out of alignment, the previously adjusted gain of the current loop will tend to make the servo oscillate so we know a new alignment is needed. The laser FM sideband used to generate

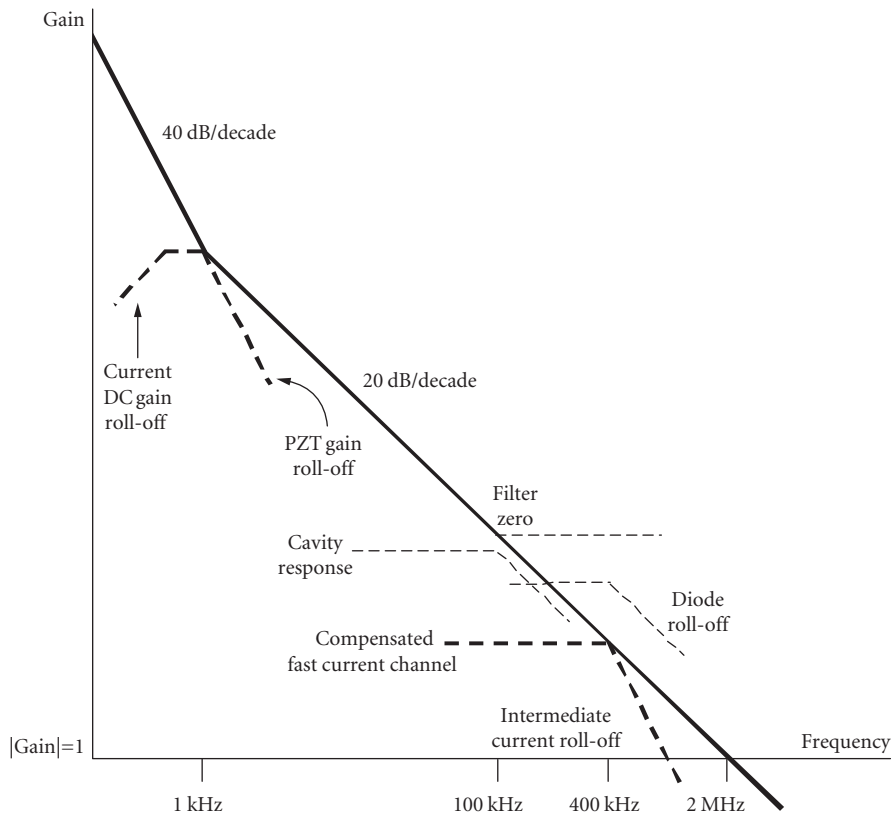


FIGURE 9 The combined loop filter function for ECDL frequency stabilization.

the locking signal is produced directly by current modulation. An electronic filter network is employed to superimpose the slow servo, fast servo, and modulation inputs to the diode. Exercise caution when accessing the diode head, as a few extra mA current increase can lead to drastic output power increase and melted laser facets, all in 1 μ s!

Since 2000 this art of locking diode lasers to cavities has been progressing rapidly and in many labs worldwide. To note just one recent measurement of subhertz linewidths, we may note the results⁶⁰ in JILA, where two diode lasers were locked to two independent vertical cavities, giving 10^{-15} frequency stability at 1-second averaging time. In another Boulder collaboration, subhertz optical beat linewidths have been measured⁶¹ between a JILA diode laser source at 689 nm and another stabilized laser used as the local oscillator for the Hg⁺ clock, at a different wavelength of 1126 nm, and located 4 km away in the NIST labs. Two optical frequency Combs were employed, along with a Nd:YAG laser whose wavelength could be transmitted by the 4-km fiber link. (This Comb technology seems to have gone quickly from being a research topic to a reliable and versatile tool!)

22.4 SUMMARY AND OUTLOOK

The technology of laser frequency stabilization has been refined and simplified over the years and has become an indispensable research tool in any modern laboratories involving optics. Research on laser stabilization has been and still is pushing the limits of measurement science. Indeed, a number of currently active research projects on fundamental physical principles benefit a great deal from stable optical sources and will need a continued progress of the laser stabilization front.⁶² Using extremely stable phase coherent optical sources, we will be entering an exciting era when the LISA interferometer will achieve picometer resolution over a five million kilometer distance in space⁶³ and a few Hertz linewidth of an ultraviolet resonance will be probed with a high S/N .^{64, 65} One has to be optimistic looking at the stabilization results of all different kinds of lasers. To list just a few examples of cw tunable lasers, we notice milliHertz linewidth stabilization (relative to a cavity) for diode-pumped solid state lasers; dozen milliHertz linewidth for Ti:Sapphire lasers; and sub-Hertz linewidths for diode and dye lasers. Long-term stability of lasers referenced to atoms and molecules⁶⁶ has reached mid 10^{-16} levels in an averaging time as short as ~ 300 s. Phase locking between different laser systems is now routinely employed, even for diode lasers that have fast frequency noise.

An important new capability comes from the accurate frequency transfer via noise-compensated optical fiber, which allows groups to collaborate in the testing of their state-of-the-art systems, thus tapping into the fruits of our neighbors' achievements in stabilizing lasers.⁶⁷ The Optical Comb makes it possible to use nearly any stable frequency as the reference. The fiber transport allowed measurement of the inaccuracy issues affecting four potential atomic clocks, based on a trapped Hg⁺ ion,⁶⁸ Al⁺ trapped ions,⁶⁸ on cold and free Ca atoms,⁶⁹ and on cold Sr atoms confined in an optical lattice.⁶⁶ Measurements confirmed that the Sr, Hg⁺, and Al⁺ systems can be expected to yield sub- 1×10^{-16} independent reproducibility (the same quality as *accuracy*, except for the authority of Cs definition as the unit of time.) Thus for the first time there is not one, but even three(!) frequency reference systems which surpass the performance of the present best standard in Physics, namely the Cs Fountain Clock. Likely there will be other optically-based systems (such as Yb in a lattice⁷⁰) which may soon be in this performance category. Clearly there will be great laser and physics fun coming in the next times! For example, an extended series of comparisons between Hg⁺ and Al⁺ was carried out at NIST under strict conditions, and provided a $< 1 \times 10^{-16}$ /year checking for possible drift in the fundamental "constants" of physics.⁶⁸ By the way, do *you* believe that the frequencies of nuclear Mössbauer transitions will stay at the same frequency, as measured by these atomic frequency sources? With the rapid progress in producing VUV Combs by Higher Harmonic Generation,^{71, 72} one can begin to dream about coherent Mössbauer spectroscopy so we can find out. In the optical comparisons, we're looking in the seventeenth decimal digit now!

Another highway toward fundamental physics involves precision laser spectroscopy of simple atoms, such as H and He⁺, but laser cooling is a challenge due to the low available power. The exciting program⁷³ at MPQ in Munich can now generate enough VUV 121.5-nm light to begin this story!

Quantum noise is the usual limit of the measurement process and therefore may be the limit of the stabilization process as well. To circumvent the quantum noise altogether is an active research field itself.⁷⁴ We, however, have not reached this quantum limit just yet. For instance, we have already stated that the Nd:YAG laser should be able to reach microHertz stability if the shot noise is the true limit. What have we done wrong? A part of the deficiency is due to the inadequacy of the measurement process, namely the lack of accuracy. This is because the signal recovery effort—modulation and demodulation process—is contaminated by spurious optical interference effects and RAM associated with the modulation frequency. Every optical surface along the beam path can be a potential time bomb to damage the modulation performance. In cases that some low contrast interference effects are not totally avoidable, we would need to have the whole system controlled in terms of the surrounding pressure and temperature. The degree to which we can exert control of course dictates the ultimate performance. The second fundamental limitation appears to be thermally generated mechanical noise of the optical coatings, which was already noted,^{41,42} but there are already schemes being discussed to buy us another decade or two.

22.5 CONCLUSIONS AND RECOMMENDATIONS

It becomes clear that there are many interlinking considerations involved in the design of laser stabilization systems, and it is difficult to present a full description in an chapter such as this. Still it is hoped that the reader will see some avenues to employ feedback control methods to the laser systems of her current interest. We are optimistic that some of this technology may become commercially available in the future, thus simplifying the user's task.

22.6 ACKNOWLEDGMENTS

The work discussed here has profited from interactions with many colleagues, postdoctoral researchers, and graduate students over many years. In particular we must thank Leo Hollberg and Miao Zhu for their earlier contributions. More recent contributors are Long Sheng Ma, and Mark Notcutt. Especially we thank Mark for his experimental work with the vertically mounted cavities, and for his useful suggestions for improving this text. The now-retired one of us (JLH) declares his joy to see the JILA laboratory prospering in the strong hands of Jun Ye. Jan also is particularly grateful to his wife Lindy for patience “beyond the call of duty” during these many years of laser research. The work at JILA has been supported over many years by the Office of Naval Research, the National Science Foundation, the Air Force Office of Scientific Research, NASA, and the National Institute of Standards and Technology, as part of its frontier research into basic standards and their applications.

22.7 REFERENCES

1. D. W. Allan, “Statistics of Atomic Frequency Standards,” *Proc. IEEE* **54**:221–230 (1966).
2. D. B. Sullivan, D. W. Allan, D. A. Howe, and F. L. Walls, (eds.), *Characterization of Clocks and Oscillators*, NIST Technical Note **1337**, U. S. Government Printing Office, Washington D.C., 1990. See also <http://tf.nist.gov/timefreq/general/generalpubs.htm> (2009).
3. D. S. Elliot, R. Roy and S. J. Smith, “Extracavity Laser Band-Shape and Bandwidth Modification,” *Phys. Rev. A* **26**:12–18 (1982).
4. J. L. Hall and M. Zhu, “An Introduction to Phase-Stable Optical Sources,” in Proc. of the *International School of Phys. “Enrico Fermi,” Course CXVIII, Laser Manipulation of Atoms and Ions*, E. Arimondo, W. D. Phillips, and F. Strumia, (eds.), North Holland, pp 671–702 (1992).
5. M. Zhu and J. L. Hall, “Stabilization of Optical-Phase Frequency of a Laser System—Application to a Commercial Dye-Laser with an External Stabilizer,” *J. Opt. Soc. Am. B* **10**: 802–816 (1993).

6. For similar issues in the microwave/rf field, see the application note *Time Keeping and Frequency Calibration*, Agilent, Palo Alto, Calif., and <http://tf.nist.gov/timefreq/general/generalpubs.htm>.
7. For general references on feedback systems, please refer to *Modern Control Systems*, 3rd ed., Richard C. Dorf, Addison-Wesley Publishing Co., Reading Mass. (1980).
8. *Feedback Control of Dynamic Systems*, 3rd ed., Gene F. Franklin, J. David Powell, and A. Emami-Naeini, Addison-Wesley Publishing Co., Reading, Mass. (1994).
9. The simulations presented here are performed with the following software by The Math Works Inc., Natick, Mass. (1996): Matlab Control Systems Toolbox: User's Guide; Matlab Simulink: User's Guide.
10. See, e.g., J. L. Hall, "Frequency Stabilized Lasers—A Parochial Review," in *Frequency-Stabilized Lasers and Their Applications*, Y. C. Chung, (ed.), *Proc. SPIE* **1837**:2–15 (1993).
11. R. L. Barger, M. S. Sorem, and J. L. Hall, "Frequency Stabilization of a Cw Dye Laser," *Appl. Phys. Lett.* **22**:573–575 (1973).
12. G. D. Houser and E. Garmire, "Balanced Detection Technique to Measure Small Changes in Transmission," *Appl. Opt.* **33**:1059–1062 (1994).
13. K. L. Haller and P. C. D. Hobbs, "Double Beam Laser Absorption Spectroscopy: Shot-Noise Limited Performance at Baseband with a Novel Electronic Noise Canceller," in *Optical Methods for Ultrasensitive Detection and Analysis: Techniques and Applications*, B. L. Fearey, (ed.), *Proc. SPIE* **1435**: 298–309 (1991). Also see <http://www.electrooptical.net> (2009).
14. J. Helmcke, S. A. Lee, and J. L. Hall, "Dye-Laser Spectrometer for Ultrahigh Spectral Resolution—Design and Performance," *Appl. Opt.* **21**:1686–1694 (1982).
15. C. E. Wieman and T. W. Hänsch, "Doppler-Free Laser Polarization Spectroscopy," *Phys. Rev. Lett.* **36**:1170–1173 (1976).
16. T. W. Hänsch and B. Couillaud, "Laser Frequency Stabilization by Polarization Spectroscopy of a Preflecting Reference Cavity," *Opt. Comm.* **35**:441–444 (1980).
17. H. Wahlquist, "Modulation Broadening of Unsaturated Lorentzian Lines," *J. Chem. Phys.* **35**:1708–1710 (1961).
18. T. J. Quinn, "Mise en Pratique of the Definition of the Metre (1992)" *Metrologia* **30**:524–541 (1994).
19. J. Hu, E. Ikonen, and K. Riski, "On the Nth Harmonic Locking of the Iodine Stabilized He-Ne-Laser," *Opt. Comm.* **120**:65–70 (1995); **121**:169 (1995).
20. M. S. Taubman and J. L. Hall, "Cancellation of laser dither modulation from optical frequency standards," *Opt. Lett.* **25**:311–313 (2000).
21. B. Smaller, *Phys. Rev.* **83**:812–820 (1951); R. V. Pound, "Electronic Stabilization of Microwave Oscillators," *Rev. Sci. Instrum.* **17**: 490–505 (1946).
22. G. C. Bjorklund, "Frequency-Modulation Spectroscopy: A New Method for Measuring Weak Absorptions and Dispersions," *Opt. Lett.* **5**:15–17 (1980).
23. J. L. Hall, L. Hollberg, T. Baer, and H. G. Robinson, "Optical Heterodyne Saturation Spectroscopy," *Appl. Phys. Lett.* **39**:680–682 (1981).
24. G. C. Bjorklund and M. D. Levenson, "Sub-Doppler Frequency-Modulation Spectroscopy of I₂," *Phys. Rev. A* **24**:166–169 (1981).
25. W. Zapka, M. D. Levenson, F. M. Schellenberg, A. C. Tam, and G. C. Bjorklund, "Continuous-Wave Doppler-Free Two-Photon Frequency Modulation Spectroscopy in Rb Vapor," *Opt. Lett.* **8**:27–29 (1983).
26. G. J. Rosasco and W. S. Hurst, "Phase-Modulated Stimulated Raman Spectroscopy," *J. Opt. Soc. Am. B* **2**:1485–1496 (1985).
27. J. H. Shirley, "Modulation Transfer Processes in Optical Heterodyne Saturation Spectroscopy," *Opt. Lett.* **7**:537–539 (1982).
28. R. W. P. Drever, J. L. Hall, F. V. Kowalski, J. Hough, G. M. Ford, A. J. Munley, and H. Ward, "Laser Phase and Frequency Stabilization Using an Optical-Resonator," *Appl. Phys. B* **31**:97–105 (1983).
29. P. Werle, "Laser Excess Noise and Interferometric Effects in Frequency-Modulated Diode-Laser Spectrometers," *Appl. Phys. B* **60**:499–506 (1995).
30. J. L. Hall, J. Ye, L.-S. Ma, K. Vogel, and T. Dinneen, "Optical Frequency Standards: Progress and Applications," in *Laser Spectroscopy XIII*, Z.-J. Wang, Z.-M. Zhang, and Y.-Z. Wang, (eds.), World Scientific, Singapore, pp. 75–80 (1998).

31. N. C. Wong and J. L. Hall, "Servo Control of Amplitude-Modulation in Frequency-Modulation Spectroscopy—Demonstration of Shot-Noise-Limited Detection," *J. Opt. Soc. Am. B* **2**:1527–1533 (1985).
32. M. S. Taubman and J. L. Hall, unpublished JILA work (1999).
33. D. R. Hjelm, S. Neegård, and E. Vartdal, "Optical Interference Fringe Reduction in Frequency-Modulation Spectroscopy Experiments," *Opt. Lett.* **20**:1731–1733 (1995).
34. G. R. Janik, C. B. Carlisle, and T. F. Gallagher, "Two-Tone Frequency-Modulation Spectroscopy," *J. Opt. Soc. Am. B* **3**:1070–1074 (1986).
35. J. L. Hall, L.-S. Ma, and G. Kramer, "Principles of Optical Phase-Locking—Application to Internal Mirror He-Ne Lasers Phase-Locked Via Fast Control of the Discharge Current," *IEEE J. Quan. Electron.* **QE-23**:427–437 (1987).
36. M. Kourogi, K. Nakagawa, and M. Ohtsu, "Wide-Span Optical Frequency Comb Generator for Accurate Optical Frequency Difference Measurement," *IEEE J. Quan. Electron.* **QE-29**:2693–2701 (1993).
37. Th. Udem, J. Reichert, R. Holzwarth, and T. W. Hänsch, "Accurate Measurement of Large Optical Frequency Differences with a Mode-Locked Laser," *Opt. Lett.* **24**:881–883 (1999).
38. S. A. Diddams, L.-S. Ma, J. Ye, and J. L. Hall, "Broadband Optical Frequency Comb Generation with a Phase-Modulated Parametric Oscillator," *Opt. Lett.* **24**:1747–1749 (1999).
39. J. L. Hall, M. Notcutt, and J. Ye, "Improving Laser Coherence," in *Laser Spectroscopy XVII*, E. Hinds, A. Ferguson, and E. Riis, (eds.), World Scientific, Singapore, pp. 3–12 (2006).
40. B. C. Young, F. C. Cruz, W. M. Itano, and J. C. Bergquist, "Visible Lasers with Subhertz Linewidths," *Phys. Rev. Lett.* **82**:3799–3802 (1999).
41. K. Numata, A. Kemery, and J. Camp, "Thermal-Noise Limit in the Frequency Stabilization of Lasers with Rigid Cavities," *Phys. Rev. Lett.* **93**:250602 (2004).
42. M. Notcutt, L. S. Ma, A. D. Ludlow, S. M. Foreman, J. Ye, J. L. Hall, "Contribution of Thermal Noise to Frequency Stability of Rigid Optical Cavity via Hertz-Linewidth Lasers," *Phys. Rev. A* **73**:031804 (R) (2006).
43. Notcutt, M. L. S. Ma, J. Ye, J. L. Hall, "Simple and Compact 1-Hz Laser System via an Improved Mounting Configuration of a Reference Cavity," *Opt. Lett.* **30**:1815–1817 (2005).
44. JILA design by Lisheng Chen, 2004. See L. Chen, J. L. Hall, J. Ye, T. Yang, E. Zang, and T. Li, "Vibration-Induced Elastic Deformation of Fabry-Perot Cavities," *Phys. Rev. A* **74**:053801 (2006).
45. Interested readers are referred to the following books for more detailed discussions. *The Quantum Physics of Atomic Frequency Standards*, vol. I and II, Jacques Vanier and Claude Audoin; Adam Hilger, Institute of Physics Publishing Ltd, Bristol, UK (1989).
46. T. M. Niebauer, J. E. Faller, H. M. Godwin, J. L. Hall, and R. L. Barger, "Frequency Stability Measurements on Polarization-Stabilized He-Ne Lasers," *Appl. Opt.* **27**:1285–1289 (1988).
47. Jun Ishikawa, NMIJ, AIST, Tsukuba, Japan, private communications (1996).
48. J. L. Hall and T. W. Hänsch, "External Dye-Laser Frequency Stabilizer," *Opt. Lett.* **9**:502–504 (1984).
49. J. Ye, L.-S. Ma, and J. L. Hall, "Ultrasensitive Detections in Atomic and Molecular Physics: Demonstration in Molecular Overtone Spectroscopy," *J. Opt. Soc. Am. B* **15**:6–15 (1998).
50. J. Ye and J. L. Hall, "Optical Phase Locking in the Microradian Domain: Potential Applications to NASA Spaceborne Optical Measurements," *Opt. Lett.* **24**:1838–1840 (1999).
51. A. Arie and R. L. Byer, "Laser Heterodyne Spectroscopy of $^{127}\text{I}_2$ Hyperfine Structure near 532 nm," *J. Opt. Soc. Am. B* **10**:1990–1997 (1993).
52. M. L. Eickhoff and J. L. Hall, "Optical Frequency Standard at 532 Nm," *IEEE Trans. Instrum. Meas.* **44**:155–158 (1995); and P. Jungner, M. Eickhoff, S. Swartz, J. Ye, J. L. Hall and S. Waltman, *IEEE Trans. Instrum. Meas.* **44**:151–154 (1995).
53. J. Ye, L. Robertsson, S. Picard, L.-S. Ma, and John L. Hall, "Absolute Frequency Atlas of Molecular I-2 Lines at 532 nm," *IEEE Trans. Instrum. Meas.* **48**:544–549 (1999).
54. J. L. Hall, L.-S. Ma, M. Taubman, B. Tiemann, F.-L. Hong, O. Pfister, and J. Ye, "Stabilization and Frequency Measurement of the I-2-Stabilized Nd:YAG laser," *IEEE Trans. Instrum. Meas.* **48**:583–586 (1999).
55. C. Ishibashi, J. Ye, and J. L. Hall, "Issues and Applications in Ultra-Sensitive Molecular Spectroscopy," *Proc. SPIE.* **4634**:58–69 (2002).
56. C. E. Wieman and L. Hollberg, "Using Diode-Lasers for Atomic Physics," *Rev. Sci. Instrum.* **62**:1–20 (1991).

57. M. G. Littman and H. J. Metcalf, "Spectrally Narrow Pulsed Dye Laser without Beam Expander," *Appl. Opt.* **17**:2224–2227 (1978).
58. B. Dahmani, L. Hollberg, and R. Drullinger, "Frequency Stabilization of Semiconductor-Lasers by Resonant Optical Feedback," *Opt. Lett.* **12**:876–878 (1987).
59. K. MacAdam, A. Steinbach, and C. E. Wieman, "A Narrow-Band Tunable Diode-Laser System with Grating Feedback, and a Saturated Absorption Spectrometer for Cs and Rb," *Am. J. Phys.* **60**:1098–1111 (1992).
60. A. D. Ludlow, X. Huang, M. Notcutt, T. Zanon-Willette, S. M. Foreman, M. M. Boyd, S. Blatt, and J. Ye, "Compact, Thermal-Noise-Limited Optical Cavity for Diode Laser Stabilization at 1×10^{-15} ," *Opt. Lett.* **32**:641–643 (2007).
61. S. M. Foreman, A. D. Ludlow, M. H. G. de Miranda, J. Stalnaker, S. A. Diddams, and J. Ye, "Coherent Optical Phase Transfer over a 32-km Fiber with 1 s Instability at 10^{-17} ," *Phys. Rev. Lett.* **99**:153601 (2007).
62. P. Fritschel, G. González, B. Lantz, P. Saha, and M. Zucker, "High Power Interferometric Phase Measurement Limited by Quantum Noise and Application to Detection of Gravitational Waves," *Phys. Rev. Lett.* **80**:3181–3185 (1998).
63. K. Danzmann, "LISA—An ESA Cornerstone Mission for a Gravitational Wave Observatory," *Class. Quantum Grav.* **14**:1399–1404 (1997).
64. J. C. Bergquist, R. J. Rafac, B. C. Young, J. A. Beall, W. M. Itano, and D. J. Wineland, "Sub-Dekahertz Spectroscopy of $^{199}\text{Hg}^+$," in *Laser Frequency Stabilization, Standards, Measurement, and Applications*, J. L. Hall and J. Ye, (eds.), *Proc. SPIE.* **4269**:1–7 (2001).
65. E. E. Eyler, D. E. Chieda, M.C. Stowe, M. J. Thorpe, T. R. Schibli, and J. Ye, "Prospects for Precision Measurements of Atomic Helium using Direct Frequency Comb Spectroscopy," *Eur. Phys. J. D.* **48**:43–55 (2008).
66. D. Ludlow, T. Zelevinsky, G. K. Campbell, S. Blatt, M. M. Boyd, M. H. G. de Miranda, M. J. Martin, J. W. Thomsen, S. M. Foreman, Jun Ye, T. M. Fortier, J. E. Stalnaker, S. A. Diddams, Y. Le Coq, Z. W. Barber, N. Poli, N. D. Lemke, K. M. Beck, and C. W. Oates, "Sr Lattice Clock at 1×10^{-16} Fractional Uncertainty by Remote Optical Evaluation with a Ca clock," *Science* **319**:1805–1808 (2008).
67. S. M. Foreman, K. W. Holman, D. D. Hudson, D. J. Jones, and J. Ye, "Remote Transfer of Ultrastable Frequency References via Fiber Networks," *Rev. Sci. Instrum.* **78**:021101 (2007).
68. T. Rosenband, D. B. Hume, P. O. Schmidt, C. W. Chou, A. Brusch, L. Lorini, W. H. Oskay, R. E. Drullinger, T. M. Fortier, J. E. Stalnaker, S. A. Diddams, W. C. Swann, N. R. Newbury, W. M. Itano, D. J. Wineland, and J. C. Bergquist, "Frequency Ratio of Al^+ and Hg^+ Single-Ion Optical Clocks; Metrology at the 17th Decimal Place," *Science* **319**:1808–1812 (2008).
69. J. E. Stalnaker, Y. Le Coq, T. M. Fortier, S. A. Diddams, C. W. Oates, and L. Hollberg, "Measurement of Excited-State Transitions in Cold Calcium Atoms by Direct Femtosecond Frequency-Comb Spectroscopy," *Phys. Rev. A* **75**:040502 (R) (2007).
70. N. Poli, Z. W. Barber, N. D. Lemke, C. W. Oates, L. S. Ma, J. E. Stalnaker, T. M. Fortier, S. A. Diddams, L. Hollberg, J. C. Bergquist, A. Brusch, S. Jefferts, T. Heavner, and T. Parker, "Frequency Evaluation of the Doubly Forbidden $^1\text{S}_0 \rightarrow ^3\text{P}_0$ Transition in Bosonic Yb-174," *Phys. Rev. A* **77**:050501 (R) (2008).
71. R. J. Jones, K. D. Moll, M. J. Thorpe, and J. Ye, "Phase-Coherent Frequency Combs in the Vacuum Ultraviolet via High-Harmonic Generation inside a Femtosecond Enhancement Cavity," *Phys. Rev. Lett.* **94**:193201 (2005).
72. C. Göhle, T. Udem, M. Herrmann, J. Rauschenberger, R. Holzwarth, H. A. Schuessler, F. Krausz, and T. W. Hänsch, "A Frequency Comb in the Extreme Ultraviolet," *Nature* **436**:234–237 (2005).
73. See <http://www.mpg.mpg.de/~haensch/antihydrogen/index.html> (2009).
74. V. B. Braginsky and F. Ya. Khalili, *Quantum Measurement*, Cambridge University Press, Cambridge, UK (1992).

This page intentionally left blank.

QUANTUM THEORY OF THE LASER

János A. Bergou^{a,b}
 Berthold-Georg Englert^{a,c,d}
 Melvin Lax^{e,*}
 Marian O. Scully^{a,c}
 Herbert Walther^{c,f,*}
 M. Suhail Zubairy^{a,g}

^a*Institute for Quantum Studies and Department of Physics
 Texas A&M University
 College Station, Texas*

^b*Department of Physics and Astronomy
 Hunter College of the City University of New York
 New York, New York*

^c*Max-Planck-Institut für Quantenoptik
 Garching bei München, Germany*

^d*Abteilung Quantenphysik der Universität Ulm
 Ulm, Germany*

^e*Department of Physics
 City College of the City University of New York
 New York, New York*

^f*Sektion Physik der Universität München
 Garching bei München, Germany*

^g*Department of Electronics
 Quaid-i-Azam University
 Islamabad, Pakistan*

23.1 GLOSSARY

Section 23.3

a_k, a_k^\dagger annihilation, creation operator for photons in the k th mode
 A Einstein coefficient for spontaneous emission
 $\mathbf{A}_k(\mathbf{r})$ k th electric mode function at position \mathbf{r}

*Deceased.

B	Einstein coefficient for absorption and stimulated emission
$\mathbf{B}(\mathbf{r}, t)$	magnetic field at position \mathbf{r} and time t
$\mathbf{B}_k(\mathbf{r})$	k th magnetic mode function at position \mathbf{r}
c	speed of light
\mathbf{e}_k	polarization unit vector of the k th mode
$\mathbf{E}_\perp(\mathbf{r}, t)$	transverse electric field at position \mathbf{r} and time t
h, \hbar	Planck's constant [$\hbar = h/(2\pi)$]
\mathbf{k}, k	wave vector, its length
$(d\mathbf{k})$	three-dimensional volume element in \mathbf{k} space
k_B	Boltzmann's constant
\mathbf{n}_k	propagation unit vector of the k th mode
N	photon number operator
N_e	number of excited-state atoms
N_g	number of ground-state atoms
\mathbf{r}	position vector
$(d\mathbf{r})$	three-dimensional volume element in \mathbf{r} space
t, dt	time, time interval
T	temperature
U_{1ph}	energy density for a one-photon state
dV	volume element
$ \text{vac}\rangle, \langle\text{vac} $	ket and bra of the photon vacuum
$\overline{w^2}$	mean square energy fluctuations
W	spectral-spatial energy density of blackbody radiation
α_k	coherent state amplitudes
$ \{\alpha\}_c\rangle$	ket of a coherent state
δ_{jk}	Kronecker's delta symbol
$\delta_\perp(\mathbf{r})$	transverse delta function at \mathbf{r} (a dyadic)
ϵ_0	dielectric constant [$\epsilon_0 = 1/(\mu_0 c^2) \approx 8.854 \times 10^{-12}$ F/m]
μ_0	permeability constant ($4\pi \times 10^{-7}$ H/m)
$\nu, d\nu$	light frequency, frequency interval
ν_k	frequency of the k th mode
ρ_{ph}	statistical operator of the photon field
ρ	probability amplitude for reflection
τ	probability amplitude for transmission
Ψ_k, Ψ_{jk}	probability amplitudes of the one-photon, two-photon states
$ \{\psi\}_1\rangle, \{\psi\}_2\rangle$	kets for one-photon, two-photon states
∇	gradient vector differential operator

Section 23.4

A	destruction operator for the field
\mathcal{A}	laser gain coefficient
a, a^\dagger	photon ladder operators
\mathcal{B}	laser saturation parameter
b_k, b_k^\dagger	destruction and creation operators for the reservoir modes
\mathcal{D}	largest eigenvalue of the laser equation

$F(t)$	field noise operator
$F_\alpha(t), F_\gamma$	noise operators associated with gain and loss, respectively
$F(1, x, y)$	hypergeometric function
g	radiant frequency stating the strength of atom-photon coupling
g_k	coupling coefficient between reservoir and field
$G^{(1)}(t_0 + t, t_0)$	field correlation function
$\mathcal{H}, \mathcal{H}_0, \mathcal{H}_1$	total, free, and interaction hamiltonians for atom-field interaction
k	wave vector for the field
K	kick operator
\mathcal{L}	superoperator for cavity damping
$M(\tau)$	superoperator describing the effect of a single inverted atom on the field
$\overline{n}, \overline{n^2}$	mean and mean squared number of photons in a laser
n_m	maximum of the photon distribution of a laser
n_{th}	mean number of photons in a thermal reservoir
N_{ex}	number of atoms traversing the cavity during the lifetime of the cavity field
$N(t_i, t, \tau)$	notch function
p	transition dipole moment
$p(n)$	photon distribution function
$P(\tau)$	distribution function for the interaction times
q	a nonnegative integer
Q	Mandel Q function
r	atom injection rate inside a laser and micromaser
$S(\omega)$	spectrum of the laser field
T	temperature
t_m	measurement time
U	time evolution operator
$U_{\text{af}}(\tau)$	atom-field time evolution operator
\mathcal{V}	interaction picture Hamiltonian
$\alpha(t, t')$	gain function of a laser
Γ	atomic decay rate via spontaneous emission
$\gamma, \gamma_a, \gamma_b$	atomic decay rates
$\delta_{nm'}$	Kronecker delta function
ν_k	frequency of the reservoir mode
σ_+, σ_-	ladder operators for a two-level atom
σ_z	atomic inversion operator
ω_0	radiant frequency of an atomic transition
λ	eigenvalue
λ_j	eigenvalues of the laser equation
$\rho(t)$	reduced density operator for the field
ρ_{at}	total density operator for the atom-field system
$\rho_{mm'}$	matrix elements of the field density operator
$\rho_n^{(k)}$	off-diagonal density matrix element
κ	cavity loss rate
χ	square of the ratio of vacuum Rabi frequency and the atomic decay rate
$\theta(t)$	step function
$\phi(t)$	phase of the field

Section 23.5

a, b, c, d	parameters appearing in Table 1
$\frac{E}{E}$	electric field strength
$\frac{E}{E}$	mean field
$F(M), F_0$	free energy of a ferromagnet, its value for $M = 0$
g_k	undefined in Eqs. (126) and (127)
$G(E), G_0$	free energy of a laser, its value for $E = 0$
H	external magnetic field
H_{n0}	heating rate
K	one-fourth the spontaneous emission rate
K_{n0}	cooling rate
n_0	number of atoms in the Bose-Einstein condensate
$\langle n_0 \rangle$	mean number of atoms in the condensate
$\langle \dot{n}_0 \rangle$	time derivative of $\langle n_0 \rangle$
$\langle n_k \rangle_{n_0}$	average number of atoms in the k th excited state, given n_0 atoms in the condensate
N	total number of Bose atoms
N', N''	normalization constants in Table 1
M	magnetization of a ferromagnet
$P(M)$	probability density for a ferromagnet
$P(E)$	probability density for a laser
$P(\alpha, \alpha^*)$	P representation for the field
S	injected signal
T_c	critical temperature
W_k	heat bath density of states
x, y	$x = \text{Re } \alpha, y = \text{Im } \alpha$
X	zero-field susceptibility of a ferromagnet
$Z(T, N)$	canonical partition function
α	eigenvalue of the coherent state $ \alpha\rangle$
ε	scaled temperature (inversion) of a ferromagnet (laser) in Fig. 9
$\zeta(3)$	Riemann's zeta function $\zeta(3) = 1.2020569 \dots$
η	scaled thermodynamical variable (in Fig. 9)
$\langle \eta_k \rangle$	average occupation number of the k th heat bath oscillator
$\Theta(\cdot)$	Heaviside's unit step function
\mathbf{k}	undefined in Eq. (15)
ξ	laser analog of X
ρ_{n_0, n_0}	probability for having n_0 atoms in the condensate
$\dot{\rho}_{n_0, n_0}$	time derivative of ρ_{n_0, n_0}
σ	population inversion
σ_t	threshold inversion
$\Phi_e(\eta)$	scaled thermodynamical potential (in Fig. 9)
Ω	trap frequency

Section 23.6

A	phase-shifted destruction operator for a free-electron laser (FEL)
c_b, c_c	probability amplitudes for atom to be in levels $ b\rangle$ and $ c\rangle$, respectively

$\mathcal{D}(\theta)$	phase diffusion function
\mathcal{E}	slowly varying field amplitude
g	coupling constant for the electron-field interaction in an FEL
j	parameter for the gain in an FEL
k	wavevector for the laser field in an FEL
\mathcal{L}	linear gain ($i = j$) and cross-coupling ($i \neq j$) Liouville operators
m	mass of electron
$O(A, A^\dagger)$	arbitrary operator containing A, A^\dagger
p	electron momentum
\bar{p}	eigenvalue of electron momentum operator
P_b, P_c	probability of atom being in states $ b\rangle$ and $ c\rangle$, respectively
P_{emission}	probability of emission of radiation
$S(T)$	time-evolution operator for the electron-photon state
T	electron-photon interaction time
\mathcal{T}	time-ordering operator
z	electron coordinate
α_i	eigenvalue of the coherent state $ \alpha_i\rangle$
α_{ij}	constants depending upon parameters of gain medium in a correlated emission laser (CEL)
β	normalized electron momentum
γ	relativistic factor for an electron
Δ	atom-field detuning in lasing without inversion (LWI)
ν_1, ν_2	frequencies of the two modes in a CEL
$\rho(t_i)$	atomic density operator at initial time t_i
$\rho_{ij}^{(0)}$	initial values of the ij th atomic matrix elements
ρ_i	classical amplitude of the i th mode
$\rho(a_1, a_1^\dagger; a_2, a_2^\dagger)$	reduced density operator for the two-mode field in a CEL
θ_i	phase of the i th field
ω_0	microwave frequency
$\omega_a, \omega_b, \omega_c$	frequencies associated with atomic levels
λ_s	wavelength of the field emitted in an FEL
λ_w	the period of the magnetic wiggler
κ	gain coefficient
$\kappa_{a \rightarrow b}, \kappa_{a \rightarrow c}$	constants depending upon the matrix elements between the relevant levels
ϕ	relative phase between atomic levels
Φ	total phase angle in two-mode CEL schemes
ψ	relative phase $\Phi + \theta_1 - \theta_2$

23.2 INTRODUCTION

Most lasers, and in particular all commercially sold ones, emit electromagnetic radiation whose properties can be accounted for quite well by a *semiclassical* description. In such a treatment, quantum aspects (level spacings, oscillator strengths, etc.) of the matter (atoms, molecules, electron-hole pairs, etc.) that constitute the gain medium are essential, but those of the electromagnetic field are disregarded. Quantum properties of the radiation are, however, of decisive importance for laser systems “at the limit” which reach fundamental bounds for the linewidth, for the regularity of photon statistics, or for other quantities of interest.

Recognition and understanding of these fundamental limitations are furnished by the quantum theory of the laser, whose foundations were laid in the 1960s. The two main approaches, the master-equation formalism and the Langevin method—equivalent in the physical contents and supplementing each other like spouses—can be roughly, and somewhat superficially, associated with the Schrödinger and the Heisenberg pictures of quantum mechanics. The master-equation method corresponds to the former, the Langevin approach to the latter. Both are reviewed in Sec. 23.4, but more room is given to the master-equation treatment. This bias originates in our intention to present a parallel exposition for both the standard laser theory and the theory of the micromaser, which in turn is traditionally and most conveniently treated by master equations.

The micromaser, in which the dynamic is dominated by the strong coupling of a single mode of the radiation field to a single atomic dipole transition, is the prototype of an open, driven quantum system. Accordingly, micromaser experiments are *the* test ground for the quantum theory of the laser; therefore, micromaser theory deserves the special attention that it receives in Sec. 23.4.

As a logical and historical preparation, we recall in Sec. 23.3, the theoretical and experimental facts that are evidence for quantum properties of electromagnetic radiation in general, and the reality of photons in particular. Some special issues are discussed in Secs. 23.5 and 23.6. In Sec. 23.5, we stress the analogy between the threshold behavior of a laser and the phase transition of a ferromagnet, and note the recent lessons about Bose-Einstein condensates taught by this analogy. Section 23.6 summarizes the most important features of some exotic lasers and masers, which exploit atomic coherences or the quantum properties of the atomic center-of-mass motion. Basics of the so-called free-electron laser are reported as well.

The quantum theory of the laser is a central topic in the field of quantum optics. An in-depth understanding of the various facets of quantum optics can be gained by studying the pertinent textbooks.^{1–18}

23.3 SOME HISTORY OF THE PHOTON CONCEPT

Early History: Einstein's Light Quanta

Planck's formula of 1900¹⁹ marks the beginning of quantum mechanics, and in particular of the quantum theory of light. It reads

$$W = \frac{8\pi\nu^2}{c^3} \frac{h\nu}{\exp(h\nu/k_B T) - 1} \quad (1)$$

and relates the spectral-spatial energy density W of blackbody radiation to the frequency ν of the radiation and the temperature T of the blackbody. Boltzmann's constant k_B and Planck's constant h are conversion factors that turn temperature and frequency into energy, and c is the speed of light. A volume dV contains electromagnetic energy of the amount $W d\nu dV$ in the frequency range $\nu \cdots \nu + d\nu$.

The first factor in Eq. (1) is the density of electromagnetic modes. It obtains as a consequence of the classical wave theory of light and owes its simplicity to an implicit short-wavelength approximation. For wavelengths of the order of magnitude set by the size of the cavity that contains the radiation, appropriate corrections have to be made that reflect the shape and size of the cavity. This is of great importance in the context of micromasers, but need not concern us presently.

The second factor in Eq. (1) is the mean energy associated with radiation of frequency ν . It is a consequence of the quantum nature of light. In the limits of very high frequencies or very low ones, it turns into the respective factors of Wien²⁰ and Rayleigh-Jeans:^{21,22}

$$\begin{aligned} & \frac{h\nu}{\exp(h\nu/k_B T) - 1} && \text{(Planck)} \\ \rightarrow & \begin{cases} h\nu \exp\left(-\frac{h\nu}{k_B T}\right) & \text{for } \nu \gg \frac{k_B T}{h} & \text{(Wien)} \\ k_B T & \text{for } \nu \ll \frac{k_B T}{h} & \text{(Rayleigh-Jeans)} \end{cases} && (2) \end{aligned}$$

The relevant frequency scale is set by $k_B T/h$; at a temperature of $T = 288$ K it is about 6×10^{12} Hz, corresponding to a wavelength of $50 \mu\text{m}$.

Ironically, Planck—whose stroke of genius was to interpolate between the two limiting forms of Eq. (2)—was not convinced of the quantum nature of electromagnetic radiation until much later. Legend has it that it was the discovery of Compton scattering in 1923 that did it.²³ We are, however, getting ahead of the story.

The true significance of Planck's formula, Eq. (1), started to emerge only after Einstein²⁴ had drawn the conclusions that led him to his famous *light-quantum hypothesis* of 1905, the *annus mirabilis*. In Pauli's words,

He immediately applied [it] to the photoelectric effect and to Stokes' law for fluorescence, later also to the generation of secondary cathode rays by X-rays and to the prediction of the high frequency limit in the *Bremsstrahlung*.²⁵

Quite a truckload, indeed.

The conflict with the well-established wave theory of light was, of course, recognized immediately, and so the introduction of light quanta also gave birth to the wave-particle duality. Upon its extension to massive objects by de Broglie in 1923 to 1924,^{26–27} it was instrumental in Schrödinger's wave mechanics.²⁸

Taylor's 1909 experiment,²⁹ in which feeble light produced interference fringes, although at most one light quantum was present in the interferometer at any time, addressed the issue of wave-particle duality from a different angle. Its findings are succinctly summarized in Dirac's dictum that “a photon interferes only with itself”^{29,30}—a statement that became the innocent victim of misunderstanding and misquotation in the course of time.³¹

Another important step was taken the same year by Einstein.³² By an ingenious application of thermodynamic ideas to Planck's formula, in particular consequences of Boltzmann's relation between entropy and statistics, he derived an expression for the mean-square energy fluctuations w^2 of the radiation in a frequency interval $\nu \cdots \nu + d\nu$ and a volume dV :

$$\overline{w^2} = \left(\frac{c^3}{8\pi\nu^2} W + h\nu \right) W d\nu dV \quad (3)$$

where W is the spectral-spatial density of Eq. (1), so that $W d\nu dV$ is the mean energy in the frequency interval and volume under consideration. The first term is what one would get if classical electrodynamics accounted for all properties of radiation. There is no room for the second term in a wave theory of light; it is analogous to the fluctuations in the number of gas molecules occupying a given volume. This second term therefore supports Einstein's particle hypothesis of 1905,²⁴ in which electromagnetic energy is envisioned as being concentrated in localized lumps that are somehow distributed over the volume occupied by the electromagnetic wave.

Wave aspects (first term) and particle aspects (second term) enter Eq. (3) on equal footing. Since the thermodynamic considerations have no bias toward either one, one must conclude that Planck's formula, Eq. (1), is unbiased as well. Electromagnetic radiation is as much a particle phenomenon as it is a wave phenomenon.

Einstein left the center stage of quantum theory for some years, returning to it after completing his monumental work on general relativity. In 1913, Bohr's highly speculative postulates³³ had suddenly led to a preliminary understanding of many features of atomic spectra (the anomalous Zeeman effect was one big exception; it remained a bewildering puzzle for another decade). In the course, “quantum theory was liberated from the restriction to such particular systems as Planck's oscillators” (Pauli²⁵).

Here was the challenge to rederive Planck's formula from Bohr's postulates, assuming that they hold for arbitrary atomic systems. Einstein's famous paper of 1917³⁴ accomplished just that, and more.

He considered radiation in thermal equilibrium with a dilute gas of atoms at temperature T . We shall here give a simplified treatment that contains the essential features without accounting for all

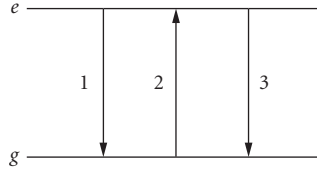


FIGURE 1 Transitions of three kinds happen between the excited level e and the ground level g : (1) spontaneous emission [see Eq. (4)], (2) absorption [see Eq. (5)], and (3) stimulated emission [see Eq. (6)].

details of lesser significance. Suppose that the energy spacing between two atomic levels e and g equals $h\nu$, so that the transition from the more energetic level e to the energetically lower level g (excited to ground state) is accompanied by the emission of a light quantum of frequency ν . According to Einstein, three processes are to be taken into account (see Fig. 1): spontaneous emission, absorption, and stimulated emission. The latter has no analog in Maxwell's electrodynamics.

Each of the three processes leads to a rate of change of the number of gas atoms in states e and g . Denoting these numbers by N_e and N_g we have the following contributions to their time derivatives. The *spontaneous emission* rate is proportional to the number of excited-state atoms:

$$\text{Spontaneous emission} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = -AN_e \quad (4)$$

where A is the first Einstein coefficient of the transition in question. The *absorption* rate is proportional to the number of ground-state atoms and to the energy density W of the radiation:

$$\text{Absorption} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = BWN_g \quad (5)$$

where B is the second Einstein coefficient. The *stimulated emission* rate is proportional to the number of excited state atoms and to the radiation energy density:

$$\text{Stimulated emission} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = BWN_e \quad (6)$$

where the same B coefficient appears as in the absorption rate.

The *detailed balance* between states e and g therefore requires

$$\text{Total} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = -AN_e + BWN_g - BWN_e = 0 \quad (7)$$

under the circumstances of thermal equilibrium. Therefore, W can be expressed in terms of the ratios A/B and N_g/N_e :

$$W = \frac{A/B}{N_g/N_e - 1} = \frac{A/B}{\exp(h\nu/k_b T) - 1} \quad (8)$$

where the second equality recognizes that Boltzmann's factor relates N_e to N_g . [The absence of additional weights here is the main simplification alluded to before; if taken into account, these weights would also require two closely related B coefficients in Eqs. (5) and (6).] In Eq. (8) we encounter the

denominator of Planck's factor from Eq. (2), and Planck's formula, Eq. (1), is recovered in full if the relation

$$A = \frac{8\pi h\nu^3}{c^3} B \quad (9)$$

is imposed on Einstein's coefficients.

The main ingredients in this derivation of Eq. (1) are the postulate of the process of stimulated emission, with a strength proportional to the density of radiation energy, and the relation between the coefficients for spontaneous emission and stimulated emission, Eq. (9).

All of this is well remembered, but there was in fact more to the 1917 paper.³⁴ It also contains a treatment of the momentum exchange between atoms and light quanta, and Einstein succeeded in demonstrating that the Maxwell velocity distribution of the atoms is consistent with the recoil they suffer when absorbing and emitting quanta. Insights gained in his study of Brownian motion (another seminal paper of 1905³⁵) were crucial for this success. When taken together, the considerations about energy balance and those concerning momentum balance are much more convincing than either one could have been alone.

The discovery of the Compton effect in 1923²³ finally convinced Bohr and other skeptics of the reality of the particlelike aspects possessed by light. But Bohr, who until then was decidedly opposed to Einstein's light-quantum hypothesis and the consequent wave-particle duality, did not give in without a last try. Together with Kramers and Slater,³⁶ he hypothesized that perhaps energy-momentum conservation does not hold for each individual scattering event, but only in a statistical sense for a large ensemble. Then one could account for Compton's data without conceding a particle nature to light in general, and X rays in particular. The refined measurements that were immediately carried out by Bothe and Geiger³⁷ showed, however, that this hypothesis is wrong: energy and momentum are conserved in each scattering event, not just statistically.

And then there was, also in 1924, Bose's seminal observation that it is possible to derive Planck's radiation law [Eq. (1)] from purely corpuscular arguments without invoking at all the wave properties of light resulting from Maxwell's field equations.^{38,39} The main ingredient in Bose's argument was the indistinguishability of the particles in question and a new way of counting them—now universally known as *Bose-Einstein statistics*—which pays careful attention to what is implied by their being indistinguishable. In the case of light quanta, an additional feature is that their number is not conserved, because light is easily emitted and absorbed. Massive particles (atoms, molecules, etc.), by contrast, are conserved; therefore, as Einstein emphasized,^{40–42} their indistinguishability has further consequences, of which the phenomenon of Bose-Einstein condensation (or should one rather say “Einstein condensation of a Bose gas”?) is the most striking one.

Quantum Electrodynamics

Theoretical studies of the quantum nature of light had a much more solid basis after Dirac's introduction of quantum electrodynamics (QED) in his seminal paper of 1927.⁴³ The basic ingredients of QED were all present in Dirac's formulation, although it is true that a consistent understanding of QED was not available until renormalized QED was developed 20 years later (see the papers reprinted in Ref. 44). In particular, the photon concept was clarified in the sense described in the following paragraphs.

The infinite number of degrees of freedom of the electromagnetic field—an operator field in QED—become manageable with the aid of a mode expansion. For the transverse part $\mathbf{E}_\perp(\mathbf{r}, t)$ of the electric field, for example, it reads

$$\mathbf{E}_\perp(\mathbf{r}, t) = \sum_k \sqrt{\frac{h\nu_k}{2\epsilon_0}} [a_k(t)\mathbf{A}_k(\mathbf{r}) + a_k^\dagger(t)\mathbf{A}_k^*(\mathbf{r})] \quad (10)$$

The mode functions $\mathbf{A}_k(\mathbf{r})$ are complex vector functions of the position vector \mathbf{r} that are eigenfunctions of the Laplace differential operator,

$$-\nabla^2 \mathbf{A}_k(\mathbf{r}) = \left(\frac{2\pi\nu_k}{c} \right)^2 \mathbf{A}_k(\mathbf{r}) \quad (11)$$

where the eigenvalue is determined by the frequency $\nu_k (>0)$ of the mode in question. The boundary conditions that the electric and magnetic field must obey at conducting surfaces imply respective boundary conditions on the $\mathbf{A}_k(\mathbf{r})$ s.

The corresponding mode expansion for the magnetic field is given by

$$\mathbf{B}(\mathbf{r}, t) = \sum_k \sqrt{\frac{\mu_0 \hbar \nu_k}{2}} [-ia_k(t)\mathbf{B}_k(\mathbf{r}) + ia_k^\dagger(t)\mathbf{B}_k^*(\mathbf{r})] \quad (12)$$

where

$$\mathbf{B}_k(\mathbf{r}) = \frac{c}{2\pi\nu_k} \nabla \times \mathbf{A}_k(\mathbf{r}); \quad \mathbf{A}_k(\mathbf{r}) = \frac{c}{2\pi\nu_k} \nabla \times \mathbf{B}_k(\mathbf{r}) \quad (13)$$

relates the two kinds of mode functions to each other.

Among others, the mode functions $\mathbf{A}_k(\mathbf{r})$ have the following important properties:

$$\text{Transverse} \quad \nabla \cdot \mathbf{A}_k(\mathbf{r}) = 0 \quad (14a)$$

$$\text{Orthonormal} \quad \int (d\mathbf{r}) \mathbf{A}_j^*(\mathbf{r}) \cdot \mathbf{A}_k(\mathbf{r}) = \delta_{jk} \quad (14b)$$

$$\text{Complete} \quad \sum_k \mathbf{A}_k(\mathbf{r}) \mathbf{A}_k^*(\mathbf{r}') = \delta_\perp(\mathbf{r} - \mathbf{r}') \quad (14c)$$

The same statements hold for the $\mathbf{B}_k(\mathbf{r})$ s as well. The property in Eq. (14a) states the radiation-gauge condition. The integration in Eq. (14b) covers the entire volume bounded by the conducting surfaces just mentioned; the eigenvalue equation Eq. (11) holds inside this volume, the so-called quantization volume. In the completeness relation in Eq. (14c), both positions \mathbf{r} and \mathbf{r}' are inside the quantization volume, and δ_\perp is the transverse delta function, a dyadic that is explicitly given by

$$\delta_\perp(\mathbf{r}) = \int \frac{(d\mathbf{k})}{(2\pi)^3} \exp(i\mathbf{k} \cdot \mathbf{r}) \left(1 - \frac{\mathbf{k}\mathbf{k}}{k^2} \right) \quad (15)$$

where 1 is the unit dyadic and $k = \sqrt{\mathbf{k} \cdot \mathbf{k}}$ is the length of the wave vector \mathbf{k} integrated over. The transverse character of $\delta_\perp(\mathbf{r})$ ensures the consistency of the properties in Eqs. (14a) and (14c).

The time dependence of $\mathbf{E}_\perp(\mathbf{r}, t)$ and $\mathbf{B}(\mathbf{r}, t)$ stems from the ladder operators $a_k(t)$ and $a_k^\dagger(t)$, which obey the bosonic equal-time commutation relations

$$[a_j, a_k] = 0 \quad [a_j, a_k^\dagger] = \delta_{jk} \quad [a_j^\dagger, a_k^\dagger] = 0 \quad (16)$$

The photon number operator

$$N = \sum_k a_k^\dagger a_k \quad (17)$$

has eigenvalues $N' = 0, 1, 2, \dots$; its eigenstates with $N' = 1$ are the one-photon states, those with $N' = 2$ are the two-photon states, and so on. The unique eigenstate with $N' = 0$ is the photon vacuum.

We denote its ket by $|\mathbf{vac}\rangle$. It is, of course, the joint eigenstate of all *annihilation operators* a_k with eigenvalue zero:

$$a_k |\mathbf{vac}\rangle = 0 \quad \text{for all } k \quad (18)$$

Application of the *creation operator* a_k^\dagger to $|\mathbf{vac}\rangle$ yields a state with one photon in the k th mode:

$$a_k^\dagger |\mathbf{vac}\rangle = \{\text{state with 1 photon of the type } k\} \quad (19)$$

More generally, the ket of a pure one-photon state is of the form

$$|\{\psi\}_1\rangle = \sum_k \psi_k a_k^\dagger |\mathbf{vac}\rangle \quad \text{with} \quad \sum_k |\psi_k|^2 = 1 \quad (20)$$

where $|\psi_k|^2$ is the probability for finding the photon in the k th mode. Similarly, the kets of pure two-photon states have the structure

$$|\{\psi\}_2\rangle = \frac{1}{\sqrt{2}} \sum_{j,k} \psi_{jk} a_j^\dagger a_k^\dagger |\mathbf{vac}\rangle \quad \text{with} \quad \psi_{jk} = \psi_{kj} \quad \text{and} \quad \sum_{j,k} |\psi_{jk}|^2 = 1 \quad (21)$$

and analogous expressions apply to pure states with 3, 4, 5, . . . photons.

Einstein's light quanta are one-photon states of a particular kind. In a manner of speaking, they are localized lumps of electromagnetic energy. In technical terms this means that the energy density

$$\begin{aligned} U_{\text{1ph}}(\mathbf{r}, t) &= \left\langle \{\psi\}_1 \left| : \left[\mathbf{E}^2 / (2\epsilon_0) + \left(\frac{\mu_0}{2} \right) \mathbf{B}^2 \right] : \right| \{\psi\}_1 \right\rangle \\ &= \frac{h}{2} \sum_{j,k} \psi_j^* \psi_k \sqrt{v_j v_k} (\mathbf{A}_k^* \cdot \mathbf{A}_j + \mathbf{B}_k^* \cdot \mathbf{B}_j) \end{aligned} \quad (22)$$

is essentially nonzero in a relatively small spatial region only. The time dependence is carried by the probability amplitudes ψ_k , the spatial dependence by the mode functions \mathbf{A}_k and \mathbf{B}_k . An arbitrarily sharp localization is not possible, but it is also not needed. The pair of colons :: symbolize the injunction to order the operator in between in the *normal* way: all creation operators a_k^\dagger to the left of all annihilation operators a_k . This normal ordering is an elementary feature of renormalized QED.

At high frequencies, or when the quantization volume is unbounded, the eigenvalues of $-\nabla^2$ in Eq. (11) are so dense that the summations in Eqs. (10), (12), and (14b) are effectively integrations, and the Kronecker delta symbol in Eq. (14b) is a Dirac delta function. Under these circumstances, it is often natural to choose plane waves

$$\mathbf{A}_k(\mathbf{r}) \sim \mathbf{e}_k \exp\left(i \frac{2\pi \nu_k}{c} \mathbf{n}_k \cdot \mathbf{r}\right) \quad \mathbf{B}_k(\mathbf{r}) \sim \mathbf{n}_k \times \mathbf{e}_k \exp\left(i \frac{2\pi \nu_k}{c} \mathbf{n}_k \cdot \mathbf{r}\right) \quad (23)$$

for the mode functions. The unit vector \mathbf{e}_k that specifies the polarization is orthogonal to the unit vector \mathbf{n}_k that specifies the direction of propagation.

With

$$\psi_k(t) \sim \exp(-i2\pi \nu_k t) \quad (24)$$

in Eq. (22), one then meets exponential factors of the form

$$\exp\left[i \frac{2\pi \nu_k}{c} (\mathbf{n}_k \cdot \mathbf{r} - ct)\right]$$

As a consequence, an einsteinian light quantum propagates without dispersion, which is the anticipated behavior.

The one-photon energy density in Eq. (22) illustrates the general feature that quantum-mechanical probabilities (the ψ_k s) with their interference properties appear together with the classical interference patterns of superposed mode functions [the $\mathbf{A}_k(\mathbf{r})$ s and $\mathbf{B}_k(\mathbf{r})$ s]. In other words, interference phenomena of two kinds are present in QED: (1) the classical interference of electromagnetic fields in the three-dimensional \mathbf{r} space, and (2) the quantum interference of alternatives in the so-called Fock space; that is, the Hilbert space spanned by the photon vacuum $|\mathbf{vac}\rangle$, the one-photon states $|\{\psi\}_1\rangle$, the two-photon states $|\{\psi\}_2\rangle$, and all multiphoton states.

In the early days of QED, this coexistence of classical interferences and quantum interferences was a research topic, to which Fermi's paper of 1929 "Sulla teoria quantistica delle frange di interferenza" is a timeless contribution.⁴⁵ He demonstrated, at the example of Lippmann fringes, a very general property of single-photon interference patterns: the photon-counting rates, as determined from quantum-mechanical probabilities, are proportional to the corresponding classical intensities.

Electromagnetic radiation is easily emitted and absorbed by antennas, processes that change the number of photons. Accordingly, the number of photons is not a conserved quantity, and therefore states of different photon numbers can be superposed. Particularly important are the *coherent states*

$$|\{\alpha\}_c\rangle = \exp\left(-\frac{1}{2}\sum_k |\alpha_k|^2 + \sum_k \alpha_k a_k^\dagger\right) |\mathbf{vac}\rangle \quad (25)$$

that are characterized by a set $\{\alpha\}_c$ of complex amplitudes α_k . As revealed in Glauber's 1963 papers,^{46–48} they play a central role in the coherence theory of light.

Since the coherent states are common eigenstates of the annihilation operators

$$a_k |\{\alpha\}_c\rangle = |\{\alpha\}_c\rangle \alpha_k \quad (26)$$

the expectation values of the electric and magnetic field operators of Eqs. (10) and (12)

$$\langle \{\alpha(t)\}_c | \mathbf{E}_\perp(\mathbf{r}, t) | \{\alpha(t)\}_c \rangle = \sum_k \sqrt{\frac{h\nu_k}{2\epsilon_0}} [\alpha_k(t) \mathbf{A}_k(\mathbf{r}) + \alpha_k^*(t) \mathbf{A}_k^*(\mathbf{r})] \quad (27)$$

$$\langle \{\alpha(t)\}_c | \mathbf{B}_\perp(\mathbf{r}, t) | \{\alpha(t)\}_c \rangle = \sum_k \sqrt{\frac{\mu_0 h\nu_k}{2}} [-i\alpha_k(t) \mathbf{B}_k(\mathbf{r}) + i\alpha_k^*(t) \mathbf{B}_k^*(\mathbf{r})]$$

have the appearance of classical Maxwell fields. In more general terms, if the statistical operator ρ_{ph} of the photonic degrees of freedom—in other words, the statistical operator of the radiation field—is a mixture of (projectors to) coherent states

$$\rho_{\text{ph}} = \sum_{\{\alpha\}_c} |\{\alpha\}_c\rangle w(\{\alpha\}_c) \langle \{\alpha\}_c| \quad (28)$$

with

$$w(\{\alpha\}_c) \geq 0 \quad \text{and} \quad \sum_{\{\alpha\}_c} w(\{\alpha\}_c) = 1 \quad (29)$$

then the electromagnetic field described by ρ_{ph} is very similar to a classical Maxwell field. Turned around, this says that whenever it is impossible to write a given ρ_{ph} in the form of Eq. (28), then some statistical properties of the radiation are decidedly nonclassical.

During the 20-year period from Dirac's paper of 1927 to the Shelter Island conference in 1947, QED remained in a preliminary state that allowed various studies—the most important ones included the Weisskopf-Wigner treatment of spontaneous emission⁴⁹ and Weisskopf's discovery that the self-energy of the photon is logarithmically divergent⁵⁰—although the not-yet-understood divergences

were very troublesome. The measurement by Lamb and Retherford⁵¹ of what is now universally known as the *Lamb shift*, first reported at the Shelter Island conference, was the crucial experimental fact that triggered the rapid development of renormalized QED by Schwinger, Feynman, and others.

Theoretical calculations of the Lamb shift rely heavily on the quantum properties of the electromagnetic field, and their marvelous agreement with the experimental data proves convincingly that these quantum properties are a physical reality. In other words, photons exist. The same remark applies to the theoretical and experimental values of the anomalous magnetic moment of the electron, one of the early triumphs of Schwinger's renormalized QED,⁵² which finally explained an anomaly in the spectra of hydrogen and deuterium that Pasternack had observed in 1938⁵³ and a discrepancy in the measurements by Millman and Kusch⁵⁴ of nuclear magnetic moments.

Photon-Photon Correlations

Interferometers that exploit not the spatial intensity variations (or, equivalently, the photon-detection probabilities) but correlations between intensities at spatially separated positions became important tools in astronomy and spectroscopy after the discovery of the Hanbury-Brown-Twiss (HB&T) effect in 1954.^{55–57} A textbook account of its classical theory is given in Sec. 4.3 of Ref. 58.

In more recent years, the availability of single-photon detectors made it possible to study the HB&T effect at the two-photon level. The essentials are depicted in Fig. 2. Two light quanta are incident on a half-transparent mirror from different directions, such that they arrive simultaneously. If their frequency contents are the same, it is fundamentally impossible to tell if an outgoing light quantum was reflected or transmitted. This indistinguishability of the two light quanta is of decisive importance in the situation where one is in each output channel. The two cases *both reflected* and *both transmitted* are then indistinguishable and, according to the laws of quantum mechanics, the corresponding probability amplitudes must be added.

Now, denoting the probability amplitudes for single-photon reflection and transmission by ρ and τ , respectively, the probability for one light quantum in each output port is given by

$$|\rho^2 + \tau^2|^2 = \left| \left(\frac{1}{\sqrt{2}} \right)^2 + \left(\frac{i}{\sqrt{2}} \right)^2 \right|^2 = 0 \quad (30)$$

where we make use of $\rho = 1/\sqrt{2}$ and $\tau = i/\sqrt{2}$, which are the values appropriate for a symmetric half-transparent mirror. Thus, the situation of one light quantum in each output port does not occur. Behind the half-transparent mirror, one always finds both light quanta in the same output port.

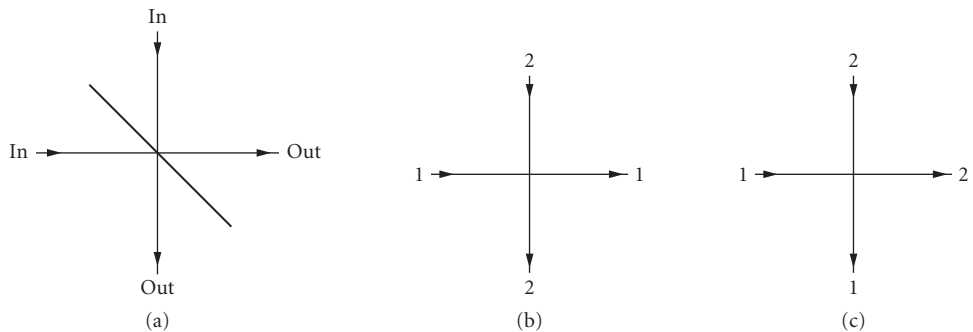


FIGURE 2 Essentials of the Hanbury-Brown-Twiss effect at the two-photon level. (a) Two light quanta are simultaneously incident at a symmetric half-transparent mirror. To obtain one quantum in each output port, the quanta must be either both transmitted (b) or both reflected (c). The probability amplitude for (b) is $(1/\sqrt{2})^2 = 1/2$, that for (c) is $(i/\sqrt{2})^2 = -1/2$. If the two cases are indistinguishable, the total amplitude is $1/2 - 1/2 = 0$.

If the light quanta arrived at very different times, rather than simultaneously, they would be distinguishable and one would have to add probabilities instead of probability amplitudes, so that

$$|\rho^2|^2 + |\tau^2|^2 = \frac{1}{2} \quad (31)$$

would replace Eq. (30). Clearly, there are intermediate stages at which the temporal separation is a fraction of the temporal coherence length and the two quanta are neither fully distinguishable nor utterly indistinguishable. The probability for one light quantum in each output port is then a function of the separation, a function that vanishes when the separation does.

Experiments that test these considerations⁵⁹ employ correlated photon pairs produced by a process known as *parametric downconversion*. Roughly speaking, inside a crystal that has no inversion symmetry a high-frequency photon is absorbed and two lower-frequency photons are emitted, whereby the conservation of energy and momentum imposes geometrical restrictions on the possible propagation directions of the three photons involved. Downconversion sources with a high luminosity are available.^{60,61}

The HB&T effect of Fig. 2 as well as closely related phenomena are crucial in many experiments in which entangled photons are a central ingredient. In particular, the recent realizations^{62,63} of schemes for quantum teleportation⁶⁴ and the experiment⁶⁵ that demonstrated the practical feasibility of quantum-dense coding⁶⁶ are worth mentioning here.

None of these exciting developments could be understood without the quantum properties of radiation. Since the photon concept, in the sense of the discussion of Eqs. (10), (12), (19), (20), and so on is an immediate consequence of these quantum properties, the existence of photons is an established experimental fact beyond reasonable doubt.

23.4 QUANTUM THEORY OF THE LASER

The quantum theory of laser radiation is a problem in nonequilibrium statistical mechanics. There are several alternative, but ultimately equivalent, approaches to the characterization of the field inside the resonator. As is customary in the Scully-Lamb quantum theory, we describe the state of the laser field by a density operator.^{67,68} In this section our main focus is on the review of the equation of motion, the so-called master equation, for this density operator as it emerges from an underlying physical model, with statistical considerations and some simplifying assumptions. The alternative procedure based on the quantum theory of noise sources introduced in Refs. 69 and 70 and summarized in Refs. 71 to 74 will also be briefly reviewed at the end of the section. For a recent, more detailed overview of the quantum Langevin point of view, we refer the reader to Ref. 75.

In general, a laser model should be based on the interaction of multimode fields with multilevel atoms as the active medium, and a detailed consideration of all possible processes among all the levels involved should be given. Decay channels and decay rates, in particular, play a crucial role in determining the threshold inversion and, thus, the necessary pumping rates. Of course, the pumping mechanisms themselves can be quite complicated. It is well established that a closed two-level model cannot exhibit inversion and, hence, lasing. In order to achieve inversion three- and four-level pumping schemes are employed routinely. On the other hand, to illustrate the essential quantum features a single mode field can serve as paradigm. The single-mode laser field inside the resonator interacts with one particular transition of the multilevel system—the lasing transition—and the role of the entire complicated level structure is to establish inversion on this transition—that is, to put more atoms in the upper level than there are in the lower one. If one is not interested in the details of how the inversion builds up, it is possible to adopt a much simpler approach than the consideration of a multilevel-multimode system. In order to understand the quantum features of the single-mode field it is sufficient to focus only on the two levels of the lasing transition and their interaction with the laser field. In this approach, the effect of pumping, decay, and so on in the multilevel structure is simply replaced by an initial condition; it is assumed that the atom is in its upper state immediately before the interaction with the laser mode begins. Since here we are primarily interested

in the quantum signatures of the laser field and not in the largely classical aspects of cavity design, pumping mechanisms, and so on, we shall follow this simpler route from the beginning. The model that accounts for the resonant interaction of a two-level atom with a single quantized mode in a cavity was introduced by Jaynes and Cummings.⁷⁶

We shall make an attempt to present the material in a tutorial way. We first derive an expression for the change of the field density operator due to the interaction with a single two-level atom, initially in its upper state, using the Jaynes-Cummings model. This expression will serve as the seed for both the laser and micromaser theories. We next briefly review how to account for cavity losses by using standard methods for modeling the linear dissipation loss of the cavity field due to mirror transmission. Then we show that with some additional assumptions the single-atom-single-mode approach can be used directly to derive what has become known as the Scully-Lamb master equation for the more traditional case of the laser and the micromaser. The additional assumptions include the Markov approximation or, equivalently, the existence of very different time scales for the atomic and field dynamics so that adiabatic elimination of the atoms and introduction of coarse-grained time evolution for the field become possible. The main difference between the laser and micromaser theories is that the interaction time of the active atoms with the field is governed by the lifetime of the atoms in the laser and by the transit time of the atoms through the cavity in the micromaser. In the laser case, the atoms decay out of the lasing levels into some far-removed other levels, and they are available for the lasing transition during their lifetime on the average. In the micromaser case, the transit time is approximately the same for all atoms in a monoenergetic pumping beam. Thus, the laser involves an extra averaging over the random interaction times. If we model the random interaction times by a Poisson distribution and average the change of the field density operator that is due to a single atom—the kick—over the distribution of the interaction times, we obtain the master equation of a laser from that of the micromaser. Historically, of course, the development was just the opposite: the master equation was derived in the context of the laser much earlier. However, it is instructive to see how the individual Rabi oscillations of single nondecaying atoms with a fixed interaction time, as in the micromaser, give rise to the saturating, nonoscillatory collective behavior of an ensemble of atoms, as in the laser, upon averaging over the interaction times. As applications of this fully quantized treatment we study the photon statistics, the linewidth, and spectral properties. Finally, we briefly discuss other approaches to the quantum theory of the laser.

Time Evolution of the Field in the Jaynes-Cummings Model

We shall consider the interaction of a single two-level atom with a single quantized mode of a resonator using the rotating-wave approximation (for a recent review of the Jaynes-Cummings model see Ref. 77). The arrangement is shown in Fig. 3.

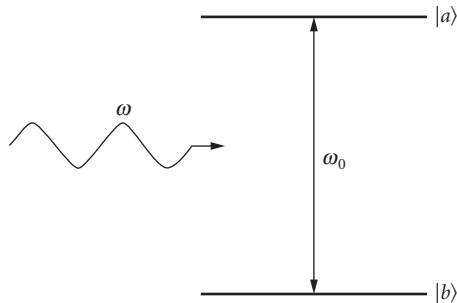


FIGURE 3 Scheme of a two-level atom interacting with a single mode quantized field. The text focuses on the resonant case, $\omega = \omega_0$.

The hamiltonian for this system is given by

$$\mathcal{H} = \mathcal{H}_0 + \mathcal{H}_1 \quad (32)$$

where

$$\mathcal{H}_0 = \frac{1}{2}\hbar\omega_0\sigma_z + \hbar\omega_0a^\dagger a \quad (33)$$

and

$$\mathcal{H}_1 = \hbar g(\sigma_+ a + a^\dagger \sigma_-) \quad (34)$$

Here a and a^\dagger are the annihilation and creation operators for the mode. The upper level of the lasing transition is denoted by $|a\rangle$ and the lower level by $|b\rangle$. The atomic lowering and raising operators are expressed in terms of the state vectors as $\sigma_- = (\sigma_+)^{\dagger} = |b\rangle\langle a|$ and the population operator as $\sigma_z = |a\rangle\langle a| - |b\rangle\langle b|$. ω_0 is the frequency of the $|a\rangle - |b\rangle$ transition. For simplicity we assume perfect resonance with the mode. Finally, g is the coupling constant between the atom and the mode. In terms of atomic and field quantities it is given by $g = p\sqrt{\omega_0/2\hbar\epsilon_0 v}$ where p is the transition dipole moment and v the quantization volume (the volume of the active medium, in our case). Again, for simplicity, p is assumed to be real.

In the interaction picture with respect to \mathcal{H}_0 , the interaction hamiltonian becomes

$$\mathcal{V} = \exp\left(\frac{i\mathcal{H}_0 t}{\hbar}\right) \mathcal{H}_1 \exp\left(-\frac{i\mathcal{H}_0 t}{\hbar}\right) = \hbar g(\sigma_+ a + a^\dagger \sigma_-) \equiv \mathcal{H}_1 \quad (35)$$

since \mathcal{H}_0 and \mathcal{H}_1 commute.

In the two-dimensional Hilbert space spanned by the state vectors $|a\rangle$ and $|b\rangle$ the interaction hamiltonian can be written as

$$\mathcal{H}_1 = \hbar g \begin{pmatrix} 0 & a \\ a^\dagger & 0 \end{pmatrix} \quad (36)$$

The time evolution operator for the coupled atom-field system satisfies the equation of motion in this picture

$$\frac{i\hbar dU}{dt} = \mathcal{V}U \quad (37)$$

and since \mathcal{H}_1 is time independent the solution is formally

$$U(\tau) = \exp\left(-\frac{i}{\hbar}\mathcal{V}\tau\right) \quad (38)$$

Using the properties of the σ_- and σ_+ matrices, it is easy to show that $U(\tau)$ can be written in the preceding 2×2 matrix representation as

$$U(\tau) = \begin{pmatrix} \cos(g\tau\sqrt{aa^\dagger}) & -i\frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}}a \\ -ia^\dagger\frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}} & \cos(g\tau\sqrt{a^\dagger a}) \end{pmatrix} \quad (39)$$

Let us now assume that initially, at t_0 , the atom is in its upper state given by the atomic density operator $\rho_{\text{at}}(t_0) = |a\rangle\langle a|$ and the field is in an arbitrary state which, in general, can be described by the density operator $\rho(t_0)$, so that the joint atom-field system is characterized by the initial density operator $\rho_{\text{af}}(t_0) = \rho_{\text{at}}(t_0) \otimes \rho(t_0)$. After the interaction, $\rho_{\text{af}}(t_0 + \tau) = U(\tau)\rho_{\text{af}}(t_0)U(\tau)^{-1}$. Our main interest here is in the evolution of the field density operator. This we obtain if we trace the atom-field density operator over the atomic states, yielding

$$\begin{aligned} \rho(t_0 + \tau) &= \text{Tr}_{\text{atom}}[\rho_{\text{af}}(t_0 + \tau)] \\ &= \cos(g\tau\sqrt{aa^\dagger})\rho(t_0)\cos(g\tau\sqrt{aa^\dagger}) + a^\dagger \frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}}\rho(t_0)\frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}}a \\ &\equiv M(\tau)\rho(t_0) \end{aligned} \quad (40)$$

Here in the last step we just introduced the superoperator M , which describes the effect of a single inverted atom on the field and is a key ingredient of laser and micromaser theory. The matrix elements in photon number representation take the form

$$(M(\tau)\rho)_{nn'} = A_{nn'}(\tau)\rho_{nn'} + B_{n-1n'-1}(\tau)\rho_{n-1n'-1} \quad (41)$$

where the coefficients are given by

$$A_{nn'}(\tau) = \cos(g\tau\sqrt{n+1})\cos(g\tau\sqrt{n'+1}) \quad (42)$$

$$B_{nn'}(\tau) = \sin(g\tau\sqrt{n+1})\sin(g\tau\sqrt{n'+1}) \quad (43)$$

For later purposes, we also introduce the change in the state of the field due to the interaction with a single inverted atom as

$$\rho(t_0 + \tau) - \rho(t_0) = M(\tau)\rho(t_0) - \rho(t_0) = (M - 1)\rho(t_0) \equiv K\rho(t_0) \quad (44)$$

The operator K , sometimes called the *kick operator*, contains all the information we will need to build the quantum theory of the single-mode laser and micromaser. In matrix representation, $[K(\tau)\rho]_{nn'} = [A_{nn'}(\tau) - \delta_{nn'}]\rho_{nn'} + B_{n-1n'-1}(\tau)\rho_{n-1n'-1}$.

For more elaborate systems (multimode lasers driven by multilevel atoms, for example) one cannot give M in such a simple analytical form, but the principle remains always the same. One should find the superoperator M or, equivalently, the kick operator $K = M - 1$ which gives the action of a single (possibly multimode) atom on the (possibly multimode) field from the general expression $\rho(t_0 + \tau) = \text{Tr}_{\text{atom}}[U_{\text{af}}(\tau)\rho_{\text{at}}(t_0) \otimes \rho(t_0)U_{\text{af}}(\tau)^{-1}] \equiv M(\tau)\rho(t_0)$. In order to determine the full time-evolution operator $U_{\text{af}}(\tau)$ of the coupled atom-field system, however, one usually needs to resort to approximation methods, such as perturbation theory, in the more complicated multilevel-multimode cases.

Derivation of the Scully-Lamb Master Equation

In 1954, Gordon, Zeiger, and Townes showed that coherent electromagnetic radiation can be generated in the radio frequency range by the maser (acronym for *microwave amplification by stimulated emission of radiation*).⁷⁸ The first maser action was observed in a beam of ammonia.⁷⁹ The maser principle was extended to the optical domain by Schawlow and Townes,⁸⁰ and also by Prokhorov

and Basov,⁸¹ thus obtaining a laser (acronym for *light amplification by stimulated emission of radiation*). A laser consists of a large ensemble of inverted atoms interacting resonantly with the electromagnetic field inside a cavity. The cavity selects only a specific set of modes corresponding to a discrete set of eigenfrequencies. The active atoms—that is, the ones that are pumped to the upper level of the laser transition—are in resonance with one of the eigenfrequencies of the cavity in the case of the single-mode laser and with a finite set of frequencies in the case of the multimode laser. As discussed in the introduction to this section, for the discussion of the essential quantum features of the radiation field of a laser it is sufficient to confine our treatment to the single-mode case, and that is what we will do for the remainder of this section. A resonant electromagnetic field gives rise to stimulated emission, and the atoms thereby transfer their energy to the radiation field. The emitted radiation is still at resonance. If the upper level is sufficiently populated, this radiation gives rise to further transitions in other atoms. In this way, all the excitation energy of the atoms is transferred to the single mode of the radiation field.

The first pulsed laser operation was demonstrated by Maiman in ruby.⁸² The first continuous wave (CW) laser, a He-Ne gas laser, was built by Javan et al.⁸³ Since then a large variety of systems have been demonstrated to exhibit lasing action. Coherent radiation has been generated this way over a frequency domain ranging from infrared to soft X rays. These include dye lasers, chemical lasers, solid-state lasers, and semiconductor lasers.

Many of the laser properties can be understood on the basis of a semiclassical theory. In such a theory the radiation field is treated classically, but the active medium is given a full quantum-mechanical treatment. Such a theory can readily explain threshold and saturation, transient dynamics, and general dependence on the external parameters (pumping and losses). It is not our aim here to give an account of the semiclassical theory; therefore, we just refer the reader to the ever instructive and wonderfully written seminal paper by Lamb⁸⁴ and a more extended version in Ref. 67. Although quantum effects play only a minor role in usual practical laser applications because of the large mean photon numbers, they are essential for the understanding of the properties of micromasers, in which excited two-level atoms interact one after the other with a single mode of the radiation field.⁸⁵ Nevertheless, the quantum properties of the laser field are of fundamental interest as well. They have been thoroughly investigated theoretically with respect to the photon statistics and the spectrum of the laser. In particular, the quantum limitation of the laser linewidth caused by the inevitably noisy contribution of spontaneous emission has attracted much attention. It gives rise to the so-called Schawlow-Townes linewidth, which is inversely proportional to the laser intensity (see Ref. 80). Because of the importance of stable coherent signals for various high-precision measurements, the problem of the intrinsic quantum-limited linewidth has gained renewed interest recently, and the investigations have been extended to cover bad-cavity lasers and several more exotic systems. In this review, however, we shall restrict ourselves to good-cavity lasers in which the cavity damping time is long compared to all other relevant time scales, and present a fully quantized theory of the most fundamental features.

Cavity Losses To account for the decay of the cavity field through the output mirror of the cavity, we simply borrow the corresponding result from reservoir theory. Its usage has become fairly standard in laser physics and quantum optics (see, for example, Refs. 67 and 68), and here we just quote the general expression without actually deriving it.

$$\left(\frac{d\varrho}{dt}\right)_{\text{loss}} = \mathcal{L}\varrho \equiv -\frac{\kappa}{2}n_{\text{th}}(aa^\dagger\varrho + \varrho aa^\dagger - 2a^\dagger\varrho a) - \frac{\kappa}{2}(n_{\text{th}}+1)(a^\dagger a\varrho + \varrho a^\dagger a - 2a\varrho a^\dagger) \quad (45)$$

This equation refers to a loss reservoir which is in thermal equilibrium at temperature T , with n_{th} being the mean number of thermal photons $n_{\text{th}} = [\exp(\hbar\omega_0/kT) - 1]^{-1}$, and κ is the cavity damping rate. For the laser case, it is sufficient to take the limiting case of a zero temperature reservoir since $\hbar\omega_0 \gg kT$ and n_{th} is exponentially small. We obtain this limit by substituting $n_{\text{th}} = 0$ into Eq. (45). For the description of most micromaser experiments, however, we need the finite temperature version, since even at very low temperatures the thermal photon number is comparable to the total number of photons in the cavity.

The Laser Master Equation After introducing the loss part of the master equation, we now turn our attention to the part that stems from the interaction with the gain reservoir. The gain reservoir is modeled by an ensemble of initially excited two-level atoms allowed to interact with the single-mode cavity field. A central role in our subsequent discussions will be played by the so-called kick operator, $K = M - 1$, describing the change of the field density operator due to the interaction with a single atom. This quantity was introduced in Eq. (44). While in the micromaser case the effect of each of the atoms can be represented by the same kick operator, since in a monoenergetic pumping beam each atom has the same interaction time with the cavity field, this is no longer the case for a laser. In a typical CW gas laser, such as the He-Ne laser, atoms are excited to the upper level of the lasing transition at random times and, more important, they can also interact with the field for a random length of time. The interaction time thus becomes a random variable. Since at any given time the number of atoms is large (about 10^6 to 10^7 active atoms in the lasing volume of a CW He-Ne laser), it is a legitimate approach to describe their effect on the field by an average kick operator. We can arrive at the interaction-time-averaged master equation quickly if we take the average of Eq. (44) with respect to the interaction time τ :

$$(M-1)\varrho(t) = \int_0^\infty d\tau P(\tau)(M(\tau)-1)\varrho(t) \quad (46)$$

where the distribution function for the interaction time $P(\tau)$ is defined as

$$P(\tau) = \gamma e^{-\gamma\tau} \quad (47)$$

This distribution function corresponds to the exponential decay law. Individual atoms can decay from the lasing levels at completely random times, but for an ensemble of atoms the probability of finding an initially excited atom still in the lasing levels in the time interval $(\tau, \tau + d\tau)$ is given by Eq. (47). With increasing τ it is increasingly likely that the atoms have decayed outside the lasing transition. Also note that our model corresponds to an open system: the atoms decay to other non-lasing levels both from the upper state $|a\rangle$ and the lower state $|b\rangle$, and, in addition, we assume that decay rate γ is the same for both levels, as indicated in Fig. 4.

Obviously, these restrictions can be relaxed and, indeed, there are various more general models available. For example, the upper level $|a\rangle$ can have two decay channels. It can decay to the lower level $|b\rangle$ and to levels outside the lasing transition. Or, in some of the most efficient lasing schemes, the lower level decays much faster than the upper one, $\gamma_b \gg \gamma_a$. In these schemes virtually no population builds up in the lower level; hence, saturation of the lasing transition occurs at much higher

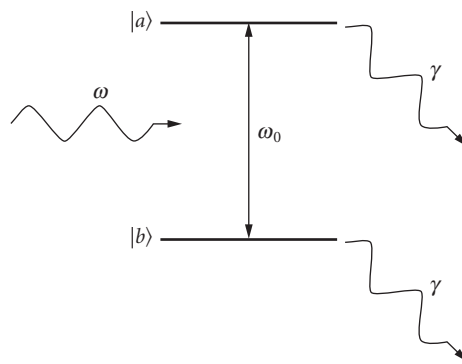


FIGURE 4 Scheme of a two-level atom, with atomic decay permitted, interacting with a single-mode quantized field.

intensities than in lasers with equal decay rates for both levels. These generalizations, however, are easily accounted for (see Ref. 67) and it is not our concern here to provide the most general treatment possible. Instead, we want to focus on the essential quantum features of the laser field and employ the simple two-level model with equal decay rates for both.

The formal averaging in Eq. (46) can be performed most easily if we transform to a particular representation of the density matrix. For our purposes, the photon number representation suffices, although other options are readily available and some of them will be summarized briefly at the end of this section. Taking the n, n' elements of Eq. (46) and using Eqs. (42) and (43), the averaging yields

$$\begin{aligned} [(M-1)\rho]_{nn'} = & -\frac{\chi(n+1+n'+1)+\chi^2(n-n')^2}{1+2\chi(n+1+n'+1)+\chi^2(n-n')^2}\rho_{nn'} \\ & +\frac{2\chi\sqrt{nn'}}{1+2\chi(n+n')+\chi^2(n-n')^2}\rho_{n-1n'-1} \end{aligned} \quad (48)$$

where the notation $\chi = g^2/\gamma^2$ is introduced. Finally, taking into account cavity losses $(\mathcal{L}\rho)_{nn'}$ from Eq. (45), with $n_{\text{th}} = 0$, we obtain the following master equation for our quantum-mechanical laser model:

$$\begin{aligned} \dot{\rho}_{nn'} = & -\frac{\mathcal{N}'_{nn'}\mathcal{A}}{1+\mathcal{N}'_{nn'}\mathcal{B}/\mathcal{A}}\rho_{nn'} + \frac{\sqrt{nn'}\mathcal{A}}{1+\mathcal{N}_{n-1n'-1}\mathcal{B}/\mathcal{A}}\rho_{n-1n'-1} \\ & -\frac{\kappa}{2}(n+n')\rho_{nn'} + \kappa\sqrt{(n+1)(n'+1)}\rho_{n+1n'+1} \end{aligned} \quad (49)$$

Here we introduced the original notations of the Scully-Lamb theory—the linear gain coefficient:

$$\mathcal{A} = 2r\chi \quad (50)$$

the self-saturation coefficient:

$$\mathcal{B} = 4\chi\mathcal{A} \quad (51)$$

and the dimensionless factors:

$$\mathcal{N}' = \frac{1}{2}(n+1+n'+1) + \frac{(n-n')^2\mathcal{B}}{8\mathcal{A}} \quad (52)$$

and

$$\mathcal{N} = \frac{1}{2}(n+1+n'+1) + \frac{(n-n')^2\mathcal{B}}{16\mathcal{A}} \quad (53)$$

Equation (49) is the Scully-Lamb master equation, which is the central equation of the quantum theory of the laser. Along with the notations introduced in Eqs. (50) to (53), it constitutes the main result of this section and serves as the starting point for our treatment of the quantum features of the laser. Among the specific problems we shall consider are the photon statistics, which is the physical information contained in the diagonal elements, and the spectrum, which is the physical information contained in the off-diagonal elements of the field density matrix.

The Micromaser Master Equation The development of the single-atom maser or micro-maser plays a particularly important role in cavity quantum electrodynamics because it realizes one of the most fundamental models, the Jaynes-Cummings hamiltonian. The experimental situation⁸⁶ is very close to the idealized case of a single two-level atom interacting with a single-mode quantized field, as previously discussed, and allows a detailed study of fundamental quantum effects in the atom-field interaction.

In the micromaser, a stream of two-level atoms is injected into a superconducting microwave cavity of very high quality Q . The injection rate r is low enough to ensure that at most only one atom is present inside the cavity at any given time and that most of the time the cavity is empty. The decay time of the high- Q cavity field is very long compared to both the interaction time τ , which is set by the transit time of atoms through the cavity, and the inverse of the single-photon Rabi frequency g^{-1} . In typical experimental situations, however, $g\tau = 1$. Therefore, a field is built up in the cavity provided the interval between atomic injections does not significantly exceed the cavity decay time. Sustained oscillation is possible with less than one atom on the average in the cavity.

In addition to the progress in constructing superconducting cavities, advances in the selective preparation of highly excited hydrogenlike atomic states, called *Rydberg states*, have made possible the realization of the micromaser. In Rydberg atoms the probability of induced transitions between adjacent states is very large, and the atoms may undergo several Rabi cycles—that is, several periodic energy exchanges between the atom and the cavity field may take place in the high- Q cavity. The lifetime of Rydberg states for spontaneous emission decay is also very long, and atomic decay can be neglected during the transit time in the cavity.

Here we set out to derive a master equation for the micromaser. For this, we consider a single-mode resonator into which two-level atoms are injected in their upper states. Due to the different time scales, the effect of cavity damping can be neglected during the interaction. Then the effect of a single atom on the field density operator, injected at t_i and interacting with the field for a time τ , is given by Eq. (44) with t_0 replaced by t_i . If several atoms are injected during a time interval Δt which is still short on the time scale governed by the cavity decay time κ^{-1} but long on the time scale of the interaction time τ , then the cumulative effect on the field is simply the sum of changes

$$\Delta \varrho = \sum_{t \leq t_i \leq t + \Delta t} (M(\tau) - 1) \varrho(t_i) = r \int_t^{t + \Delta t} (M(\tau) - 1) \varrho(t_i) dt_i \quad (54)$$

where in the last step we turned the sum into an integral by using the injection rate r as the number of excited-state atoms entering the cavity per unit time.

At this point we introduce some of the most important approximations of laser physics, or reservoir theory in general—the so-called Markov or adiabatic approximation and coarse time graining. These approximations are based on the existence of three very different time scales in the problem. First, the interaction time τ for individual atoms is of the order of the inverse single-photon Rabi frequency g^{-1} , and is the shortest of all. In fact, on the time scale set by the other relevant parameters, it appears as a delta-function-like kick to the state of the field with the kick operator K of Eq. (44). The second time scale is set by r^{-1} , the average time separation between atomic injections. It is supposed to be long compared to τ but short compared to the cavity-damping time κ^{-1} . Thus, we have the following hierarchy of timescales: $\tau \ll r^{-1} \ll \kappa^{-1}$. When we turned the sum in Eq. (54) into an integral we already tacitly assumed that there is a time scale on which the injection appears to be quasicontinuous. We now see that it is the time scale set by κ^{-1} . This is the time scale that governs the time evolution of the cavity field. In the evaluation of the integral in Eq. (54) we assume that Δt is an intermediate time interval such that $r^{-1} < \Delta t < \kappa^{-1}$. Then during this interval the state of the field does not change appreciably, and we can replace $\varrho(t_i)$ on the right-hand side of Eq. (54) by $\varrho(t)$. This is the essential step in transforming the integral equation into a differential equation, and it constitutes what is called the *Markov approximation*. It is also called the *adiabatic approximation* since the field changes very slowly (adiabatically) on the time scale set by the atoms. As a result, $\varrho(t)$ can now be taken out of the integral, and the integration in Eq. (54) yields $\Delta \varrho = r \Delta t (M - 1) \varrho(t)$. Dividing both sides by Δt , we obtain

$$\frac{\Delta \varrho}{\Delta t} = r(M(\tau) - 1) \varrho \quad (55)$$

The left-hand side is not a true derivative; it only appears to be one on the time scale of the cavity decay time. However, if we are interested in the large-scale dynamics of the field, we can still regard it as a good approximation to a time derivative. It is called the *coarse-grained derivative* and the Eq. (55)

now properly describes the time rate of change of the field due to the interaction with an ensemble of active atoms, the gain reservoir.

Equation (55) gives the time rate of change of the field density operator due to the gain reservoir $(d\varrho/dt)_{\text{gain}}$. To this, we add the time rate of change of the density operator due to the cavity losses by hand. For the parameters of the Garching micromaser experiment, $T = 0.5$ K and $\omega_0/2\pi = 21.5$ GHz, yielding $n_{\text{th}} = 0.15$. The thermal background cannot be neglected since, as we shall see, the steady-state field contains but a few photons. The complete master equation for the micromaser, including both gain and loss, is then simply the sum of Eqs. (55) and (45):

$$\frac{d\varrho}{dt} = \left(\frac{d\varrho}{dt}\right)_{\text{gain}} + \left(\frac{d\varrho}{dt}\right)_{\text{loss}} = r(M(\tau)-1)\varrho + \mathcal{L}\varrho \quad (56)$$

For later purposes, we also give the master equation in matrix representation:

$$\begin{aligned} \frac{d\varrho_{nn'}}{dt} = & r[(A_{nn'}(\tau)-1)\varrho_{nn'} + B_{n-1n'-1}(\tau)\varrho_{n-1n'-1}] \\ & - \frac{\kappa}{2}n_{\text{th}}[(n+n'+2)\varrho_{nn'} - 2\sqrt{nn'}\varrho_{n-1n'-1}] \\ & - \frac{\kappa}{2}(n_{\text{th}}+1)[(n+n')\varrho_{nn'} - 2\sqrt{(n+1)(n'+1)}\varrho_{n+1n'+1}] \end{aligned} \quad (57)$$

Here $A_{nn'}(\tau)$ and $B_{nn'}(\tau)$ are given by Eqs. (42) and (43). Equation (57) is identical to the one obtained by more standard methods⁸⁷ and employed widely in the context of micromasers. It forms the basis for most studies (with a few notable exceptions, as discussed at the end of this section) on the quantum statistical properties of the micromaser and, naturally, it will be our starting point as well.

Physics on the Main Diagonal: Photon Statistics

To begin to bring to light the physical consequences of the laser and maser master equations, we shall first focus on the diagonal elements of the field density matrix ϱ_{nn} , which give us the photon-number distribution since $\varrho_{nn} = p(n)$ is the probability of finding n photons in the cavity mode. The case of the laser is sufficiently different from that of the micromaser that we shall deal with them separately.

Laser Photon Statistics Taking the diagonal $n = n'$ elements in Eq. (49) and regrouping the terms, we obtain the following equation for the photon-number probabilities:

$$\dot{p}(n) = -\frac{(n+1)\mathcal{A}}{1+(n+1)\mathcal{B}/\mathcal{A}}p(n) + \kappa(n+1)p(n+1) + \frac{n\mathcal{A}}{1+n\mathcal{B}/\mathcal{A}}p(n-1) - \kappa np(n) \quad (58)$$

Here the overdot stands for the time derivative. Note that diagonal elements are coupled only to diagonal elements. This holds quite generally; Eq. (49) describes coupling along the same diagonal only. For example, elements on the first side diagonal are coupled to other elements on the first side diagonal, and so on, and in general only elements with the same difference $n - n'$ are coupled. The quantity $k = n - n'$ corresponds to elements on the k th side diagonal. Elements on different diagonals are not coupled, which greatly simplifies the solution of laser-related problems.

Before we begin the solution of Eq. (58), we want to give a simple intuitive physical picture of the processes it describes in terms of a probability flow diagram, shown in Fig. 5.

The left-hand side is the rate of change of the probability of finding n photons in the cavity. The right-hand side contains the physical processes that contribute to the change. Each process is

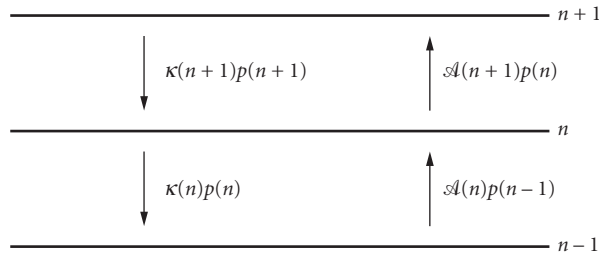


FIGURE 5 Probability flow diagram for the laser.

represented by an arrow in the diagram. The processes are proportional to the probability of the state they are starting from and this will be the starting point of the arrow. The tip of the arrow points to the state the process is leading to. There are two kinds of elementary processes: emission of a photon into the cavity mode (upward arrows) and loss of a photon from the cavity through the output mirror (downward arrows). Furthermore, the processes that start from a given state and end in a different one will decrease the probability of the state they are starting from and increase the probability of the state they are ending at. For example, the first term on the right-hand side describes emission of a photon into the cavity mode provided there are n photons already present before the emission event takes place. Since there will be $n+1$ photons after the emission, this process decreases the probability of finding n photons, hence the minus sign. The total emission rate $\mathcal{A}(n+1)$ has a contribution from stimulated emission, the $\mathcal{A}(n)$ term, and another one from spontaneous emission, the \mathcal{A} term. The third term on the right-hand side corresponds similarly to emission, conditioned on the presence of $n-1$ photons in the cavity initially. After the emission, there will be n photons, hence the plus sign. The second term describes the loss of a photon through the cavity mirror, provided there are n photons initially. After the escape of a photon, there will be $n-1$ photons; therefore, this term decreases $p(n)$. Finally, the last term corresponds similarly to a loss process, with initially $n+1$ photons in the cavity. After the escape of a photon there will be n photons left, so this term increases the probability $p(n)$ of finding n photons in the cavity.

After this brief discussion of the meaning of the individual terms, we now turn our attention to the solution of the equation. Although it is possible to obtain a rather general time-dependent solution to Eq. (58), our main interest here is in the steady-state properties of the field. To obtain the steady-state photon statistics, we replace the time derivative with zero. Note that the right-hand side of the equation is of the form $F(n+1) - F(n)$, where

$$F(n) = \kappa n p(n) - \frac{n\mathcal{A}}{1+n\mathcal{B}/\mathcal{A}} p(n-1) \quad (59)$$

simply meaning that in steady-state $F(n+1) = F(n)$. In other words, $F(n)$ is independent of n and is, therefore, a constant c . Furthermore, the equation $F(n) = c$ has normalizable solution only for $c = 0$. From Eq. (59) we then immediately obtain

$$p(n) = \frac{\mathcal{A}/\kappa}{1+n\mathcal{B}/\mathcal{A}} p(n-1) \quad (60)$$

which is a very simple two-term recurrence relation to determine the photon-number distribution. Before we present the solution, a remark is called for here. The fact that $F(n) = 0$ and $F(n+1) = 0$ hold separately is called the *condition of detailed balance*. As a consequence, we do not need to deal with all four processes affecting $p(n)$. It is sufficient to balance the processes connecting a pair of adjacent levels in Fig. 5, and instead of solving the general three-term recurrence relation resulting from the steady-state version of Eq. (58), it is enough to solve the much simpler two-term recursion, Eq. (60).

It is instructive to investigate the photon statistics in some limiting cases before discussing the general solution. Below threshold the linear approximation holds. Since only very small n states are

populated appreciably, the denominator on the right-hand side of Eq. (60) can be replaced by unity in view of $n\mathcal{B}/\mathcal{A} \ll 1$. Then

$$p(n) = p(0) \left(\frac{\mathcal{A}}{\kappa} \right)^n \quad (61)$$

The normalization condition $\sum_{n=0}^{\infty} p(n) = 1$ determines the constant $p(0)$, yielding $p(0) = (1 - \mathcal{A}/\kappa)$. Finally,

$$p(n) = \left(1 - \frac{\mathcal{A}}{\kappa} \right) \left(\frac{\mathcal{A}}{\kappa} \right)^n \quad (62)$$

Clearly, the condition of existence for this type of solution is $\mathcal{A} < \kappa$. Therefore, $\mathcal{A} = \kappa$ is the threshold condition for the laser. At threshold, the photon statistics change qualitatively and very rapidly in a narrow region of the pumping parameter. It should also be noted that below threshold the distribution function Eq. (62) is essentially of thermal character. If we introduce an effective temperature T defined by

$$\exp\left(-\frac{\hbar\omega_0}{kT}\right) = \frac{\mathcal{A}}{\kappa} \quad (63)$$

we can cast Eq. (62) to the form

$$p(n) = \left[1 - \exp\left(-\frac{\hbar\omega_0}{kT}\right) \right] \exp\left(-\frac{n\hbar\omega_0}{kT}\right) \quad (64)$$

This is just the photon-number distribution of a single mode in thermal equilibrium with a thermal reservoir at temperature T . The inclusion of a finite temperature-loss reservoir to represent cavity losses will not alter this conclusion about the region below threshold.

There is no really good analytical approximation for the region around threshold, although the lowest-order expansion of the denominator in Eq. (60) yields some insight. The solution with this condition is given by

$$p(n) = p(0) \left(\frac{\mathcal{A}}{\kappa} \right)^n \prod_{k=0}^{n-1} \left(1 - \frac{k\mathcal{B}}{\mathcal{A}} \right) \quad (65)$$

This equation clearly breaks down for $n > \mathcal{A}/\mathcal{B} = n_{\max}$, where $p(n)$ becomes negative. The resulting distribution is quite broad, exhibiting a long plateau and a rapid cutoff at n_{\max} . The broad plateau means that many values of n are approximately equally likely; therefore, the intensity fluctuations are large around threshold. The most likely value of $n = n_{\text{opt}}$ can be obtained from the condition $p(n_{\text{opt}} - 1) = p(n_{\text{opt}})$ since $p(n)$ is increasing before $n = n_{\text{opt}}$ and decreasing afterwards. This condition yields $n_{\text{opt}} = (\mathcal{A} - \kappa)/\mathcal{B}$, which is smaller by the factor κ/\mathcal{A} than the value obtained from the full nonlinear equation [cf. Eq. (70) following].

The third region of special interest is the one far above threshold. In this region, $\mathcal{A}/\kappa \gg 1$ and the n values contributing the most to the distribution function are the ones for which $n \gg \mathcal{A}/\mathcal{B}$. We can then neglect 1 in the denominator of Eq. (60), yielding

$$p(n) = \exp\left(-\bar{n} \frac{\bar{n}^n}{n!}\right) \quad (66)$$

with $\bar{n} = \mathcal{A}^2/(\kappa\mathcal{B})$. Thus, the photon statistics far above threshold are poissonian, the same as for a coherent state. This, however, does not mean that far above threshold the laser is in a coherent state. As we shall see later, the off-diagonal elements of the density matrix remain different from those of a coherent state for all regimes of operation.

After developing an intuitive understanding of the three characteristically different regimes of operation, we give the general solution of Eq. (60), valid in all three regimes:

$$p(n) = p(0) \prod_{k=1}^n \frac{(\mathcal{A}/\kappa)}{(1+k\mathcal{B}/\mathcal{A})} \quad (67)$$

The normalization constant $p(0)$ may be expressed in terms of the confluent hypergeometric function

$$p(0) = \left[\sum_{n=0}^{\infty} \frac{\left(\frac{\mathcal{A}}{\mathcal{B}}\right)! \left(\frac{\mathcal{A}^2}{\mathcal{B}\kappa}\right)^n}{\left(n + \frac{\mathcal{A}}{\mathcal{B}}\right)!} \right]^{-1} \\ = \left[F\left(1; \frac{\mathcal{A}}{\mathcal{B}} + 1; \frac{\mathcal{A}^2}{\mathcal{B}\kappa}\right) \right]^{-1} \quad (68)$$

In Fig. 6, the photon-number distribution is displayed for various regimes of operation.

It is interesting to note that $p(n)$ is a product of n factors of the form $(\mathcal{A}/\kappa)/(1+k\mathcal{B}/\mathcal{A})$. This is an increasing function of k as long as the factors are larger than 1 and decreasing afterward. The maximum of the distribution function can be found from the condition

$$\frac{(\mathcal{A}/\kappa)}{(1+n_m\mathcal{B}/\mathcal{A})} = 1 \quad (69)$$

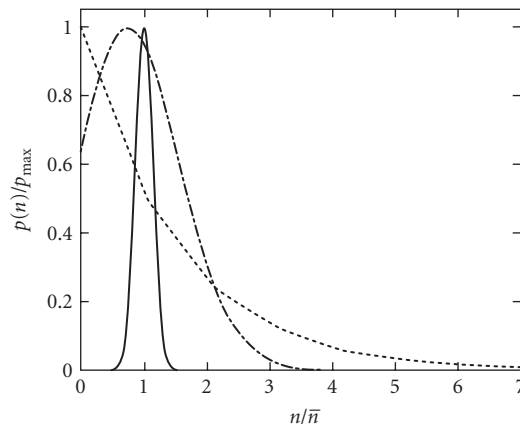


FIGURE 6 Photon statistics of the laser for various regimes of operation. Dotted line: laser 50 percent below threshold ($\mathcal{A}/\kappa = 1/2$), distribution thermal in character, Mandel $Q_M = 1$. Dot-dashed line: laser 10 percent above threshold, ($\mathcal{A}/\kappa = 1.1$), distribution is broad, $Q_M = 5$. Solid line: laser 100 percent above threshold ($\mathcal{A}/\kappa = 2$), distribution super-Poissonian, $Q_M = 1$. Since in the various regimes the actual photon numbers and $p(n)$ values differ by several orders of magnitude, in order to comparatively display the statistics in one figure, we plot $p(n)/p_{\max}$ versus $n/\sqrt{\bar{n}}$. This way the maximum of each curve is 1 and unity on the horizontal axis corresponds to the average photon number. $\mathcal{A}/\mathcal{B} = 50$ is used for all plots, for illustrative purposes only. For more realistic values, the above-threshold distributions are much narrower on this scale.

or

$$n_m = \frac{\mathcal{A}}{\mathcal{B}} \frac{\mathcal{A} - \kappa}{\kappa} \quad (70)$$

Clearly, for $\mathcal{A} < \kappa$ the maximum is at $n = 0$ and the distribution is monotonically decreasing, which is characteristic of a thermal distribution. This is in agreement with the previous findings for the below-threshold region. Around $\mathcal{A} = \kappa$ the distribution quickly changes its character. The factor \mathcal{A}/\mathcal{B} governs the magnitude of the photon number while $(\mathcal{A} - \kappa)/\kappa$ is a measure of how far away from threshold the laser is operating. Typical values for CW gas lasers (the He-Ne laser, for example) are $\mathcal{A}/\mathcal{B} \approx 10^6$ and $\kappa \approx 10^6$ Hz. Around threshold, $\mathcal{A} = \kappa$ and the factors appearing in $p(n)$, given by Eq. (69), are effectively unity for a broad range of n . For example, for $\mathcal{A}/\kappa = 1.001$ (i.e., one-tenth of a percent above threshold), the factors are slightly above 1 for $1 < n < 1000$. So, in the threshold region, the distribution very quickly changes from a thermal one, dominated by the vacuum state, to a broad distribution with large intensity fluctuations. Farther above threshold the distribution becomes more and more peaked around n_m and becomes essentially poissonian for $\mathcal{A}/\kappa > 2$.

It is easy to obtain the mean photon number \bar{n} from Eq. (67):

$$\bar{n} = \sum_{n=0}^{\infty} n p(n) \frac{\mathcal{A}}{\mathcal{B}} = \frac{\mathcal{A} - \kappa}{\kappa} + \frac{\mathcal{A}}{\mathcal{B}} p(0) \quad (71)$$

Above threshold, $p(0) \ll 1$ and the last term becomes quickly negligible. Then \bar{n} coincides with n_m , the maximum of the distribution. We can obtain \bar{n}^2 similarly. The result is

$$\bar{n}^2 = \bar{n} + \frac{\mathcal{A}^2}{\mathcal{B}\kappa} \quad (72)$$

Using Eq. (71), the variance can be expressed as

$$\sigma^2 = \bar{n}^2 - \bar{n} = \bar{n} \frac{\mathcal{A}}{\mathcal{A} - \kappa} \quad (73)$$

From here we see that the variance always exceeds that of a poissonian distribution ($\sigma^2 > \bar{n}$), but it approaches one far above threshold. A good characterization of the photon-number distribution is given by the Mandel Q_M parameter:

$$Q_M = \frac{\bar{n}^2 - \bar{n}^2}{\bar{n}} - 1 \quad (74)$$

For our case, it is given by

$$Q_M = \frac{\kappa}{\mathcal{A} - \kappa} \quad (75)$$

Since $Q_M > 0$ above threshold, the field is superpoissonian. Very far above threshold, when $\mathcal{A} \gg \kappa$, Q_M approaches zero, which is characteristic of a poissonian distribution, again in agreement with our discussion of the far-above-threshold region.

Micromaser Photon Statistics As a first application of the micromaser master equation, Eq. (56), we shall study the steady-state photon statistics arising from it. To this end we take the diagonal $n = n'$ elements, and after regrouping the terms we obtain:

$$\begin{aligned} \dot{p}(n) = & -N_{\text{ex}} \sin^2(g\tau\sqrt{n+1})p(n) + (n_{\text{th}} + 1)(n+1)p(n+1) - n_{\text{th}}(n+1)p(n) \\ & + N_{\text{ex}} \sin^2(g\tau\sqrt{n})p(n-1) - (n_{\text{th}} + 1)np(n) + n_{\text{th}}np(n-1) \end{aligned} \quad (76)$$

Here the overdot stands for derivative with respect to the scaled time $t' = \kappa t$. $N_{\text{ex}} = r/\kappa$ is the number of atoms traversing the cavity during the lifetime of the cavity field, and the diagonal matrix elements of the density operator $p(n) = \rho_{nn}$ are the probabilities of finding n photons in the cavity. The various processes in the right-hand side of this equation are again visualized in Fig. 5. They have the structure $F(n+1) - F(n)$ where $F(n+1)$ corresponds to the processes connecting $p(n+1)$ to $p(n)$. In the steady state the left-hand side is zero and $F(n+1) = F(n)$, yielding $F(n) = \text{constant}$. The only normalizable solution to the photon statistics arises when this constant is zero, $F(n) = 0$. Once again, this is the condition of detailed balance because the probability flows between adjacent levels are separately balanced. More explicitly, this leads to the following recurrence relation for the photon-number probabilities:

$$p(n) = \frac{N_{\text{ex}} \sin^2(g\tau\sqrt{n})/n + n_{\text{th}}}{n_{\text{th}} + 1} p(n-1) \quad (77)$$

The solution to this simple recurrence relation is straightforward:

$$p(n) = p(0) \prod_{i=1}^n \frac{N_{\text{ex}} \sin^2(g\tau\sqrt{i})/i + n_{\text{th}}}{n_{\text{th}} + 1} \quad (78)$$

where $p(0)$ is determined from the normalization condition $\sum_{n=0}^{\infty} p(n) = 1$. The photon-number distribution $p(n)$ versus n can be multip peaked in certain parameter regimes. This can be easily understood on the basis of Fig. 5. The gain processes (upward arrows) balance the loss (downward arrows). Since the gain is an oscillatory function of n , several individual peaks [with the property $p(n+i) = p(n)$] will develop for those values of n where the gain perfectly balances the loss. The resulting mean photon number (first moment of the distribution) and photon-number fluctuations (second moment Q_M) versus the scaled interaction parameter $\theta = g\tau\sqrt{N_{\text{ex}}}$ are displayed in Fig. 7.

The mean photon number is an oscillatory function of the scaled interaction time. The oscillations correspond to subsequent Rabi cycles the atoms are undergoing in the cavity as function of the interaction time. The first threshold occurs at $\theta = 1$; the higher ones occur where θ is approximately an integer multiple of 2π . Around the thresholds the micromaser field is superpoissonian; in the parameter region between the thresholds it is subpoissonian, which is a signature of its nonclassicality.

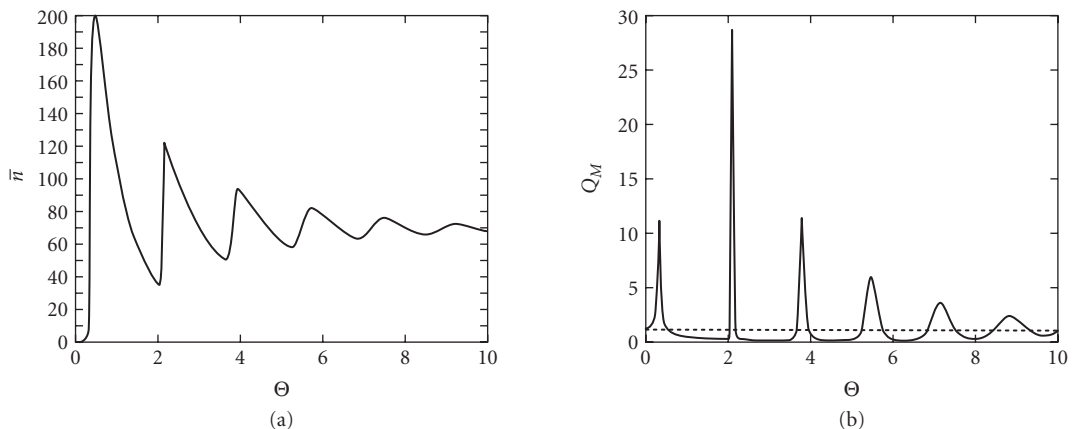


FIGURE 7 (a) Mean photon number and (b) Mandel Q_M parameter versus interaction parameter for the micromaser, for $N_{\text{ex}} = 200$.

Physics off the Main Diagonal: Spectrum and Linewidth

In the subsection on laser photon statistics we have shown that the laser has poissonian photon statistics far above threshold, just as a coherent state would. However, it is erroneous to conclude from this that the laser field is in a coherent state far above threshold. Only the intensity of the laser becomes stabilized due to a delicate balance of the nonlinear gain and loss. Any deviation from the steady-state intensity induces a change that tries to restore the steady-state value; as a result, the intensity is locked to this value. The phase of the field, on the other hand, evolves freely and is not locked to any particular value. In fact, it performs a random walk due to the separate spontaneous emission events. Each such event contributes a small random phase change to the instantaneous phase of the field. The mean of these changes averages to zero, but the mean of the square of these changes remains finite. As a result, the phase undergoes a diffusionlike process and will become uniformly distributed over the 2π interval. Any information contained in the instantaneous phase will be erased in this process. The time scale for the decay of the phase information is given by the rate of the phase diffusion. In the following, we shall determine this characteristic time scale, the so-called phase diffusion constant for the laser and micromaser.

Spectral Properties of the Laser Field The decay of the phase information can be directly read out from the temporal behavior of the two-time correlation function of the field amplitude:

$$g^{(1)}(t_0 + t, t_0) = \frac{\langle a^\dagger(t_0 + t)a(t_0) \rangle}{\langle a^\dagger(t_0)a(t_0) \rangle} \quad (79)$$

With increasing time difference between their amplitudes, the fields become less and less correlated since spontaneous emission randomizes their phases. At steady state, the two-time correlation function [Eq. (79)] depends only on the time difference t and is independent of the choice of the initial time t_0 .

A quantum regression “theorem,”⁸⁸ based on a system-reservoir factorization of the density matrix, was developed to permit the time evolution of the two-time correlation function at steady state to be calculated from the time evolution of the single-time correlation function for a Markov process. The unusual success of this procedure (see Ref. 89 for a comparison between experiment and theory for the phase linewidth, the intensity linewidth, the photo-count distribution, and the spectral moments) requires additional explanation. This was supplied in a proof that regression is valid when the system is markovian.⁹⁰ In the quantum case, the system is only approximately markovian. But this assumption has already been made in all cases for which solutions have been found. Therefore, it is sufficient for us to study the time evolution of the amplitude itself:

$$\langle a^\dagger(t) \rangle = \sum_{n=0}^{\infty} \sqrt{n+1} \rho_{nn+1} \quad (80)$$

At this point it is useful to define a column vector with the components $\rho_n^{(k)} \equiv \rho_{nn+k}$. This way, we simply arrange the elements of the k th diagonal in the form of a vector. For example, elements on the first side diagonal correspond to $k = 1$, and so on. Let us note that the equation of motion for the off-diagonal matrix elements of the density operator can now be written in a simple matrix form

$$\dot{\rho}_n^{(k)} = A_{nn'} \rho_{n'}^{(k)} \quad (81)$$

where summation is implied over repeated indexes and the matrix elements $A_{nn'}$ can be read out from Eq. (49). They are given by

$$A_{nn'} = - \left[\frac{\mathcal{N}'_{nn+k} \mathcal{A}}{1 + \mathcal{N}_{nn+k} \mathcal{B}/\mathcal{A}} + \kappa \left(n + \frac{k}{2} \right) \right] \delta_{nn'} + \frac{\sqrt{n(n+k)} \mathcal{A}}{1 + \mathcal{N}_{n-1+n+k-1} \mathcal{B}/\mathcal{A}} \delta_{nn'+1} + \kappa \sqrt{(n+1)(n+k+1)} \delta_{nn'-1} \quad (82)$$

Clearly, $A_{nn'}$ is a tridiagonal matrix. Due to its linearity, we can look for the solution of Eq. (81) by the simple exponential Ansatz, $\mathbf{e}_n^{(k)}(t) = e^{-\lambda t} \mathbf{e}_n^{(k)}(0)$. With this substitution, Eq. (81) can be written in the form of an eigenvalue equation to determine λ ,

$$\lambda \mathbf{e}_n^{(k)} = - \sum_{n-1}^{n+1} A_{nn'} \mathbf{e}_{n'}^{(k)} \quad (83)$$

We restrict the following treatment to $k = 1$ because that is what we need for the calculation of the $g^{(1)}$ correlation function. Higher-order correlation functions, $g^{(k)}$ with $k > 1$, are related to $\mathbf{e}_n^{(k)}$ with $k > 1$. From the structure of the $-A$ matrix, one can show that all eigenvalues are positive. There is a smallest eigenvalue, which we denote by D . This eigenvalue will dominate the longtime behavior of the field amplitude, as can easily be seen from the following considerations. Let us denote the set of eigenvalues by $\{D, \lambda_j\}$ with $j = 1, 2, 3, \dots$. Then $\mathbf{e}_n^{(1)}(t)$ can be written as

$$\mathbf{e}_n^{(1)}(t) = \mathbf{e}_{n0}^{(1)}(0) \exp(Dt) + \sum_{j=1}^{\infty} \mathbf{e}_{nj}^{(1)}(0) \exp(-\lambda_j t) \quad (84)$$

From this we see that, indeed, the longtime behavior of the off-diagonal elements will be governed by the first term, since the other terms decay faster according to our assumption of D being the smallest positive eigenvalue. Therefore, our task is reduced to the determination of D . In order to obtain an analytical insight, we can proceed as follows. First, let us note that in the longtime limit all elements of the vector $\mathbf{e}_n^{(1)}(t)$ decay the same way—they are proportional to $\exp(-Dt)$. Therefore, the sum of the elements also decays with the same rate, D , in this limit. It is quite easy to obtain an equation of motion for the sum of the elements. Starting from Eq. (81) and using Eq. (82) for the case $k = 1$, we immediately obtain

$$\dot{\mathbf{e}}^{(1)} = - \sum_{n=0}^{\infty} \left[\frac{n+3/2 - \sqrt{(n+1)(n+2)}}{1 + (n+3/2)\mathcal{B}/\mathcal{A}} \mathcal{A} + \kappa(n+1/2 - \sqrt{n(n+1)}) \right] \mathbf{e}_n^{(1)} \quad (85)$$

Here we introduced the notation $\sum_{n=0}^{\infty} \mathbf{e}_n^{(1)} = \mathbf{e}^{(1)}$ and used the fact that $\mathcal{N}'_{nn+1} = n+3/2 + \mathcal{B}/(8\mathcal{A}) \approx n+3/2$ and $\mathcal{N}_{nn+1} = n+3/2 + \mathcal{B}/(16\mathcal{A}) \approx n+3/2$ since $\mathcal{B}/\mathcal{A} \approx 10^{-6}$ and can therefore safely be neglected next to $3/2$. In the longtime limit the time derivative on the left-hand side can simply be replaced by $-D$ due to Eq. (84). It is also plausible to assume that in the same limit those values of n will contribute the most that lie in the vicinity of \bar{n} . Then we can expand the coefficients around the steady-state value of the photon number. This is certainly a good approximation in some region above threshold. The key point is that after the expansion the coefficients of $\mathbf{e}_n^{(1)}$ on the right-hand side become independent of the summation index n and can be factored out from the sum. Then, after the summation, the quantity $\mathbf{e}^{(1)}$ also appears on the right-hand side of the equation:

$$D \mathbf{e}^{(1)} = \frac{\mathcal{A} + \kappa}{8\bar{n}} \mathbf{e}^{(1)} \quad (86)$$

From this we can simply read out the decay rate:

$$D = \frac{\mathcal{A} + \kappa}{8\bar{n}} \quad (87)$$

This quantity, called the *phase diffusion coefficient*, plays a crucial role in determining the transient behavior of the laser as well as its spectral properties. It exhibits the characteristic line narrowing for high intensity, first found by Schawlow and Townes.⁸⁰ The mean amplitude, Eq. (80), can now be written as

$$\langle a^\dagger(t) \rangle = e^{-Dt} \langle a^\dagger(0) \rangle \quad (88)$$

The decay of any initial coherent component of the laser field is governed by the phase diffusion constant, due to the randomization of the initial phase information. The randomization is due to two separate processes, as can be read out from the analytical expression, Eq. (87) of the phase diffusion constant. The part proportional to the spontaneous emission rate \mathcal{A} is due to the random addition of photons to the field via spontaneous emission; the part proportional to the cavity decay rate is due to leaking of vacuum fluctuations into the cavity through the output mirrors. Both processes randomize the phase of the initial field; as a result, the phase performs a random walk with a diffusion rate given by Eq. (87). Ultimately, of course, vacuum fluctuations are also responsible for spontaneous emission.

The phase diffusion constant also determines the linewidth of the spectrum of the laser field. Using the quantum regression theorem, we immediately find that the (nonnormalized) steady-state field correlation function is given by

$$g^{(1)}(t_0+t, t_0) = \langle a^\dagger(t_0+t)a(t_0) \rangle = \bar{n} \exp(i\omega_0 t - Dt) \quad (89)$$

where ω_0 denotes the operating frequency of the laser, as before. The power spectrum is given by the Fourier transform of the field correlation function:

$$S(\omega) = \frac{1}{\pi} \operatorname{Re} \int_0^\infty g^{(1)}(t_0+t, t_0) e^{-i\omega t} dt = \frac{\bar{n}}{\pi} \frac{D}{(\omega - \omega_0)^2 + D^2} \quad (90)$$

This is a lorentzian spectrum centered around the operating frequency, $\omega = \omega_0$. The full width at half-maximum (FWHM) is given by $2D$. Figure 8 depicts the normalized spectrum $S(\omega)/S(\omega_0)$ versus the detuning $\Delta = (\omega - \omega_0)/D$.

It should be emphasized that our method of obtaining the preceding analytical approximations is justified only if the mean photon number is large, the photon-number distribution consists of a single large peak, and cross-coupling between intensity and phase, arising from the nonlinearity of the gain very far above threshold, is negligible. These conditions are met for a laser in some region above threshold. Near the threshold, however, the intensity fluctuations cannot be neglected. From

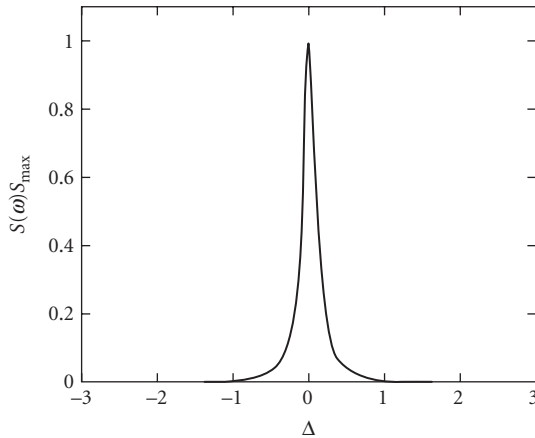


FIGURE 8 Laser spectrum $S(\omega)/S_{\max}$ versus detuning $\Delta = (\omega - \omega_0)/\kappa$, for a laser 100 percent above threshold (note that detuning is in units of bare-cavity linewidth). For our parameters the spectral width D is about 1 percent of the bare-cavity linewidth, an example of the Schawlow-Townes narrowing

numerical studies⁹¹ it was concluded that their contribution to the linewidth is approximately equal to the phase diffusion constant, and the linewidth in this region is about twice what is predicted by Eq. (90); i.e., it is $4D$; and far above threshold (for $\mathcal{A}\kappa > 2$), the linewidth is smaller than the prediction of Eq. (90). These numerical findings are confirmed by a recent analytical approach to the problem.⁹² One of the earliest quantum calculations of the laser linewidth⁶⁹ was based on the quantum theory of noise sources.⁷⁰ The laser linewidth above threshold was established to be due mostly to phase noise,⁶⁹ and it was explained that the factor of 2, too large, was obtained by many previous authors because they had permitted amplitude fluctuations, which were valid near threshold but not far above. These ideas were confirmed by analytic calculations below and above threshold⁹³ and by a numerical solution of the associated Fokker-Planck equation.^{94,95} Moreover, the effects of laser detuning on the linewidth were determined without assuming that the light field decays much slower than the atomic decay rates. It was found that the effective linewidth was a harmonic mean between the field and electronic decay rates, as shown in Eq. (35) of Ref. 69.

The calculations of the phase linewidth done in this section are equivalent to the quasilinear approximation employed in the Langevin noise source procedure in Ref. 93 and shown in Ref. 94 to be valid for dimensionless pump rates $p > 10$. Below threshold, $p < -10$, the components of the electromagnetic field can be treated as uncoupled gaussian variables, leading to the Schawlow-Townes formula. Only the region in the immediate vicinity of threshold requires careful analysis. In that region, it is shown in footnote 10 in Ref. 93 that the coefficients in the Fokker-Planck equation can be expanded in powers of n , retaining only the first nonlinear term. The result then reduces to the rotating-wave Van der Pol oscillator. One advantage of this reduction is that the equation can be scaled in time and in amplitude, leaving the dimensionless pump parameter p as the only remaining parameter. This greatly reduces any subsequent numerical calculations. The ability to retain only one nonlinear term is based only on the requirement that we are dealing with a good oscillator—that is, one whose signal-to-noise ratio is large. Equivalent approximations can be introduced in the density matrix treatment, as well.

The assumption made here that the phase linewidth comes predominantly from the smallest decay eigenvalue was established in Ref. 94, where the actual line shape is shown to be a sum of lorentzians, and the percentage from each lorentzian is calculated. The intensity fluctuations are also expressed as a sum of lorentzians but several modes contribute, not just the lowest. The percentage in each is given in Ref. 94. Part of the reason for this is that the modes approach degeneracy in pairs as one moves above threshold. The lifting of such a degeneracy was shown in Ref. 96 to be associated with a phase transition. Since our system is finite, the phase change occurs gradually and can be observed. The view of lasing as a phase transition will be explored in more detail in Sec. 23.5.

Spectral Properties of the Micromaser Field If we take the $n, n+1$ matrix elements of Eq. (57), then, following the methods of the previous subsection, it is straightforward to derive the diffusion constant. In particular, elements of the first side diagonal can again be arranged in the form of a column vector, and this vector again satisfies an equation of motion similar to Eq. (81), with an appropriate redefinition of the tridiagonal matrix appearing in the equation:

$$\begin{aligned}
 A_{nn'} = & -\{r[1 - \cos(g\tau\sqrt{n+1})\cos(g\tau\sqrt{n+2})] + \kappa(2n_{\text{th}}(n+1) + n+1/2)\}\delta_{nn'} \\
 & + (r\sin(g\tau\sqrt{n})\sin(g\tau\sqrt{N+1}) + \kappa n_{\text{th}}\sqrt{n(n+1)})\delta_{nn'+1} \\
 & + \kappa(n_{\text{th}}+1)\sqrt{(n+1)(n+2)}\delta_{nn'-1}
 \end{aligned} \tag{91}$$

The exponential Ansatz for the decay of the column vector again turns the equation of motion into an eigenvalue equation for the matrix A and the longtime behavior will again be governed by the smallest eigenvalue, which we denote by D . Summing over all elements of the resulting eigenvalue

equation and replacing the coefficients on the right-hand side by their longtime value—that is, expanding the coefficients around $n = \bar{n}$ —will finally yield

$$D = 2r \sin^2 \left(\frac{g\tau}{4\sqrt{\bar{n}}} \right) + \kappa \frac{(2n_{\text{th}} + 1)}{8\bar{n}} \quad (92)$$

for the phase diffusion constant of the micromaser. Here \bar{n} is the mean photon number of the (single-peaked) distribution. This expression was first derived by Scully et al.⁹⁷ It leads to a lorentzian spectrum similar to Eq. (90) but with D appropriately replaced by that of the micromaser. It should be noted, however, that the analytical formula has a more restricted validity than in the case of the laser. Namely, in the case of the micromaser, the photon-number distribution can be multi-peaked, and the simple expansion around the \bar{n} value corresponding to a single dominating peak may not hold. For this more general case several numerical approaches have been developed (see Ref. 98). When Eq. (92) is valid it coincides with the results of the numerical calculations.

For small values of the argument in the sine function, Eq. (92) can be cast to a form which is formally identical to the usual laser phase diffusion constant⁹⁹[cf. Eq. (87)]:

$$\begin{aligned} D &= r \frac{g^2 \tau^2}{8\bar{n}} + \kappa \frac{2n_{\text{th}} + 1}{8\bar{n}} \\ &= \frac{\mathcal{A} + \kappa(2n_{\text{th}} + 1)}{8\bar{n}} \end{aligned} \quad (93)$$

Here we introduced the small signal gain $\mathcal{A} \equiv rg^2 \tau^2$, in analogy to the laser gain. However, for other values of the argument the first term in Eq. (92) may dominate, and the phase diffusion constant can exceed the bare-cavity linewidth, which is a unique quantum feature of the micromaser and makes it distinctly different from the classical Schawlow-Townes-type behavior.

Without going into the specifics, we just mention a few other aspects of the quantum theory of the micromaser. As we have just seen a steady state is reached due to the equilibrium between the gain and loss processes. In some cases, however, a steady state can be reached even in a lossless cavity. This happens if the probability flow in Fig. 5 due to the gain process is interrupted for some value of the interaction parameter. The upward flow is interrupted when

$$g\tau \sqrt{n_q + 1} = q\pi \quad (94)$$

and the downward flow is interrupted when

$$g\tau \sqrt{n_q} = q\pi \quad (95)$$

with $q = 1, 2, \dots$ in both cases. They are called the *upward* and *downward trapping state*, respectively.¹⁰⁰ In such cases the state of the field is a number state. The signature of the number state is a large maximum in the linewidth at the corresponding interaction parameter. It can easily be understood qualitatively: since the state of the field is a number state, it cannot have any phase information. Therefore, the phase randomizes on a very rapid time scale; in other words, the phase correlations decay very rapidly, and a large phase diffusion constant ensues. Further, if the atoms are injected in a coherent superposition of their upper and lower levels in the cavity then, under certain conditions, the so-called tangent and cotangent states of the field may develop.¹⁰¹ Finally, it should be mentioned that the master equations Eqs. (49) and (56) hold only for the case when the time interval between consecutive atomic injections is completely random. Other arrival times statistics, including the case of regular injections, have been investigated by a number of authors.¹⁰² A closely related area of recent theoretical studies pertains to the detection of the statistical properties of the (experimentally inaccessible) intracavity field via monitoring the statistics of the outgoing atoms for poissonian¹⁰³ as well as nonpoissonian¹⁰⁴ pumping. For regular (subpoissonian) pumping, transient oscillations in field correlation function and a corresponding multi-peaked spectrum were predicted.¹⁰⁵

Most of the predictions of this theory have been confirmed by experiments. For example, trapping states have been observed in recent experiments by the Garching group.¹⁰⁶ Without providing an exhaustive list, we just refer the reader to recent progress in the experimental department.¹⁰⁷

Other Approaches

So far we have considered laser theory based on a density-operator approach. An equivalent approach to a laser theory can be formulated using a Heisenberg-Langevin approach. In this approach explicit equations of motion are derived for the field operator.

The quantum-noise operator formalism was presented in essentially its final form by Lax at the 1965 meeting on the *Physics of Quantum Electronics*.⁶⁹ There, it was applied to the laser, calculating, the laser linewidth in its most general form. For the general theory, see the 1966 Brandeis lecture notes of Lax.⁷¹

In the present formulation, our goals are more specific—namely, the Heisenberg-picture quantum theory of the laser. To that end, we will here give a quantum-noise treatment along the lines of that followed in the previous subsection—namely, a single laser mode damped at a rate κ by a (dissipative) reservoir and driven by atoms injected into the laser cavity at random times t_j .

First we discuss the simple example of damping of the field by a reservoir and derive the quantum Langevin equation for the field operator. We then discuss the gain noise in a laser and derive the laser linewidth.

Damping of Field by Reservoir We consider the damping of a single-mode field of frequency ν interacting with a reservoir consisting of simple harmonic oscillators. The system describes, for example, the damping of a single-mode field inside a cavity with lossy mirrors. The reservoir in this case consists of a large number of phononlike modes in the mirror.

The field is described by the creation and destruction operators a^\dagger and a , whereas the harmonic oscillators of frequency $\nu_k = ck$ are described by the operators b_k^\dagger and b_k . The field-reservoir system evolves in time under the influence of the total hamiltonian:

$$\mathcal{H} = \hbar\omega_0\left(a^\dagger a + \frac{1}{2}\right) + \hbar\sum_k \nu_k \left(b_k^\dagger b_k + \frac{1}{2}\right) + \hbar\sum_k g_k (ab_k^\dagger + b_k a^\dagger) \quad (96)$$

Here g_k are the coupling constants and we have made the rotating-wave approximation. We are interested in the evolution of the field operator a . The Heisenberg equations of motion for the field and reservoir are

$$\dot{a}(t) = -i\omega_0 a(t) - i\sum_k g_k b_k(t) \quad (97)$$

$$\dot{b}_k(t) = -i\nu_k b_k(t) - ig_k a(t) \quad (98)$$

The equation for b_k can be formally solved, and the resulting expression is substituted in Eq. (97). In the Weisskopf-Wigner approximation, the annihilation operator in the interaction picture $a = a(t) \exp[i\omega_0(t-t_0)]$ satisfies a Langevin equation

$$\dot{a} = -\frac{\kappa}{2}a + F(t) \quad (99)$$

where

$$F(t) = -i\sum_k g_k b_k(0) \exp[-i(\nu_k - \omega_0)(t-t_0)] \quad (100)$$

is a noise operator. Equation (99) clearly indicates that the damping of the field (represented by the term $-\kappa a/2$) is accompanied by noise.

For the damping of the single-mode field inside a cavity via transmission losses, the damping constant κ is related to the quality factor Q of the cavity via $\kappa = \omega_0/Q$.

Atomic (Gain) Noise and Laser Linewidth As discussed earlier, the natural linewidth of the laser arises due to spontaneous emission by the atoms. In the density-operator approach, a fully nonlinear treatment was followed. Here, we present a simple linear analysis to calculate the laser linewidth in the Heisenberg-Langevin approach. We assume that the atoms are long lived, and that they interact with the cavity field for a time τ . This treatment allows us to include the memory effects inside a laser, and is one of the simplest examples of a nonmarkovian process.^{108,109}

We start with the hamiltonian describing the atom-field interaction:

$$\mathcal{H} = \mathcal{H}_F + \mathcal{H}_{\text{atom}} + \hbar g \sum_i \{ \sigma^i a^\dagger N(t_i, t, \tau) + Hc \} \quad (101)$$

where \mathcal{H}_F and $\mathcal{H}_{\text{atom}}$ describe the field and atoms, respectively; g is the atom-field coupling constant; and σ_i is the lowering operator for the i th atom; Hc is the injunction to add the Hermitian conjugate. The operators a and a^\dagger represent the annihilation and creation operators, and $N(t_i, t, \tau)$ is a notch function which has the value

$$N(t_i, t, \tau) = \begin{cases} 1 & \text{for } t_i \leq t < \tau \\ 0 & \text{otherwise} \end{cases} \quad (102)$$

Using this hamiltonian, we write the equations for the atom-field operators in the interaction picture as

$$\dot{a} = -ig \sum_i \sigma^i N(t_i, t, \tau) - \frac{1}{2} \kappa a(t) + F_\kappa(t) \quad (103)$$

$$\dot{\sigma}^i = ig N(t_i, t, \tau) \sigma_z^i a(t)$$

where the effects of cavity damping are determined by the cavity decay rate κ and the associated Langevin noise source F_κ . Integrating the equation for the atom operator and substituting it into that for the field operator, we obtain

$$\dot{a}(t) = \int_{-\infty}^t dt' \alpha(t, t') a(t') - \frac{1}{2} \gamma a + F_\alpha(t) + F_\kappa(t) \quad (104)$$

where

$$\alpha(t, t') = g^2 \sum_i N(t_i, t, \tau) N(t_i, t', \tau) \sigma_z^i(t') \quad (105)$$

$$F_\alpha(t) = -ig \sum_i N(t_i, t, \tau) \sigma^i(t_i) \quad (106)$$

Here, the noise operator Eq. (106) may be seen to have the moments

$$\langle F_\alpha(t) \rangle = 0 \quad (107)$$

$$\langle F_\alpha^\dagger(t) F_\alpha(t') \rangle = g^2 \sum_{ij} N(t_i, t, \tau) N(t_j, t', \tau) \langle \sigma^{\dagger i}(t_i) \sigma^j(t_j) \rangle \quad (108)$$

Because we are injecting our lasing atoms in the upper state, the atomic average is given by $\langle \sigma^{\dagger i}(t_i) \sigma^j(t_j) \rangle = \delta_{ij}$. After replacing the sum upon i in Eq. (108) by an integration over injection times t_j , we find

$$\langle F_\alpha^\dagger(t) F_\alpha(t') \rangle = r g^2 \{ N(t' - \tau, t, \tau) [t - (t' - \tau)] - N(t', t, \tau) [t - (t' + \tau)] \} \quad (109)$$

where r is the atomic injection rate. The phase variance can then be calculated through the noise operator product:

$$\langle \phi^2(t) \rangle = -\frac{1}{2\pi} \int_0^t dt' \int_0^{t'} dt'' \langle F^\dagger(t') F(t'') \exp\{i[\phi(t') - \phi(t'')]\} \rangle \quad (110)$$

On insertion of Eq. (109) into Eq. (110), the expression for the generalized maser phase diffusion noise $\langle \phi^2(t) \rangle$ is found to be

$$\langle \phi^2(t) \rangle = \left(\frac{\mathcal{A}}{2\bar{n}} \right) \left[\left(\frac{t^2}{\tau} - \frac{t^3}{3\tau^2} \right) \theta(\tau - t) + \left(t - \frac{\tau}{3} \right) \theta(t - \tau) \right] \quad (111)$$

Here $\mathcal{A} = rg^2\tau^2$ is the small-signal gain of the maser [cf. Eq. (93), with $n_{th} = 0$ and using that in steady state $\mathcal{A} = \kappa$]. In the case involving atoms which are injected at random times t_i but which decay via spontaneous emission to far-removed ground states at a rate γ , a similar but more complicated analysis can be carried out. The result in this case is given by

$$\langle \phi^2(t) \rangle = \left(\frac{\mathcal{A}}{2\bar{n}} \right) [t + \gamma^{-1}(e^{-\gamma} - 1)] \quad (112)$$

Here $\mathcal{A} = 2rg^2/\gamma^2$ is the small-signal gain of the laser [cf. Eq. (50)]. In both of the preceding cases, we find that for times $t = t_m$ small compared to the atomic lifetime, the phase diffusion is quadratic in the measurement time t_m ; that is, we now have a phase error which goes as

$$\Delta\phi^2 = \left(\frac{\mathcal{A}t_m}{2\bar{n}} \right) \left(\frac{\gamma t_m}{2} \right) \quad (113)$$

Therefore, we see that the quantum noise due to spontaneous emission is reduced from the Schawlow-Townes linewidth $2D = \mathcal{A}/2\bar{n}$ by the factor $\gamma t_m/2$, which can be a significant reduction for short measurement times. For times long compared to the atomic lifetime, however, the Schawlow-Townes result is obtained from both Eqs. (111) and (112) as expected.

23.5 THE LASER PHASE-TRANSITION ANALOGY

Considerations involving the analogies between phase transitions in ferromagnets, superfluids, and superconductors have emphasized the similarities between these systems near their critical temperatures.¹¹⁰

A natural comparison can be made between second-order phase transitions, such as the order-disorder transitions of ferromagnetic and ferroelectric materials or the vapor-liquid transition of a pure fluid, and the laser threshold. As we have discussed in Sec. 23.4, the state of a laser changes abruptly upon passing through the threshold point. This point is characterized by a threshold population inversion.

The physical basis for this similarity becomes evident when it is recalled that the usual treatments of laser behavior are self-consistent theories. In the laser analysis we assume that each atom evolves in a radiation field due to all the other atoms, and then calculate the field produced by many such evolving atoms in a self-consistent fashion. In this way the laser problem is similar to that of a ferromagnet, in which each spin sees a mean magnetic field due to all the other spins and aligns itself accordingly, thus contributing to the average magnetic field.

Following this point of view, we can discuss the laser theory using the language of second-order phase transitions.

The density matrix of the laser field obeys Eq. (49). The time dependence of the expectation value \bar{E} of the electric field operator $E=(a+a^\dagger)$ is there given by the following equation:

$$\dot{\bar{E}} = \frac{1}{2}(\mathcal{A}-\kappa)\bar{E} - \frac{\mathcal{B}}{2}\bar{E}^3 \quad (114)$$

Here we have assumed that the laser is operating close to threshold ($\mathcal{B}\bar{n}/\mathcal{A} \ll 1$) so that we retain only the terms proportional to \mathcal{B} . In addition, we assume $\bar{E} \gg 1$. We can then replace \bar{E}^3 by \bar{E}^3 and Eq. (114) becomes the well-known result of Lamb's semiclassical theory. The steady-state properties of the laser oscillator are described by the following equation of state:

$$(\mathcal{A}-\kappa)\bar{E} - \mathcal{B}\bar{E}^3 = 0 \quad (115)$$

The threshold condition is given by $\mathcal{A}=\kappa$ as before. Upon putting $\mathcal{A}=a\sigma$, $\mathcal{B}=b\sigma$, and $\kappa=a\sigma_t$ where σ_t is the threshold population inversion, the steady-state solution of Eq. (115) is

$$\begin{aligned} \bar{E} &= 0 && \text{if } \sigma - \sigma_t < 0 \text{ (below threshold)} \\ \bar{E} &= \left[\frac{a}{b} \left(\frac{\sigma - \sigma_t}{\sigma} \right)^{1/2} \right] && \text{if } \sigma - \sigma_t > 0 \text{ (above threshold)} \end{aligned} \quad (116)$$

Equation (116) is formally identical to the equation for a ferromagnet in the Weiss mean-field theory. The electric field E corresponds to the static magnetization M , which is the order parameter in the ferromagnetic transition. The quadratic polarization $P=(\mathcal{A}\bar{E}-\mathcal{B}\bar{E}^3)/2$ in Eq. (115) corresponds to the magnetic field H generated by a magnetization M , and the term $\kappa\bar{E}/2$ corresponds to a local magnetic field which is assumed proportional to M in the mean-field theory. Furthermore, the steady-state points depend $\sigma - \sigma_t$ in the same way that M in the ferromagnetic case depends on $T - T_c$, where T_c is the critical temperature. Therefore, σ and σ_t correspond to T and T_c , respectively. The similarity between these two systems is summarized in Table 1 and illustrated in Fig. 9.

We recall that the probability density $P(M)$ for a ferromagnetic system with magnetization M near a phase transition is given by, in thermal equilibrium,

$$P(M) = N'' \exp\left(-\frac{F(M)}{k_b T}\right) \quad (117)$$

where

$$F(M) = \frac{1}{2}c(T - T_c)M^2 + \frac{1}{4}dT M^4 \quad (118)$$

is the free energy. In the corresponding laser analysis, the probability density for the electromagnetic field $P(E)$ is derived in the form

$$P(E) = N' \exp\left(-\frac{G(E)}{k_b \sigma}\right) \quad (119)$$

For this purpose we transform the laser equation for the density matrix for the field [Eq. (49)] into an equivalent equation in terms of the $P(\alpha, \alpha^*)$ representation defined by

$$\rho = \int d^2\alpha P(\alpha, \alpha^*) |\alpha\rangle\langle\alpha| \quad (120)$$

where $|\alpha\rangle$ is an eigenstate of the annihilation operator a with eigenvalue α . The P representation allows us to evaluate any normally ordered correlation function of the field operators using the

TABLE 1 Summary of Comparison between the Laser and a Ferromagnetic System Treated in a Mean-Field Approximation

Parameter	Ferromagnet	Laser
Order parameter	Magnetization M	Electric field strength E
Reservoir variable	Temperature T	Population inversion σ Threshold inversion σ_t
Coexistence curve*	$M = \Theta(T_c - T) \left[\frac{c}{d} \frac{T - T_c}{T} \right]^{1/2}$	$E = \Theta(\sigma - \sigma_t) \left[\frac{a}{b} \frac{\sigma - \sigma_t}{\sigma} \right]^{1/2}$
Symmetry breaking mechanism	External field H	Injected signal S
Critical isotherm†	$M = \left[\frac{H}{dT_c} \right]^{1/2}$	$E = \left[\frac{2S}{b\sigma_t} \right]^{1/2}$
Zero field susceptibility*	$X \equiv \left(\frac{\partial M}{\partial H} \right) \Big _{H=0}$ $= \Theta(T_c - T) [2c(T_c - T)]^{-1}$ $+ \Theta(T - T_c) [c(T - T_c)]^{-1}$	$\xi \equiv \left(\frac{\partial E}{\partial S} \right) \Big _{S=0}$ $= \Theta(\sigma_t - \sigma) \left[\frac{a(\sigma_t - \sigma)}{2} \right]^{-1}$ $+ \Theta(\sigma - \sigma_t) [a(\sigma - \sigma_t)]^{-1}$
Thermo-dynamic potential	$F(M) = \frac{1}{2}c(T - T_c)M^2$ $+ \frac{1}{4}dT M^4$ $- HM + F_0$	$G(E) = -\frac{a}{4}(\sigma - \sigma_t)E^2$ $+ \frac{1}{8}b\sigma E^4$ $- SE + G_0$
Statistical distribution	$P(M) = \mathcal{N}'' \exp\left(-\frac{F(M)}{k_B T}\right)$	$P(E) = \mathcal{N}' \exp\left(-\frac{G(E)}{k_B \sigma}\right)$

* $\Theta(\cdot)$ is Heaviside's unit step function.

†Value of order parameter at critical point.

methods of classical statistical mechanics. The quantity $P(\alpha, \alpha^*)$ represents the probability density for finding the electric field corresponding to α .

Near threshold, $P(\alpha, \alpha^*)$ obeys the following Fokker-Planck equation:

$$\begin{aligned} \frac{\partial P}{\partial t} = & -\frac{\partial}{\partial \alpha} \left[\frac{1}{2}(\mathcal{A} - \kappa)\alpha P - \frac{1}{2}\mathcal{B}|\alpha|^2 \alpha P \right] \\ & - \frac{\partial}{\partial \alpha^*} \left[\frac{1}{2}(\mathcal{A} - \kappa)\alpha^* P - \frac{1}{2}\mathcal{B}|\alpha|^2 \alpha^* P \right] + \mathcal{A} \frac{\partial^2 P}{\partial \alpha \partial \alpha^*} \end{aligned} \quad (121)$$

The steady-state solution of this equation is given by

$$P(\alpha, \alpha^*) = \mathcal{N} \exp \left[\frac{(\mathcal{A} - k)|\alpha|^2 - \mathcal{B}|\alpha|^4/2}{2\mathcal{A}} \right] \quad (122)$$

where \mathcal{N} is a normalization constant. The P representation can be rewritten in terms of the variables $x = \text{Re } \alpha$ and $y = \text{Im } \alpha$ as

$$P(x, y) = \mathcal{N} \exp \left[-\frac{G(x, y)}{K\sigma} \right] \quad (123)$$

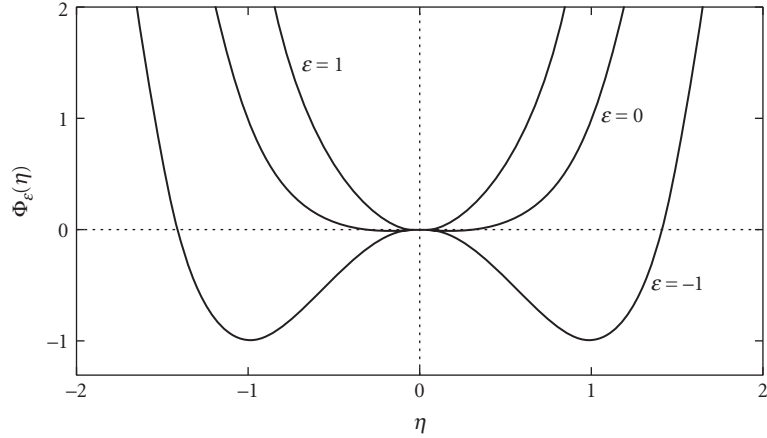


FIGURE 9 Scaled thermodynamical potentials. The $H = 0$ version of $F(M)$ and the $S = 0$ version of $G(E)$ can be expressed in terms of $\Phi_\varepsilon(\eta) = 2\varepsilon\eta^2 + \eta^4$ because

$$\Phi_\varepsilon(\eta) = \begin{cases} \frac{[F(M) - F_0]}{\frac{c^2 T_c}{4d}} & \text{with } M = (c/d)^{1/2} (T/T_c)^{1/4} \eta \\ & \text{and } \varepsilon = (T - T_c)/(TT_c)^{1/2} \\ \frac{[G(E) - G_0]}{\frac{a^2 \sigma_t}{8b}} & \text{with } E = (a/b)^{1/2} (\sigma_t/\sigma)^{1/4} \eta \\ & \text{and } \varepsilon = (\sigma_t - \sigma)/(\sigma\sigma_t)^{1/2} \end{cases}$$

are equivalent to the respective entries in Table 1. The plot shows $\Phi_\varepsilon(\eta)$ for $\varepsilon=1$ (ferromagnet above, T_c , laser below threshold), $\varepsilon=0$ (ferromagnet at T_c , laser at threshold), and $\varepsilon=-1$ (ferromagnet below T_c , laser above threshold).

with

$$G(x, y) = -\frac{1}{4}a(\sigma - \sigma_t)(x^2 + y^2) + \frac{1}{8}b\sigma(x^2 + y^2)^2 \quad (124)$$

Here $K = a/4$ is one-fourth of the gain of one atom, $a(\sigma - \sigma_t) = \mathcal{A} - \kappa$, and $b\sigma = \mathcal{B}$.

We can see that the steady-state situation of the laser corresponds to the minimum value of G , i.e., $\partial G/\partial x = \partial G/\partial y = 0$. These solutions are $x = y = 0$ and $|\alpha|^2 = (x^2 + y^2) = a(\sigma - \sigma_t)/b\sigma$. Thus, for $(\sigma - \sigma_t) < 0$, the only allowed solution is $x = y = 0$. However, for $(\sigma - \sigma_t) > 0$, $x = y = 0$ is an unstable solution as the second derivative of G with respect to x and y is positive. This is seen clearly in Fig. 9, where we have plotted G versus $x = E$ for $y = 0$.

We thus see that G behaves in essentially the same way as the free energy of a thermodynamic system.

It should be emphasized that in the thermodynamic treatment of the ferromagnetic order-disorder transition, there are three variables required: (1) magnetization M , (2) external magnetic field H , and (3) temperature T . In order to have a complete analogy, it is important to realize that in addition to the electric-field-magnetization, population inversion-temperature correspondences, there must exist a further correspondence between the external magnetic field and a corresponding symmetry-breaking mechanism in the laser analysis. As shown in Ref. 111 and illustrated in Fig. 10, this symmetry breaking mechanism in the laser problem corresponds to an injected classical signal S . This leads to a skewed effective free energy.

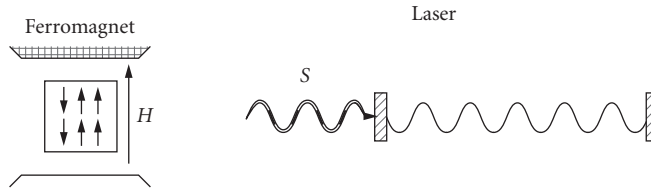


FIGURE 10 Figure depicting the broken symmetry mode of operation for both a ferromagnet and a laser.

An example of how the analogy can provide us with deeper insight is contained in the fact that we are able to *guess* correctly the $P(E)$ for a laser influenced by an injected signal, by analogy with the corresponding magnetic problem in the broken symmetry mode of operation.

More recently we have been turning the tables and using the quantum laser theory to learn about Bose-Einstein condensation (BEC). Recent experiments on BEC in laser-cooled gases,^{112–114} and in He^4 in porous gel,¹¹⁵ have stimulated a wealth of theoretical work^{116–118} on the equilibrium and non-equilibrium properties of confined quantum degenerate gases.^{119–124} Presently the partition function, critical temperature, and other such quantities are of interest for N bosons in a box below T_c . But the canonical ensemble is difficult to use in practice, because the state sums must be restricted to N particles. Indeed, the canonical partition function for a Bose gas of N particles at temperature T has not been so widely studied as one might have thought. To quote Herzog and Olshanii,

To our knowledge there is no simple analytic expression for the canonical partition function in [the case of N bosons in a three-dimensional trap].¹²¹

Furthermore, there are questions of principle concerning the critical temperature and the validity of using phase-transition concepts in a mesoscopic sample having a small number of particles ($N \approx 10^3$). In fact, Uhlenbeck pointed out to Einstein many years ago that BEC rigorously occurs only in the limit of infinite particle number.¹²⁵ Indeed, for a finite number of atoms there is no sharp “critical point” or critical temperature T_c . But the same can be said for the laser threshold. There is a *gradual* transition from disorder to order in both cases. However, as discussed later, even when fluctuations are present, T_c for a Bose gas and the laser threshold inversion are well defined.

Motivated by the preceding, we extend the laser-phase transition analogy to include BEC. We present a new approach to the problem of N bosons in thermal equilibrium below T_c . We emphasize that the present work provides another example¹²⁶ in which steady-state (detailed balance) solutions to nonequilibrium equations of motion provide a supplementary approach to conventional statistical mechanics (e.g., partition-function calculations). The present approach lends itself to different approximations; yielding, among other things, a simple (approximate) analytic expression for the ground-state density matrix for N trapped bosons and the partition function for same.

Thus, we seek a nonequilibrium equation of motion for the ground state of an ideal Bose gas in a three-dimensional harmonic trap coupled to the thermal reservoir, as shown elsewhere.¹²⁷

$$\begin{aligned} \dot{\rho}_{n_0, n_0} = & -K_{n_0} (n_0 + 1) \rho_{n_0, n_0} + K_{n_0-1} n_0 \rho_{n_0-1, n_0-1} \\ & - H_{n_0} n_0 \rho_{n_0, n_0} + H_{n_0+1} (n_0 + 1) \rho_{n_0+1, n_0+1} \end{aligned} \quad (125)$$

The cooling and heating coefficients K_{n_0} and H_{n_0} are given by

$$K_{n_0} = \sum_k 2\pi W_k g_k^2 \langle \eta_k + 1 \rangle \langle n_k \rangle_{n_0} \quad (126)$$

and

$$H_{n_0} = \sum_k 2\pi W_k g_k^2 \langle \eta_k \rangle \langle n_k + 1 \rangle_{n_0} \quad (127)$$

where W_k is the heat-bath density of states, $\langle \eta_k \rangle$ is the average occupation number of the k th heat-bath oscillator, and $\langle \eta_k \rangle_{n_0}$ is the average number of atoms in the k th excited state, given n_0 atoms in the condensate. Here the coefficient K_{n_0} denotes the cooling rate from the excited states to the ground state, and similarly H_{n_0} stands for the heating rate for the ground state.

The heating term is approximately

$$H_{n_0} = \kappa \sum_k \langle \eta(\varepsilon_k) \rangle = \kappa \sum_{\ell, m, n} \left[\exp\left(\frac{\hbar\Omega}{k_B T}\right) (\ell + n + m) - 1 \right]^{-1} \quad (128)$$

In the weak trap limit, this yields

$$H_{n_0} = \kappa \left(\frac{k_B T}{\hbar\Omega} \right)^3 \zeta(3) \quad (129)$$

where $\zeta(3)$ is the Riemann zeta function and Ω is the trap frequency. Likewise, the cooling term in Eq. (125) is governed by the total number of excited state bosons,

$$K_{n_0} = \kappa \sum_k \langle n_k \rangle_{n_0} = \kappa(N - n_0) \quad (130)$$

By writing the equation of motion for $\langle n_0 \rangle$ from Eq. (125), using H_{n_0} in the weak trap limit, and Eq. (130) for K_{n_0} , we find

$$\langle \dot{n}_0 \rangle = \kappa \left[(N+1) \langle n_0 \rangle - \langle (n_0+1)^2 \rangle - \zeta(3) \left(\frac{k_B T}{\hbar\Omega} \right)^3 \langle n_0 \rangle \right] + \kappa(N+1) \quad (131)$$

Noting that near T_c , $\langle n_0 \rangle = N$, we may neglect $\langle (n_0+1)^2 \rangle$ compared to $N \langle n_0 \rangle$, and neglecting the spontaneous emission term $\kappa(N+1)$, Eq. (131) becomes

$$\langle \dot{n}_0 \rangle = \kappa \left[N - \zeta(3) \left(\frac{k_B T}{\hbar\Omega} \right)^3 \right] \langle n_0 \rangle \quad (132)$$

We now define the critical temperature (in analogy with the laser threshold) such that cooling (gain) equals heating (loss) and $\langle \dot{n}_0 \rangle = 0$ at $T = T_c$; this yields

$$T_c = \left(\frac{\hbar\Omega}{k_B} \right) \left[\frac{N}{\zeta(3)} \right]^{1/3} \quad (133)$$

Thus, by defining the critical temperature as that temperature at which the rate of removal of atoms from the ground state equals the rate of addition, we arrive at the usual definition for the critical temperature, even for mesoscopic systems.

23.6 EXOTIC MASERS AND LASERS

Lasing without Inversion

For a long time, it was considered that population inversion was necessary for laser action to take place. Recently, it has been shown both theoretically^{128–130} and experimentally^{131–134} that it is also

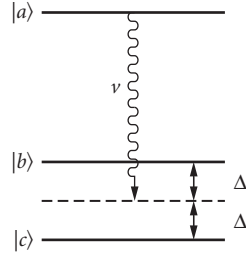


FIGURE 11 Level diagram for lasing without inversion.

possible to achieve lasing without inversion (LWI). In LWI, the essential idea is the cancellation of absorption by atomic coherence and interference.

Consider a system of three-level atoms interacting with a laser field in a cavity. The simple model we will focus on is that of Fig. 11. The atoms have one upper level $|a\rangle$ and two lower levels $|b\rangle$ and $|c\rangle$, with energies $\hbar\omega_a$, $\hbar\omega_b$, and $\hbar\omega_c$, respectively. The cavity field of frequency ν can be detuned from the atomic transition, as shown in the figure. The transitions $|a\rangle \rightarrow |b\rangle$ and $|a\rangle \rightarrow |c\rangle$ are now induced by one classical light field of frequency ν . The transition $|b\rangle \rightarrow |c\rangle$ is dipole forbidden. The atoms are pumped at a rate r_a in a coherent superposition of states

$$\rho(t_i) = \rho_{aa}^{(0)} |a\rangle\langle a| + \rho_{bb}^{(0)} |b\rangle\langle b| + \rho_{cc}^{(0)} |c\rangle\langle c| + \rho_{bc}^{(0)} |b\rangle\langle c| + \rho_{cb}^{(0)} |c\rangle\langle b| \quad (134)$$

Here $\rho_{\alpha\alpha}^{(0)}$ ($\alpha = a, b, c$) are the level populations and $\rho_{\alpha\alpha'}^{(0)}$ ($\alpha \neq \alpha'$) are the atomic coherences. We give a simple argument to show how cancellation of absorption can lead to lasing without inversion in this scheme.

As the levels $|b\rangle$ and $|c\rangle$ are independent, the probability of emission is given by

$$\begin{aligned} P_{\text{emission}} &= P_b + P_c \\ &= (\kappa_{a \rightarrow b} |2\mathcal{E}|^2 + |\kappa_{a \rightarrow c}| 2\mathcal{E}^2) \rho_{aa}^{(0)} \end{aligned} \quad (135)$$

where $\kappa_{a \rightarrow b}$ and $\kappa_{a \rightarrow c}$ are constants which depend on the matrix element between the relevant levels and the coupling of the atom with the field. On the other hand, the absorption probability is given by

$$\begin{aligned} P_{\text{absorption}} &= \kappa |C_b + C_c| 2\mathcal{E}^2 \\ &= \kappa (\rho_{bb}^{(0)} + \rho_{cc}^{(0)} + \rho_{bc}^{(0)} + \rho_{cb}^{(0)}) \mathcal{E}^2 \end{aligned} \quad (136)$$

where c_a and c_b are the probability amplitudes for the states $|b\rangle$ and $|c\rangle$. Therefore, the rate of growth of the laser field amplitude, under appropriate conditions, becomes

$$\dot{\mathcal{E}} = \frac{\mathcal{A}}{2} (\rho_{aa}^{(0)} - \rho_{bb}^{(0)} - \rho_{cc}^{(0)} - \rho_{bc}^{(0)} - \rho_{cb}^{(0)}) \mathcal{E} \quad (137)$$

Here \mathcal{A} is a constant. Thus, if the terms $\rho_{bc}^{(0)}$ and $\rho_{cb}^{(0)}$ cancel $\rho_{bb}^{(0)}$ and $\rho_{cc}^{(0)}$, we have

$$\dot{\mathcal{E}} = \frac{\mathcal{A}}{2} \rho_{aa}^{(0)} \mathcal{E} \quad (138)$$

and we can have lasing even if only a small fraction of atoms is in the excited state $|a\rangle$, that is, even if $\rho_{aa} < (\rho_{bb} + \rho_{cc})$.

Physically, the lack of absorption in the three-level system considered here is a manifestation of quantum coherence phenomena. When an atom makes a transition from the upper level to the two lower levels, the total transition probability is the sum of $|a\rangle \rightarrow |b\rangle$ and $|a\rangle \rightarrow |c\rangle$ probabilities. However, the transition probability from the two lower levels to the single upper level is obtained by squaring the sum of the two probability amplitudes. When there is coherence between the two lower levels, this can lead to interference terms yielding a null in the transition probability corresponding to photon absorption.

Correlated (Spontaneous) Emission Laser

As discussed earlier, the fundamental source of noise in a laser is spontaneous emission. A simple pictorial model for the origin of the laser linewidth envisions it as being due to the random phase diffusion process arising from the addition of spontaneously emitted photons with random phases to the laser field. Here we show that the quantum noise leading to the laser linewidth can be suppressed below the standard Schawlow-Townes limit by preparing the atomic systems in a coherent superposition of states as in the Hanle-effect and quantum-beat experiments. In such coherently prepared atoms, the spontaneous emission is said to be *correlated*. Lasers operating via such a phase-coherent atomic ensemble are known as *correlated emission lasers* (CELs).¹³⁵

An interesting aspect of the CEL is that it is possible to eliminate the spontaneous emission quantum noise in the relative linewidths by correlating the two spontaneous emission noise events.

A number of schemes exist in which quantum noise quenching below the standard limit can be achieved. In two-mode schemes a correlation between the spontaneous emission events in two different modes of the radiation field is established via atomic coherence so that the relative phase between them does not diffuse or fluctuate. In a Hanle laser¹³⁶ and a quantum-beat laser¹³⁷ this is achieved by pumping the atoms coherently such that every spontaneously emitting atom contributes equally to the two modes of the radiation, leading to a reduction and even vanishing of the noise in the phase difference. In a two-photon CEL, a cascade transition involving three-level atoms is coupled to only one mode of the radiation field.¹³⁸ A well-defined coherence between the upper and lower levels $|a\rangle$ and $|c\rangle$ leads to a correlation between the light emitted by an $|a\rangle \rightarrow |b\rangle$ and a subsequent $|b\rangle \rightarrow |c\rangle$ transition.

The quantum theory of quantum-beat or Hanle-effect lasers may be conveniently cast in terms of the equation of motion for the density matrix describing the laser radiation field $\rho(a_1, a_1^\dagger; a_2, a_2^\dagger)$; that is,

$$\dot{\rho} = \sum_{ij} \mathcal{L}_{ij} \rho \quad (139)$$

where the linear gain and cross-coupling Liouville operators are given by

$$\mathcal{L}_{ii} \rho = -\frac{1}{2} [\alpha_{ii} \rho a_i a_i^\dagger + \alpha_{ii}^* a_i a_i^\dagger \rho - (\alpha_{ii} + \alpha_{ii}^*) a_i^\dagger \rho a_i] \quad (140)$$

$$\mathcal{L}_{12} \rho = -\frac{1}{2} [\alpha_{12} \rho a_2 a_1^\dagger + \alpha_{21}^* a_2 a_1^\dagger \rho - (\alpha_{12} + \alpha_{21}^*) a_1^\dagger \rho a_2] e^{i\Phi} \quad (141)$$

$$\mathcal{L}_{21} \rho = -\frac{1}{2} [\alpha_{21} \rho a_1 a_2^\dagger + \alpha_i^* a_1 a_2^\dagger \rho - (\alpha_i + \alpha_i^*) a_2^\dagger \rho a_1] e^{i\Phi} \quad (142)$$

Here α_{ij} are constants that depend on the parameters of the gain medium such as detunings, Rabi frequency of the driving field, and so on. When the coherent mixing of levels $|a\rangle$ and $|b\rangle$ is produced via a microwave signal having frequency ω_0 , the phase angle ϕ is given by $\Phi(t) = (\nu_1 - \nu_2 - \omega_0)t - \phi$ where ϕ is the (microwave determined) atomic-phase difference $\phi_a - \phi_b$. In the case of polarization-induced

coherent mixing, the phase angle is $\Phi(t) = (v_2 - v_1)t - \phi$, where ϕ is again the relative phase between levels $|a\rangle$ and $|b\rangle$ but determined this time by the state of elliptic polarization of the pump light used to excite the atoms.

The Liouville equation (31) for the reduced-density operator for the field can be converted into an equivalent Fokker-Planck equation by introducing coherent state representation for a_1 and a_2 and the P representation $P(\alpha, \alpha^*)$ for ρ . If we define the coherent states as

$$a_i |\alpha_1, \alpha_2\rangle = \alpha_i |\alpha_1, \alpha_2\rangle; i = 1, 2 \quad (143)$$

where α is an arbitrary complex number, and we represent α_i as

$$\alpha_i = \rho_i \exp(i\theta_i); i = 1, 2 \quad (144)$$

then the Fokker-Planck equation in terms of ρ_i, θ_i will contain a term which describes diffusion of relative phase $\theta = \theta_1 - \theta_2$ as

$$\dot{P} = \frac{\partial^2}{\partial \theta^2} [\mathcal{D}(0)P] \quad (145)$$

with

$$\mathcal{D} = \frac{1}{16} \left\{ \left(\frac{\alpha_{11}}{\rho_1^2} + \frac{\alpha_{22}}{\rho_2^2} \right) - \frac{(\alpha_{12} + \alpha_{21}^*)e^{-i\psi}}{\rho_1 \rho_2} \right\} \quad (146)$$

and $\psi = \Phi + \theta_1 - \theta_2$ with θ_i being the phase of the i th field. The diffusion constant \mathcal{D} for the relative phase vanishes for $\psi = 0$, $\rho_1 = \rho_2$, and $\alpha_{11} = \alpha_{22} = \alpha_{12} = \alpha_{21}^*$, thus leading to CEL action.

Free-Electron Laser

A coherent emission of radiation in a free-electron laser (FEL) is due to the bunching of a relativistic electron beam propagating along a periodic magnetic structure. The electrons experience a Lorentz force and thus follow oscillating orbits and radiate. This spontaneous emission coupled with the periodic magnetic structure give rise to a periodic ponderomotive potential. The electrons bunch together and radiate coherently.¹³⁹ The spontaneous emission pattern of a relativistic electron of energy $E = \gamma mc^2$ with $\gamma \gg 1$ is mostly in the forward direction. For a magnetic wiggler of period λ_w , the spectrum in the forward direction is symmetric about the wavelength $\lambda_s \cong \lambda_w / 2\gamma^2$. Thus, a change of the periodicity of the wiggler λ_w can be used to tune the coherent light emitted by the FEL over a very wide range.

Many interesting features of FEL can be understood classically. However, the quantum-statistical properties of radiation emitted by FEL exhibit many interesting features such as squeezing and sub-poissonian statistics.¹⁴⁰⁻¹⁴²

Here we describe a free-electron amplifier in the small-signal noncollective regime. Such an FEL can be described by the one-electron nonrelativistic Bambini-Renieri hamiltonian which refers to a moving frame, where the laser and the wiggler frequencies coincide with $\omega = ck/2$.¹⁴³ In this frame, resonance occurs when the electron is at rest; therefore, the electron can be treated nonrelativistically. The hamiltonian is given by

$$\mathcal{H} = \frac{p^2}{2m} + \hbar \omega A^\dagger A + i\hbar g(A - A^\dagger) \quad (147)$$

with $A = a \exp(ikz)$. Here a is the annihilation operator of the laser field, p and z are the electron's momentum and coordinate with $[z, p] = i\hbar$, $[A, A^\dagger] = 1$, $[p, A] = \hbar k A$, m is the effective mass of the electron, and

$$g = \left(\frac{e^2 B}{mk} \right) \left(\frac{2}{V \epsilon_0 \hbar \omega} \right)^{1/2} \quad (148)$$

with V the quantization volume and B the magnetic-field strength of the wiggler field in the moving frame. In Eq. (147) we have already taken the classical limit of the wiggler field. By transforming to the interaction picture we obtain

$$\mathcal{H}_I = ig\hbar \left\{ \exp \left[-\frac{it(\hbar k^2 + 2kp)}{2m} \right] A^\dagger - Hc \right\} \quad (149)$$

We now consider an initial state made up by an electron with momentum p and the field vacuum, i.e., $|\text{in}\rangle = |\bar{p}, 0\rangle$

$$p|\bar{p}, 0\rangle = \bar{p}|\bar{p}, 0\rangle \quad (150)$$

$$A|\bar{p}, 0\rangle = 0 \quad (151)$$

$$A^\dagger|\bar{p}, 0\rangle = |\bar{p} - \hbar k, 1\rangle \quad (152)$$

The final-state expectation value of any operator $O(A, A^\dagger)$ is then

$$\langle \text{out} | O | \text{out} \rangle = \langle \bar{p}, 0 | \bar{s}^\dagger(T) O S(T) | \bar{p}, 0 \rangle \quad (153)$$

where

$$S(T) = T \exp \left[-\frac{i}{\hbar} \int_{-T/2}^{T/2} dt H_I(t) \right] \quad (154)$$

is the time-evolution operator for the electron-photon state.

The evaluation of Eq. (153) is straightforward in the small-signal limit along the lines given in Ref. 141, and we obtain

$$(\Delta A_1)^2 = \frac{1}{4} - \frac{\hbar k^2}{2m} j \frac{\partial j}{\partial \beta} \quad (155a)$$

$$(\Delta A_2)^2 = \frac{1}{4} + \frac{\hbar k^2}{2m} j \frac{\partial j}{\partial \beta} \quad (155b)$$

$$(\Delta A_1)(\Delta A_2) = \frac{1}{4} \quad (155c)$$

$$\Delta n^2 - \langle n \rangle = -\frac{2\hbar k^2}{m} j^3 \frac{\partial j}{\partial \beta} \quad (155d)$$

where

$$j = \left(\frac{2g}{\beta} \right) \sin \left(\frac{\beta T}{2} \right) \quad (156)$$

$$\beta = \frac{k\bar{p}}{m} \quad (157)$$

In our notation, the gain of the free-electron laser is proportional to $-j\partial j/\partial\beta$. Hence, Eqs. (155a) and (155b) show that, depending on the sign of the gain, either A_1 or A_2 is squeezed while, because of Eq. (155c), minimum uncertainty is maintained. Here, we have defined squeezing with respect to the operator A instead of the annihilation operator a of the radiation field. This must be so because we employ electron-photon states, and the annihilation of a photon always comes up to increasing the momentum of the electron by $\hbar k$. Finally, Eq. (155d) shows that we have subpoissonian, poissonian, or superpoissonian statistics if the electron momentum is below resonance ($\beta < 0$), at resonance ($\beta = 0$), or below resonance ($\beta > 0$), respectively.

Exploiting the Quantized Center-of-Mass Motion of Atoms

In the treatment of the interaction of a two-level atom with photons of a single, dynamically privileged mode by the Jaynes-Cummings model, as discussed in Sec. 23.4, the center-of-mass motion of the atom is regarded as classical. This is a well-justified approximation, since the atom's kinetic energy of typically $\sim 10^{-2}$ eV is many orders of magnitude larger than the interaction energy of typically $\sim 10^{-11}$ eV if the atom belongs to a thermal beam. For ultracold atoms, however, matters can be quite different, and the quantum properties of the center-of-mass motion must be taken into account.

Early studies showed that very slow atoms can be reflected at the entry port of a resonator¹⁴⁴ or trapped inside.¹⁴⁵ The reflection probability is considerable even if the photon lifetime is not short as compared with the relatively long interaction time.¹⁴⁶

Whereas Refs. 144 to 146 deal mainly with the mechanical effects on the center-of-mass motion of the atom, the modifications in the maser action are addressed in Refs. 147 to 150. For thermal atoms, the emission probability displays the usual Rabi oscillations (see Sec. 23.4) as a function of the interaction *time*. For very slow atoms, however, the emission probability is a function of the interaction *length* and exhibits resonances such as the ones observed in the intensity transmitted by a Fabry-Perot resonator. The resonances occur when the resonator length is an integer multiple of half the de Broglie wavelength of the atom inside the cavity.

A detailed calculation¹⁴⁷ shows that the emission probability is 50 percent at a resonance, irrespective of the number of photons that are present initially. Owing to this unusual emission probability, a beam of ultracold atoms can produce unusual photon distributions, such as a shifted thermal distribution. In the trilogy (Refs. 148 to 150) this *microwave amplification by z-motion-induced emission of radiation* (mazer) is studied in great detail.

In order to see the mazer resonances for atoms with a certain velocity spread, the interaction length has to be small. Therefore, micromaser cavities of the usual cylindrical shape, for which the smallest cavity length is given by half the wavelength of the microwaves, cannot be used for this purpose. But cavities of the reentrant type (familiar as components of klystrons) allow for an interaction length that is much smaller than the wavelength. With such a device, an experiment with realistic parameters seems possible.¹⁴⁹ As a potential application, we mention that a working mazer could be used as a velocity filter for atoms.¹⁵¹

23.7 ACKNOWLEDGMENTS

The authors gratefully acknowledge the generous support of the Office of Naval Research over the many years spent on completing the works reviewed here. It is also a pleasure to acknowledge the Max-Planck-Institute for Quantum Optics (Garching, Germany) for providing the excellent working atmosphere and the intellectual stimulus for the most interesting works in this field.

23.8 REFERENCES

1. J. R. Klauder and E. C. G. Sudarshan, *Fundamentals of Quantum Optics* (W. A. Benjamin, New York, 1970).
2. R. Loudon, *The Quantum Theory of Light* (Oxford University Press, New York, 1973).
3. W. H. Louisell, *Quantum Statistical Properties of Radiation* (John Wiley, New York, 1973).
4. H. M. Nussenzveig, *Introduction to Quantum Optics* (Gordon and Breach, New York, 1974).
5. M. Sargent III, M. O. Scully, and W. E. Lamb, Jr., *Laser Physics* (Addison-Wesley, Reading, Mass., 1974).
6. L. Allen and J. H. Eberly, *Optical Resonance and Two-Level Atoms* (John Wiley, New York, 1975).
7. H. Haken, *Light*, Vols. I and II (North-Holland, Amsterdam, 1981).
8. P. L. Knight and L. Allen, *Concepts of Quantum Optics* (Pergamon Press, Oxford, 1983).
9. P. Meystre and M. Sargent III, *Elements of Quantum Optics* (Springer-Verlag, Berlin, 1990).
10. C. W. Gardiner, *Quantum Noise* (Springer-Verlag, Berlin, 1991).
11. C. Cohen-Tannoudji, J. Dupont-Roc, and G. Grynberg, *Atom-Photon Interactions* (John Wiley, New York, 1992).
12. H. Carmichael, *An Open Systems Approach to Quantum Optics* (Springer-Verlag, Berlin, 1993).
13. W. Vogel and D.-G. Welsch, *Lectures on Quantum Optics* (Akademie Verlag, Berlin, 1994).
14. J. Peřina, Z. Hradil, and B. Jurčo, *Quantum Optics and Fundamentals of Physics* (Kluwer, Dordrecht, Netherlands, 1994).
15. D. F. Walls and G. J. Milburn, *Quantum Optics* (Springer-Verlag, Berlin, 1994).
16. L. Mandel and E. Wolf, *Optical Coherence and Quantum Optics* (Cambridge University Press, London, 1995).
17. E. R. Pike and S. Sarkar, *Quantum Theory of Radiation* (Cambridge University Press, London, 1995).
18. M. O. Scully and M. S. Zubairy, *Quantum Optics* (Cambridge University Press, London, 1997).
19. M. Planck, *Verh. Phys. Ges.* **2**:202, 237 (1900).
20. W. Wien, *Ann. Physik* **58**:662 (1896).
21. Lord Rayleigh, *Phil. Mag.* **49**:539 (1900).
22. J. H. Jeans, *Phil. Mag.* **10**:91 (1905).
23. A. H. Compton, *Phys. Rev.* **21**:483 (1923).
24. A. Einstein, *Ann. Physik* **17**:132 (1905).
25. W. Pauli, "Einstein's Contributions to Quantum Theory," in *Albert Einstein: Philosopher-Scientist*, P. A. Schilpp (ed.) (Library of Living Philosophers, Evanston, Ill., 1949).
26. L. de Broglie, *C. R. Acad. Sci. Paris* **177**:517 (1923).
27. L. de Broglie, *These* (Masson et Cie., Paris, 1924).
28. E. Schrödinger, *Ann. Physik* **79**:361 (1926).
29. G. I. Taylor, *Proc. Camb. Phil. Soc.* **15**:114 (1909).
30. P. A. M. Dirac, *The Principles of Quantum Mechanics*, 4th ed. (Oxford University Press, Oxford, 1958).
31. R. J. Glauber, *Am. J. Phys.* **63**:12 (1995).
32. A. Einstein, *Phys. Z.* **10**:185 (1909).
33. N. Bohr, *Phil. Mag.* **26**:1,476, 857 (1913).
34. A. Einstein, *Phys. Z.* **18**:121 (1917).
35. A. Einstein, *Ann. Physik* **17**:549 (1905).
36. N. Bohr, H. A. Kramers, and I. C. Slater, *Phil. Mag.* **47**:785 (1924).
37. W. Bothe and H. Geiger, *Z. Phys.* **26**:44 (1924).
38. S. N. Bose, *Z. Phys.* **26**:178 (1924).
39. S. N. Bose, *Z. Phys.* **27**:384 (1924).
40. A. Einstein, "Quantentheorie des einatomigen Gases," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1924, p. 261.

41. A. Einstein, "Quantentheorie des einatomigen Gases. 2. Abhandlung," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1925, p. 3.
42. A. Einstein, "Quantentheorie des idealen Gases," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1925, p. 18.
43. P. A. M. Dirac, *Proc. Roy. Soc. London* **A114**:710 (1927).
44. J. Schwinger (ed), *Quantum Electrodynamics* (Dover, New York, 1958).
45. E. Fermi, *Rend. Lincei* **10**:72 (1929).
46. R. J. Glauber, *Phys. Rev.* **130**:2529 (1963).
47. R. J. Glauber, *Phys. Rev.* **131**:2766 (1963).
48. R. J. Glauber, *Phys. Rev. Lett.* **10**:84 (1963).
49. V. F. Weisskopf and E. P. Wigner, *Z. Phys.* **63**:47 (1930).
50. V. F. Weisskopf, *Phys. Rev.* **56**:72 (1939).
51. W. E. Lamb, Jr., and R. C. Retherford, *Phys. Rev.* **72**:241 (1947).
52. J. Schwinger, *Phys. Rev.* **73**:416 (1948).
53. S. Pasternack, *Phys. Rev.* **54**:1113 (1938).
54. S. Millman and P. Kusch, *Phys. Rev.* **57**:438 (1940).
55. H. Hanbury-Brown and R. Q. Twiss, *Phil. Mag.* **45**:663 (1954).
56. H. Hanbury-Brown and R. Q. Twiss, *Nature* (London) **178**:1046 (1956).
57. H. Hanbury-Brown and R. Q. Twiss, *Proc. Roy. Soc.* **A242**:300 (1957).
58. R. G. Newton, *Scattering Theory of Waves and Particles* (McGraw-Hill, New York, 1966).
59. C. K. Hong, Z. Y. Ou, and L. Mandel, *Phys. Rev. Lett.* **59**:2044 (1987).
60. P. G. Kwiat, K. Mattle, H. Weinfurter, and A. Zeilinger, *Phys. Rev. Lett.* **75**:4337 (1995).
61. P. G. Kwiat, E. Waks, A. G. White, I. Appelbaum, and P. H. Eberhard, *Phys. Rev. A* **60**:773 (1999).
62. D. Bouwmeester, J.-W. Pan, K. Mattle, M. Eibl, H. Weinfurter, and A. Zeilinger, *Nature* (London) **390**:575 (1997).
63. D. Boschi, S. Branca, F. De Martini, L. Hardy, and S. Popescu, *Phys. Rev. Lett.* **80**:1121 (1998).
64. C. H. Bennett, G. Brassard, C. Crepeau, R. Josza, A. Peres, and W. Wootters, *Phys. Rev. Lett.* **70**:1895 (1993).
65. K. Mattle, H. Weinfurter, P. G. Kwiat, and A. Zeilinger, *Phys. Rev. Lett.* **76**:4656 (1996).
66. C. H. Bennett and S. J. Wiesner, *Phys. Rev. Lett.* **69**:2881 (1992).
67. M. Sargent, M. O. Scully, and W. E. Lamb, *Laser Physics* (Addison-Wesley, Reading, Mass., 1974).
68. M. O. Scully and M. S. Zubairy, *Quantum Optics* (Cambridge University Press, Cambridge, 1997).
69. M. Lax, "Phase Noise in a Homogeneously Broadened Maser," in *Physics of Quantum Electronics*, P. L. Kelley, B. Lax, and P. E. Tannenwald (eds.) (McGraw-Hill, New York, 1966).
70. M. Lax, *Phys. Rev.* **145**:110 (1966).
71. M. Lax, "Fluctuations and Coherence Phenomena in Classical and Quantum Physics," in *1966 Brandeis Summer Lecture Series on Statistical Physics*, Vol. 2, M. Chretien, E. P. Gross, and S. Dreser (eds.) (Gordon and Breach, New York, 1968; Mir, Moscow, 1975).
72. W. H. Louisell, *Quantum Statistical Properties of Radiation* (John Wiley, New York, 1973).
73. H. Haken, "Laser Theory," in *Encyclopedia of Physics*, Vol. 35/c, S. Flügge (ed.) (Springer, Berlin, 1970).
74. H. Risken, *The Fokker Planck Equation* (Springer, Heidelberg, 1984).
75. M. Lax, "The Theory of Laser Noise," keynote address, 1990 Conference on Laser Science and Optics. Applications, *Proc. SPIE* **1376**:2 (1991).
76. E. T. Jaynes and F. W. Cummings, *Proc. IEEE* **51**:89 (1963).
77. B. W. Shore and P. L. Knight, *J. Mod. Opt.* **40**:1195 (1993).
78. J. P. Gordon, H. J. Zeiger, and C. H. Townes, *Phys. Rev.* **95**:282L (1955).
79. J. P. Gordon, H. J. Zeiger, and C. H. Townes, *Phys. Rev.* **99**:1264 (1955).
80. A. L. Schawlow and C. H. Townes, *Phys. Rev. A* **112**:1940 (1958).
81. N. G. Basov and A. M. Prokhorov, *Dokl. Ak. Nauk* **101**:47 (1955).

82. T. H. Maiman, *Nature* (London) **187**:493 (1960).
83. A. Javan, W. R. Bennett, and D. R. Herriott, *Phys. Rev. Lett.* **6**:106 (1961).
84. W. E. Lamb, *Phys. Rev.* **134**:A1429 (1964).
85. For recent reviews, see, e.g., B.-G. Englert, M. Löffler, O. Benson, B. Varcoe, M. Weidinger, and H. Walther, *Fortschr. Phys.* **46**:897 (1998); G. Raithel, C. Wagner, H. Walther, L. M. Narducci, and M. O. Scully, in *Advances in Molecular and Optical Physics*, P. Berman (ed.) (Academic, New York, 1994), Suppl. 2.
86. D. Meschede, H. Walther, and G. Müller, *Phys. Rev. Lett.* **54**:551 (1985).
87. P. Filipowicz, J. Javanainen, and P. Meystre, *Phys. Rev. A* **34**:3077 (1986).
88. M. Lax, *Phys. Rev.* **129**:2342 (1963).
89. M. Lax and M. Zwanziger, *Phys. Rev. A* **7**:750 (1973).
90. M. Lax, *Phys. Rev.* **172**:350 (1968).
91. N. Lu, *Phys. Rev. A* **47**:4322 (1993).
92. U. Herzog and J. Bergou, *Phys. Rev. A* **62** (2000). In press.
93. M. Lax, *Phys. Rev.* **160**:290 (1967).
94. R. D. Hempstead and M. Lax, *Phys. Rev.* **161**:350 (1967).
95. H. Risken and H. D. Vollmer, *Z. Physik* **191**:301 (1967).
96. M. Kac, in 1966 *Brandeis Summer Lecture Series on Statistical Physics*, Vol. 1, M. Chretien, E. P. Gross, and S. Dreser (eds.) (Gordon and Breach, New York, 1968).
97. M. O. Scully, H. Walther, G. S. Agarwal, T. Quang, and W. Schleich, *Phys. Rev. A* **44**:5992 (1991).
98. N. Lu, *Phys. Rev. Lett.* **70**:912 (1993); N. Lu, *Phys. Rev. A* **47**:1347 (1993); T. Quang, G. S. Agarwal, J. Bergou, M. O. Scully, H. Walther, K. Vogel, and W. P. Schleich, *Phys. Rev. A* **48**:803 (1993); K. Vogel, W. P. Schleich, M. O. Scully, and H. Walther, *Phys. Rev. A* **48**:813 (1993); R. McGowan and W. Schieve, *Phys. Rev. A* **55**:3813 (1997).
99. S. Qamar and M. S. Zubairy, *Phys. Rev. A* **44**:7804 (1991).
100. P. Filipowicz, J. Javanainen, and P. Meystre, *J. Opt. Soc. Am. B* **3**:154 (1986).
101. J. J. Slosser, P. Meystre, and S. L. Braunstein, *Phys. Rev. Lett.* **63**:934 (1989).
102. J. Bergou, L. Davidovich, M. Orszag, C. Benkert, M. Hillery, and M. O. Scully, *Phys. Rev. A* **40**:7121 (1989); J. D. Cresser, *Phys. Rev. A* **46**:5913 (1992); U. Herzog, *Phys. Rev. A* **52**:602 (1995).
103. H. J. Briegel, B.-G. Englert, N. Sterpi, and H. Walther, *Phys. Rev. A* **49**:2962 (1994); U. Herzog, *Phys. Rev. A* **50**:783 (1994); J. D. Cresser and S. M. Pickless, *Phys. Rev. A* **50**:R925 (1994); U. Herzog, *Appl. Phys. B* **60**:S21 (1995).
104. H.-J. Briegel, B.-G. Englert, Ch. Ginzel, and A. Schenzle, *Phys. Rev. A* **49**:5019 (1994).
105. H.-J. Briegel, B.-G. Englert, *Phys. Rev. A* **52**:2361 (1995); J. Bergou, *Quantum and Semiclass. Optics* **7**:327 (1995); U. Herzog and J. Bergou, *Phys. Rev. A* **54**:5334 (1996); *ibid.* **55**:1385 (1997).
106. M. Weidinger, B. T. H. Varcoe, R. Heerlein, and H. Walther, *Phys. Rev. Lett.* **82**:3795 (1999).
107. H. Walther, *Phys. Rep.* **219**:263 (1992).
108. M. O. Scully, G. Süssmann, and C. Benkert, *Phys. Rev. Lett.* **60**:1014 (1988).
109. M. O. Scully, M. S. Zubairy, and K. Wódkiewicz, *Opt. Commun.* **65**:440 (1988).
110. H. E. Stanley, *Introduction to Phase Transitions and Critical Phenomena* (Oxford University Press, Oxford, 1971).
111. V. DeGiorgio and M. O. Scully, *Phys. Rev. A* **2**:1170 (1970).
112. M. Anderson, J. Ensher, M. Matthews, C. Wieman, and E. Cornell, *Science* **269**:198 (1995).
113. C. Bradley, C. Sackett, J. Tollett, and R. Hulet, *Phys. Rev. Lett.* **75**:1687 (1995).
114. K. Davis, M. Mewes, M. Andrews, N. van Druten, D. Durfee, D. Kurn, and W. Ketterle, *Phys. Rev. Lett.* **75**:3969 (1995).
115. M. H. W. Chen, K. I. Blum, S. Q. Murphy, G. K. S. Wong, and J. D. Reppy, *Phys. Rev. Lett.* **61**:1950 (1998).
116. K. Huang, *Statistical Mechanics* (John Wiley, New York, 1987).
117. R. Arnowitt and M. Girardeau, *Phys. Rev.* **113**:745 (1959).

118. G. Baym and C. Pethick, *Phys. Rev. Lett.* **76**:6 (1996).
119. H. Politzer, *Phys. Rev. A* **54**:5048 (1996).
120. S. Grossmann and M. Holthaus, *Phys. Rev. E* **54**:3495 (1996).
121. C. Herzog and M. Olshanii, *Phys. Rev. A* **55**:3254 (1997).
122. P. Navez, D. Bitouk, M. Gajda, Z. Idziaszek, and K. Rzążewski, *Phys. Rev. Lett.* **79**:1789 (1997).
123. M. Wilkens and C. Weiss, *J. Mod. Opt.* **44**:1801 (1997).
124. S. Grossmann and M. Holthaus, *Phys. Rev. Lett.* **79**:3557 (1997).
125. H. Woolf (ed.), *Some Strangeness in the Proportion: A Centennial Symposium to Celebrate the Achievements of Albert Einstein* (Addison-Wesley, Reading, Mass., 1980), p. 524.
126. J. Goldstein, M. O. Scully, and P. Lee, *Phys. Lett.* **A35**:317 (1971).
127. M. O. Scully, *Phys. Rev. Lett.* **82**:3927 (1999).
128. O. Kocharovskaya and Ya I. Khanin, *JETP Lett.* **48**:630 (1988).
129. S. E. Harris, *Phys. Rev. Lett.* **62**:1033 (1989).
130. M. O. Scully, S.-Y. Zhu, and A. Gavrieleides, *Phys. Rev. Lett.* **62**:2813 (1989).
131. A. Nottelmann, C. Peters, and W. Lange, *Phys. Rev. Lett.* **70**:1783 (1993).
132. E. S. Fry, X. Li, D. Nikonov, G. G. Padmabandu, M. O. Scully, A. V. Smith, F. K. Tittel, C. Wang, S. R. Wilkinson, and S.-Y. Zhu, *Phys. Rev. Lett.* **70**:3235 (1993).
133. W. E. van der Veer, R. J. J. van Diest, A. Dönszelmann, and H. B. van Linden van den Heuvell, *Phys. Rev. Lett.* **70**:3243 (1993).
134. A. S. Zibrov, M. D. Lukin, D. E. Nikonov, L. W. Hollberg, M. O. Scully, V. L. Velichansky, and H. G. Robinson, *Phys. Rev. Lett.* **75**:1499 (1995).
135. M. O. Scully, *Phys. Rev. Lett.* **55**:2802 (1985).
136. J. Bergou, M. Orszag, and M. O. Scully, *Phys. Rev. A* **38**:768 (1988).
137. M. O. Scully and M. S. Zubairy, *Phys. Rev. A* **35**:752 (1987).
138. M. O. Scully, K. Wodkiewicz, M. S. Zubairy, J. Bergou, N. Lu, and J. Meyer ter Vehn, *Phys. Rev. Lett.* **60**:1832 (1988).
139. J. M. J. Madey, *J. Appl. Phys.* **42**:1906 (1971).
140. W. Becker, M. O. Scully, and M. S. Zubairy, *Phys. Rev. Lett.* **48**:475 (1985).
141. W. Becker and M. S. Zubairy, *Phys. Rev. A* **25**:2200 (1982).
142. R. Bonifacio, *Opt. Commun.* **32**:440 (1980).
143. A. Bambini and A. Renieri, *Opt. Commun.* **29**:244 (1978).
144. B.-G. Englert, J. Schwinger, A. O. Barut, and M. O. Scully, *Europhys. Lett.* **14**:25 (1991).
145. S. Haroche, M. Brune, and J.-M. Raimond, *Europhys. Lett.* **14**:19 (1991).
146. M. Battocletti and B.-G. Englert, *J. Phys. II France* **4**:1939 (1994).
147. M. O. Scully, G. M. Meyer, and H. Walther, *Phys. Rev. Lett.* **76**:4144 (1996).
148. G. M. Meyer, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4142 (1997).
149. M. Löffler, G. M. Meyer, M. Schröder, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4153 (1997).
150. M. Schröder, K. Vogel, W. P. Schleich, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4164 (1997).
151. M. Löffler, G. M. Meyer, and H. Walther, *Lett.* **41**:593 (1998).

This page intentionally left blank.

PART

5

DETECTORS

This page intentionally left blank.

Paul R. Norton

*U.S. Army Night Vision and Electronics Directorate
Fort Belvoir, Virginia*

*Second revision and update from an article by Stephen F. Jacobs**

24.1 SCOPE

The purpose of this chapter is to describe the range of detectors commercially available for sensing optical radiation. Optical radiation over the range from vacuum ultraviolet to the far-infrared or submillimeter wavelength (25 nm to 1000 μm) is considered. We will refer to the following spectral ranges:

25–200 nm	vacuum ultraviolet	VUV
200–400 nm	ultraviolet	UV
400–700 nm	visible	VIS
700–1000 nm	near infrared	NIR
1–3 μm	short-wavelength infrared	SWIR
3–5 μm	medium-wavelength infrared	MWIR
5–14 μm	long-wavelength infrared	LWIR
14–30 μm	very long wavelength infrared	VLWIR
30–100 μm	far-infrared	FIR
100–1000 μm	submillimeter	SubMM

We begin by giving a brief description of the photosensitive mechanism for each type of detector. The usefulness and limitations of each detector type are also briefly described. Definitions of the technical terms associated with the detection process are listed. The concept of sensitivity is defined, and D^* (D^*) is presented as a measure of ideal performance. Examples are then given of the limiting cases for D^* under several conditions. In addition, other detector performance parameters are described which may be of critical interest in a specific application, including spectral response, responsivity, quantum efficiency, noise, uniformity, speed, and stability. Finally, manufacturers' specifications for a range of available detectors are compiled and a list of manufacturers is included for each type of detector.

*In *Handbook of Optics*, first edition, McGraw-Hill, 1978. Section 4, "Nonimaging Detectors," by Stephen F. Jacobs, Optical Sciences Center, University of Arizona, Tucson, Arizona.

The sensitivity of many detectors has reached the limit set by room-temperature background photon fluctuations (radiation shot noise). For these detectors, sensitivity may be enhanced by providing additional cooling, while restricting their spatial field of view and/or spectral bandwidth. At some point, other factors such as amplifier noise may limit the improvement.

Techniques for evaluating detector performance are not covered in this treatment but can be found in Refs. 1–6.

24.2 THERMAL DETECTORS

Thermal detectors sense the change in temperature produced by absorption of incident radiation. Their spectral response can therefore be as flat as the absorption spectrum of their blackened coating and window* will allow. This makes them useful for spectroscopy and radiometry. These detectors are generally operated at room temperature, where their sensitivity is limited by thermodynamic considerations^{7,8} to 3 pW for 1-s measurement time and 1-mm² sensitive area. This limit has been very nearly reached in practice, whereas cooled bolometers have been made to reach the background photon noise limit. Figure 1 illustrates the basic structural elements of a thermal detector.

Construction of the detector seeks to minimize both the thermal mass of the sensitive element and the heat loss from either conductive or convective mechanisms. Heat loss may ideally be dominated by radiation. This allows the incident photon flux to give a maximum temperature rise (maximum signal), but results in a correspondingly slow response time for this class of detectors. The response time τ of thermal detectors is generally slower than 1 ms, depending on thermal capacity C and heat loss per second per degree G , through the relation

$$\tau = C/G$$

A short time constant requires a small C . However, for room-temperature operation, ultimate sensitivity is limited by the mean spontaneous temperature fluctuation:

$$\Delta T = T \sqrt{\frac{k}{C}}$$

where k is Boltzmann's constant. There is thus a trade-off between time constant and ultimate sensitivity.

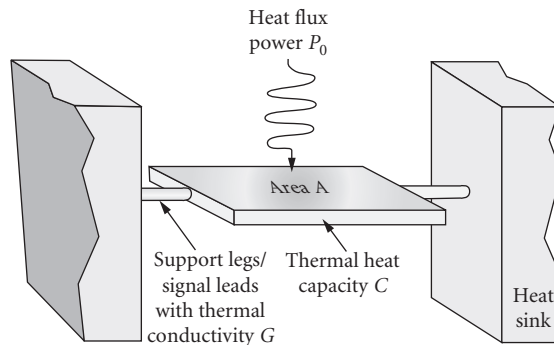


FIGURE 1 Structural elements of a thermal detector. Sensitive area A with a thermal heat capacity of C , supported by leads having thermal conductivity G , and with a heat flux of P_0 incident on the pixel.

*No windows exist without absorption bands anywhere between the visible and millimeter region. Some useful window materials for the far-infrared are diamond, silicon, polyethylene, quartz, and CsI.

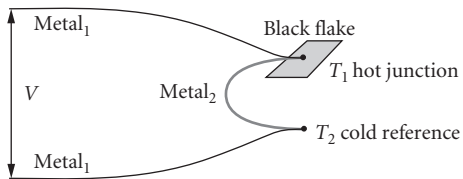


FIGURE 2 Thermocouple detector structure.

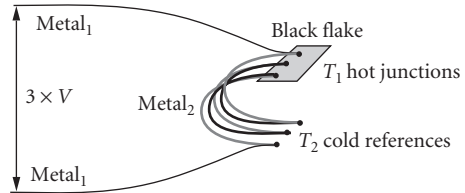


FIGURE 3 Thermopile detector structure.

Thermocouple/Thermopile

The thermocouple receiver, illustrated in Fig. 2, is a thin, blackened flake connected thermally to the junction of two dissimilar metals or semiconductors. Heat absorbed by the flake causes a temperature rise of the junction, and hence a thermoelectric emf is developed which can be measured, for example, with a voltmeter.

Thermocouples are limited in sensitivity by thermal (Johnson-Nyquist) noise but are nevertheless respectably sensitive. Their usefulness lies in the convenience of room temperature operation, wide spectral response, and their rugged construction. Thermocouples are widely used in spectroscopy.

Thermopiles consist of thin-film arrays of thermocouples in series, as illustrated in Fig. 3. This device multiplies the thermocouple signal corresponding to the number of junctions in series. The device may be constructed with half the thermocouples acting as reference detectors attached to a heat sink.

Bolometer/Thermistor

The receiver is a thin, blackened flake or slab, whose impedance is highly temperature dependent—see Fig. 4. The impedance change may be sensed using a bridge circuit with a reference element in the series or parallel arm. Alternatively, a single bolometer element in series with a load and voltage source may be used.

Most bolometers in use today are of the thermistor type made from oxides of manganese, cobalt, or nickel. Their sensitivity closely approaches that of the thermocouple for frequencies higher than 25 Hz. At lower frequencies there may be excess or $1/f$ noise. Construction can be very rugged for systems applications. Some extremely sensitive low-temperature semiconductor and superconductor bolometers are available commercially.

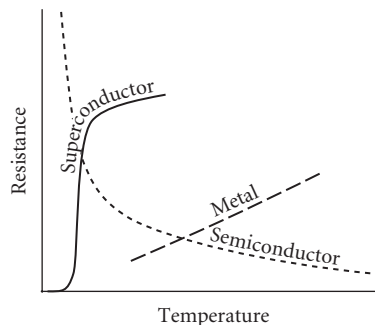


FIGURE 4 Temperature dependence characteristics of three bolometer material types.

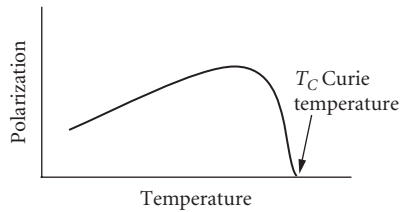


FIGURE 5 Ferroelectric materials exhibit residual polarization with no applied bias.

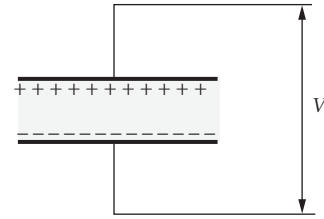


FIGURE 6 The pyroelectric effect produces a surface charge when the temperature changes.

Pyroelectric

Ferroelectric material exhibits a residual polarization in the absence of any electric field, as illustrated in Fig. 5. Dipole moments, initially aligned by applying an external field, result in a surface charge which is normally slowly neutralized by leakage. This polarization is temperature-dependent (pyroelectric effect), and when incident radiation heats an electroded sample, there is a change in surface charge (open-circuit voltage) which is proportional to the incident radiation power—see Fig. 6. Electrically, the device behaves like a capacitor, requiring no bias and therefore exhibiting no current noise. The signal, however, must be chopped or modulated. Sensitivity is limited either by amplifier noise or by loss-tangent noise. Response speed can be engineered, with a proportional decrease in sensitivity, making pyroelectric detectors useful for moderately fast laser pulse detection. Other common applications include power meters. Microphonic noise in applications associated with vibrations can be a problem with some of these devices.

24.3 QUANTUM DETECTORS

Photon detectors respond in proportion to incident photon rates (quanta) rather than to photon energies (heat). Thus, the spectral response of an ideal photon detector is flat on an incident-photon-rate basis but linearly rising with wavelength on an incident-power (per watt) basis. The sensitivity of efficient quantum detectors can approach the limits of photon noise fluctuations provided that the detector temperature is sufficiently low for photon-induced mechanisms to dominate thermally induced mechanisms in the detector. Quantum detectors generally have sub-microsecond time constants. Their main disadvantage is the associated cooling required for optimum sensitivity. (These remarks do not apply to photographic detection, which measures cumulative photon numbers.)

Photoemissive

The radiation is absorbed by a photosensitive surface which usually contains alkali metals (cesium, sodium, or potassium). Incident quanta release photoelectrons (Fig. 7), via the photoelectric effect, which are collected by a positively biased anode. This is called a diode phototube; it can be made the basis for the multiplier phototube (photomultiplier phototube, or photomultiplier) by the addition of a series of biased dynodes which serve as secondary emission multipliers.

In spectral regions where quantum efficiency is high ($\lambda < 550$ nm), the photoemissive detector is very nearly ideal. Sensitivity is high enough to count individual photons. Amplification does not degrade the signal-to-noise ratio. The sensitive area is conveniently large. Photomultiplier signal response time (transit spread time) can be made as short as 0.1 ns. Since the sensitivity in red-sensitive tubes is limited by thermally generated electrons, sensitivity can be improved by cooling.

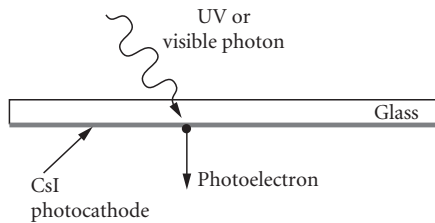


FIGURE 7 Electrons are ejected from a photoemissive surface when excited by photons.

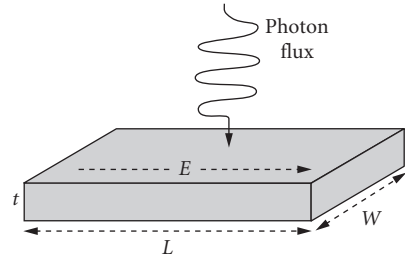


FIGURE 8 Photoconductor device structure.

Photoconductive

The radiation is absorbed by a photoconductive (PC) material, generally a semiconductor, either in thin-film or bulk form, as illustrated in Fig. 8. Each incident quantum may release an electron-hole pair or an impurity-bound charge carrier, thereby increasing the electrical conductivity. The devices are operated in series with a bias voltage and load resistor. Very low impedance photoconductors may be operated with a transformer as the load. Since the impedance of photoconductors varies with device type and operating conditions from less than $50\ \Omega$ to more than $10^{14}\ \Omega$, the load resistor and preamplifier must be chosen appropriately. Figure 9 shows the current-voltage characteristics of a photoconductor.

Photoconductors that utilize excitation of an electron from the valence to conduction band are called *intrinsic* detectors. Those which operate by exciting electrons into the conduction band or holes into the valence band from states within the band—impurity-bound states, quantum wells, or quantum dots—are called *extrinsic* detectors. Figure 10 illustrates these two mechanism types. Intrinsic detectors are most common at the short wavelengths, out to about $20\ \mu\text{m}$. Extrinsic detectors are most common at longer wavelengths. A key difference between intrinsic and extrinsic detectors is that intrinsic detectors do not require as much cooling to achieve high sensitivity at a given spectral response cutoff as extrinsic detectors. Thus, intrinsic photoconductors such as HgCdTe will operate out to 15 to $20\ \mu\text{m}$ at $77\ \text{K}$, while comparable extrinsic detectors with similar cutoff must be cooled below 30 to $40\ \text{K}$.

A further distinction may be made by whether the semiconductor material has a direct or indirect bandgap. This difference shows up near the long-wavelength limit of the spectral response

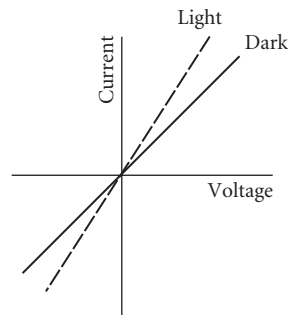


FIGURE 9 Photoconductor current-voltage characteristics in the dark and in the light.

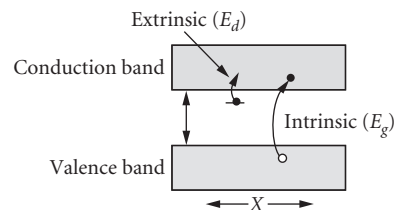


FIGURE 10 Intrinsic detectors excite electrons between the valence and conduction band. Extrinsic detectors excite electrons (or holes) from states within the band to the conduction (valence) band.

where detectors made from direct bandgap materials such as InGaAs, InSb, or HgCdTe have a sharper spectral cutoff than indirect bandgap materials such as silicon and germanium.

Photoconductors can have high quantum efficiency from the visible region out to the far infrared but lack the nearly ideal high amplification of photomultipliers. They are therefore most commonly used in the spectral region beyond 1 μm , where efficient photoemitters are unavailable. Photoconductors do, however, provide current gain which is equal to the recombination time divided by the majority-carrier transit time. This current gain leads to higher responsivity than is possible with (nonavalanching) photovoltaic (PV) detectors. For applications where photovoltaic detection would be amplifier noise limited, the larger photoconductive responsivity makes it possible to realize greater sensitivity with the photoconductor. In general, lower-temperature operation is associated with longer-wavelength sensitivity in order to suppress noise due to thermally induced transitions between close-lying energy levels. Ideally, photoconductors are limited by generation-recombination noise in the photon-generated carriers. Response time can be shorter than 1 μs and in some cases response times can be shorter than 1 ns for small elements. Response across a photoconductive element can be nonuniform due to recombination mechanisms at the electrical contacts, and this effect may vary with electrical bias.

Photovoltaic*

The most widely used photovoltaic detector is the *pn* junction type (Fig. 11), where a strong internal electric field exists across the junction even in the absence of radiation. Photons incident on the junction of this film or bulk material produce free hole-electron pairs which are separated by the internal electric field across the junction, causing a change in voltage across the open-circuit cell or a current to flow in the short-circuited case.

As with the photoconductor, quantum efficiency can be high from the visible to the very long-wavelength infrared, generally about 20 to 25 μm . The limiting noise level can ideally be $\sqrt{2}$ times lower than that of the photoconductor, thanks to the absence of recombination noise. Lower temperatures are associated with longer-wavelength operation. Response times can be less than a nanosecond,

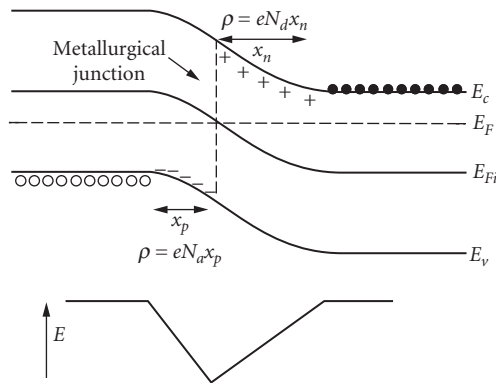


FIGURE 11 *pn* junction showing how the band bends across the junction. The junction width is given by $W = x_p + x_n$.

*The use of a photovoltaic detector at other than zero bias is often referred to as its photoconductive mode of operation because the circuit then is similar to the standard photoconductor circuit. This terminology is confusing with regard to detection mechanism and will not be used here.

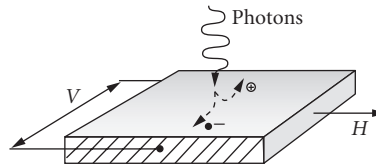


FIGURE 12 Photoelectromagnetic detector configuration—incident photons create electron-hole pairs that diffuse away from the surface. A magnetic field perpendicular to the diffusion separates the charges toward opposite sides, creating a voltage.

and are generally limited by device capacitance and detection-circuit resistance. The *pin* diode has been developed to minimize capacitance for high-bandwidth applications. The advantages of nearly ideal internal amplification have now become available in avalanche photodiodes sensitive out to $1.55\ \mu\text{m}$. This internal gain is most important for high-frequency operation, where external load resistance must be kept small and would otherwise introduce limiting thermal noise, and for situations involving low signal flux where amplifier noise is otherwise dominant.

Photoelectromagnetic

A thin slab of photoconductive material is oriented with radiation incident on a large face and a magnetic field perpendicular to it, as illustrated in Fig. 12. Electron-hole pairs generated by the incident photons diffuse through the material and are separated by the magnetic field, causing a potential difference at opposite ends of the detector.

These detectors require no cooling or biasing electric field but do require a (permanent) magnetic field. Photoelectromagnetic InSb at room temperature has response up to $7.5\ \mu\text{m}$, where it is as sensitive as a thermocouple of equal size, and has a response time less than $1\ \mu\text{s}$. Another competing uncooled detector is InAs, which is far more sensitive out to $3.5\ \mu\text{m}$. HgCdTe is also available in the photoelectromagnetic (PEM) configuration out to the LWIR spectral region. Cooled infrared detectors are one to two orders of magnitude more sensitive.

Photographic

The receiver is an emulsion containing silver halide crystals. Incident photons are absorbed by the halide ion, which subsequently loses its electron. This electron eventually recombines with a silver ion and reduces it to a neutral silver atom. As more photons are absorbed, this process is repeated until a small but stable cluster of reduced silver atoms is formed within the crystal (latent image). Internal amplification is provided by introduction of an electron donor (photographic) developer, which, using the latent image as a catalytic center, reduces all the remaining silver ions within the exposed crystal to neutral silver atoms. The density of reduced crystals is a measure of the total radiation exposure.

The spectral region of sensitivity for photographic detection coincides rather closely with that of the photoemissive detector. For $\lambda > 1.2\ \mu\text{m}$ ($\sim 1\ \text{eV}$) there is too little energy in each photon to form a stable latent image. The basic detection process for both detectors operates well for higher energy, shorter wavelength radiation. The problem in ultraviolet and x-ray operation is one of eliminating nonessential materials, for example, the emulsion which absorbs these wavelengths.

The photographic process is an integrating one in that the output (emulsion density) measures the cumulative effect of all the radiation incident during the exposure time. The efficiency of the photographic process can be very high, but it depends upon photon energy; for example, in the

visible region it takes only 10 to 100 photons to form a stable latent image (developable grain). The photographic process enjoys a large and efficient internal amplification ability (development) wherein the very small energy of the photons' interaction is converted into readily observed macroscopic changes. An extensive discussion of photographic detection is found in Chap. 29, "Photographic Films."

Photoionization

The radiation is absorbed by a gas. If the photon energy exceeds the gas-ionization threshold, ion pairs can be produced with very high efficiency. They are collected by means of an applied voltage. Operating in a dc mode, these detectors are known as ionization and gas-gain chambers. When a pulse mode is used, the detectors are known as proportional and photon (Geiger) counters.

Photoionization detectors have a high sensitivity and a low noise level. They may also be quite selective spectrally since the choice of window and gas independently set upper and lower limits on detectable photon energies. Manufacturers' specifications are not discussed for these detectors as applications are still few enough to be treated as individual cases.⁹

24.4 DEFINITIONS

The following definitions will be used:

Avalanche photodiode (APD) A photodiode designed to operate in strong reverse bias where electron and/or hole impact ionization produces multiplication of photogenerated carriers.

Background temperature The effective temperature of all radiation sources viewed by the detector exclusive of the signal source.

Blackbody D star D_{BB}^* (cm Hz^{1/2}/W also called "Jones") Similar to $D^*(\lambda)$ or $D^*(T_B f)$ except that the source is a blackbody whose temperature must be specified.

Blackbody detectivity D_{BB} (W⁻¹) A measure of detector sensitivity, defined as $DBB (NEP_{BB})^{-1}$.

Blackbody noise-equivalent power NEP_{BB} Same as spectral NEP, except that the source has blackbody spectral character whose temperature must be specified, for example, NEP (500 K, 1, 800) means 500-K blackbody radiation incident, 1-Hz electrical bandwidth, and 800-Hz chopping frequency.

Blackbody responsivity R_{BB} Same as spectral sensitivity except that the incident signal radiation has a blackbody spectrum (whose temperature must be specified).

Blip detector or blip condition Originally meaning background-limited impurity photoconductor, this term has come to mean performance of any detector where the limiting noise is due to fluctuations in the arrival rate of background photons.

Cutoff wavelength λ_c The wavelength at which the detectivity has degraded to one-half its peak value.

Dark current The output current which flows even when input radiation is absent or negligible. (Note that although this current can be subtracted out for the dc measurements, the shot noise on the dark current can become the limiting noise.)

Detective quantum efficiency The square of the ratio of measured detectivity to the theoretical limit of detectivity.

Detective time constant $\tau_d = 1/2\pi f_d$ where f_d is the frequency at which D^* drops to 0.707 ($1/\sqrt{2}$) times its maximum value. A physics convention defines it as $1/e$, or 0.368 of the maximum value.

Dewar A container (cryostat) for holding detector coolant.

Equivalent noise input (ENI) A term meaning nearly the same thing as NEP_{BB} (287 K, 1, f). The difference is that the peak-to-peak value of square-wave chopped input flux is used, rather than the rms value of the sinusoidally chopped input flux. (See recommendation IRE.²)

Excess noise A term usually referring to noises other than generation-recombination, shot, or thermal.

Extrinsic semiconductor transition Incident photons produce a free electron in the conduction band and bound hole at a donor impurity site or a bound electron at an acceptor impurity site and a free hole in the valence band by excitation of an impurity level.

Field of view (FOV) The solid angle from which the detector receives radiation.

Flicker noise See Modulation noise.

Generation noise Noise produced by the statistical fluctuation in the rate of production of photoelectrons.

Generation-recombination (GR) noise Charge carriers are generated both by (optical) photons and (thermal) phonons. Fluctuations in these generation rates cause noise; fluctuations in carrier-recombination times cause recombination noise. The phonon contribution can be removed by cooling. The remaining photon contribution is indistinguishable from radiation shot noise (photon noise). With photovoltaic pn junctions, carriers are swept away before recombination, so that recombination noise is absent.

Guard ring An electrically biased field plate or surrounding diode used in some photodiodes, usually used to control surface recombination effects and thus reduce the leakage current in the detection circuit.

Intrinsic semiconductor transition Incident photons produce a free electron-free hole pair by direct excitation across the forbidden energy gap (valence to conduction band).

Johnson noise Same as Thermal noise.

Jones Unit of measure for D^* $\text{cm Hz}^{1/2}/\text{W}$.

Maximized D star, $D^(\lambda_{pk}, f_o)$, $\text{cm Hz}^{1/2}/\text{W}$ or *Jones** The value of $D^*(\lambda_{pk}, f)$ corresponding to wavelength λ_{pk} and chopping frequency of maximum D^* .

Modulation (or $1/f$) noise A consensus regarding the origin(s) of the mechanism has not been established, and although a quantum theory has been proposed, other mechanisms may dominate. As its name implies, it is characterized by a $1/f^n$ noise power spectrum, where $0.8 < n < 2$. This type of noise is prominent in thermal detectors and may dominate the low-frequency noise characteristics of photoconductive and photovoltaic quantum detectors as well as other electronic devices such as transistors and resistors.

Multiplier phototube or multiplier photodiode Phototube with built-in amplification via secondary emission from electrically biased dynodes.

NEI photons/($\text{cm}^2 \text{sec}$) noise equivalent irradiance is the signal flux level at which the signal produces the same output as the noise present in the detector. This unit is useful because it directly gives the photon flux above which the detector will be photon noise limited. See also Spectral noise equivalent power (NEP).

Noise spectrum The electrical power spectral density of the noise.

Nyquist noise Same as Thermal noise.

Photo cell See Photodiode.

Photoconductive gain The ratio of carrier lifetime divided by carrier transit time in a biased photoconductor.

Photodiode The term photodiode has been applied both to vacuum- or gas-filled photoemissive detectors (diode phototubes, or photo cells) and to photovoltaic detectors (semiconductor pn junction devices).

Photomultiplier Same as Multiplier phototube.

Photon counting Digital counting of individual photons from the photoelectrons produced in the detector in contrast to averaging of the photocurrent. This technique leads to very great sensitivity but can be used only for quite low light levels.

Quantum efficiency The ratio of the number of countable output events to the number of incident photons, for example, photoelectrons per photon, usually referred to as a percentage value.

RMS noise $V_{n,rms}$ That component of the electrical output which is not coherent with the radiation signal (generally measured with the signal radiation removed).

RMS signal $V_{s,rms}$ That component of the electrical output which is coherent with the input signal radiation.

Response time τ Same as Time constant.

Responsive quantum efficiency See Quantum efficiency.

Sensitivity Degree to which detector can sense small amounts of radiation.

Shot noise This current fluctuation results from the random arrival of charge carriers, as in a photodiode. Its magnitude is set by the size of the unit charge.

$$i_{n,rms} = (2ei_{dc}\Delta f)^{1/2}$$

Spectral D-double-star $D^{**}(\lambda, f)$ A normalization of D^* to account for detector field of view. It is used only when the detector is background-noise-limited. If the FOV is 2π sr, $D^{**} = D^*$.

Spectral D-star $D^*(\lambda, f)$ ($cm\ Hz^{1/2}/W$ or *Jones*) A normalization of spectral detectivity to take into account the area and electrical bandwidth dependence, for example, $D^*(1\ \mu m, 800\ Hz)$ means D^* at $\lambda = 1\ \mu m$ and chopping frequency 800 Hz; unit area and electrical bandwidth are implied. For background-noise-limited detectors the FOV and the background characteristics must be specified. For many types of detectors this normalization is not valid, so that care should be exercised in using D^* .

Spectral detectivity $D(\lambda)$ (W^{-1}) A measure of detector sensitivity, defined as $D(\lambda) = (NEP_{\lambda})^{-1}$. As with NEP, the chopping frequency electrical bandwidth, sensitive area, and, sometimes, background characteristics should be specified.

Spectral noise equivalent power NEP_{λ} The rms value of sinusoidally modulated monochromatic radiant power incident upon a detector which gives rise to an rms signal voltage equal to the rms noise voltage from the detector in a 1-Hz bandwidth. The chopping frequency, electrical bandwidth, detector area, and, sometimes, the background for characteristics should be specified. $NEP(1\ \mu m, 800\ Hz)$ means noise equivalent power at 1- μm wavelength, 1-Hz electrical bandwidth, and 800-Hz chopping rate. Specification of electrical bandwidth is often simplified by expressing NEP in units of $W/Hz^{1/2}$.

Spectral responsivity $R(\lambda)$ The ratio between rms signal output (voltage or current) and the rms value of the monochromatic incident signal power or photon flux. This is usually determined by taking the ratio between a sample detector and a thermocouple detector. The results are given as relative response/watt or relative response/photon, respectively.

Temperature noise Fluctuations in the temperature of the sensitive element, due either to radiative exchange with the background or conductive exchange with a heat sink, produce a fluctuation in signal voltage. For thermal detectors, if the temperature noise is due to the former, the detector is said to be at its theoretical limit. For thermal detectors:

$$\overline{(\Delta T)^2} = \frac{4kT^2G\Delta f}{K^2 + 4\pi^2 f^2 C^2}$$

where $\overline{(\Delta T)^2}$ = mean square temperature fluctuations

K = thermal conductance

C = heat capacity

Thermal noise (also known as Johnson or Nyquist noise) Noise due to the random motion of charge carriers in a resistive element:

$$V_{n,rms} = (4kTR\Delta f)^{1/2} \quad k = \text{Boltzmann's constant}$$

Thermopile A number of thermocouples mounted in series in such a way that their thermojunctions lie adjacent to one another in the plane of irradiation.

Time constant τ (see also *detective time constant*) A measure of the detector's speed of response. $\tau = 1/(2\pi f_c)$, where f_c is that chopping frequency at which the responsivity has fallen to 0.707 ($1/\sqrt{2}$) times its maximum value. Sometimes a physics convention defines it as $1/e$, or 0.368 of the maximum value:

$$R(f) = \frac{R_0}{(1 + 4\pi^2 f^2 \tau^2)^{1/2}}$$

24.5 DETECTOR PERFORMANCE AND SENSITIVITY

D^*

A figure of merit defined by Jones in 1958 is used to compare the sensitivity of detectors.¹⁰ It is called D^* . Although the units of measure are $\text{cm Hz}^{1/2}/\text{W}$, this unit is now referred to as a *Jones*. D^* is the signal-to-noise (S/N) ratio of a detector measured under specified test conditions, referenced to a 1-Hz bandwidth and normalized by the square root of the detector area (A) in square centimeters. Specified test conditions usually consist of the blackbody signal source temperature, often 500 K for infrared detectors, and the signal chopping frequency. If the background temperature is other than room temperature (295 or 300 K in round numbers), then that should be noted.

By normalizing the measured S/N ratio by the square root of the detector area, the D^* figure of merit recognizes that the statistical fluctuations of the background photon flux incident on the detector (photon noise) are dependent upon the square root of the number of photons and thus increase as the square root of the detector area, while the signal will increase in proportion to the detector area itself. This figure of merit therefore provides a valid comparison of detector types that may have been made and tested in different sizes.

The ultimate limit in S/N ratio for any radiation power detector is set by the statistical fluctuation in photon arrival times. For ideal detectors which are photon-noise-limited, and where only generation noise is present, we shall discuss limiting detectivity for three cases:

1. Photon detector where arrival rate of signal photons far exceeds that of background photons (all other noise being negligible)
2. Photon detector where background photon arrival rate exceeds signal photon rate (all other noise being negligible)
3. Thermal detector, background limited

The rate of signal-carrier generation is

$$n = \eta AN_s \tag{1}$$

where η = detector quantum efficiency and AN_s = average rate of arrival of signal photons.

It can be shown¹¹ that in a bandwidth Δf , the rms fluctuation in carrier-generation rate is

$$\delta n_{rms} = (2P_N \Delta f)^{1/2} \tag{2}$$

where P_N is the frequency dependence of the mean square fluctuations in the rate of carrier generation, that is,

$$P_N = A \int_0^\infty \eta(\nu) (\Delta N)^2 d\nu \quad (3)$$

where $(\Delta N)^2$ is the mean square deviation in the total rate of photon arrivals per unit area and frequency interval including signal and background photons. For thermally produced photons of frequency ν (see Ref. 12).

$$(\Delta N)^2 = \bar{N} \frac{e^{h\nu/kT}}{e^{h\nu/kT} - 1} = \frac{2\pi\nu^2}{c^2} \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} \quad (4)$$

where \bar{N} is the average rate of photon arrivals per unit area and frequency interval. Then, for the special case of $h\nu \gg kT$

$$\delta n_{\text{rms}} = (2A\eta\bar{N}\Delta f)^{1/2} \quad (5)$$

This is also the case for a laser well above threshold. Here the photon statistics become Poisson, and $(\Delta N)^2 = \bar{N}$ even when $h\nu$ is not greater than kT .

Photon Detector, Strong-Signal Case This is generally a good approximation for visible and higher photon energy detectors since the background radiation is often weak or negligible. When signal photons arrive at a much faster rate than background photons

$$\delta n_{\text{rms}} = (2A\eta\bar{N}_s\Delta f)^{1/2} \quad (6)$$

then

$$\text{NEP} = \frac{N_s A h \nu}{n / \delta n_{\text{rms}} (\Delta f)^{1/2}} = \left(\frac{2N_s A}{\eta} \right)^{1/2} h \nu \quad (7)$$

or the noise-equivalent quantum rate is

$$\text{NEQ} = \left(\frac{2 \times \text{incident photon rate}}{\text{quantum efficiency}} \right)^{1/2} \quad (8)$$

Photon Detector, Background-Limited Case This is usually a good approximation for detecting low signal levels in the infrared where background flux levels exceed signal flux levels in many applications. When the background photon noise rate N_B exceeds the signal photon rate ($N_B \gg N_s$)

$$\delta n_{\text{rms}} \approx (2A\eta\bar{N}_B\Delta f)^{1/2} \quad (9)$$

the noise-equivalent power is

$$\text{NEP} = \frac{N_s A h \nu}{(n / \delta n_{\text{rms}}) (\Delta f)^{1/2}} = \left(\frac{2N_B A}{\eta} \right)^{1/2} h \nu \quad (10)$$

The noise-equivalent quantum rate is

$$\text{NEQ} = \left(\frac{2 \times \text{incident background photon rate}}{\text{quantum efficiency}} \right)^{1/2} \quad (11)$$

or

$$D^* = \frac{A^{1/2}}{\text{NEP}} = \left(\frac{\eta}{2N_B} \right)^{1/2} \frac{1}{h\nu} \quad (12)$$

or

$$\text{Area-normalized quantum detectivity} = \frac{A^{1/2}}{\text{NEQ}} = \left(\frac{\eta}{2N_B} \right)^{1/2} \quad (13)$$

For the general case of a detector with area A seeing 2π sr of blackbody background at temperature T , $(\Delta N)^2$ is that in Eq. (4)

$$\overline{(\delta n)^2} = 2A\Delta f \int_0^\infty \eta(\nu) (\Delta N)^2 d\nu = 4\pi A\Delta f \int_0^\infty \eta(\nu) \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \quad (14)$$

Then for $\Delta f = 1$ Hz,

$$\text{NEP} = \frac{h\nu}{c\eta} \left[4\pi A \int_0^\infty \eta(\nu) \frac{\nu^2 e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{1/2} \quad (15)$$

or

$$D^*(T, \lambda) = \frac{c\eta}{2\pi^{1/2}h\nu} \left[\int_0^\infty \eta(\nu) \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{-1/2} \quad (16)$$

Assuming $\eta(\nu)$ is independent of frequency but falls back to zero for $\nu < \nu_c$

$$D^*(T, \lambda) = \frac{c\eta^{1/2}}{2\pi^{1/2}h\nu} \left[\int_0^\infty \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{-1/2} \quad (17)$$

Figure 13 shows photon-noise-limited D^* versus cutoff wavelength λ_c for various thermal-background temperatures.¹³ Note that these curves are not independent. $D^*(T, \lambda_c)$ is related to $D^*(T', \lambda'_c)$ by the formula

$$D^*(T, \lambda_c) = \left(\frac{T'}{T} \right)^{5/2} D^*(T', \lambda'_c) \quad \text{where } \lambda'_c = \frac{T}{T'} \lambda_c \quad (18)$$

This relation is useful for determining values of $D^*(T, \lambda_c)$, which do not appear in Fig. 13, in terms of a value of $D^*(T', \lambda'_c)$, which does appear. For example, to find $D^*(1000 \text{ K}, 4 \mu\text{m})$ from the 500-K curve

$$D^*(1000, 4) = \left(\frac{500}{1000} \right)^{5/2} D^* \left(500, 4 \times \frac{1000}{500} \right) = 2.3 \times 10^9 \text{ Jones} \quad (19)$$

If higher accuracy is desired than can be determined from Fig. 13, one can use the preceding formula in combination with Table 1, which gives explicit values of $D^*(\lambda_c)$ versus λ_c for $T = 295 \text{ K}$.

The effect on D^* of using a narrow bandwidth detection system is illustrated in Fig. 14. Such a system may be configured with a cold narrow bandwidth filter, or with a narrow bandwidth amplifier—in order to limit the background flux noise or the electrical bandwidth noise, respectively. Q refers to the factor of the reduction provided.

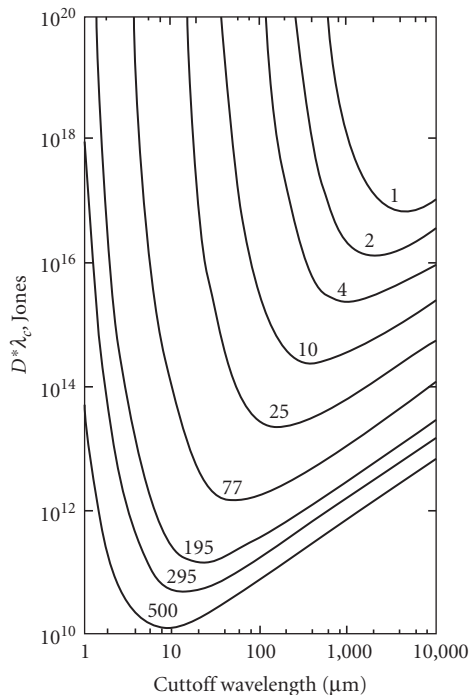


FIGURE 13 Photon-noise-limited D^* at peak wavelength—assumed to be cut-off wavelength—for background temperatures 1, 2, 4, 10, 25, 77, 195, 295, and 500 K (assumes 2π FOV and $\eta = 1$). (Reprinted from Ref. 13.)

TABLE 1 D^* versus λ_c for $T = 295$ K

λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$
1	2.19×10^{13}	10	5.35×10^{10}	100	1.67×10^{11}	1000	1.55×10^{12}
2	4.34×10^{13}	20	5.12×10^{10}	200	3.20×10^{11}	2000	3.10×10^{12}
3	1.64×10^{12}	30	6.29×10^{10}	300	4.74×10^{11}	3000	4.64×10^{12}
4	3.75×10^{11}	40	7.68×10^{10}	400	6.28×10^{11}	4000	6.19×10^{12}
5	1.70×10^{11}	50	9.13×10^{10}	500	7.82×10^{11}	5000	7.73×10^{12}
6	1.06×10^{11}	60	1.06×10^{11}	600	9.36×10^{11}	6000	9.28×10^{12}
7	7.93×10^{10}	70	1.21×10^{11}	700	1.09×10^{12}	7000	1.08×10^{13}
8	6.57×10^{10}	80	1.36×10^{11}	800	1.24×10^{12}	8000	1.24×10^{13}
9	5.80×10^{10}	90	1.52×10^{11}	900	1.40×10^{12}	9000	1.39×10^{13}

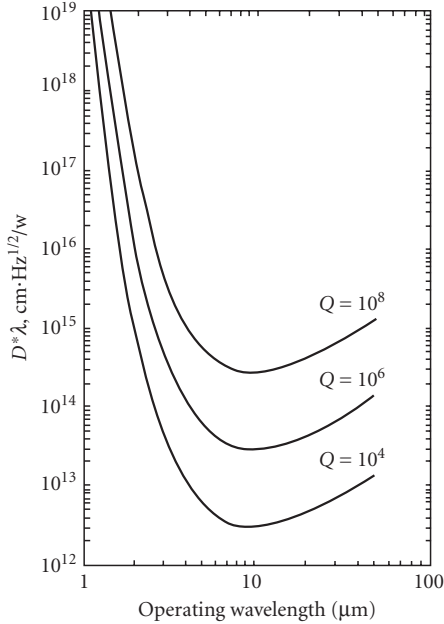


FIGURE 14 Photon noise limit of a narrow-band quantum counter as a function of operating wavelength for a 290-K background, 2π FOV, and $\eta = 1$. (From Ref. 14.)

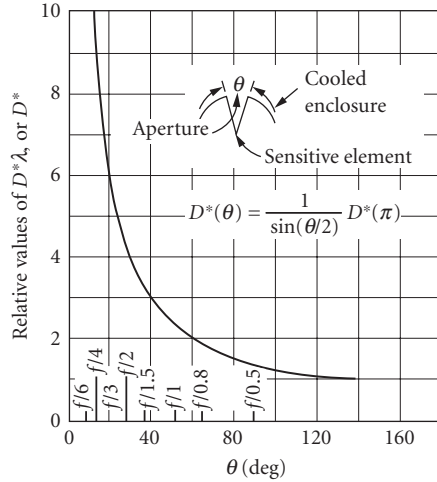


FIGURE 15 Relative increase in photon-noise-limited $D^*(\lambda_{pk})$ or D^* achieved by using a cooled aperture in front of lambertian detector. (From Ref. 14.)

Figure 15 shows the relative increase in photon-noise-limited D^* achievable by limiting the FOV through use of a cooled aperture.

The photon-noise-limited sensitivity shown in Fig. 13 and Table 1 apply to photovoltaic and photoemissive detectors. Figure 14 is for photoconductors. For photoconductors, recombination noise results in a $\sqrt{2}$ reduction in D^* at all wavelengths.

Thermal Detectors Limiting sensitivity of an ideal thermal detector has been discussed previously.^{12,14,15} Assuming no shortwave- or long-wavelength cutoffs exist,

$$D^* = \frac{\varepsilon}{[8\varepsilon\sigma k(T_1^5 + T_2^5)]^{1/2}} = \frac{4 \times 10^{16} \varepsilon^{1/2}}{(T_1^5 + T_2^5)^{1/2}} \text{ Jones} \quad (20)$$

where T_1 = detector temperature
 T_2 = background temperature
 ε = detector emissivity
 σ = Stefan-Boltzmann constant
 k = Boltzmann constant

D^* versus T_2 is plotted for various T_1 in Fig. 16. Figure 17 shows the effect of both short- and long-wavelength cutoffs on bolometer sensitivity,¹⁶ with the ideal photoconductor curve for reference. D^* can be seen to increase rapidly when the cutoff is set to avoid the high flux density from the 300-K background that peaks around 10 μm .

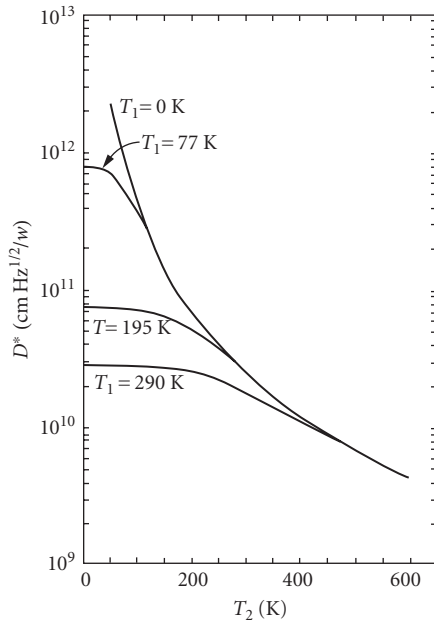


FIGURE 16 Photon-noise-limited D^* for thermal detectors as a function of detector temperature T_1 and background temperature T_2 (2π FOV: $\eta = 1$). (From Ref. 14.)

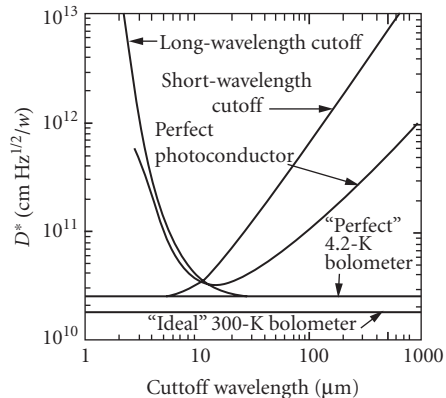


FIGURE 17 The detectivity of a “perfect” bolometer plotted as a function of both short- and long-wavelength cutoffs. Plots for a perfect photoconductor and two other cases are included for comparison. The background temperature is 300 K. (From Ref. 16.)

24.6 OTHER PERFORMANCE PARAMETERS

Spectral Response

Spectral response provides key information regarding how the detector will respond as a function of wavelength or photon energy. Spectral response may be limited by the intrinsic detector material properties, a coating on the detector, or by a window through which the radiation must pass. Relative response is the spectral response ratioed against a detector with a nominally wavelength independent response, such as a thermocouple having a spectrally-broad black coating. Relative response is plotted as a function of wavelength with either a vertical scale of W^{-1} or photon^{-1} . Thermal detectors tend to be spectrally flat in the first case while quantum detectors are generally flat in the second case. The curves are typically shown with the peak value of the spectral response normalized to a value of 1. The spectral response curve can be used together with the blackbody D^* to calculate D^* as a function of wavelength, which is shown in Fig. 18 for selected detectors.

Responsivity and Quantum Efficiency

Responsivity and quantum efficiency are important figures of merit relating to the detector signal output. Responsivity is a measure of the transfer function between the input signal photon power or flux and the detector electrical signal output. Thermal detectors will typically give this responsivity in volts/watt. Photoconductors will usually quote the same units, but will also frequently reference the value to the peak value of relative response per watt from the spectral response curve. This value

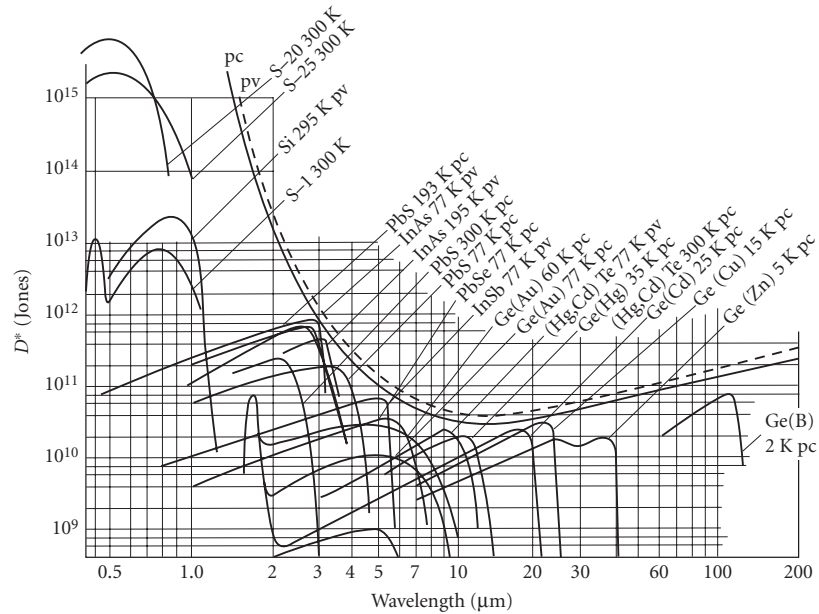


FIGURE 18 D^* versus λ for selected detectors.

is actually realized if the detector bias circuit load resistor is significantly larger than the detector resistance. Photoconductor responsivity is given by:

$$\text{Responsivity}_{\text{peak}} = \frac{\eta q R_d E \tau (\mu_n + \mu_p)}{h \nu \ell} \quad (21)$$

where η is the quantum efficiency, q is the electronic charge, R_d is the detector resistance, E is the electric field, τ is time constant, μ are the mobilities for electrons (n) and holes (p), $h\nu$ is the photon energy, and ℓ is the device length. Photomultiplier tubes and photovoltaic detectors will usually reference the responsivity in amperes per watt, again referenced to peak spectral response.

Detector response performance is also conveyed from the detector quantum efficiency. In the case of photovoltaic detectors which, in the absence of avalanche operation have a gain of unity, quantum efficiency is essentially the current per photon. For a blip photovoltaic detector, the quantum efficiency also determines the D^* . Quantum efficiency is not readily measured for photoconductors and photomultiplier tubes unless the internal gain is carefully calibrated. It is sometimes inferred from the measured D^* for photoconductive devices which are blip—see definition of detective quantum efficiency.

Noise, Impedance, Dark and Leakage Current

Noise has a number of potential origins. Background photon flux-limited detectors have noise dominated by the square root of the number of background photons striking the detector per second [see Eq. (9)]. Other noise sources may contribute or dominate. Among these are

- Johnson, Nyquist, or thermal noise which is defined by the detector temperature and impedance
- Modulation or $1/f$ noise which may dominate at lower frequencies
- Amplifier noise
- Shot noise from dark or leakage current

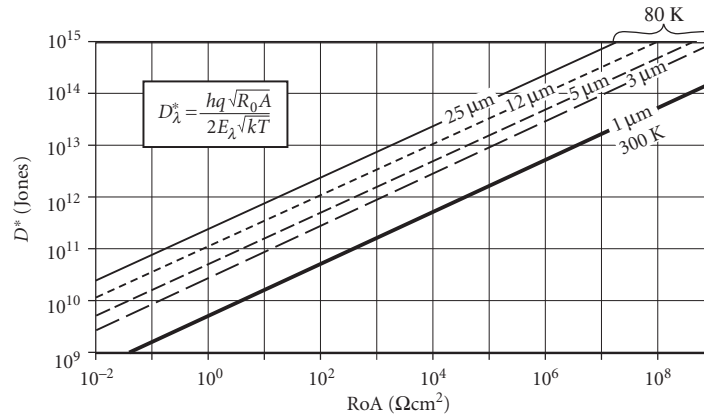


FIGURE 19 Zero-bias impedance-area product (R_oA or shunt resistance per unit area) of a photodiode can limit the D^* as shown. The limiting D^* depends on R_oA , temperature, and photon energy or wavelength. Examples are illustrated for a 1- μm cutoff diode at 300 K, and for 3-, 5-, 12-, and 25- μm cutoff devices at 80 K. Other noise mechanisms, such as photon noise, typically limit D^* to lower values than the highest values shown here.

The impedance of a photodiode may limit performance, depending upon the detector operating conditions. Figure 19 illustrates the diode impedance per unit area (R_oA) limiting value of D^* for silicon detectors at room temperature and longer-wavelength infrared photodiodes at 80 K.

Measurement of the noise as a function of frequency can be valuable for characterizing the relevant noise sources. Selection of an appropriate preamplifier is also critical, particularly for detectors having very low or very high impedance. Integration of preamplifiers together with detectors has significantly improved the overall performance of many detectors. The use of phase-sensitive lock-in amplifiers in combination with a modulated signal can also improve the signal-to-noise ratio.

Uniformity

One cannot assume that the response of a detector will be uniform across its sensitive area. Material inhomogeneity and defects and/or fabrication variables can give rise to nonuniformity. Lateral collection from near the perimeter of a photodiode may give a gradual response decrease away from the edge—this effect will typically be accompanied by a change in response speed as well. Recombination at the electrical contacts to a photoconductor can limit the lifetime, and hence the photoconductive gain, for carriers generated near the contact, a phenomenon called sweep-out. Recombination may be enhanced at surfaces and edges also. Laser spot scanning is useful to check the detector spatial uniformity, although laser sources may not be readily available at all the wavelengths of interest. An alternative method is to move the detector around under a fixed small aperture in conjunction with a light source.

Speed

Detector response speed is often related inversely to detector sensitivity. Thermal detectors often show this characteristic because the signal is proportional to the inverse of response speed, while the noise is amplifier or Johnson limited. Excluding detectors with internal carrier multiplication mechanisms, the best detectors from broad experience seem limited to a D^*f^* product of a few times 10^{17} Jones Hz. D^*f^* may be proportionally higher for devices with gain, since speed can be increased to a greater extent by using a lower value of load resistance without becoming Johnson-noise-limited. The user should be aware that with many detectors it is possible to operate them in a circuit to maximize

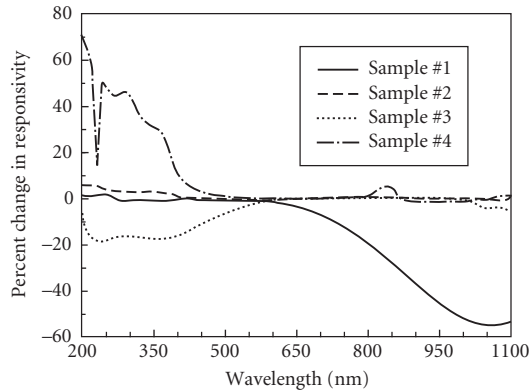


FIGURE 20 After only 3 hours of UV irradiation, these silicon detectors showed great variations in responsivity. (Reprinted from the September 1993 issue of *Photonics Spectra*, © Laurin Publishing Co., Inc.)

sensitivity or speed, but not both at the same time. Speed may vary across the sensitive area of the detector and with temperature, wavelength, and electrical bias.

Stability

Detector performance may change or drift with time. Changes in operating temperature, humidity, and exposure to elevated temperatures as well as to visible, ultraviolet, and high-energy radiation can affect device operation. These effects arise from the temperature dependence of electronic properties in solids, as well as from the critical role played by electrical charge conditions near the surface of many device types. Sensitivity changes in a sample of silicon detectors from four vendors illustrate this point. Wide variations in responsivity change after UV exposure, as shown in Fig. 20. In applications where stability is of significant concern, these effects must be carefully reviewed along with the detector supplier's experience in these matters.

24.7 DETECTOR PERFORMANCE

Manufacturers' Specifications

Table 2 lists the detector materials covered in Sec. 24.7.

TABLE 2 Detector Materials Covered in Sec. 24.7

Thermocouple	GaAsP	Ge:Au
Thermopile	CdS, CdSe	HgCdTe
Thermistor bolometer	CdTe	PbSnTe
Pyroelectric	GaAs	Ge:Hg
InSb hot electron bolometer	Si	Si:Ga
Ge:Ga bolometer	InGaAs	Si:B
Photoemissive	Ge	Ge:Cu
GaN, InGaN	PbS	Ge:Zn
SiC	InAs	Ge:Ga
TiO ₂	PbSe	Photographic
GaP	InSb	

Detector sensitivity can be the determining factor in the design and performance of a sensor system. Detector performance is subject to the development of improved materials, fabrication techniques, and the ingenuity of device engineers and inventors. The descriptions given here may improve with time, and consultation with manufacturers and users is recommended. Today, the internet can be the quickest and most up-to-date source of currently available manufacturers and specifications for the devices they offer. Many suppliers noted in this section may have gone out of business—a search on the internet is the best choice for finding active vendors. Other than a general Web search, some collections of device suppliers can be found at

<http://www.photonics.com/bgHome.aspx>

<http://laserfocusworld.365media.com/laserfocusworld/search.asp>

<http://www.physicstoday.org/ptbg/search.jsp>

Thermocouple The thermocouple offers broad uniform spectral response, a high degree of stability, and moderate sensitivity. Its slow response and relative fragility have limited its use to laboratory instruments, such as spectrometers.

Compared with thermistors, thermocouples are slower, require no bias, and have higher stability but much lower impedance and responsivity. This increases the amplification required for the thermocouple; however, the only voltage appearing is the signal voltage, so that the serious thermistor problem of bridge-circuit bias fluctuations is avoided. With proper design, performance should not be amplifier-limited but limited instead by the Johnson noise of the thermocouple. Thermocouples perform stably in dc operation, although the instability of dc amplifiers usually favors ac operation.

The inherent dc stability of thermocouples is attractive for applications requiring no moving parts, and recently a relatively rugged solid-backed evaporated thermocouple has been developed whose sensitivity approaches that of the thermistor bolometer.

Sensitivity: $D^* 1 \times 10^9$ Jones for 20-ms response time; spectral response depends on black coating (usually gold black) (see Fig. 21)

Noise: White Johnson noise, falling off with responsivity (see Fig. 22)

Resistance: 5 to 15 Ω typical

Responsivity: 5 V/W (typical), 20 to 25 V/W (selected)

Time constant: 10 to 20 ms (typical)

Operating temperature: Normally ambient

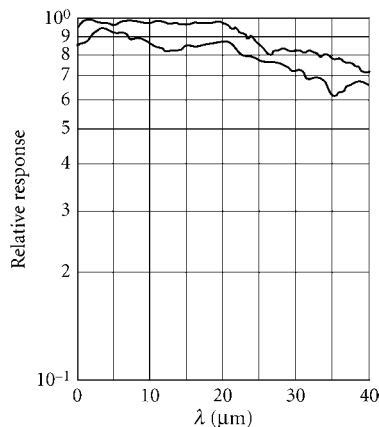


FIGURE 21 Typical thermocouple spectral response curves (CsI window) for two different manufactures.

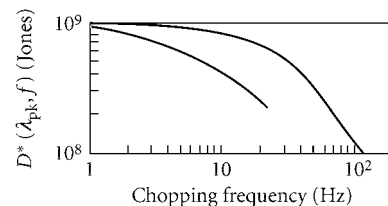


FIGURE 22 Typical thermocouple D^* (noise) frequency response for two manufactures. (From Ref. 1.)

Sensitive area: 0.1×1 to 0.3×3 or 0.6×2 mm (typical)

Linearity: 0.1 percent in region investigated (6×10^{10} to 6×10^8 W incident)

Recommended circuit: Transformer coupled into low-noise (bipolar or JFET) amplifier with good low-frequency noise characteristics

Manufacturers: Perkin-Elmer, Charles Reeder, Beckman Instruments, Farrand, Eppley Laboratory

Thermopile Thermopiles are made by evaporating an array of metal junctions, such as chromel-constantan or manganin-constantan, onto a substrate. The thin-film construction is rugged, but the Coblenz-type may be quite delicate. Wire-wound thermopile arrays are also available which are very robust. Devices with arrays of semiconductor silicon junctions are also available. The array may typically be round, square, or rectangular (for matching a spectrometer slit) and consist of 10 to 100 junctions. Configuration options include matched pairs of junction arrays or compensated arrays to provide an unilluminated reference element. A black coating, such as 3M black or lampblack, is used to provide high absorption over a broad spectral range, as illustrated in Figs. 23 and 24. The housing window may limit the spectral range of sensitivity. Typical applications include power meters and radiometers.

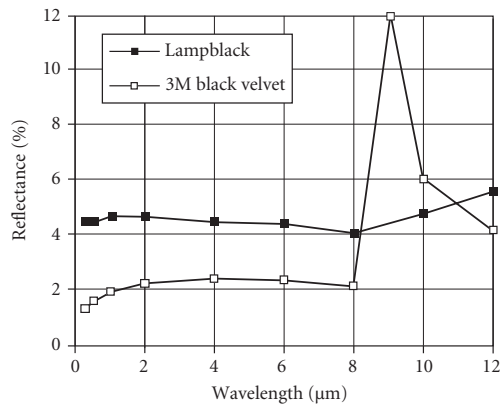


FIGURE 23 Spectral reflectance of two black coatings used in the construction of thermopile detectors. (From Eppley Laboratory studies.)

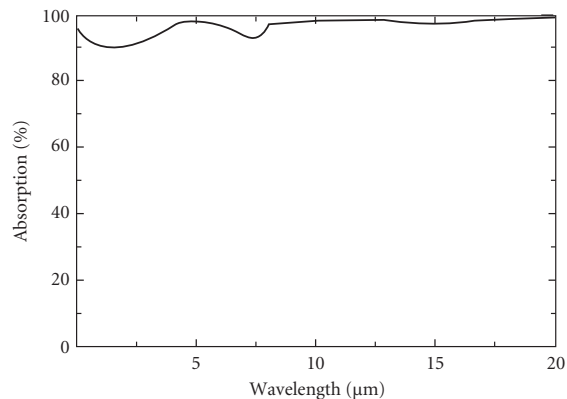


FIGURE 24 Spectral absorption of a thermopile detector coating made from a black metal oxide. (Oriol Corporation.)

Sensitivity: D^* 0.5 to 4×10^8 Jones for 30-ms typical response time. D^* may be dependent upon sensitive area; spectral response depends on black coating and window (see Fig. 23 and Fig. 24)

Noise: White Johnson noise, falling off with responsivity, typical range is 5 to 30 nV/Hz^{1/2}

Resistance: 2Ω to 60 kΩ typical

Responsivity: 4 to 250 V/W (typical) depends on the number of junctions and time constant

Time constant: 10 ms to 2 s (typical)

Operating temperature: Normally ambient

Sensitive area: 0.5 to 6-mm diameter, 0.025×0.025 to 3×3 mm, various rectangular, 0.4×3 , 0.6×2 , 0.6×4 mm (typical)

Recommended circuit: Low noise ($0.5 V_{p-p}$, dc to 1 Hz), low drift with voltage gain of 1000 and input impedance of 1 MΩ

Manufacturers: Armtech, Beckman Instruments, Concept Engineering, Dexter Research Center, Edinburgh Instruments, Eppley Laboratory, Farrand, Gentech, Molectron Detector, Ophir Optronics, Oriol, Scientech, Scitec, Swan Associates

Thermistor Bolometer Thermistors offer reliability, moderate sensitivity, and broad spectral coverage without cooling. Construction is rugged and highly resistant to vibration, shock, and other extreme environments. Response is slower than 1 ms, and trade-off exists between speed and sensitivity.

Thermistor elements are made of polycrystalline Mn, Ni, and Co oxides. In their final form they are semiconductor flakes 10 μm thick, which undergo a temperature resistance change of ~4 percent per Kelvin. Since thermistor resistance changes with ambient temperature enough to alter the biasing significantly, it is usually operated in a bridge circuit, with a nearly identical thermistor shielded from signal radiation and used for a balance resistor.

Sensitivity:

$$\text{NEP} = 8.9 \times 10^{-10} \sqrt{\frac{A(\text{mm}^2)}{\tau_{\text{rms}}}} \text{ W}$$

$$D^* = 1.1 \times 10^9 \sqrt{\tau_{\text{rms}}} \text{ Jones}$$

Spectral response: Depends on coating (usually Zapon lacquer); see Fig. 25.

Quantum efficiency: Depends on blackening coating, typically 80 percent.

Noise: Thermal-noise-limited above 20 Hz ($V_{\text{noise}} = \sqrt{4kTR\Delta f}$); below that, 1/f type noise—see Fig. 26. Used in balanced-bridge circuit (two flakes in parallel); limiting noise due to thermal noise in both flakes.

Resistance: For standard 10-μm-thick flakes, two different resistivities are available: 2.5 MΩ/sq or 250 kΩ/sq. Note that in a bolometer bridge, the resistance between the output connection and ground is half that of single flake.

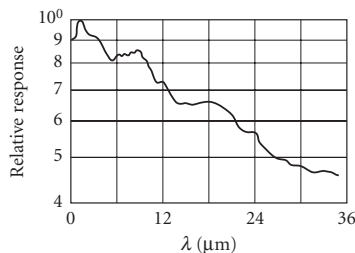


FIGURE 25 Typical thermistor spectral response (no window).

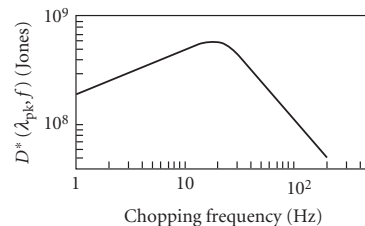


FIGURE 26 Typical thermistor D^* (noise) frequency spectrum. (From Ref. 1.)

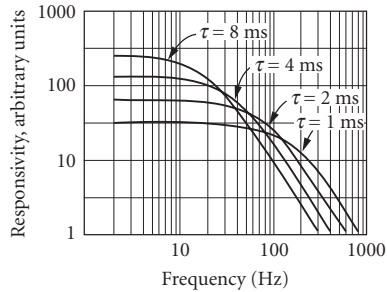


FIGURE 27 Relative response vs. frequency for various time-constant thermistor detectors ($A = 1 \times 1$ mm). (From Barnes Engineering, Bull. 2-100.)

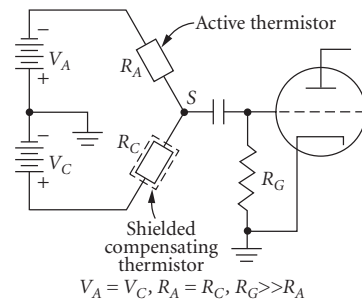


FIGURE 28 Bolometer electrical circuit. (From Barnes Engineering, Bull. 2-100.)

Time constant: τ is 1- to 20-ms standard for nonimmersed detectors and 2- to 10-ms standard for immersed detectors.

Sensitive area: 0.1×0.1 mm, 5×5 mm standard.

Operating temperature: Normally ambient, 285–370 K.

Responsivity: Depends on bias, resistance, area, and time constant $\mathcal{R} \propto \sqrt{R\tau/A} \approx 10^3$ V/W for 0.1×0.1 mm area, 250 k Ω resistance, and $\tau = 4$ ms (see Fig. 27 for frequency-time-constant dependence with given area). Output voltage (responsivity) can be increased to a limited degree by raising bias voltage. Figure 29 shows the deviation from Ohm's law due to heating. Bias should be held below 60 percent of peak voltage. Listed responsivity is that of active flake. In the bridge circuit, responsivity is half this value.

Sensitivity profile: Approximately 10 percent for 10- μ m scan diameter over a 1×1 mm cell.

Linearity: ± 5 percent 10^{-6} to 10^{-1} W/cm 2 .

Recommended circuit: See Figs. 28 and 29.

Manufacturers: Servo Corporation of America, Thermometrics.

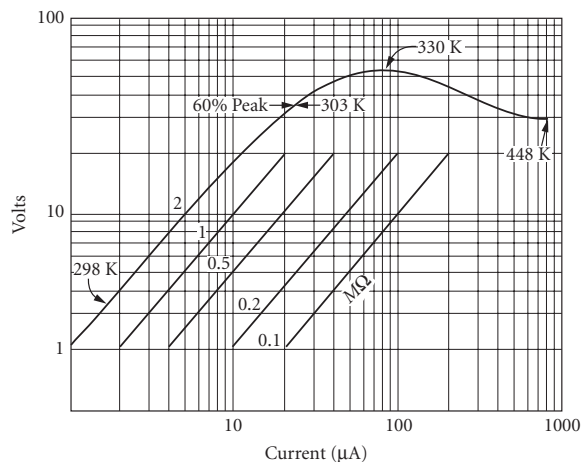


FIGURE 29 Thermistor voltage-current characteristics, showing typical flake temperatures under different conditions. (From Barnes Engineering, Bull. 2-100.)

Pyroelectric Lithium tantalate (LiTaO₃), triglycine sulfate (TGS), and other pyroelectric materials provide an uncooled thermal detector with good sensitivity. The devices are capable of fast response, limited inversely by the preamplifier feedback resistance, but with responsivity and D^* traded for speed. This detector's principle of operation is the pyroelectric effect, which is the change of electric polarization with temperature. Pyroelectric detectors offer rugged construction and the absence of $1/f$ noise because no bias is involved.

Lack of $1/f$ noise, combined with the ability to easily trade off speed and sensitivity, makes pyroelectric detectors useful for scanning applications and energy measurement of pulsed optical sources. In addition, the NEP is independent of area at low frequencies (10 Hz), so that these detectors are useful for large-area applications (preamplifier $1/f$ noise may limit, however). Pyroelectric detectors are useful for calorimetry since the pyroelectric effect is an integrated volume effect and the output signal is unaffected by spatial or temporal distribution of the radiation, up to damage threshold or depolarizing temperature. For higher damage thresholds, lead zirconate titanate ceramic (Clevite PZT-5) exhibits a much smaller pyroelectric effect than TGS, but its high Curie temperature of 638 K makes it more useful than TGS for high-energy applications.

Sensitivity: Sensitive from ultraviolet to millimeter wavelengths. For $\lambda < 2 \mu\text{m}$, TGS must be blackened, which slows response. Normally ($\lambda > 2 \mu\text{m}$) a transparent electrode is used, since TGS absorption is high from 2 to 300 μm . Beyond 300 μm , poor absorption and increased reflectivity reduce sensitivity. Spectral response depends largely on coating. See Fig. 30 for spectral response with modified 3M black. Figure 31 illustrates the relative spectral response for a LiTaO₃ device.

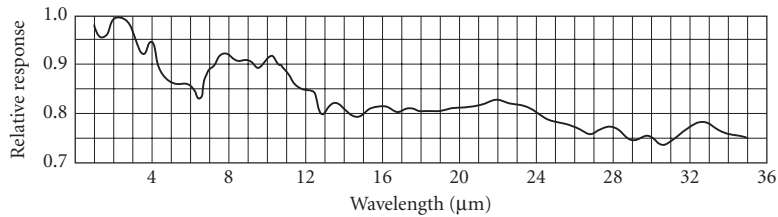


FIGURE 30 Relative spectral response of TGS detectors with modified 3M black. (From Barnes Engineering, Bull. 2-100.)

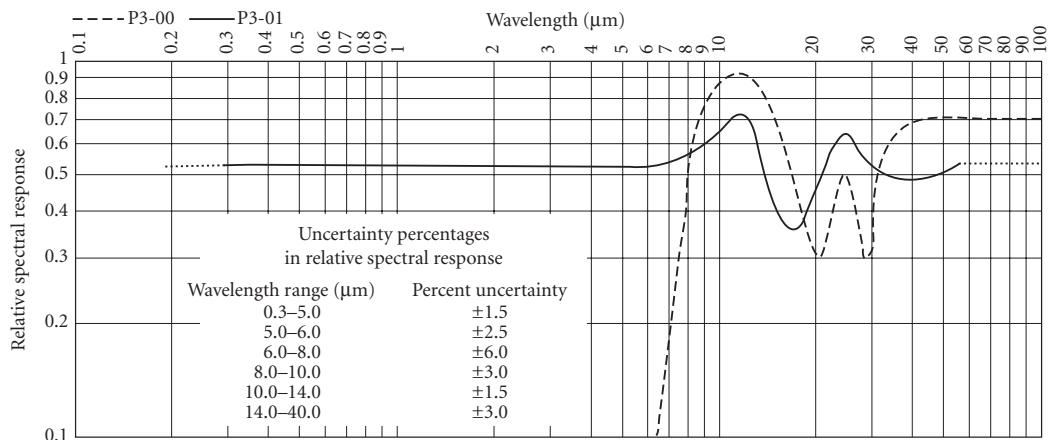


FIGURE 31 Relative spectral response of LiTaO₃ pyroelectric detectors showing both a black spectral coating and an optional coating tuned to the 8 to 14- μm LWIR band. (From Molectron Detector, Inc.)

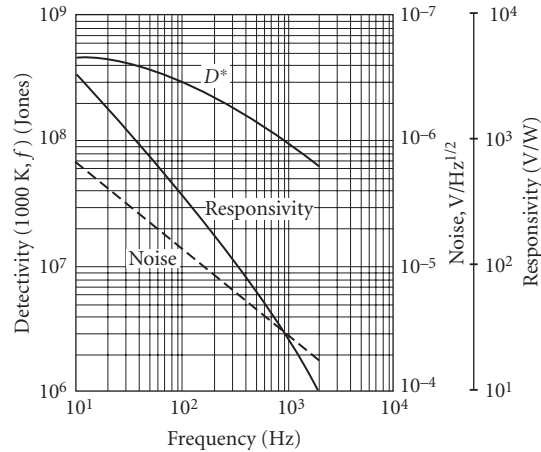


FIGURE 32 Typical D^* , responsivity, and noise versus frequency for TGS ($A = 1 \times 1$ mm; $T = 296$ K). (From Barnes Engineering, Bull. 2-220A.)

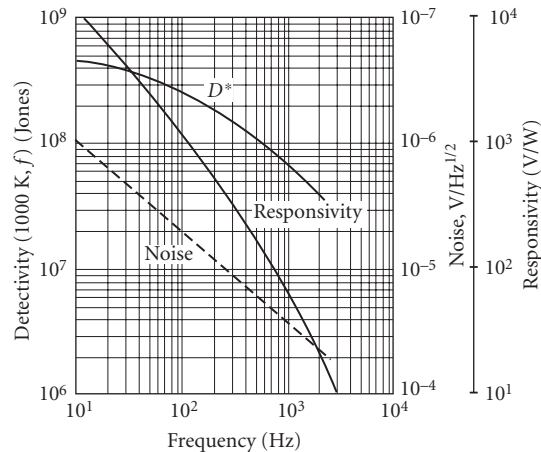


FIGURE 33 Typical D^* , responsivity, and noise versus frequency for TGS ($A = 1 \times 1$ mm; $T = 296$ K). (From Barnes Engineering, Bull. 2-220A.)

D^* is independent of A at low frequencies (10 Hz) (see Figs. 32 and 33). Figure 34 shows NEP versus A for various frequencies.

Quantum efficiency: Depends on coating absorptivity (for 3M black typically $\eta > 75$ percent).

Noise: (See Figs. 32 and 33). Limited by loss-tangent noise up to frequencies that become limited by amplifier short-circuit noise (see Fig. 35).

Operating temperature: Ambient, up to 315 K. Can be repolarized if $T > T_{\text{curie}} = 322$ K for TGS. Irreversible damage at $T = 394$ K (see Fig. 36). Other pyroelectric materials have significantly higher Curie temperatures (398 to 883 K).

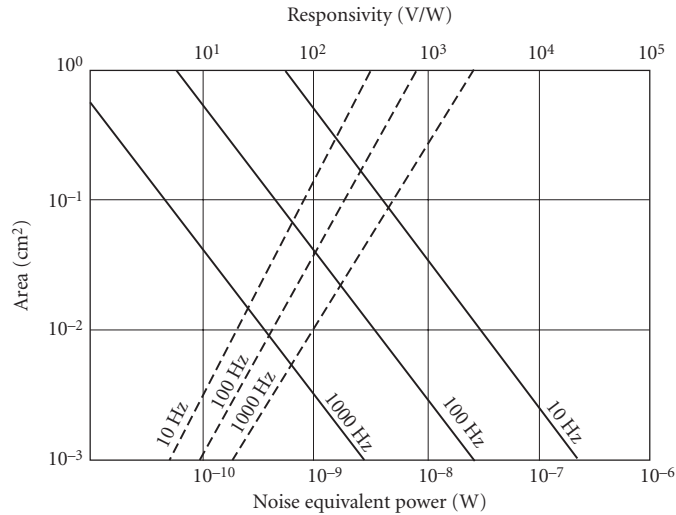


FIGURE 34 Noise equivalent power in watts (broken lines) and responsivity (solid lines) versus TGS detector area for various frequencies. (From Barnes Engineering, Bull. 2-220B.)

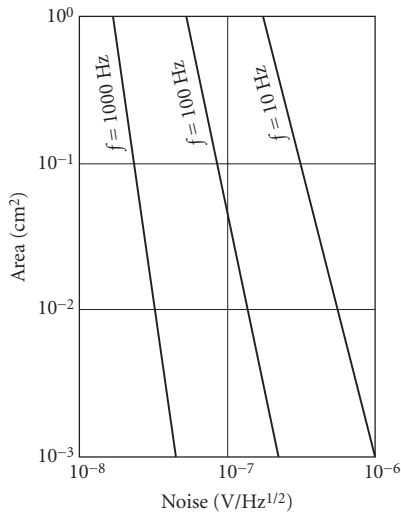


FIGURE 35 TGS noise versus detector area for various operating frequencies. (From Barnes Engineering, Bull. 2-100.)

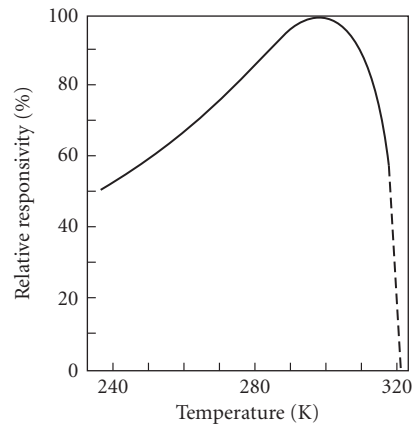


FIGURE 36 Relative responsivity versus temperature for TGS. (From Barnes Engineering, Bull. 2-220B.)

Output impedance: 50 Ω to 10 KΩ, set by built-in amplifier (see Fig. 37).

Responsivity: See Figs. 32 to 34, 36, and 38.

Capacitance: 5 pF for 0.5 × 0.5 mm; 20 pF for 1 × 1 mm; 100 pF for 5 × 5 mm.

Sensitive area: 2 to 50-mm diameter round, 0.5 × 0.5-mm to 10 × 10-mm square, typical.

Time constant: Not pertinent, response speed set by the preamplifier feedback resistor (see Fig. 38).

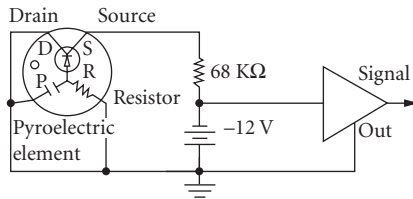


FIGURE 37 Pyroelectric detector amplifier circuit. (From Barnes Engineering, Bull. 2-220A.)

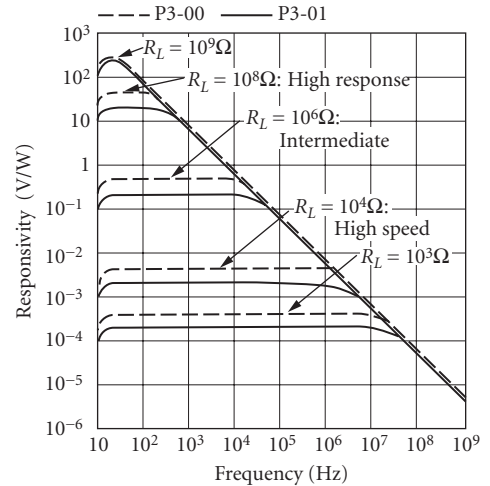


FIGURE 38 Responsivity vs. frequency for two pyroelectric models using various feedback resistors. (From Molecron Detector, Inc.)

Linearity: 5 percent between 10^{-6} and 10^{-1} W/cm².

Sensitivity profile: Depends on coating or transparent electrode; 5 to 7 percent across 12×12 -mm; spot size $< 250 \mu\text{m}$.

Recommended circuit: See Figs. 37 and 39. Field effect transistor (FET) amplification stage usually built in. Since $\mathcal{R} \propto 1/f$, use of an amplifier with $1/f$ noise and gain $\propto f$ is recommended. Then output signal and signal-to-noise ratio are independent of frequency.

Manufacturers: Alrad Instruments, Belov Technology, CSK Optronics, Delta Developments, EG&G Heimann, Electro-Optical Systems, Eltec, Gentec, Graseby, International Light, Laser Precision, Molecron Detector, Oriel, Phillips Infrared Defence Components, Sensor-Physics, Servo Corporation of America, Spiricon, Thermometrics.

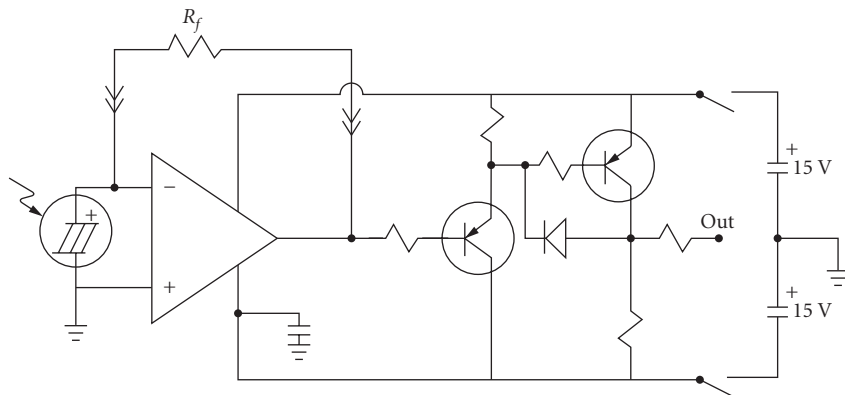


FIGURE 39 Resistive feedback circuit used with LiTaO₃ detectors. (From Molecron Detector, Inc.)

InSb Hot-Electron Bolometer At temperatures of liquid helium and lower, free carriers in indium antimonide (InSb) can absorb radiation in the far-infrared and submillimeter spectral region. Because the mobility of the electrons varies as $T_e^{3/2}$ under these conditions, the conductivity of the material is modulated. This mechanism offers submicrosecond response and broad far-infrared coverage out to millimeter wavelengths but requires liquid-helium cooling and very sophisticated receiver design.

Technically, these devices may be classed as bolometers, since incident radiation power produces a heating effect which causes a change in free-charge mobility. In the normal bolometer, the crystal lattice absorbs energy and transfers it to the free carriers via collisions. However, in InSb bolometers incident radiation power is absorbed directly by free carriers, the crystal lattice temperature remaining essentially constant. Hence the name electron bolometer. Note that this mechanism differs from photoconductivity in that free-electron mobility rather than electron number is altered by incident light (hence there is no photoconductive gain).

Sensitivity: $D^*(2 \text{ mm}, 900) = 4 \times 10^{11}$ Jones (see Fig. 40).

Noise: See Figs. 41 and 42.

Responsivity: 1000 V/W.

Time constant: 250 ns.

Sensitive area: 5×5 mm typical.

Operating temperature: 1.5 to 4.2 K.

Impedance: Without bias, 200Ω ; optimum bias, 150Ω , depends on bias (see Fig. 29).

Recommended circuit: Optimum bias 0.5 mA (see Fig. 43).

Manufacturer: Infrared Laboratories.

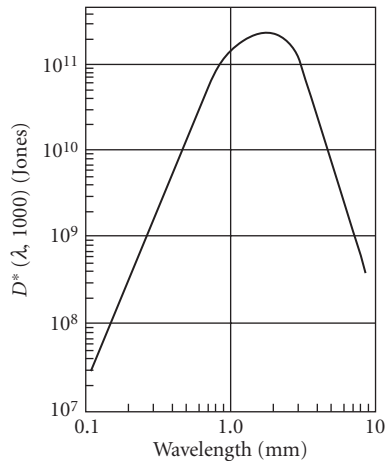


FIGURE 40 D^* versus λ for InSb electron bolometer ($H = 0$). (From Raytheon, *IR Millimeter Wave Detector*, 1967.)

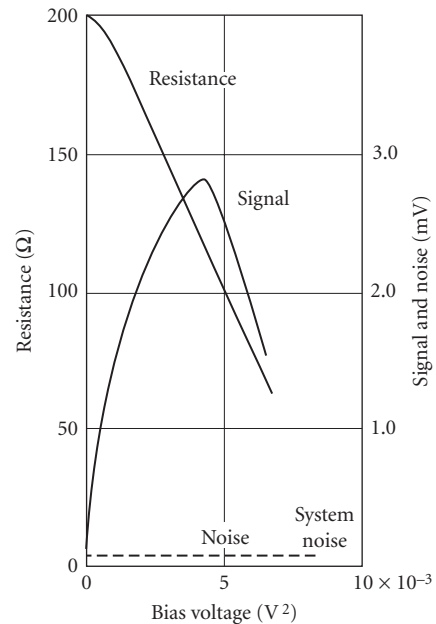


FIGURE 41 InSb electron bolometer, typical resistance, signal, and noise versus bias voltage squared ($T = 5 \text{ K}$; $R_L = 200 \Omega$; gain $= 2.4 \times 10^4$; $F = 1100 \text{ Hz}$). (From Santa Barbara Research Center, *Prelim. Res. Rep.*, 1967.)

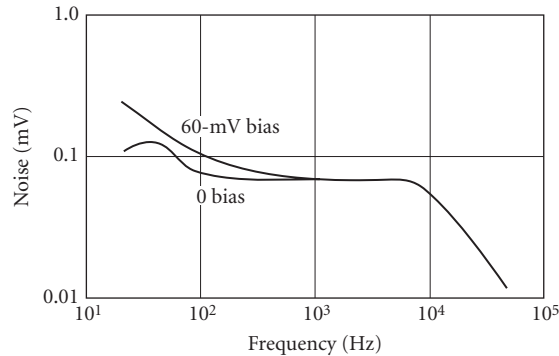


FIGURE 42 InSb electron bolometer, typical noise spectrum ($T = 5$ K; $R_L = 200 \Omega$; gain = 2.4×10^4 ; $\Delta f = 5.6$ Hz). (From Santa Barbara Research Center, Prelim. Res. Rep., 1967.)

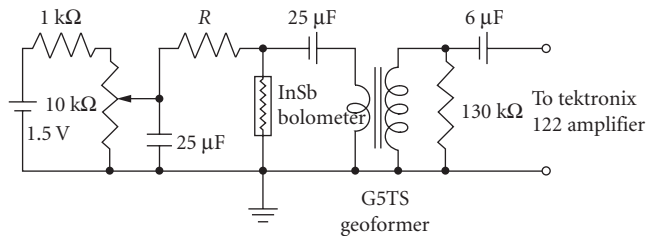


FIGURE 43 Electron bolometer biasing circuit. (From Santa Barbara Research Center, Prelim. Res. Rep., 1967.)

Ge(Ga) Low-temperature Bolometer The Ge(Ga) bolometer offers very high sensitivity and broad spectral coverage in the region 1.7 to 2000 μm . Liquid-helium cooling is required. A trade-off exists between response time (seconds) and sensitivity (10^{-14} -W NEP). Operation at 1000 Hz can be achieved still maintaining 2×10^{-13} W NEP.

Sensitivity: Depends on thermal conductance G (see Figs. 44 and 45). $\text{NEP}(\lambda, 10 \text{ Hz}) = V_n/S = 3 \times 10^{-14}$ W for $G = 1 \mu\text{W/K}$; $A = 1 \text{ mm}^2$, $D^*(\lambda, 10 \text{ Hz}) = 3 \times 10^{13}$ Jones ($Q < 0.2 \mu\text{W}$). For this detector, $\text{NEP} \approx 4T(kG)^{1/2}$ and does not vary with $A^{1/2}$ (T is heat-sink temperature, and k is Boltzmann's constant). Thus D^* cannot be used as a valid means of comparison with other detectors; 300-K background-limited performance is achievable for 2π FOV when the bolometer is operated at 4.25 K with $G = 10^{-3}$ W/K. (For $A = 0.1 \text{ cm}^2$, the time constant is 50 s.)

Responsivity: Typically, responsivity = 2.5×10^5 V/W = $0.7(R/TG)^{1/2}$, where R = resistance. Responsivity, and hence NEP, depends on thermal conductance G , which in turn is set by background power. G ranges from 0.4 to 1000 $\mu\text{W/K}$.

Thermal conductance: Typically $G = 1 \mu\text{W/K}$ for background $Q < 0.2 \mu\text{W}$; note that $Q < 1/2$ (optimum bias power P).

Sensitive area: 0.25 \times 0.25 to 10 \times 10 mm.

Resistance: 0.5 M Ω .

Operating temperature: 2 K (see Fig. 44). In applications where radiation noise can be eliminated there is much to be gained by operating at the lowest possible temperature. Figure 45 shows the theoretical NEP and time constant at 0.5 K, assuming that current noise remains unimportant.

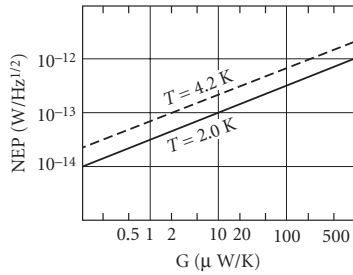


FIGURE 44 Germanium bolometer NEP versus conductance. (*Infrared Laboratories, Inc.*)

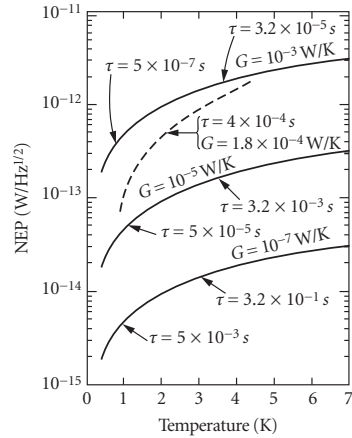


FIGURE 45 Germanium bolometer NEP versus temperature. Solid curves are theoretical: $NEP \approx 4T(kG)^{1/2}$. (*From Ref. 17.*)

Quantum efficiency: Depends on blackened coating and window. For $\lambda < 100 \mu\text{m}$, absorptivity exceeds 95 percent. For $\lambda > 100 \mu\text{m}$, efficiency varies with geometry.

Time constant: Response time constant is proportional to G^{-1} . Therefore, if G must be increased to accommodate larger background, the time constant is decreased proportionally. Responsivity and NEP, however, are degraded as $G^{-1/2}$.

Noise: $V_n = 1 \times 10^{-8} \text{ V/Hz}^{1/2}$; thermal noise is due to R and R_L .

Recommended circuit: Standard photoconductive circuit, with load resistor, grid resistor, and blocking capacitor at low temperature (see Fig. 46). See Fig. 47 for typical electrical characteristics. Bias power $P = 0.1 \text{ TG}$.

Manufacturer: Infrared Laboratories, Inc.

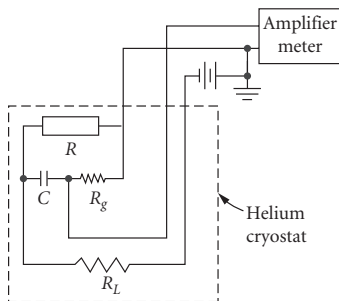


FIGURE 46 Germanium bolometer circuit and cryostat.

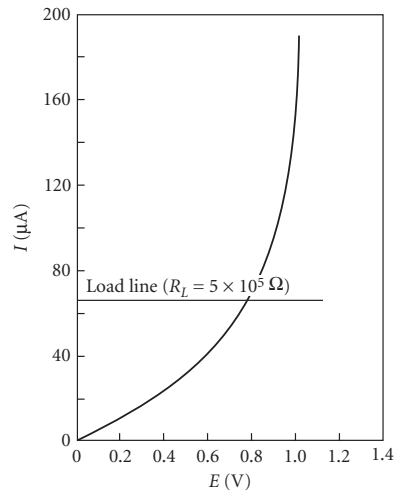


FIGURE 47 Load curve for a typical germanium bolometer ($T = 2.15 \text{ K}$) with load line, showing optimum operating point. (*From Ref. 17.*)

Photoemissive Detectors Photoemissive detector is generally the detector of choice in the UV, visible, and near-IR where high quantum efficiency is available. In the spectral region $\lambda < 600$ nm, the photomultiplier, or multiplier phototube, has close to ideal sensitivity; that is, selected photomultiplier tubes (PMT) are capable of detecting single photon arrivals (but at best only with about 30 percent quantum efficiency) and amplifying the photocurrent (pulse) enormously without seriously degrading the signal-to-noise ratio. Time resolution can be as short as 0.1 ns. Only very specialized limitations have precluded their use for $\lambda < 800$ nm, for example, cost, ruggedness, uniformity of manufacture, or need for still faster response. Recently these limitations have all been met individually but generally not collectively. Where adequate light is available, the simple phototube has advantages over the multiplier phototube in that high voltages are not required, the output level is not sensitive to applied voltage, and dynode fatigue is eliminated.

Microchannel plate tubes (MCPT) are a variant of the photomultiplier tube where the current amplifying dynode structure is replaced by an array of miniature tubes in which the photocathode current is amplified. MCP tubes are more compact than PMTs and are reliable in operating conditions of high environmental stress. The same range of photocathode materials is available in MCPTs as PMTs. MCPTs can provide a wide range of electron gain as available depending upon whether a single MCP or a stack of MCPs is used. The structure of a PMT and MCPT are compared in Fig. 48.

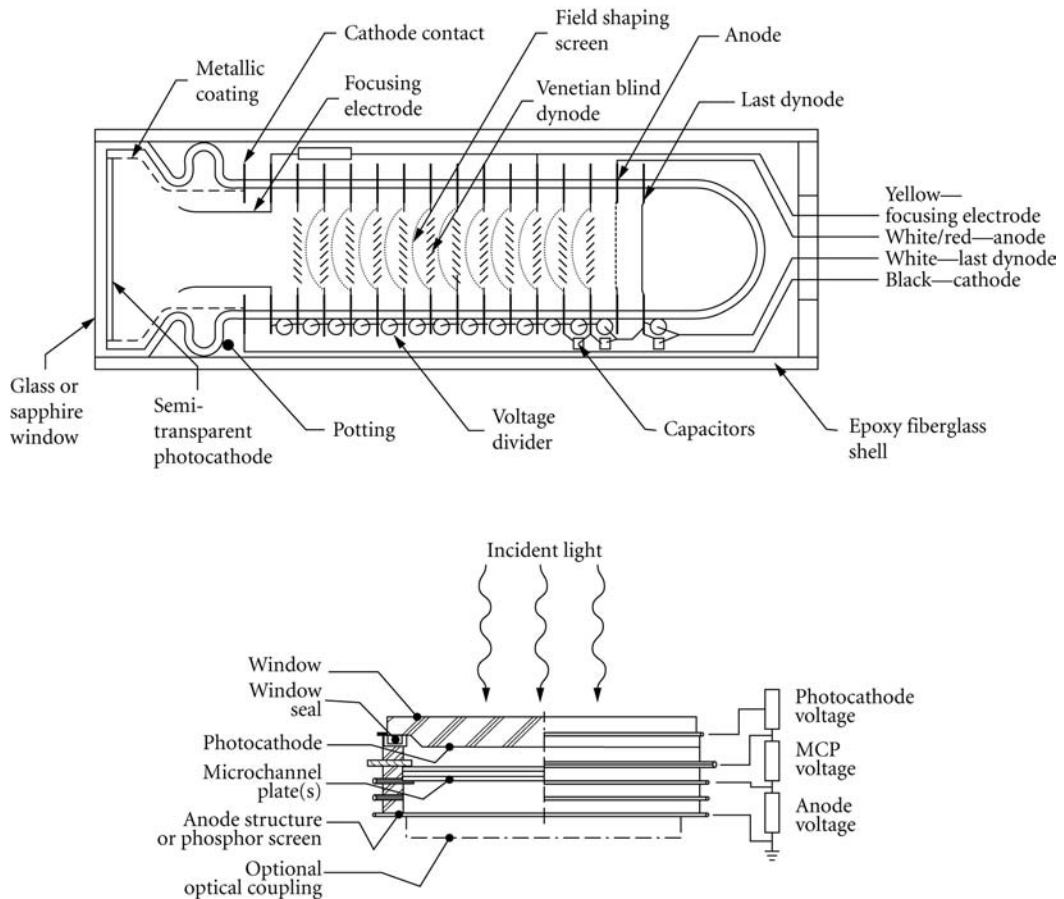


FIGURE 48 Comparison of photomultiplier tube (PMT) and microchannel plate tube (MCPT) construction. (*EMR Photoelectric.*)

Sensitivity In modern phototubes, shot noise due to the cathode dark current is by far the most important noise source. The most common descriptions of phototube sensitivity list both current responsivity (amperes per watt), dark current, and dark noise. Several useful measures of sensitivity are noise equivalent input (NEI) (see Sec. 24.4 “Definitions”), noise equivalent power (NEP), or its reciprocal $D \equiv 1/\text{NEP}$. NEP and NEI in the range 10^{-14} to 10^{-17} W/Hz^{1/2} are not uncommon.

Detectivity is generally limited by dark-current shot noise. Dark current depends on photocathode material, area, and temperature. Thus the best detectivity is obtained with small effective sensitive area. Cooling is especially useful for red-sensitive and near-IR tubes and is generally not worthwhile for others (see “Operating temperature”). Special tube housings which can provide thermoelectric cooling are available.

The spectral response curves shown in Fig. 49 are for the combination of photocathode and window. Historically, this method of description gave rise to the S-response designations, most of which are now obsolete. It is often desirable to separate photocathode response from window transmission. Thus, Fig. 50 shows the quantum efficiency (electrons per incident photon) of a number of photocathodes without window losses. For $\lambda < 400$ nm, each photocathode should maintain its peak quantum efficiency, up to photon energies where multiple photoemissions take place. In Fig. 51, the spectral dependence of quantum efficiency for a variety of modern photocathode/window combinations is illustrated.

D^* is a meaningful figure of merit for phototubes whose sensitivity is limited by dark noise (shot noise on the dark current) and whose emitting photocathode area is clearly defined, but D^* must be used with caution because, although modern phototubes are generally dark-noise-limited devices, they are often limited by noise in signal, that is, the noise content of the signal itself.¹⁸ Serious errors in predicting the detection capability of phototubes will arise if noise in signal is ignored and D^* is presumed to be the important limiting parameter [see Eqs. (6) to (8)]. Very little reliable data are presently available on D^* for photoemitters. However D^* curves for S-1, S-20, and S-25 are shown in Figs. 52 (300 K) and 53 (PMT cooled to 200 K).

Short-wavelength considerations Window considerations are as follows:

For $\lambda > 200$ nm: Windows are essential, as all useful photocathode materials are oxidized and performance would otherwise be destroyed.

$200 \text{ nm} > \lambda > 105$ nm: Photocathode materials are not oxidized by dry air (moisture degrades performance). Windows are optional. LiF windows have shortest known cutoff, 105 nm. For $\lambda < 180$ nm, it is generally advisable to flush with dry nitrogen.

$\lambda < 105$ nm: No windows are available.

Since air absorbs radiation in the region 0.2 to 200 nm (ozone absorbs 200 to 300 nm), it is necessary to include the (windowless) detector in a vacuum enclosure with the source.

A useful technique for avoiding the far-ultraviolet window-absorption problem (provided $\lambda >$ vacuum ultraviolet) is to coat the outside of the window of a conventional PMT with an efficient fluorescent material, for example, sodium salicylate, which absorbs in the ultraviolet and reemits in the blue, and is efficiently detected by most photocathodes.⁹

Solar-blind considerations are as follows—although most photocathodes have high quantum efficiency at short wavelengths, background-noise considerations often preclude their use at very short wavelengths, and very wide bandgap semiconductor photocathodes such as CsI, KBr, Cs₂Te, and Rb₂Te (having peak quantum efficiency a little greater than 10 percent) often give better signal-to-noise ratio. This sacrifice in quantum efficiency to obtain insensitivity to wavelengths greater than those of interest would not be necessary if suitable short-wavelength pass filters were readily available.

For applications where it is desirable that the detector not see much solar radiation, one can use photocathodes whose high work function precludes photoemission for photons of too low an energy. Figure 54 shows quantum efficiency versus λ for three such photocathodes, tungsten, CsI, and Cs₂Te, compared with Cs₃Sb and GaAs(Cs), which are not solar blind.

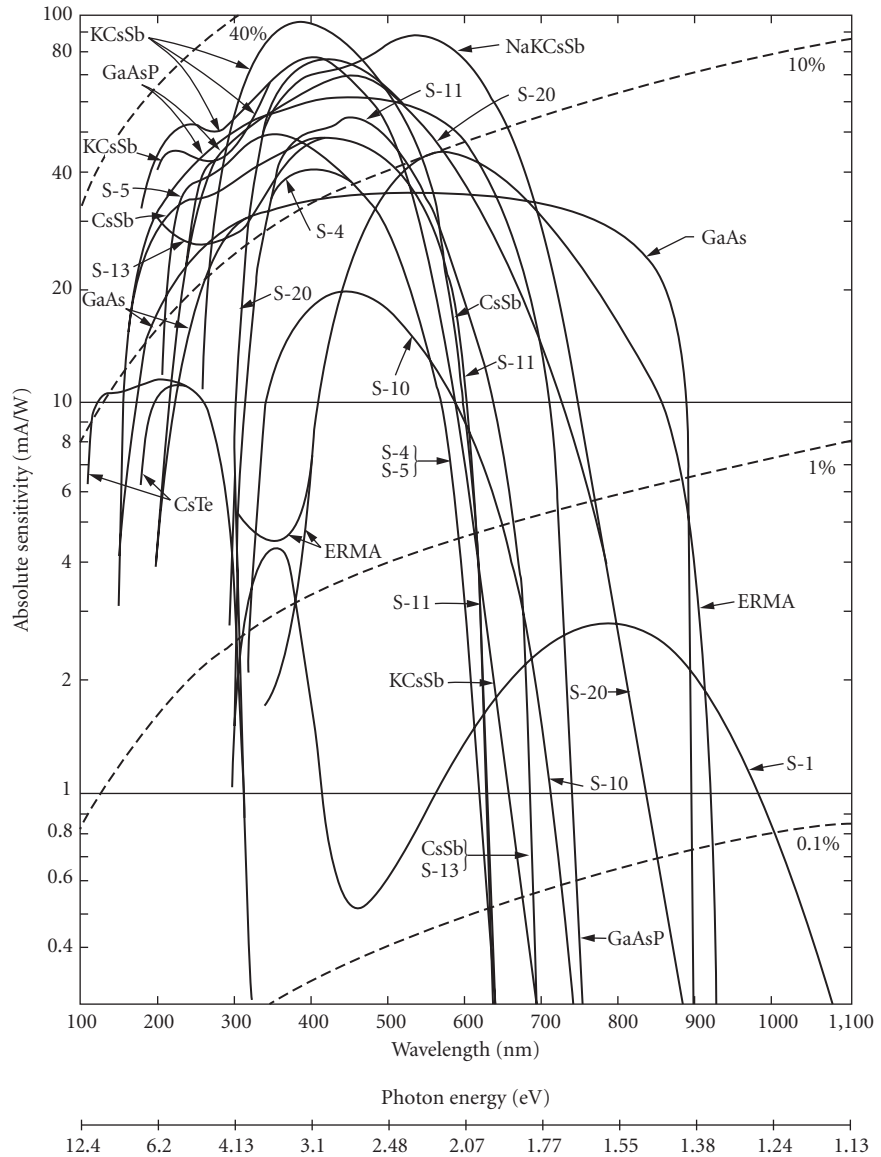


FIGURE 49 Spectral sensitivity of various photoemitters. Dotted lines indicate photocathode quantum efficiency. Chemical formulas are abbreviated to conserve space. S-1 = AgOCs with lime or borosilicate crown-glass window; S-4 = Cs_3Sb with lime or borosilicate crown-glass window (opaque photocathode); S-5 = Ss_3Sb with ultraviolet-transmitting glass window; S-8 = Cs_3Bi with lime or borosilicate crown-glass window; S-10 = AgBiOCs with lime or borosilicate crown-glass window; S-11 = Cs_3Sb with lime or borosilicate crown-glass window (semitransparent photocathode); S-13 = Cs_3Sb with fused-silica window (semitransparent photocathode); S-19 = Cs_3Sb with fused-silica window (opaque semicathode); S-20 = $\text{Na}_2 \text{KCsSb}$ with lime or borosilicate glass window. ERMA = extended red multialkali (RCA; ITT uses MA for multialkali). This curve is representative of several manufacturers' products. Many variations of this response are available, for example, trade-offs between short- and long-wavelength response. (From RCA Electronic Components, chart. PIT-701 B.)

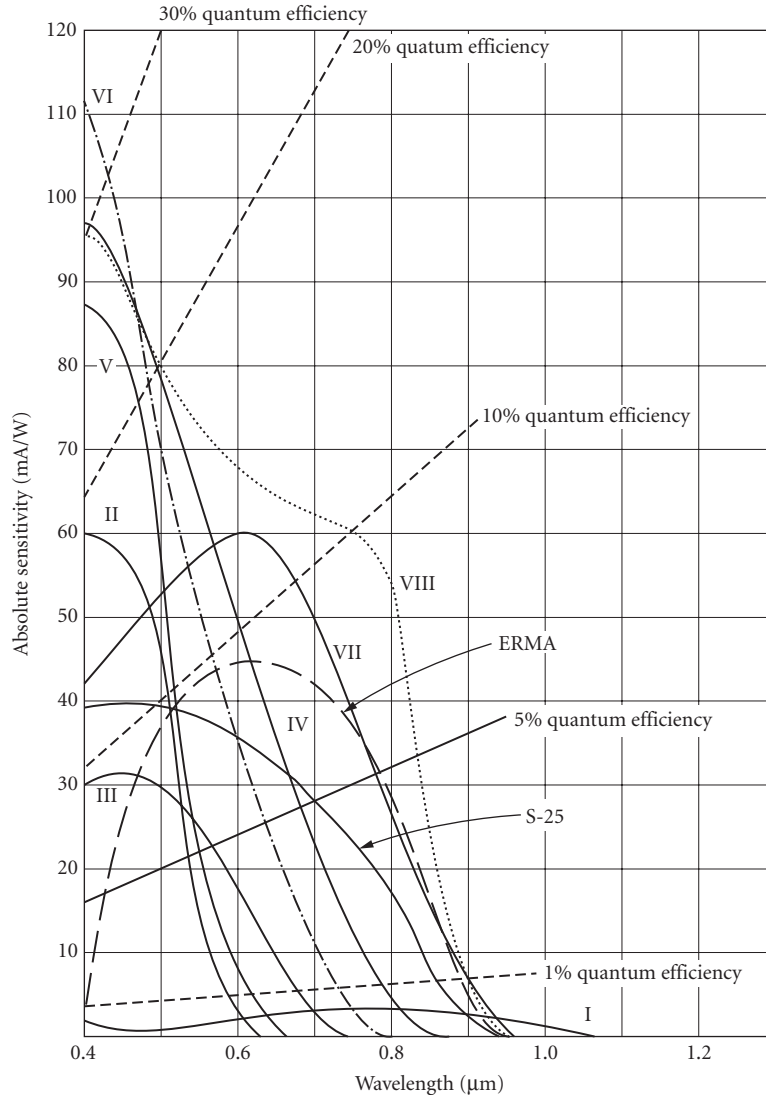
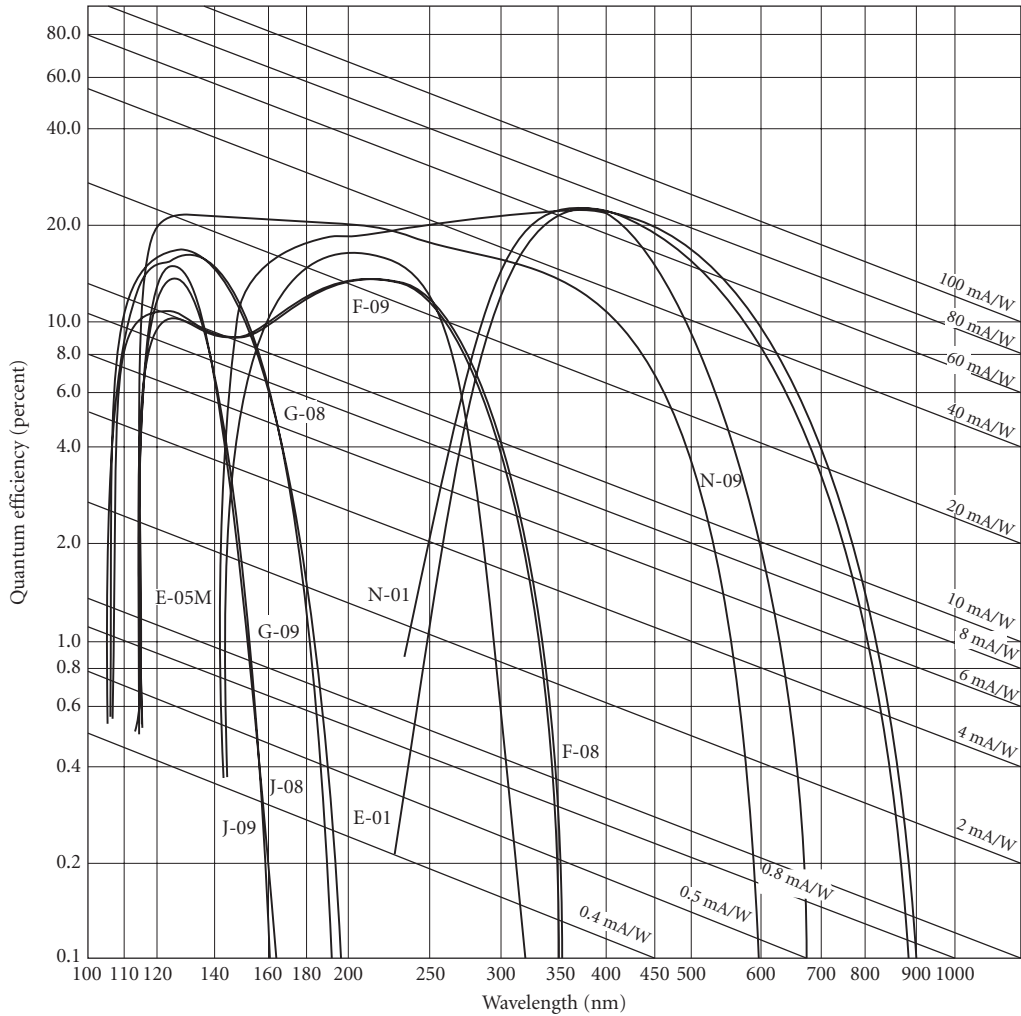


FIGURE 50 Photocathode responsivity and quantum efficiency versus wavelength (no windows). For abbreviations, see Fig. 49. I = S-1, II = S-11, III = S-10, IV = S-20, V = K₂CsSb, VI = K₂CsSb(O), VII = NaKCsSb, VIII = GaAs(Cs). (Based on material from RCA.)

Quantum efficiency Figures 49, 50, 51, and 54 show photocathode spectral quantum efficiency (probability that one photoelectron is emitted when a single photon is incident). Note that there are fairly few basic photocathode materials and that the window often determines effective quantum efficiency at short wavelengths.

For $\lambda < 40$ nm, a wide variety of photocathode materials are available with high quantum efficiency. Many of these materials, such as tungsten, are not destroyed by being subjected to air, so that open structures can be used, consisting of a photocathode multiplier chain without window. The complete windowless structure is then placed in a vacuum enclosure with the source of radiation.



Photocathode key

Key letter	Description	Long-wavelength cutoff (Note 1)	Long-wavelength sensitivity (Note 2)
E	Tri-alkali (S-20)	850 nm	780 nm
F	Cesium telluride	355 nm	340 nm
G	Cesium iodide	195 nm	185 nm
J	Potassium bromide	165 nm	150 nm
N	High-temperature Bi-alkali	690 nm	640 nm
Q	Rubidium telluride	320 nm	300 nm

Window material key

Key no.	Description	Short-wavelength cutoff*
01	Borosilicate Glass	270 nm
05	UV Grade Sapphire	145 nm
08	UV Grade Lithium Fluoride	105 nm
09	Magnesium Fluoride	115 nm

Note 1—Point at which QE becomes 1% (typical) of peak QE.

Note 2—Point at which QE is 1% (typical).

*10% Energy transmission

FIGURE 51 Quantum efficiency of photocathode/window combinations as a function of wavelength. (EMR Photoelectric)

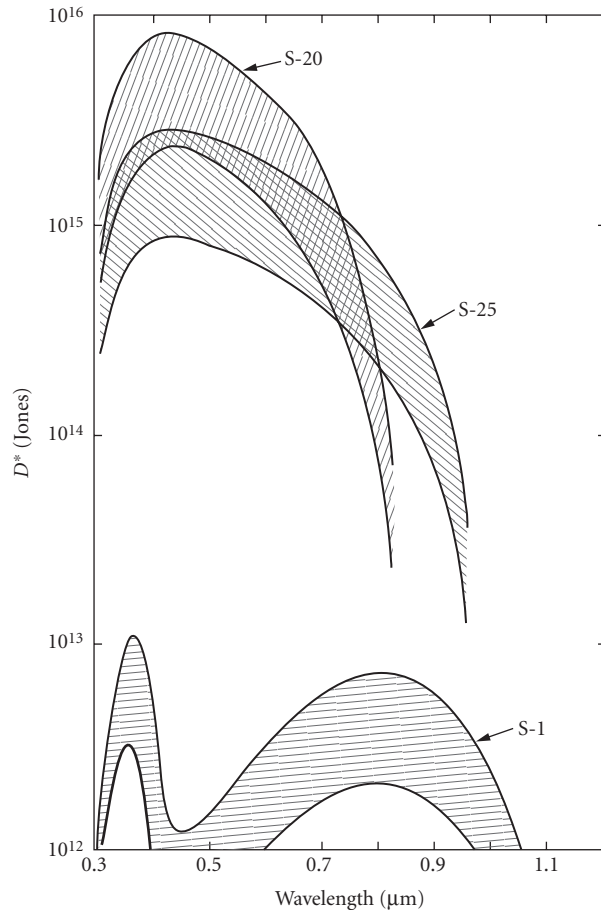


FIGURE 52 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on material from RCA)

The quantum efficiency at any wavelength can be calculated from the formula

$$\eta = \frac{\mathcal{R} \times 1239.5}{\lambda} \quad (22)$$

where \mathcal{R} = photocathode response, A/W, and λ = wavelength, nm.

A useful technique for improving quantum efficiency, reported by Livingston¹⁹ and Gunter,²⁰ involves multipassing the photocathode by trapping the light inside the photocathode using a prism.

Responsivity PMT responsivity depends upon photocathode quantum efficiency and subsequent dynode gain. For most purposes, the dynode gain in a well-designed PMT introduces no significant degradation in the photocathode signal-to-noise ratio. Figure 49 shows photocathode response expressed in photocurrent (amperes) per incident radiation power (watts).

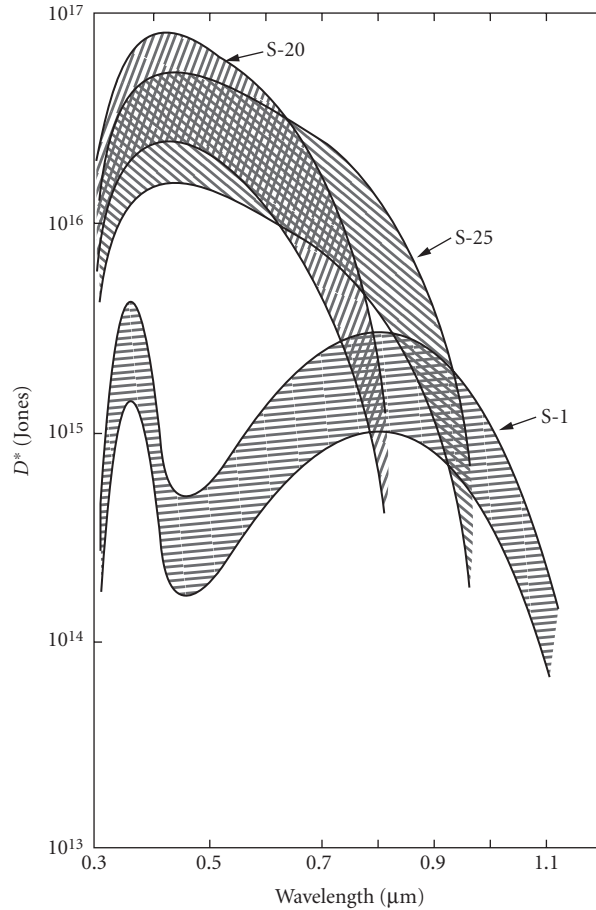


FIGURE 53 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on material from RCA) 2 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors", ITT Industrial Labs, 1969.)

Noise The limiting noise in a PMT depends on the level of illumination. For low-level detection, limiting noise is the shot noise on the dark current,

$$i_n = (2e j_{\text{dark}} \Delta f)^{1/2} \quad (23)$$

For high illumination levels the shot noise on the signal photocurrent

$$i_n = (2e i_{\text{signal}} \Delta f)^{1/2} \quad (24)$$

far exceeds that on the dark current. Manufacturers usually express noise as photocathode dark current or anode dark current for given gain, which is therefore traceable to photocathode dark current.

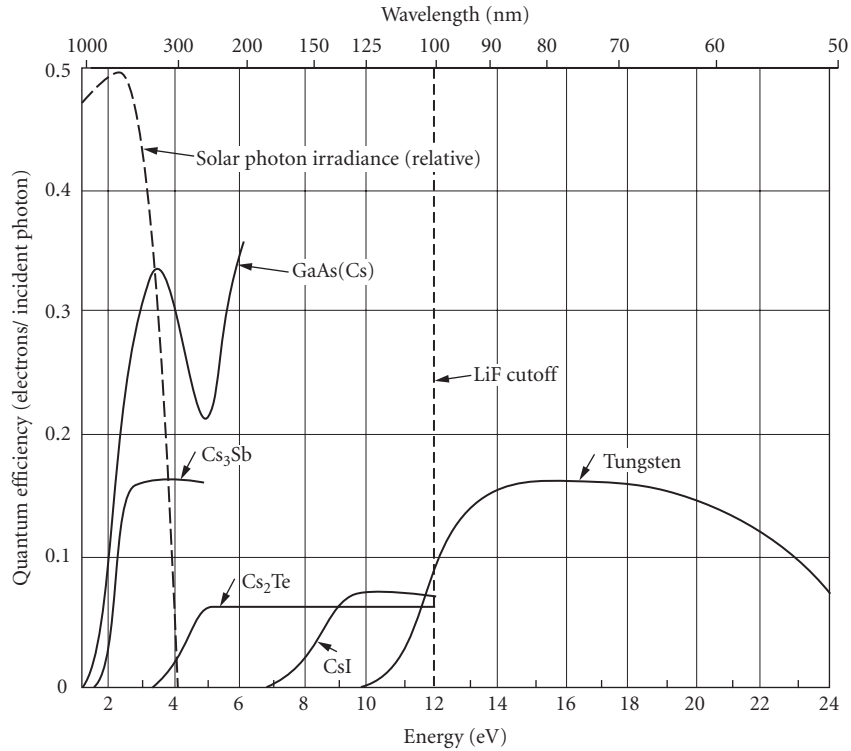


FIGURE 54 Quantum efficiency versus wavelength (photon energy) for several photoemitters.

Photocathode dark current is approximately proportional to photocathode area so that small photocathode effective areas can be expected to have reduced noise. Figure 55 shows how anode dark current and gain increase with applied voltage for a typical PMT.

Minimum detectable power is related to limiting noise through responsivity via

$$\text{NEP} = i_n / \mathcal{R} \quad (25)$$

where \mathcal{R} is in amperes per watt.

Operating temperature Dark current due to thermionic emission, usually greater in red-sensitive tubes, can be reduced by cooling (see Ref. 21). The trialkali (S-20) performance does not benefit from cooling below 255 K. Maximum beneficial cooling (three to four dark counts per second) for AgOCs (S-1), (Cs)Na₂KSb (S-20), and Cs₃Sb (S-11) is 195, 255, and 239 K, respectively. Most photocathodes become noisier as temperature rises above ambient because of increased thermionic emission. Because its thermionic emission starts at a very low value, (Cs)Na₂KSb is a useful photocathode up to temperatures of approximately 373 K.

Response time The rise time of photomultiplier tubes depends chiefly on the spread in transit time during the multiplication process. For photomultiplier tubes, this spread is about 10 ns. Some tubes with specially designed electron optics can give spread as low as 1 ns. The crossed-field PMT makes possible a spread as small as 0.1 ns. Microchannel plate tubes have response times of a nanosecond or less.

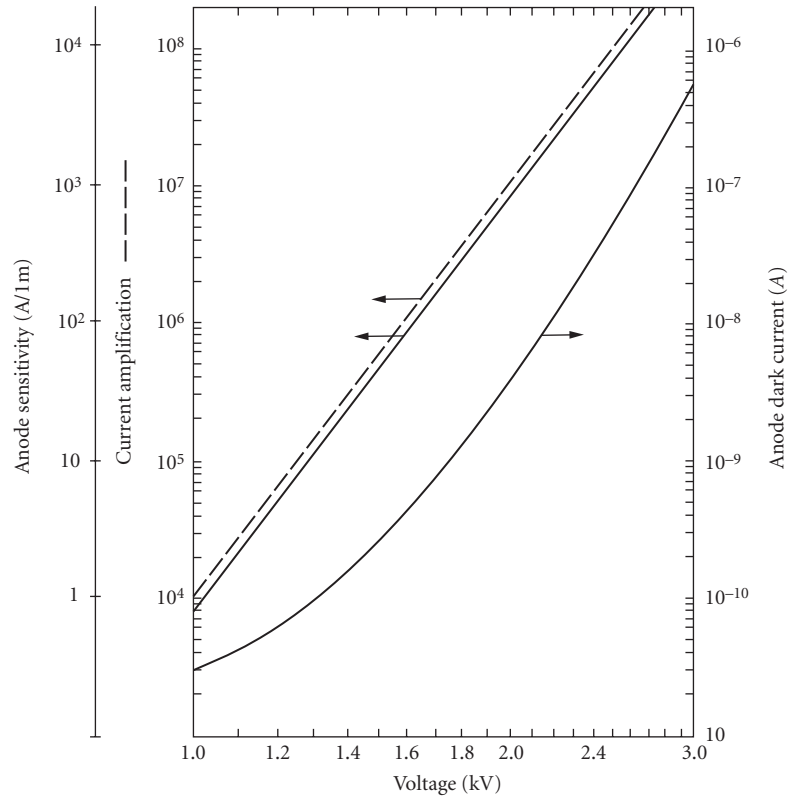


FIGURE 55 Typical current amplification and anode dark current as a function of applied voltage. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors," ITT Industrial Labs, 1969.)

For high-speed work (<1 ns rise time), good transmission-line technique must be used to obtain impedance match and to avoid reflection. The bandwidth of the output circuit will depend upon the total capacitance (PMT circuit plus stray capacitances) and the value of the load resistance.

Linearity Photomultiplier tubes are nearly all linear to about 1 percent for cathode currents of 0.1 μA or less. Some tubes may be linear to better than 0.1 percent but must be individually selected.²² Probably most of the nonlinearity results from the dynode structure.

Sensitive area No fundamental limitation. Only recently available with very small effective areas for extremely low dark current. Magnetic focusing has been used so that only a small fraction of the photocathode is used electron-optically.

Sensitivity profile Usually uniform within 20 to 50 percent. Microchannel plate detectors may have uniformity of ± 5 percent.

Stability PMTs are subject to short- and long-term drift which can depend upon anode current, changes in anode current, storage times, and aging or anode life. They are also subject to change if exposed to magnetic fields or changes in temperature. Vibration of the tube may modulate the signal (microphonic effect).

Recommended circuit See Fig. 56.

1. Since a PMT is a current generator, increasing output resistance R_1 increases output voltage. An upper limit to R_1 may be imposed either by the time-constant limitation or by nonlinearity, which results from a space charge produced near the anode when the anode is left nearly floating electrically.
2. The rated photocathode current (referred to anode current through gain) should not be exceeded.
3. Care should be taken not to destroy the photocathode with light (heating).
4. When large currents are drawn, it may affect later dynode interstage voltage and hence gain, causing nonlinearity; for example, in Fig. 44, if the photocurrent from DY 10 to anode becomes comparable to the biasing current, through R_{11} , the gain of the final stage is reduced. This can be avoided by biasing the dynodes with constant-voltage sources.
5. To avoid dynode damage, final dynode current must not exceed the value suggested by the manufacturer.

Photon counting At the photocathode, the shot-noise-limited signal-to-noise ratio (with negligible dark current) is

$$\frac{i_s}{i_n} \frac{i_s}{(2e i_s \Delta f)^{1/2}} \left(\frac{i_s}{2e \Delta f} \right)^{1/2} = \left(\frac{N_s}{2 \Delta f} \right)^{1/2} \quad (26)$$

where N_s is the photoelectron rate at the photocathode. Thus, for extremely low levels of illumination, the ideal signal-to-noise ratio becomes very poor. At this point there is much to be gained by abandoning attempts to measure the height of the fluctuating signal (Fig. 57a) in favor of digitally recording the presence or absence of individual pulses (Fig. 57e).

Single photoelectron counting can be achieved by using a pulse amplifier (see Fig. 58), which suppresses spurious dark-noise pulses not identical in amplitude and shape to those produced by photoelectrons.

An upper practical limit for (random) photon counting is set by convenient amplifier bandwidths at about 10^5 s^{-1} . For reasonable (1 percent) statistical accuracy, this implies a 10-MHz bandwidth.

Gallium phosphide dynodes The development of GaP dynodes for increased secondary-electron production^{23,24} makes possible unambiguous discrimination of small numbers of individual photoelectron counts which was not previously possible with lower dynode gains. This is shown in Fig. 59, where the spread in number of secondary electrons ($N \times \text{gain}$) is just $(N \times \text{gain})^{1/2}$.

In addition to the aforementioned fundamental advantage of high dynode gain, the large gain per stage in the first dynodes also helps discriminate against noise introduced by later stages of amplification. Also, fewer stages of amplification are required.

Manufacturers ADIT, EMR Photoelectric, Bicon, Burle, Edinburgh Instruments, Galileo Electro-Optics, Hamamatsu, K and M Electronics, id Quantique, International Light, Optometrics USA, Oriol, Phillips Components, Photek, Photon Technology, Photonis, Penta Laboratories, Thorn EMI, Varo.

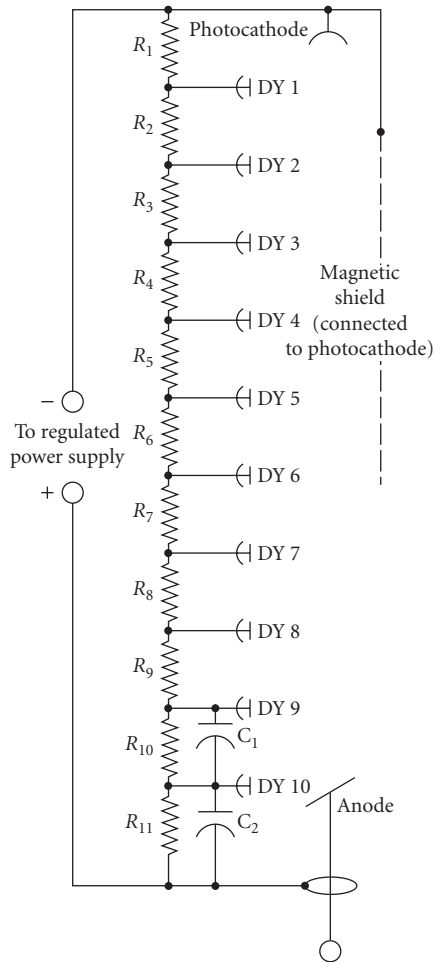


FIGURE 56 $C_1 = 68 \text{ pF} \pm 10 \text{ percent}$, 500 V (dc working); $C_2 = 270 \text{ pF} \pm 10 \text{ percent}$, 500 V (dc working); $R_1 = 220 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; $R_2 = 240 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; $R_3 = 330 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; R_4 to $R_{11} = 220 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors," ITT Industrial Labs, 1969.)

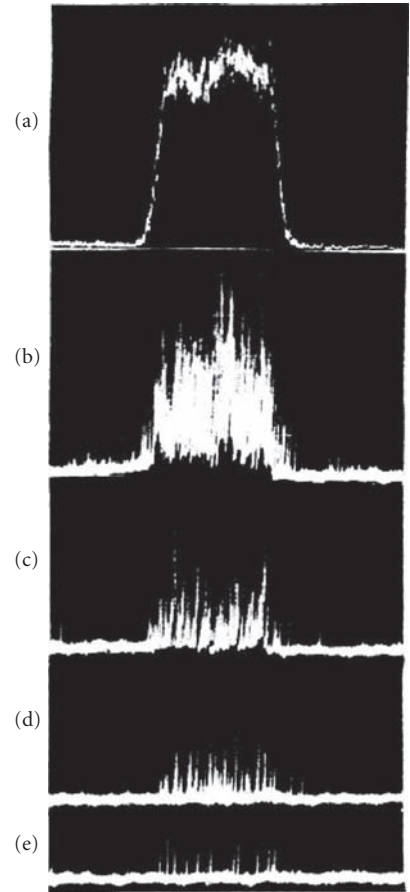


FIGURE 57 Oscilloscope presentation of PMT output when reviewing square-wave chopped light pulse. In (a) to (e) the intensity is reduced and gain is increased commensurately. (Courtesy of E.H. Eberhardt.)

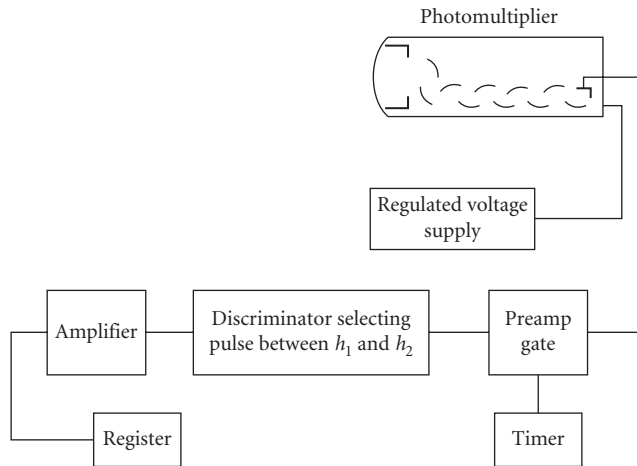


FIGURE 58 Photomultiplier and associated circuits for photon counting. (ITT Report E5.)

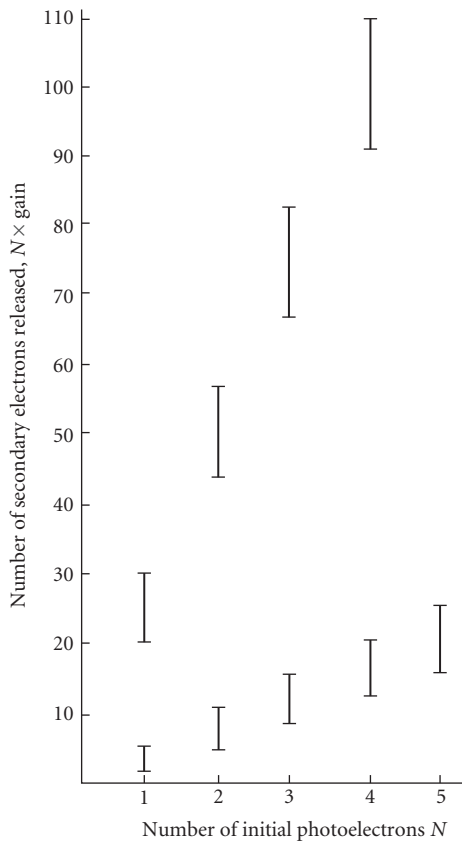


FIGURE 59 Spread in number of secondary electrons for various phototube gains.

GaN and AlGaN Gallium nitride photovoltaic detectors, with a bandgap of 3.39 eV have spectral response in the ultraviolet (UV) from 200 to 365 nm, as illustrated in Fig. 60. By using aluminum-gallium nitride—an alloy mixture of AlN and GaN—the spectral response can be tuned to shorter wavelength cutoffs. Spectral response examples are shown in Figs. 61 to 64 to compare GaN with one particular AlGaN alloy. Some devices may be tailored to custom UV bands, such as UVA (320 to 400 nm), UVB (280 to 320 nm), or UVC (100 to 280).

Response at visible wavelengths is low or absent, so that no special filtering may be required to detect UV in the presence of visible lighting or solar radiation—but see the logarithmic spectral Figs. 62 and 64 to see the degree of longer wavelength response. These solid-state devices are potentially useful for operation at elevated temperatures, in high-vibration environments, and in other environments unsuitable for photomultiplier tubes.

The photoconductive GaN devices use interdigitated contact electrodes because of the very high impedance of the GaN films, but currently there may not be any available commercially.

Response: Photovoltaic 0.1 A/W

Dark current: 0.05-nA photovoltaic

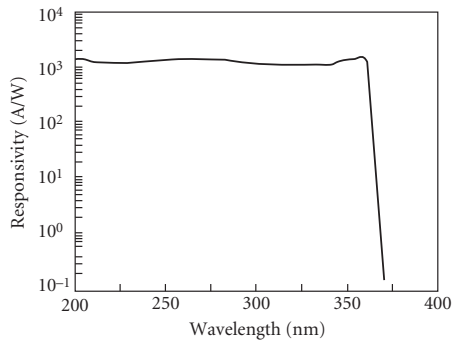


FIGURE 60 Response in amperes per watt for a GaN detector. (Reprinted from *Appl. Phys. Lett.*, vol. 60, no. 23, 1992, p. 2918.)

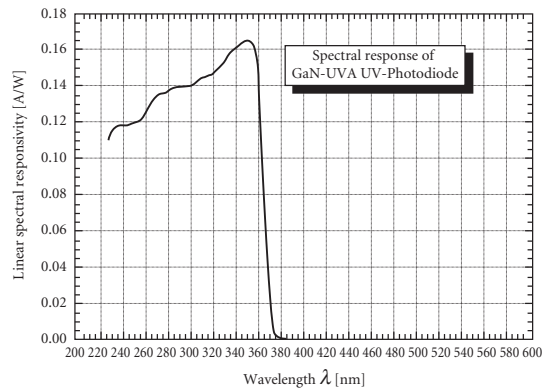


FIGURE 61 Spectral response of a UVA GaN detector shown on a linear vertical scale of amps/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaNall.pdf>)

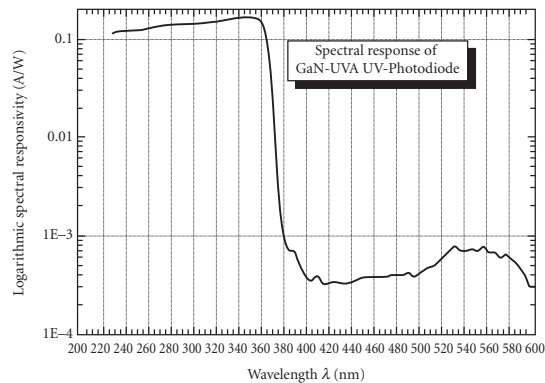


FIGURE 62 Spectral response of a UVA GaN detector shown on a logarithmic vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaNall.pdf>)

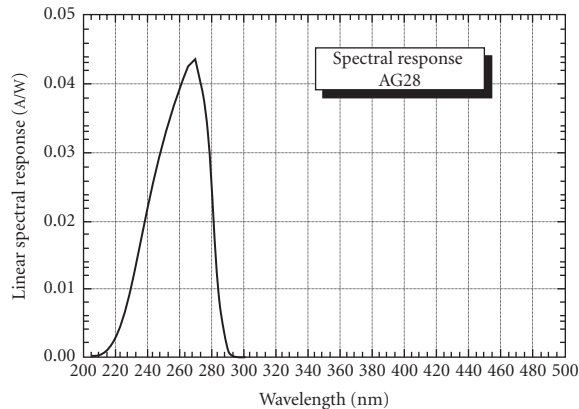


FIGURE 63 Spectral response of a UVC AlGaIn detector shown on a linear vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaIn.pdf>.)

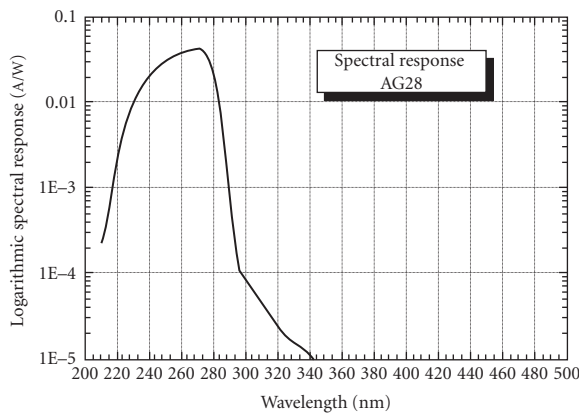


FIGURE 64 Spectral response of a UVC GaN detector shown on a logarithmic vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaIn.pdf>.)

Capacitance: 24 pF photovoltaic at 0-V bias

Time constant: Photovoltaic 0.10 ns

Size: 0.076 mm²

Devices with AlGaIn alloys have wider bandgaps and generally lower leakage currents.

Response: Photovoltaic 0.045 A/W.

Dark current: 0.1-pA photovoltaic at 0.1-V reverse bias

Capacitance: 24-pF photovoltaic at 0-V bias

Size: 0.076 mm²

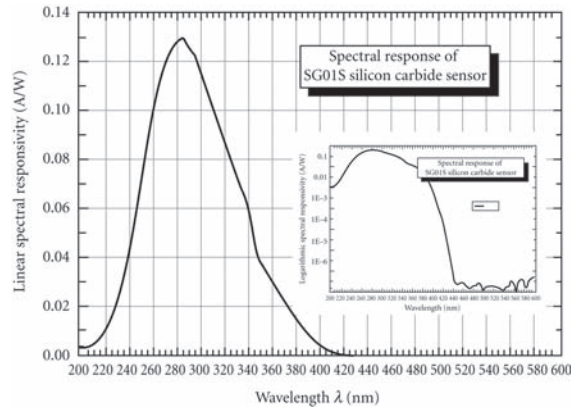


FIGURE 65 Spectral response of an unfiltered, broadband SiC detector shown on a linear vertical scale of amps/watt versus wavelength. The inset shows the same data on a logarithmic scale. (<http://www.boselec.com/products/documents/UVPhotodetectors2-08WWW.pdf>.)

Manufacturers: Advanced Photonix: http://www.advancedphotonix.com/ap_products/standard_GaN.asp?from=leftnav, Boston Electronics: <http://www.boselec.com/products/detuv.html>, Orion Semiconductor: <http://www.orion-semi.com>, SVT Associates: <http://www.svta.com/products/uv/uv.htm>.

SiC Silicon carbide UV detectors are available in photovoltaic structures. The 3-eV bandgap of SiC is slightly narrower than GaN, thereby giving a response may extend to slightly longer wavelength. However, because the bandgap of SiC is indirect, unlike GaN which is direct, the response cut-on is more gradual in SiC, peaking at a wavelength much shorter than the wavelength corresponding to 3 eV (413 nm)—see Fig. 65. SiC detectors with integrated filters are available.

Response: 0.13 A/W peak

Dark current: 1 fA for a 1 × 1 mm device

Capacitance: 195 pF for a 1 × 1 mm device

Sizes: 0.25 × 0.25 mm, 0.5 × 0.5 mm, 1 × 1 mm

Manufacturers: Boston Electronics: <http://www.boselec.com/products/detuv.html>, Electro Optical Components: http://www.eoc-inc.com/UV_detectors_silicon_carbide_photodiodes.htm

TiO₂ Detectors With a bandgap of 3.2 eV, TiO₂ is another UV photodetector. Photovoltaic devices are made with Schottky diodes. An unfiltered spectral response is shown in Fig. 66. TiO₂ detectors with integrated filters are available.

Response: 0.021 A/W peak

Dark current: 100 pA for a 5.4 × 2.9 mm device

Sizes: 2.2 × 1.9 mm, 5.4 × 2.9 mm

Manufacturer: Boston Electronics: <http://www.boselec.com/products/detuv.html>

GaP Gallium phosphide can provide Schottky photodiodes that cover the UV to mid-visible spectral region as shown in Fig. 67. The bandgap of GaP is 2.26 eV and is indirect, leading to a soft spectral

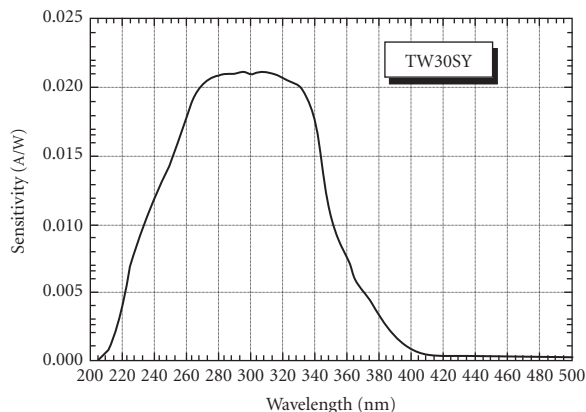


FIGURE 66 Spectral response of a TiO_2 Schottky photodiode detector shown on a linear vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/UVPhotodetectors2-08WWW.pdf>.)

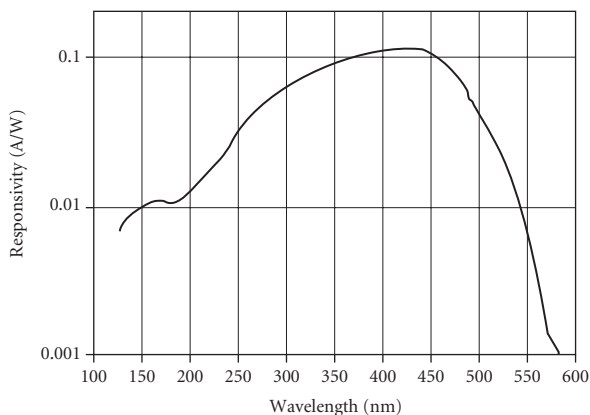


FIGURE 67 Spectral response of a GaP Schottky photodiode detector shown on a logarithmic vertical scale of amps/watt versus wavelength. (<http://www.thorlabs.com/Thorcat/12100/12174-S01.PDF>.)

turn-on and a peak quite far from the wavelength corresponding to the bandgap (549 nm). GaP devices with integrated filters to restrict the response to the UV region are also available.

Response: 0.12 A/W peak

Dark current: 1 nA max for a 2.5×2.5 mm device

NEP @ 440 nm: 1×10^{-14} W/ $\sqrt{\text{Hz}}$ @ 5 V bias

Rise time: 1 nsec @ 5 V bias for a 2.5×2.5 mm device

Fall time: 140 nsec @ 5 V bias for a 2.5×2.5 mm device

Sizes: 1.1×1.1 mm, 2.3×2.3 mm, 2.5×2.5 mm, 4.6×4.6 mm

Manufacturers: Hamamatsu: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd144/index_en.html, Electro Optical Components: <http://www.eoc-inc.com/ifw/EPD-365-0-2-5.pdf>, Thor Labs: <http://thorlabs.com/thorProduct.cfm?partNumber=FGAP71>

GaAsP Gallium arsenide phosphide alloys can provide photodiodes that cover from the UV to the near-infrared spectral region. The bandgap of GaP is 2.26 eV and is indirect, while that of GaAs is 1.43 eV and is direct. GaAsP alloys from 0 to ~50% GaP are direct bandgap materials while those with higher percentages of GaP are indirect*. A variety of alloys are available, covering the following spectral bands:

Spectral Band (nm)	λ Peak (nm)	Response at Peak (A/W)	Sizes (mm)
400–760	710	0.4	1.3 × 1.3, 2.7 × 2.7, 5.6 × 5.6
300–680	640	0.3	1.3 × 1.3, 2.7 × 2.7, 5.6 × 5.6
300–580	470	0.25	0.8 × 0.8
280–580	470	0.2	0.8 × 0.8
260–400	370	0.06	0.8 × 0.8
190–760	710	0.22	2.3 × 2.3, 4.6 × 4.6
190–680	610	0.18	10.1 × 10.1

Spectral responses of these alloys are shown in Figs. 68 to 73 (ref: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd143/index_en.html?sort=WAVE_LENGTH4&desc=1&style=F1 for all six figures).

Manufacturer: Hamamatsu: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd143/index_en.html

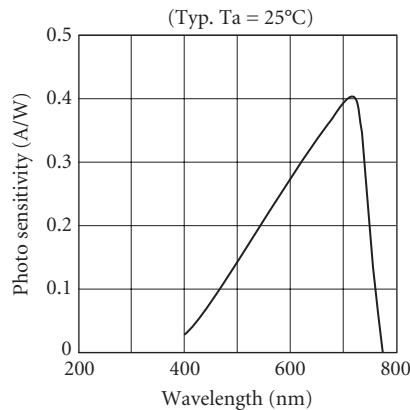


FIGURE 68 Spectral response of a 400 to 760-nm GaAsP photodiode detector with a vertical scale of amps/watt versus wavelength.

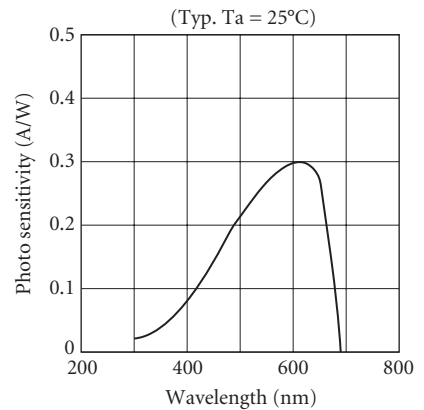


FIGURE 69 Spectral response of a 300 to 680-nm GaAsP photodiode detector with a vertical scale of amps/watt versus wavelength.

*<http://www.iue.tuwien.ac.at/phd/palankovski/node37.html>.

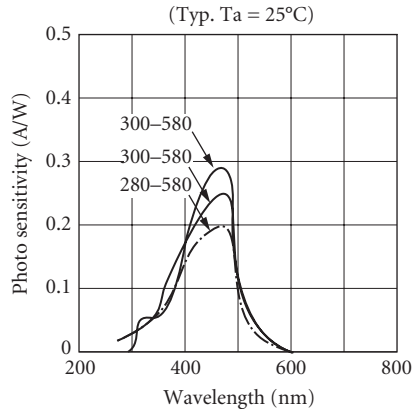


FIGURE 70 Spectral response of 300 to 580 and 280- to 580-nm GaAsP photodiode detectors with a vertical scale of amperes/watt versus wavelength.

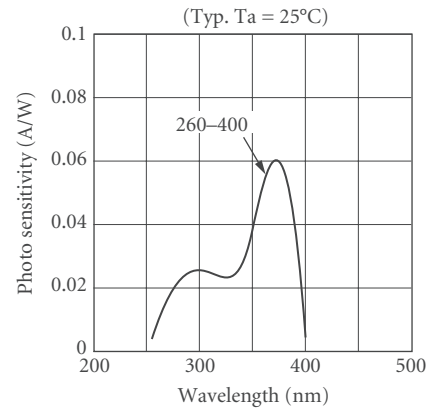


FIGURE 71 Spectral response of a 260- to 400-nm GaAsP photodiode detector with a vertical scale of amperes/watt versus wavelength.

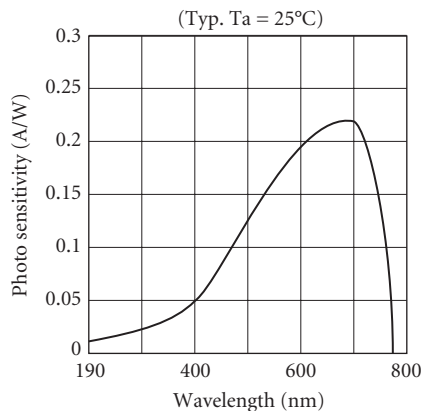


FIGURE 72 Spectral response of a 190- to 760-nm GaAsP photodiode detector vertical scale of amps/watt versus wavelength.

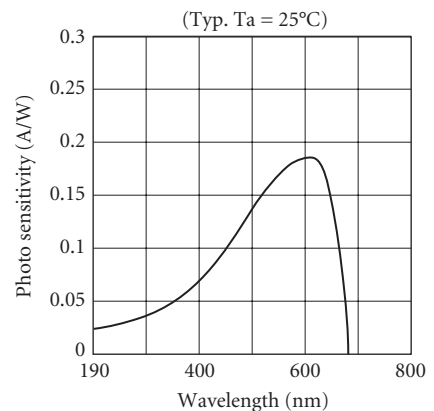


FIGURE 73 Spectral response of a 190- to 680-nm GaAsP photodiode detector vertical scale of amps/watt versus wavelength.

CdS and CdSe Cadmium sulfide and cadmium selenide photoconductors are available for detection of visible light out to 700 to 800 nm. CdS and CdSe films have sheet resistivity in the range of 20 m Ω per square at an illumination level of 2 footcandles. The devices are typically made in a linear or serpentine configuration consisting of 2 to 500 squares to maximize the length-to-width ratio. A variety of material “types” are available, offering unique spectral curves for various applications, depending upon the source color. CdS and CdSe are typically slow detectors, with response times of 5 to 100 ms, with speed improving at higher light levels. These devices exhibit “memory” or “history” effects, where the response is dependent upon the storage condition preceding use—the length of storage and time in use, and differences between the storage light level and the light level during use. These history effects may amount to changes in resistance from less than 10 percent to over 500 percent. CdSe has comparably greater memory effect than CdS.

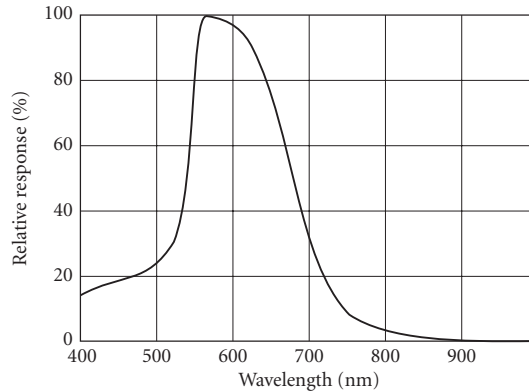


FIGURE 74 Relative spectral response of a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

CdS and CdSe are useful for a variety of commercial applications, both analog and digital, such as camera exposure control, automatic focus and brightness controls, densitometers, night light controls, etc. They are comparatively inexpensive and are available in a wide range of packages and resistance values, including dual cell configurations.

Spectral response: See Figs. 74 to 76.

Resistance and sensitivity: See Figs. 77 to 78.

Temperature coefficient of resistance: See Figs. 79, 80.

Light history effects: See Table 3.

Detector size: 4×4 mm to 12×12 mm approximate, dual elements available.

Manufacturers: In the previous edition, the listed supplier was EG&G VACTEK. Their product line has been acquired by Perkin Elmer. In this transition, all but two of the detector “types” have

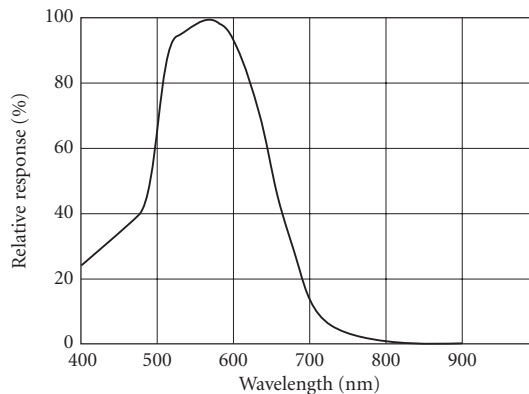


FIGURE 75 Relative spectral response of a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

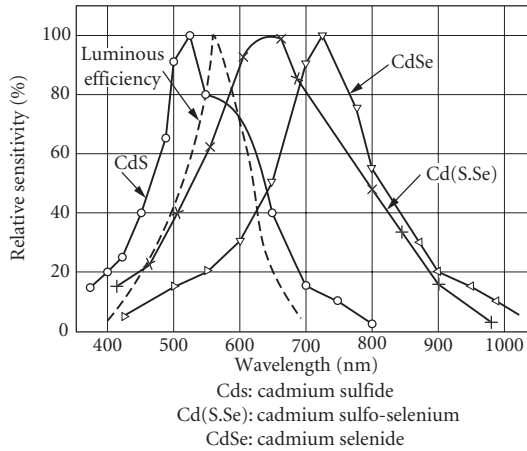


FIGURE 76 Spectral response of CdS, CdSSe, and CdSe photocell detectors together with the human eye response or luminous efficiency. (http://www.selcoproducts.com/CFM/photocells/photozell_PDF/Selco_PhotoCells_Construct.pdf.)

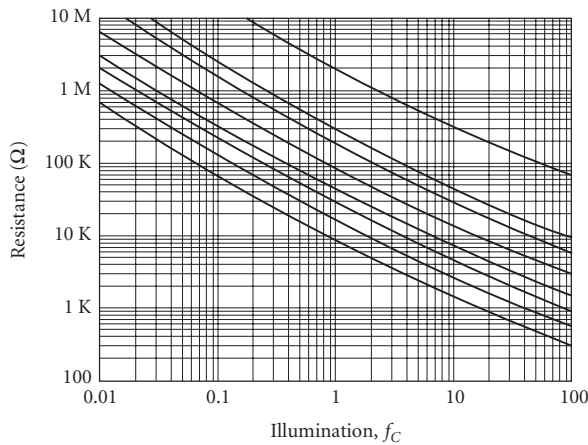


FIGURE 77 Resistance as a function of illumination for a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptois.pdf.)

been dropped—Perkin Elmer now only sells types 0 and 3. Other manufacturers offer comparable varieties and charts from at least one other producer are included.

Jameco Electronics: <http://www.jameco.com/webapp/wcs/stores/servlet/CategoryDisplay?storeId=10001&catalogId=10001&langId=-1&categoryId=151080>, Perkin Elmer: <http://optoelectronics.perkinelmer.com/catalog/Category.aspx?CategoryName=Photocells>, Selco Products: http://www.selcoproducts.com/CFM/photozell_toc.cfm, Silonex: <http://www1.silonex.com/optoelectronics/optophotoc.html>

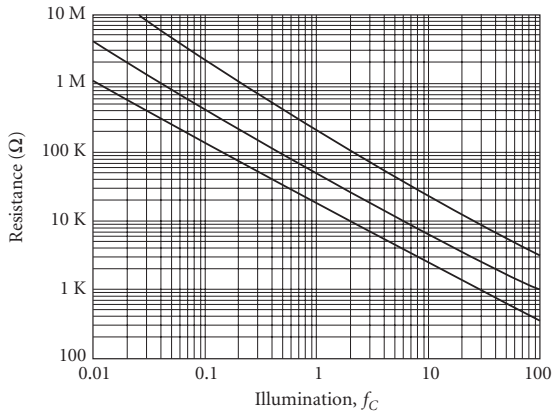


FIGURE 78 Resistance as a function of illumination for a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

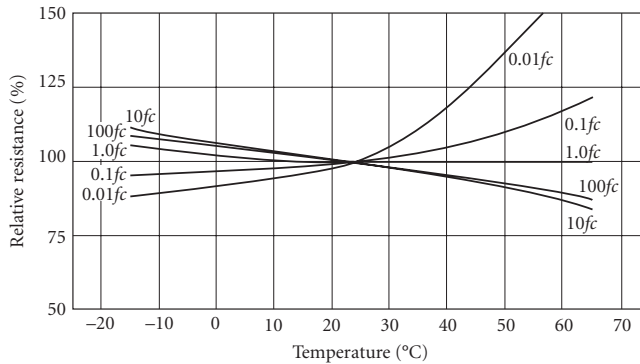


FIGURE 79 Relative resistance as a function of temperature for a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

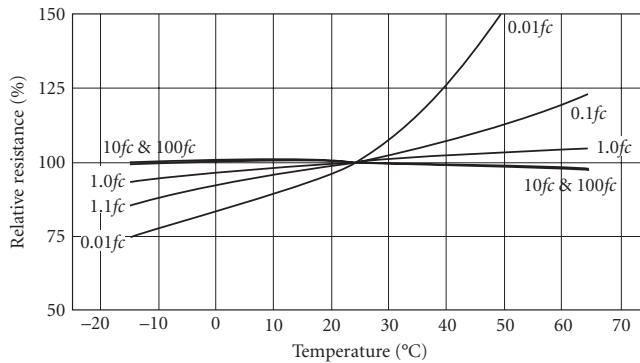


FIGURE 80 Relative resistance as a function of temperature for a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

TABLE 3 Typical Variation of Resistance with Light History Expressed as a Ratio R_{LH}/R_{DH} at Various Test Illumination Levels in Foot Candles.

Illumination (foot candles)	0.01	0.1	1.0	10	100
R_{LH}/R_{DH} ratio	1.55	1.35	1.20	1.10	1.10

(http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptois.pdf.)

R_{LH} is the resistance after “infinite” exposure to light, while R_{DH} is the resistance after “infinite” exposure to a dark environment. Infinite may be approximated by 24 hours.

CdTe Cadmium telluride and cadmium zinc telluride detectors are chemical group II-VI materials having an energy bandgap of about 1.6 eV, corresponding to a spectral cutoff in the vicinity of 775 nm. These devices, however, are principally used for gamma ray detection because of their high z number which translates into a high absorption coefficient for gamma rays. The principal advantage of CdTe in this application is its ability to operate at room temperature, in comparison with Ge gamma ray detectors which must typically be cooled to 77 K. Figure 81 illustrates the absorption of CdTe as a function of gamma ray energy out to 300 keV.

Sensitivity: See Fig. 81.

Standard sizes: Wafers in 10- and 16-mm diameter; rods $7 \times 2 \times 2$ mm; cubes $2 \times 2 \times 2$ mm.

Standard thickness: 1 and 2 mm.

Bias voltage: 150 to 300 V/cm.

Operating temperature range: -10 to $+55^\circ\text{C}$

Leakage current: 10 to 300 nA

Capacitance: 10 pF

Response time: $< 1 \mu\text{s}$.

Manufacturers: Acrorad: <http://www.acrorad.co.jp/us/index.html>, Aurora, II-VI eV Products: <http://www.evproducts.com/>, Perkin Elmer, Radiation Monitoring Devices: <http://www.rmdinc.com/products/p007.html>

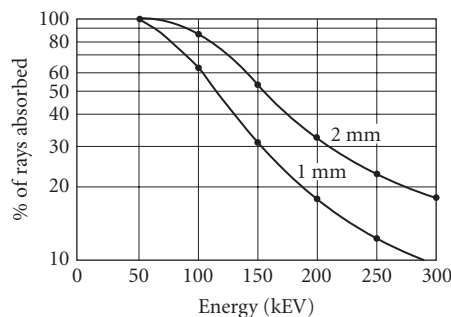


FIGURE 81 The high percentage of rays absorbed by CdTe makes these detectors highly sensitive. At 100 keV, a 2-mm-thick detector absorbs 85 percent of the rays. (*Radiation monitoring devices, Cadmium Telluride brochure.*)

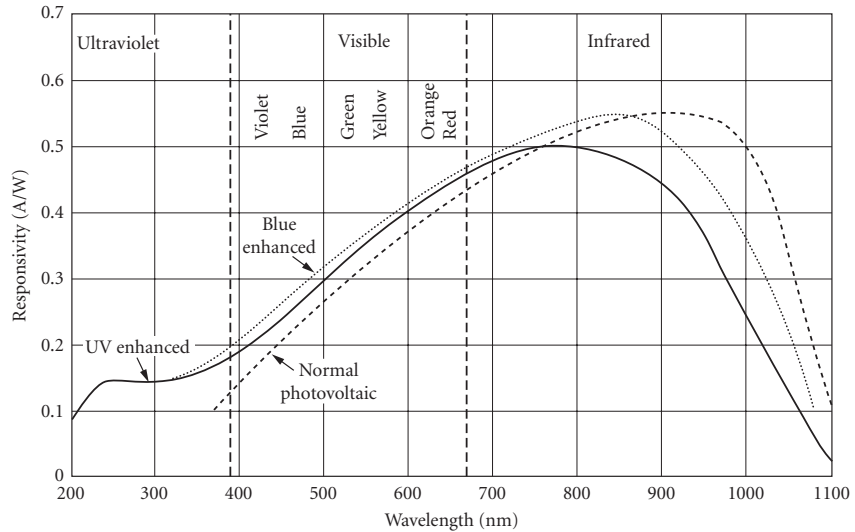


FIGURE 82 Typical spectral response for pn junction, blue-enhanced, and UV-enhanced silicon photodiodes. (UDT Sensors, Inc., *Optoelectronic Components Catalog*.)

Si Silicon photovoltaic detectors are widely available. They are useful at wavelengths shorter than about $1.1\ \mu\text{m}$ and can even be used for x-ray and gamma-ray detection. There are four main silicon detector types:

- pn junction photodiodes, generally formed by diffusion, but ion implantation can also be used.
- pin junctions, which have lower capacitance and hence higher speed, and because of a thicker active region have enhanced near-IR spectral response.
- UV- and blue-enhanced photodiodes
- Avalanche photodiodes with significant internal gain, combining high speed and sensitivity

The main parameters of interest are spectral response (see Fig. 82), time constant, and zero-bias resistance or reverse-bias leakage current. Silicon material has an indirect bandgap and hence the spectral cutoff is not very sharp near its long-wavelength limit as shown in Fig. 82. The effective time constant of pn junction silicon detectors is generally limited by resistance-capacitance (RC) considerations rather than by the inherent speed of the detection mechanism (drift and/or diffusion). High reverse bias may or may not shorten charge collection time, but it generally reduces cell capacitance, and therefore the RC product, therefore, reverse bias usually results in faster response.

On the other hand, increased reverse bias causes increased noise, so that a trade-off exists between speed and sensitivity. For high-frequency applications, load resistance should be made small, although this makes Johnson (thermal) noise comparatively larger, which limits sensitivity (see Fig. 83). In order to keep sensitivity high when using these devices at high frequency, operational (current-mode) amplifiers, which can be built into the detector package, and avalanche photodiodes, which incorporate built-in gain before the load resistor is encountered, have been developed. Very careful regulation of the detector bias is required for stable operation of avalanche photodiodes.

Silicon pn junction photodiodes These are general purpose when high sensitivity is required and time constants on the order of a microsecond are permissible. The device construction is illustrated in Fig. 84. These devices are typically operated in a photovoltaic mode at zero bias, but can be used in a photoconductive mode in which the device is reverse biased.

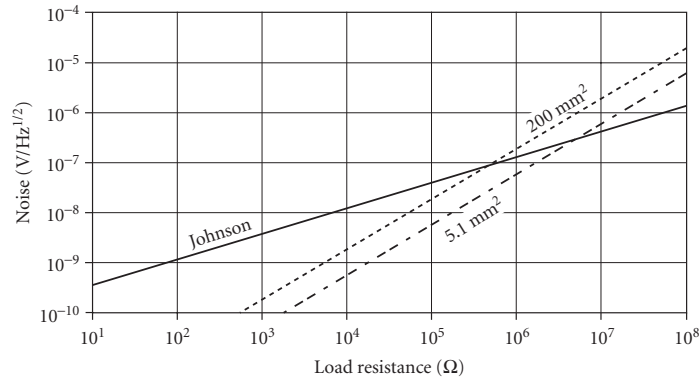


FIGURE 83 Output noise as a function of circuit load resistance for *pin* silicon photodiodes with areas of 5.1 and 200 mm², compared with the Johnson noise of the load resistor. Dark current measured at 10-V reverse bias for the detector with area of 5.1 mm² is 10 nA, and 100 nA for the detector with an area of 200 mm². Note that good preamplifiers have a noise level of about 1 nV/Hz^{1/2}, depending upon the bandwidth. (Detector data from UDT Sensors, *Optoelectronic Components Catalog*.)

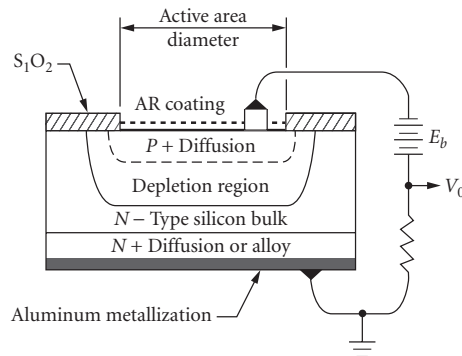


FIGURE 84 Planar diffused silicon photodiode construction. (UDT Sensors, Inc., *Optoelectronic Components Catalog*.)

Sensitivity: $D^*(\lambda_{pk}) \approx \text{mid-}10^{12} \text{ to } 10^{13}$ Jones, $D^*(2800 \text{ K}) \approx 2 \times 10^4$ Jones, becoming amplifier-limited for small-area detectors (see Figs. 85 and 86). D^* can also be estimated from the R_0A product (detector zero-bias resistance or shunt resistance diode area), which is illustrated in Fig. 87, in combination with Fig. 19, which illustrates the dependence of D^* on R_0A product.

Noise: See Figs. 88 (noise vs. bias) and 89 (noise vs. temperature); as T drops, impedance rises, so that decreasing noise current produces increasing noise voltage. However, the signal increases even faster, yielding an improved signal-to-noise ratio with cooling. Figure 90 (noise vs. frequency) shows the dependence on bias.

Capacitance: Capacitance is proportional to area and increases slightly with temperature (see Fig. 91).

Responsivity: See Figs. 82 and 88.

Quantum efficiency: >90 percent quantum efficiency achievable with antireflection coating.

Sensitive area: 0.2 to 600 mm² areas are readily available.

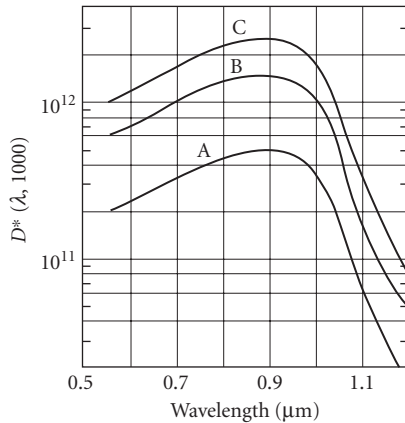


FIGURE 85 D^* versus λ for small-area junction silicon photodiodes: curves A, B, and C correspond to areas of 0.02, 0.2, and 1 cm^2 . The lower D^* for smaller area detector performance is due to amplifier limitations rather than intrinsically poorer D^* , for small-area detectors. (*Texas Instruments, Infrared Devices, SC-8385-366. Reprinted by permission of Texas Instruments.*)

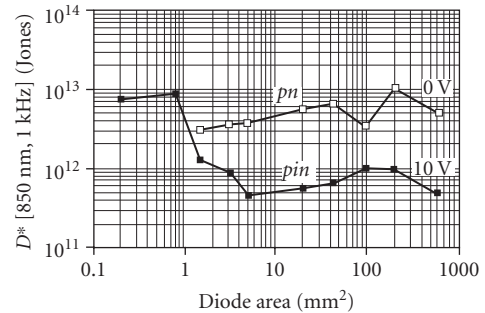


FIGURE 86 D^* as a function of diode area for pn junction silicon photodiodes operated in the photovoltaic mode (0 V) and pin junction diodes operated in the photoconductive or reverse-bias mode (10 V). (*UDT Sensors, Optoelectronic Components Catalog.*)

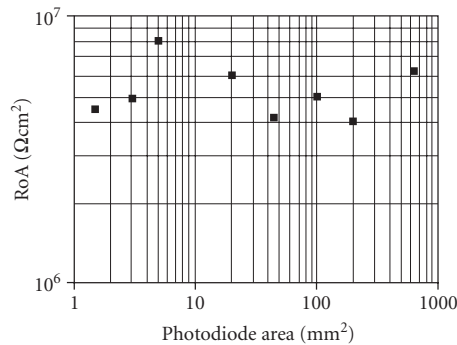


FIGURE 87 The resistance-area product (RoA) at zero bias and 295 K of silicon pn junction photodiodes. The lack of area dependence is evidence that intrinsic properties of the junction, rather than surface effects, are dominant in these devices.

Time constant: Inherently slow for high-sensitivity applications, generally limited by RC (depends directly on device area), but can be limited by carrier diffusion outside the depletion region or by trapping of carriers in deep impurity centers. Typical data for a circuit using a 50- Ω load resistor is illustrated in Fig. 92.

Operating temperature: Ambient, but noise (leakage current) can be reduced by operating at lower temperatures (see Fig. 76 for typical signal and noise vs. temperature).

Uniformity: Typically ± 8 percent across a diode area with a 40- μm focused light spot.

Linearity: 5 percent or better over 10 orders of magnitude flux from 10^{-13} to 10^{-3} W/cm^2 .

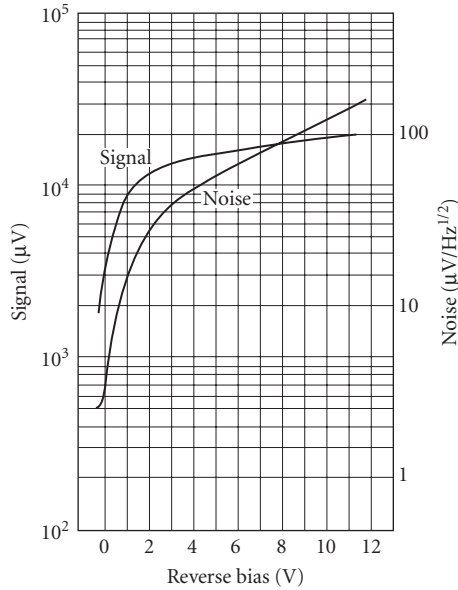


FIGURE 88 Typical *pn* junction signal and noise versus reverse bias ($R_L = 10\text{ M}\Omega$). (Electronuclear Laboratories, Bull. 1053, 1966.)

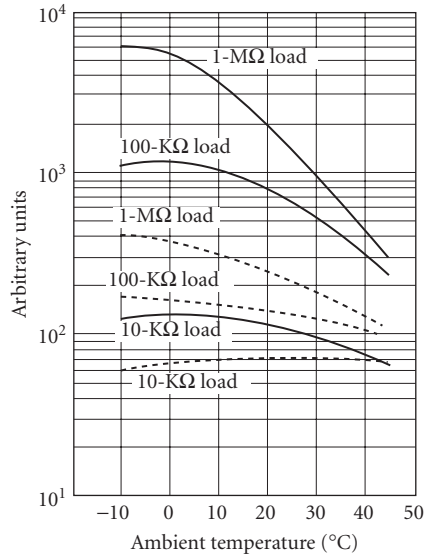


FIGURE 89 Relative signal and noise versus temperature for *pn* junction silicon photodiode at zero bias; — = signal; --- = noise. (Electronuclear Laboratories, Bull. 1052, 1966.)

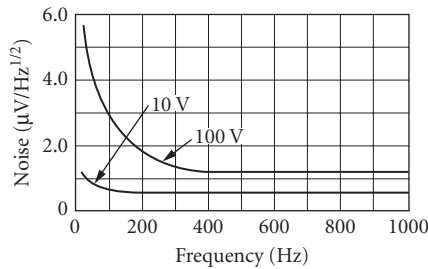


FIGURE 90 Typical *pn* junction and *pin* junction noise-frequency spectrum for different reverse bias ($A = 1 \times 1\text{ mm}$; $R_L = 1\text{ M}\Omega$). (Electronuclear Laboratories, Bull. 1078, 1966.)

Recommended circuit: See Fig. 93. High-impedance FET current-mode amplifier to supply fixed bias voltage, regardless of current.

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: Advanced Photonix, EG&G Canada, EG&G Heimann, Edmund Optics, Electro Optical Systems, Electro-Optics Technology, International Radiation Detectors, Janos Technology, Laser Precision Corp, Laser Systems Devices, Melles Griot, Newport/Klinger, Ophir Optronics, Optical Signature, Opto-Electronics, Optometrics, Oriel, Photonic Detectors, RMD, Sapidyne, Scientific Instruments, SEMICOA, Silonex, Spire, UDT Sensors.

Silicon *pin* junction photodiodes The *pin* junction detector is faster but is also somewhat less sensitive than conventional *pn* junction detectors. *PIN* photodiodes have slightly extended red

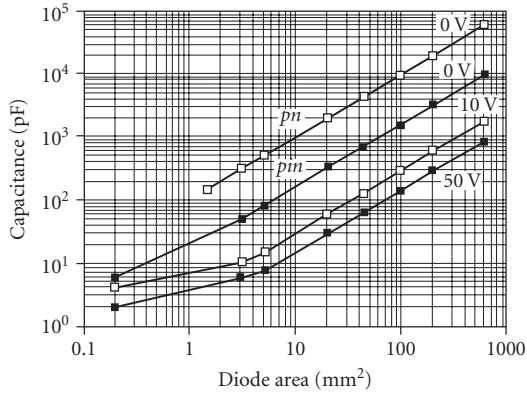


FIGURE 91 Capacitance as a function of detector area for *pn* junction silicon photodiodes operated in the photovoltaic mode (0 V) and *pin* junction photodiodes at 0-, 10-, and 50-V reverse bias. The larger depletion width, which is a consequence of the lightly doped “i” region in the *pin* device, gives *pin* diodes lower capacitance for the same device area. (UDT Sensors, *Optoelectronic Components Catalog*.)

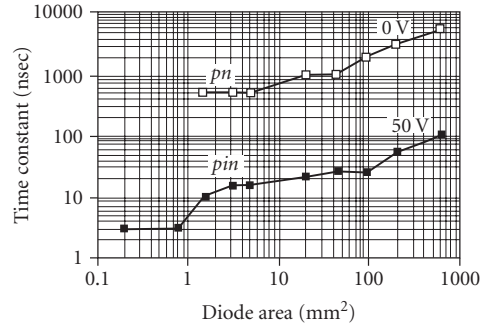
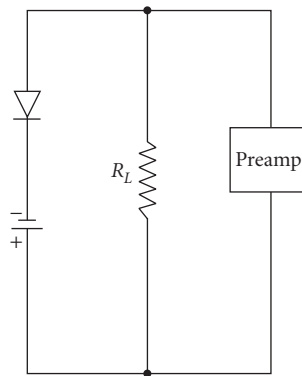
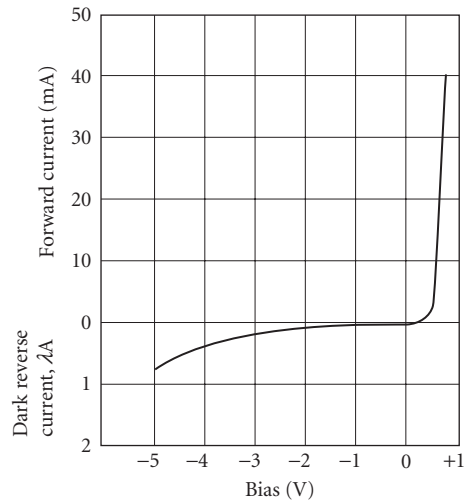


FIGURE 92 Time constant for *pn* junction silicon photodiodes operated in the photovoltaic mode (0 V) and a *pin* junction detector in the photoconductive or reverse-bias mode (50 V). A 50- Ω load was used in both cases, which limits sensitivity because of Johnson noise (see also Fig. 83). (UDT Sensors, *Optoelectronic Components Catalog*.)



(a)



(b)

FIGURE 93 *pn* junction silicon photodiode: (a) recommended circuit; (b) typical electrical characteristics. (Texas Instruments, Bull. SC-8385-366. Reprinted by permission of Texas Instruments.)

response. In the normal *pn* junction, charge-collection time has a slow and a fast component. The fast component is due to photons absorbed in the depletion layer of the *pn* junction. Since the electric field in the depletion region is strong, carriers are quickly separated by drift through the electric field across the depletion region. However, photons absorbed deeper in the material, beyond the depletion region, produce carriers which must diffuse to the junction before they are collected, and diffusion times are on the order of a microsecond. This component becomes more significant near the

long-wavelength limit of the spectral response. Application of reverse bias in an ordinary pn junction detector reduces the capacitance, shortening the RC time constant, and increases the width of the depletion layer thereby increasing the fraction of photons absorbed within the high field region and proportionally increasing the fraction of the fast component of the response.

However, the doping level of the ordinary pn junction limits the extent of the depletion layer increase to only 5 to 10 μm at a reverse bias of 50 V (this assumes an abrupt junction with a concentration of $1 \times 10^{15} \text{ cm}^{-3}$). pin detectors incorporate a very lightly doped region between the p - and n -regions that allows a modest reverse bias to form a depletion region the full thickness of the material (500 μm for a typical silicon wafer). Extended red response in a pin device is a consequence of the extended depletion layer width, since longer wavelength photons will be absorbed in the active device region. Unfortunately, the higher dark current collected from generation within the wider depletion layer results in lower sensitivity. Generation of carriers can be minimized by minimizing the concentration of deep-level impurity centers in the detector with careful manufacturing. Operation at lower temperature will also reduce the dark current.

Sensitivity: $D_{\text{pk}}^* 1 \times 10^{12}$ Jones for 2-mm² area (depends slightly on bias, see Figs. 86 and 94). For high-speed operation, detectivity is lower (see Fig. 83).

Noise: Depends upon diode area and circuit load resistance. Johnson noise will dominate at low values of load resistance when circuit is optimized for fast response. Preamp noise may also limit. See Fig. 83.

Responsivity: Similar to pn junction. See Fig. 82.

Quantum efficiency: 90 percent quantum efficiency achievable with antireflection coating.

Capacitance: Proportional to detector area. See Fig. 91.

Operating temperature: See Fig. 95.

Time constant: Varies with capacitance (device area); see Fig. 92.

Sensitive area: 0.2 to 600 mm² readily available.

Recommended circuit: Same as for pn junction photodiodes. See Figs. 93 and 96.

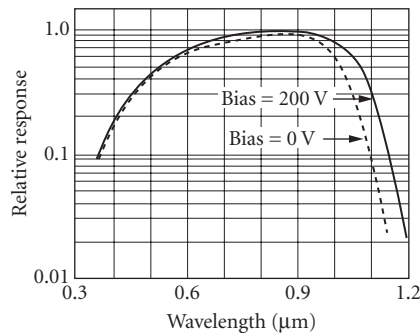


FIGURE 94 Dependence of spectral response on bias for silicon photodiodes. (Electronuclear Laboratories, Bull. 1076, 1968.)

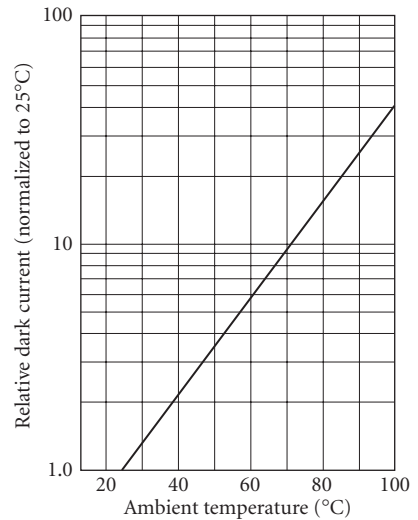


FIGURE 95 Relative dark current versus temperature for pin junction silicon photodiodes—bias 100 V. (Electronuclear Laboratories, Bull. 1076, 1969.)

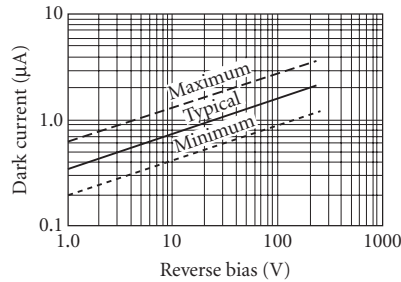


FIGURE 96 Typical dark current versus bias for *pin* silicon photodiode ($A = 1 \times 1 \text{ mm}$). (*Electronuclear Laboratories, Bull. 1078, 1968.*)

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: Same as for *pn* junction photodiodes.

UV- and blue-enhanced photodiodes Blue- and UV-enhanced photodiodes may improve the quantum efficiency by 50 to 100 percent over standard photodiodes in the blue and UV spectral region. The quantum efficiency of ordinary *pn* and *pin* junction photodiodes degrades rapidly in the blue and UV spectral regions. This is because the high absorption coefficient of silicon at these wavelengths causes the photocarriers to be generated within the heavily doped *p*- (or *n*-) type contact surface where the lifetime is short due to the high doping and/or surface recombination. Blue- and UV-enhanced photodiodes optimize the response at short wavelengths by minimizing near-surface carrier recombination. This can be achieved by using very thin and highly graded *p* (or *n* or metal Schottky) contacts, by using lateral collection to minimize the percentage of the surface area which is heavily doped, and/or passivating the surface with a fixed surface charge to repel minority carriers from the surface. These devices typically have quartz windows or UV-transmissive glass, compatible with good transmission into the UV spectrum. The user should be aware that UV and higher energy radiation in particular can alter the fixed charge conditions in the surface region of silicon and other detectors (typically in the surface oxide) which can cause the detector performance to drift and/or be unstable (see Fig. 20).

Sensitivity: See Figs. 82 and 97; $D^*_{pk} \approx 3\text{--}5 \times 10^{12}$ Jones for diodes with areas of 1–100 mm² at $V_R = 0$, $R_L = 40 \text{ M}\Omega$.

Quantum efficiency: Same as *pn* junction photodiodes, but enhanced in the UV and blue regions by 50 percent or more (see Fig. 82).

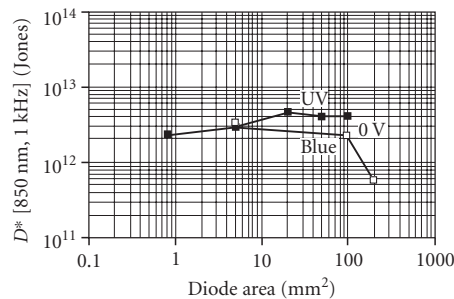


FIGURE 97 D^* as a function of diode area for blue- and UV-enhanced silicon photodiodes operating in the photovoltaic mode (0-V bias). (*UDT Sensors, Optoelectronic Components Catalog.*)

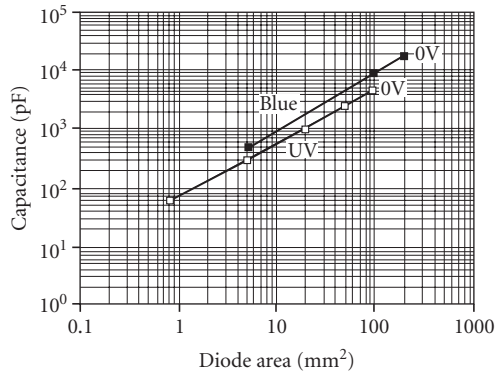


FIGURE 98 Capacitance as a function of detector area for blue- and UV-enhanced silicon photodiodes operated in the photovoltaic mode (0 V). The capacitance per unit area is close to that of *pn* junction photodiodes shown in Fig. 91. (UDT Sensors, *Opto electronic Components Catalog*.)

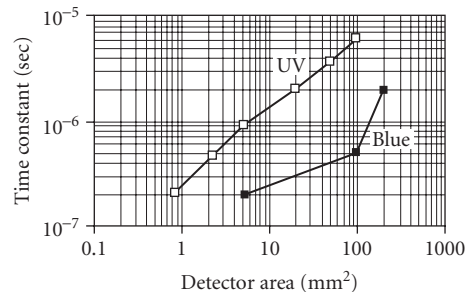


FIGURE 99 Time constant for UV- and blue-enhanced silicon photodiodes as a function of detector area at zero-volts bias. A 50- Ω load was used in both cases. (UDT Sensors, *Optoelectronic Components Catalog*.)

Responsivity: see Fig. 82.

Capacitance: Comparable to *pn* junction photodiodes (See Fig. 98).

Operating temperature: Ambient.

Time constant: Dependent upon device type. Increases with device area; 200 ns to 6 μ s for areas 1 to 100 mm² (see Fig. 99).

Sensitive area: 1 to 200 mm² readily available.

Recommended circuit: Same as *pn* junction photodiodes.

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: See list for silicon *pn* junction photodiodes.

Silicon avalanche photodiodes The avalanche photodiode, is especially useful where both fast response and high sensitivity are required. Whereas normal photodiodes become Johnson- or thermal-noise-limited when used with a low-impedance load resistor for fast response, avalanche photodiodes make use of internal multiplication, associated with reverse breakdown in the *pn* junction in order to keep the detector noise above the Johnson-noise level. [Because the response time is usually RC-limited, small load resistors (often 50 Ω) are used to achieve fast signal response.] However, as the load resistor is decreased, the detector noise voltage decreases in direct proportion, whereas the Johnson noise decreases only as the square root of the load resistor. Thus, the detector noise voltage can become lower than the Johnson noise for load resistance values smaller than a critical value. With an APD device, lower values of load resistors can be used without reaching the critical value because the internal gain boosts the detector noise voltage.

Stable avalanche or multiplication is made possible by a guard-ring construction using n^+pp^+ , Schottky- nn^+ , or $n^+p\pi p^+$ structure; beveled *pin* structure (see Fig. 100); mesa structures; or other structures which prevent surface breakdown.²⁵ However, very careful bias control is essential for stable performance. An optimum gain exists below which the system is limited by receiver noise and above which shot noise dominates receiver noise and the overall noise increases faster than the signal (Fig. 101).

In addition to fast-response applications, avalanche photodiodes are useful whenever amplifier noise is limiting, for example, small-area devices. Signal-to-noise-ratio improvements of one to two orders of magnitude over a nonavalanche detector can be achieved.

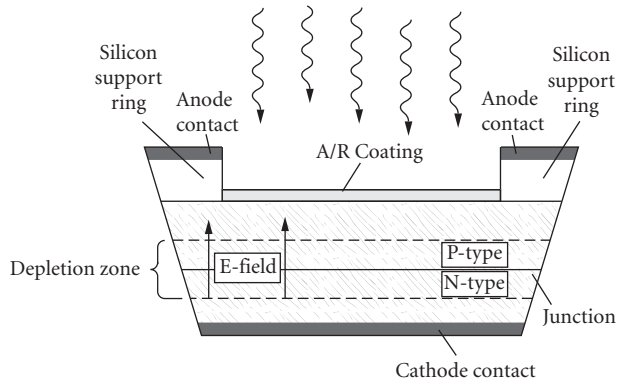


FIGURE 100 Cross section of a beveled-edge silicon avalanche photodiode. The beveled edge prevents early breakdown. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

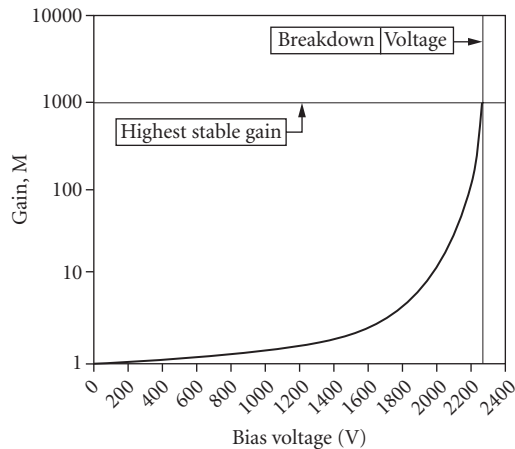


FIGURE 101 Gain as a function of reverse bias can reach 1000. This operating point is very close to breakdown and requires careful bias control. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

APDs are sometimes used in combination with scintillator crystals such as CsI to detect high-energy radiation in the range of 10 to 1000 keV.

Sensitivity: $D^* 3\text{--}5 \times 10^{13}$ Jones (see Fig. 102).

Noise: Function of detector area. As gain increases, noise (dark current) increases (see Fig. 103). Optimum gain is where avalanche noise equals system noise. Thus, optimum gain is a function of system noise.

Responsivity: Photocurrent is the product of the incident optical power in watts, wavelength in micrometers, and quantum efficiency (η) divided by 1.24 and multiplied by the avalanche gain M .

$$I_{\text{photo}} = M(P \eta \lambda / 1.24) \quad (\text{See Figs. 101 and 104})$$

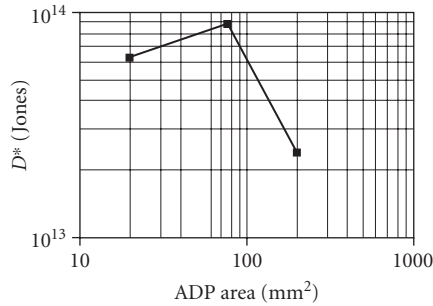


FIGURE 102 D^* for three silicon avalanche photodiodes shown as a function of diode area. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

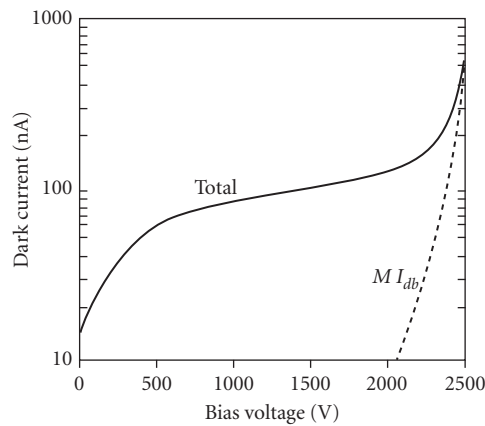


FIGURE 103 Dark current as a function of reverse bias for a 16-mm-diameter APD. At low-bias surface dark current dominates, but avalanche-multiplied bulk dark current increases rapidly as the gain increases. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

Quantum efficiency: Typically 85-percent peak (see Fig. 104).

Capacitance: Depends on bias and area (see Fig. 105).

Sensitive area: 20 to 200 mm².

Series resistance: Depends on area; typical values are 40 Ω for 5-mm diameter to 5 Ω for 16-mm diameter

Time constant: See Fig. 106.

Recommended circuit: Requires a filtered high-voltage dc supply that itself must have very low noise and a load resistor. The output may be ac- or dc-coupled. (See Fig. 107.)

Operating temperature: 40 to 45°C.

Stability: Exposure to UV or high-energy radiation may affect dark current. See Fig. 20 and section relating to stability. Check with manufacturer.

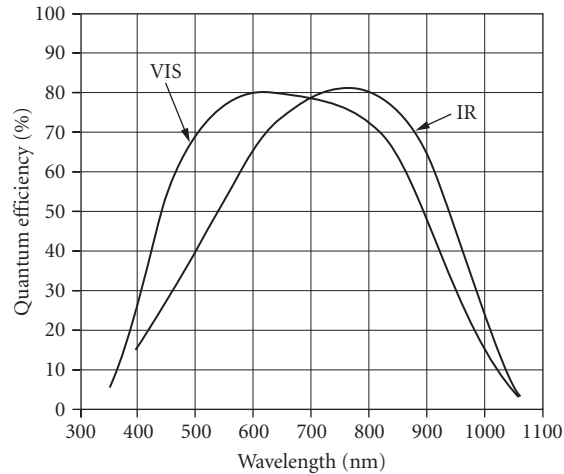


FIGURE 104 APD quantum efficiency at high gain. Adjustment of the oxide deposited on the surface produces two different curves. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

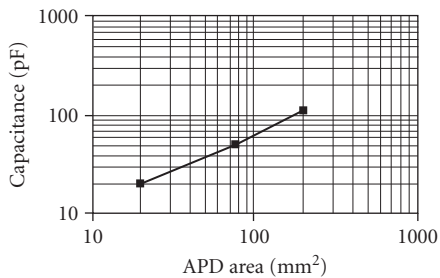


FIGURE 105 Capacitance for three silicon avalanche photodiodes shown as a function of diode area. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

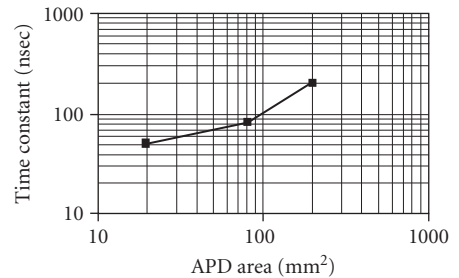


FIGURE 106 Time constant for three silicon avalanche photodiodes shown as a function of diode area. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

Manufacturers: Advanced Photonix, Devar, EG&G Judson, EG&G Vactec, EG&G Canada, Edmund Optics, Electro-Optical Systems, Hamamatsu, Janos Technology, Newport/Klinger, Opto-Electronics (Ontario), Oriel, Photonic Detectors, Photonic Packaging Technologies, RMD, Texas Optoelectronics, Thorn EMI Electron Tubes.

InGaAs Indium gallium arsenide detectors have been developed for optimum performance with fiber-optic communications at 1.3 and 1.55 μm . This detector material has a direct bandgap and represents one of several compound semiconductor alloy systems specially developed for photodetectors. In the case of this alloy of two group III-V chemical compound semiconductors, the ratio of InAs to GaAs controls the spectral cutoff, allowing the detector to be optimized for a particular wavelength. InGaAs detectors have generally been specialized for high-speed applications with

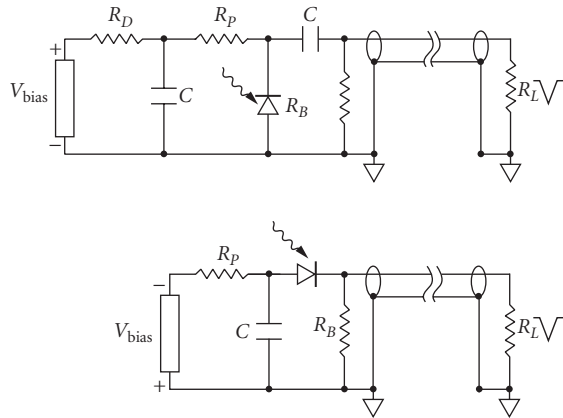


FIGURE 107 AC- and DC-coupled APD circuits with a filter following the power supply and a coaxial cable on the output. (Advanced Photonix, Avalanche Photodiode Catalog.)

optimum sensitivity since these performance factors can drive a fiber-optic system throughput and cost. For this reason, available devices include

- PIN photodiodes
- Avalanche photodiodes

The significance of these devices to fiber-optic applications is reflected in the number of vendors who sell integrated packages of InGaAs photodetectors combined with preamplifiers and fiber-optic pigtailed.

The InGaAs alloy system allows the spectral response to be tailored to longer wavelengths than the quartz fiber-optic bands and devices with cutoffs of 2.2 and 2.6 μm are also available.

InGaAs pin photodiode

Sensitivity: D^* mid- 10^{12} Jones for 1.67- μm cutoff; $D^* \approx 1 \times 10^{12}$ Jones for 2.2- μm cutoff, and $D^* \approx 5 \times 10^{11}$ Jones for 2.6- μm cutoff.

Responsivity: 0.85 to 0.95 A/W in the range of 1.3 to 1.55 μm .

Quantum efficiency: For 1.67- μm cutoff, see Fig. 108.

Dark current: See Fig. 109.

Capacitance: 0.7 to 1.2×10^4 pF/cm for 1.7- μm cutoff; 2.5×10^4 pF/cm for 1.85- μm cutoff; 3×10^4 pF/cm for 2.15- μm cutoff; 5×10^4 pF/cm for 2.65- μm cutoff. See also Fig. 110 for bias dependence.

Time constant: Varies with resistance-capacitance time (see Fig. 111). Since capacitance depends upon reverse bias, the time constant varies proportionally (see Fig. 110 for dependence of capacitance on bias).

Size: 0.05 to 3-mm diameter.

Recommended circuit: Standard photodiode options; zero bias for best sensitivity, reverse bias for maximum speed.

InGaAs avalanche photodiode

Sensitivity: $D^* \approx 5 \times 10^{11}$ Jones for 1.7- μm cutoff. In the fiber-optics industry, the sensitivity is also given in power units of dBm. Figure 112 compares InGaAs *pin*, APD, and Ge APD sensitivities.

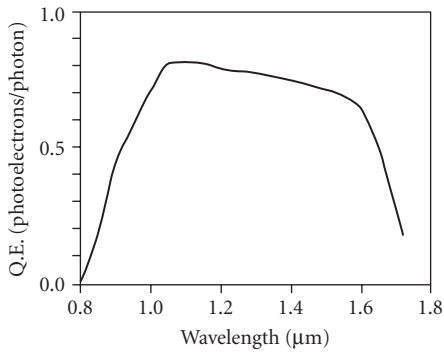


FIGURE 108 Spectral dependence of quantum efficiency for an InGaAs detector having a cutoff of 1.67 μm . (*Sensors Unlimited, data sheet.*)

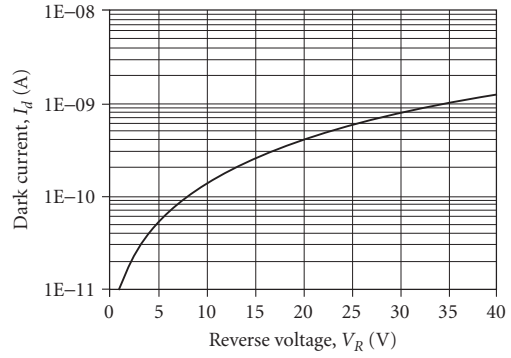


FIGURE 109 Dark current as a function of reverse-bias voltage for a 60- μm -diameter InGaAs detector having a cutoff of 1.67 μm . (*Fermionics, InGaAs Photodiodes.*)

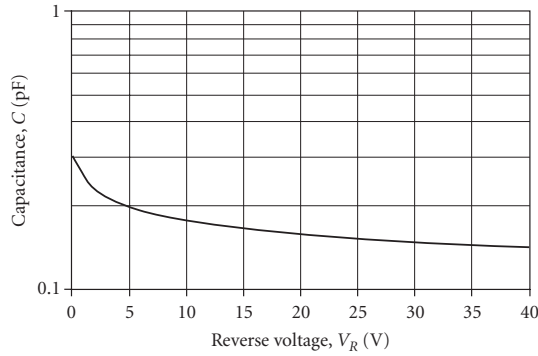


FIGURE 110 Capacitance as a function of reverse-bias voltage for a 60- μm -diameter InGaAs detector having a cutoff of 1.67 μm . (*Fermionics, InGaAs Photodiodes.*)

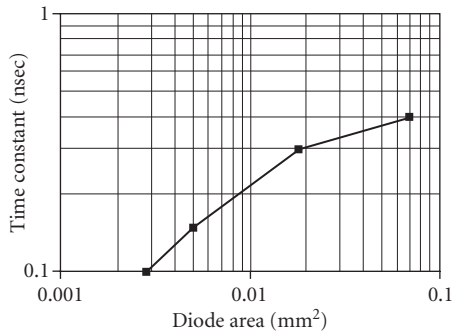


FIGURE 111 Time constant for small InGaAs *pin* photodiodes as a function of diode area measured with a 50- Ω load.

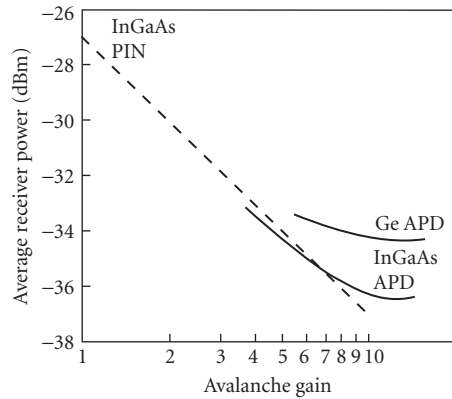
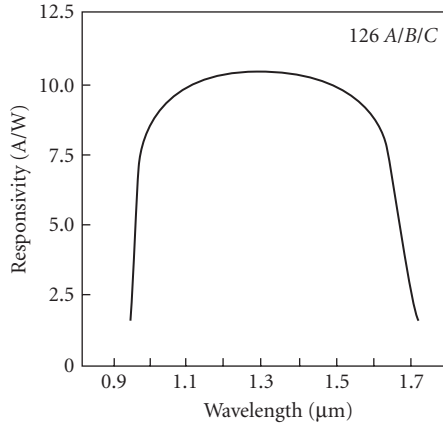


FIGURE 112 APD receiver sensitivity. Typical receiver sensitivity at a receiver rate of 1.7 Gbit/s and $\lambda = 1.3 \mu\text{m}$ for an InGaAs *pin*, Ge APD, and an InGaAs APD. (*AT&T, 126A/B, CASTROTEC InGaAs.*)



Note: Responsivity = (chip quantum efficiency) × gain × λ (μm)/1.24.
The minimum chip quantum efficiency is 80%, and the minimum pigtail coupling efficiency is 90%.

FIGURE 113 Responsivity for an InGaAs avalanche photodiode versus wavelength for avalanche gain of 12. (AT&T, 126A/B, C ASTROTEC InGaAs.)

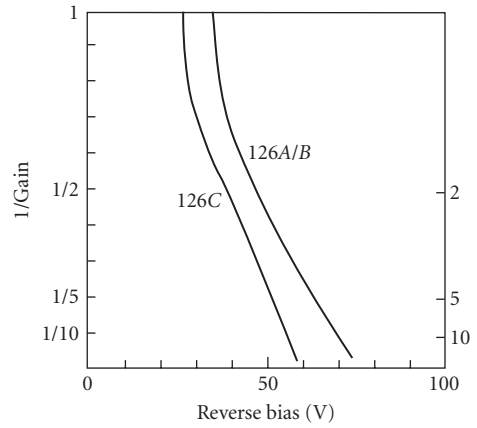


FIGURE 114 Inverse of avalanche gain for an InGaAs avalanche photodiode versus reverse bias. (AT&T, 126A/B, C ASTROTEC InGaAs.)

Spectral response: 1.0 to 1.65 μm ; see Fig. 113.

Responsivity: 8 to 10 A/W typical.

Avalanche gain: Critically depends upon reverse bias, see Fig. 114.

Capacitance: At gain of 12, 7.5×10^4 pF/cm².

Bandwidth: Up to 3 GHz; see Figs. 115 and 116.

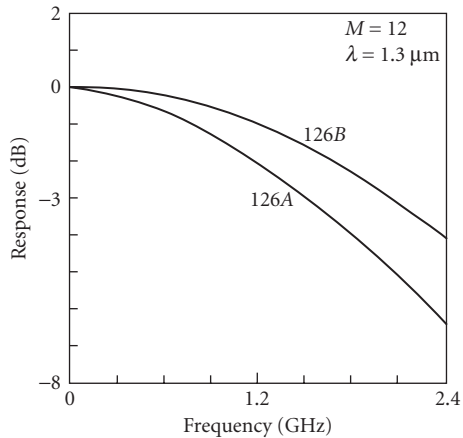


FIGURE 115 Frequency response of an InGaAs APD (126A/B). (AT&T, 126A/B, C ASTROTEC InGaAs.)

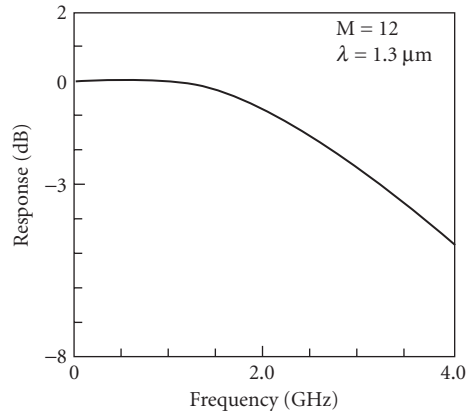


FIGURE 116 Frequency response of InGaAs APD (126C). (AT&T, 126A/B, C ASTROTEC InGaAs.)

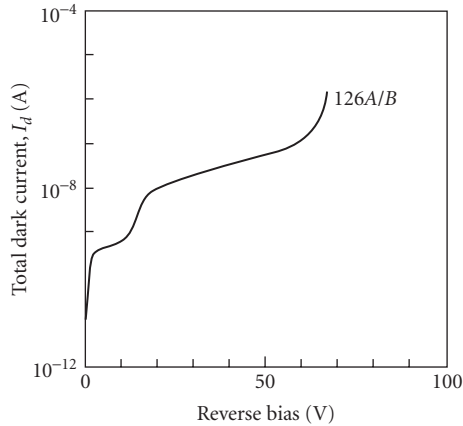


FIGURE 117 Dark current versus reverse bias of InGaAs APD (126A/B). (AT&T, 126A/B, C ASTROTEC InGaAs.)

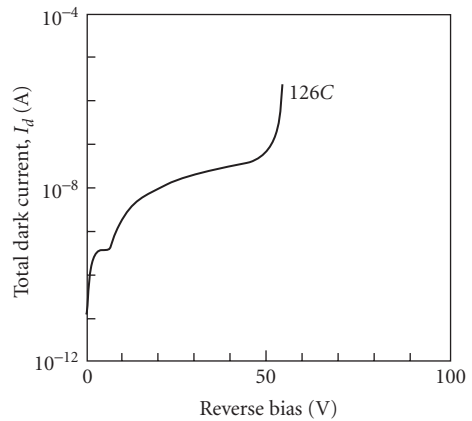


FIGURE 118 Dark current versus reverse bias of InGaAs APD (126C). (AT&T, 126A/B, C ASTROTEC InGaAs.)

Size: 0.04 to 0.5-mm diameter.

Dark current: Dependent upon reverse bias, device structure, and temperature. See Figs. 117, 118, and 119.

Recommended circuit: See Figs. 107 and 120.

Manufacturers: Advanced Photonix, AT&T, EG&G Canada, Edinburgh Instruments, Edmund Optics, Electro-Optical Systems, Electro-Optics Technology, Emcore, Epitaxx, Fermionics, GCA Electronics, Germanium Power Devices, Hamamatsu, New England Photoconductor, New Focus,

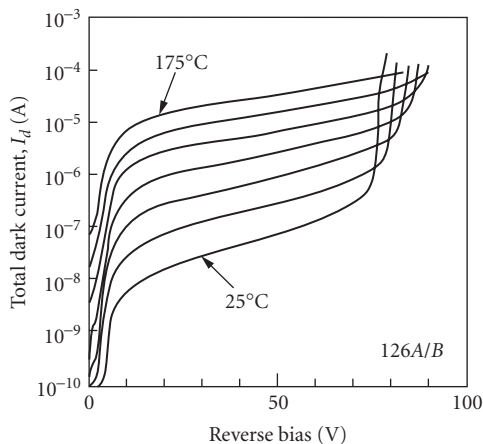


FIGURE 119 Dark current versus voltage of InGaAs APD as a function of temperature at 25°C increments. Note: The temperature dependence of the 126C dark current is the same as the 126A/B. (AT&T, 126A/B, C ASTROTEC InGaAs.)

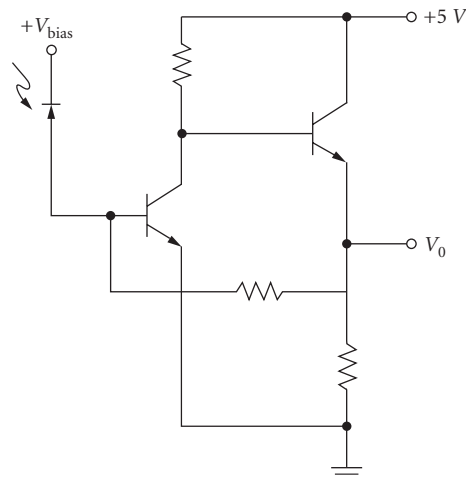


FIGURE 120 Bipolar transimpedance amplifier for InGaAs avalanche photodiode. (AT&T, 126A/B, C ASTROTEC InGaAs; *Optical Fiber Communications*, John M. Senior, © 1985, Prentice-Hall; ISBN-0-13-638248-7.)

Newport, North Coast Scientific, Opto-Electronics, Ortel, Photonic Detectors, Goodrich Sensors Unlimited, Spire, Swan Associates, Telcom Devices, Teledyne Judson Technologies, UDT Sensors.

Ge Germanium intrinsic photodetectors are similar to intrinsic silicon detectors but offer spectral response out to 1.5 to 1.8 μm . *PN* junction photodiodes offer submicrosecond response or high sensitivity from the visible region to 1.8 μm . Zero bias is generally used for high sensitivity and large reverse bias for high speed. As in the case of silicon, germanium has an indirect bandgap and soft spectral cutoff. The previous discussion on silicon detectors applies in general, with the exception that blue- and UV-enhanced devices are not relevant to germanium detectors. Germanium detectors, because of their narrower bandgap, have higher leakage currents at room temperature, compared to silicon detectors. Detector impedance increases about an order of magnitude by cooling 20°C below room temperature. Thus, performance can improve significantly with thermoelectric cooling or cooling to liquid nitrogen temperature.

As with silicon, the device structure and bias configuration can affect spectral response and rise time. Three detector types are available:

- *pn* junction
- *pin* junction
- Avalanche photodiode

Germanium *pn* and *pin*

Sensitivity: D^* (peak, 300 Hz, room temperature) $>2 \times 10^{11}$ Jones, increases significantly with cooling by thermoelectric cooler or liquid nitrogen. (See Figs. 121, 122, and 123.)

Quantum efficiency: >50 percent with antireflection coating.

Noise: See Figs. 124 and 125.

Responsivity: 0.9 A/W at peak wavelength. See Fig. 121.

Capacitance: Lower for *pin* structure compared with *pn* diode. See Fig. 126.

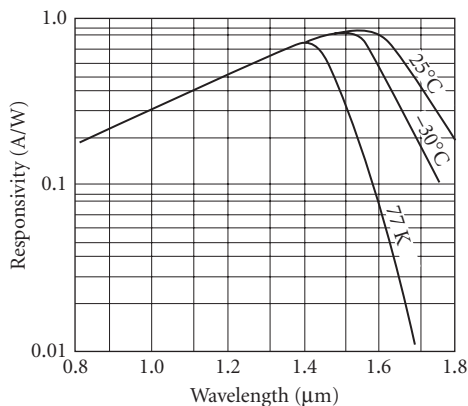


FIGURE 121 Spectral response for a germanium *pn* junction photodiode at three temperatures. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

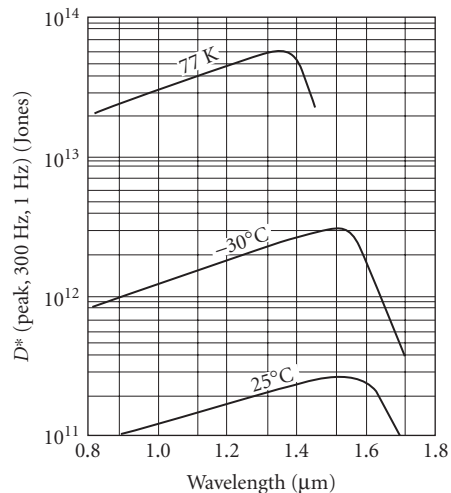


FIGURE 122 D^* as a function of wavelength for a germanium *pn* junction photodiode at three temperatures. (EG&G Judson, *Infrared Detectors*, 1994.)

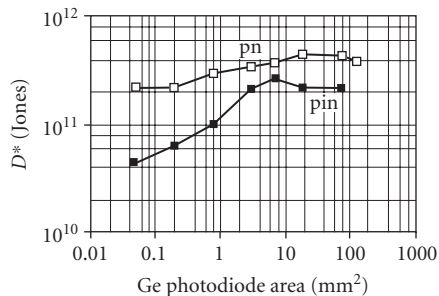


FIGURE 123 D^* for germanium pn and pin photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

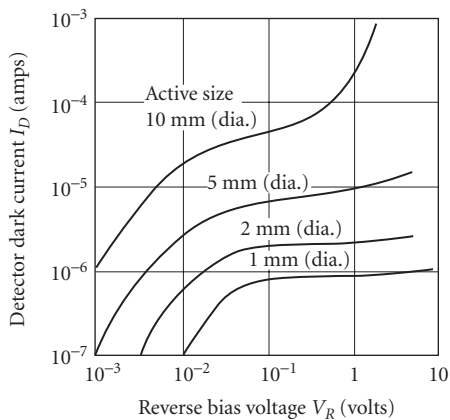


FIGURE 125 Dark current as a function of reverse bias for germanium pn junction photodiodes of different diameters at 25°C. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

Time constant: PIN diodes provide faster response. See Fig. 127.

Sensitive area: 0.25 to 13-mm diameter standard.

Operating temperature: Ambient, TE-cooled, or liquid nitrogen.

Profile: ± 2 percent across active area at 1.3 μm .

Linearity: Excellent over 10 orders of magnitude. See Fig. 128.

Recommended circuit: See previous section on silicon photodiodes.

Manufacturers: Edinburgh Instruments, Electro-Optical Systems, Judson, Electro-Optical Systems, Fastpulse Technology, Germanium Power Devices, Infrared Associates, Newport, North Coast Scientific, Opto-Electronics, Oxford Instruments, Scientific Instruments.

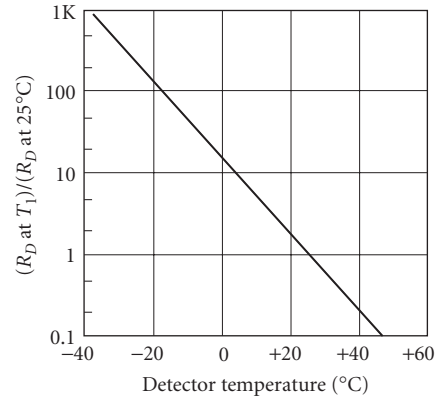


FIGURE 124 Ratio of resistance at temperature T to the resistance at 25°C for a germanium pn junction photodiode. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

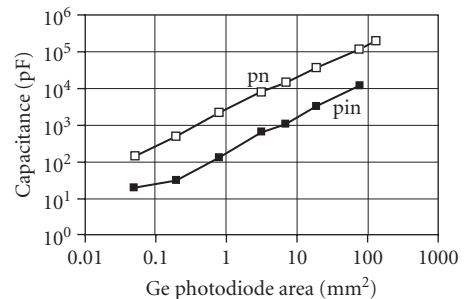


FIGURE 126 Capacitance for germanium pn and pin photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

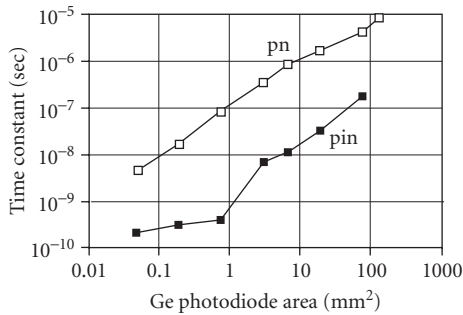


FIGURE 127 Time constant for germanium *pn* and *pin* photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

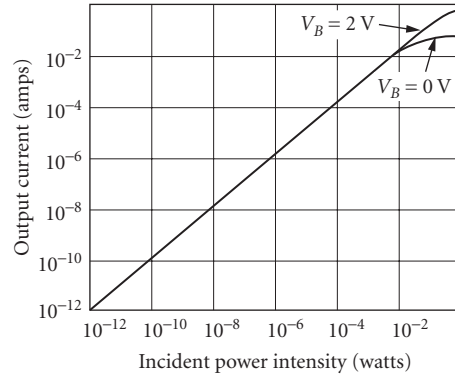


FIGURE 128 Linearity of a germanium *pn* junction photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

Germanium avalanche photodiode The germanium avalanche photodiode is similar to the silicon APD but has lower optimum gain, longer cut-off wavelength (1.7 μm), and higher leakage current. Germanium APDs combine the sensitivity of a Ge *pn* photodiode and the speed of a *pin* Ge photodiode.

Sensitivity: $D^* 2 \times 10^{11}$ Jones at 30 MHz for a diode with area of 5×10^{-2} mm^2 or about the same as for a *pn* Ge diode with the same area, and about a factor of 4 higher than a *pin* Ge photodiode (compare with Figs. 122 and 123). D^* depends on gain.

Gain: See Fig. 129.

Dark current: See Fig. 130.

Capacitance: 2 pF at 20-V reverse bias for 100- μm diameter, 8 pF at 20-V reverse bias for 300- μm diameter.

Quantum efficiency: 60 to 70 percent at 1.3 μm .

Responsivity: Photocurrent is the product of the incident optical power in watts, wavelength in micrometers, and quantum efficiency (η) divided by 1.24 and multiplied by the avalanche gain M . $I_{\text{photo}} = M(P\lambda\eta/1.24)$. (See Figs. 121 and 129).

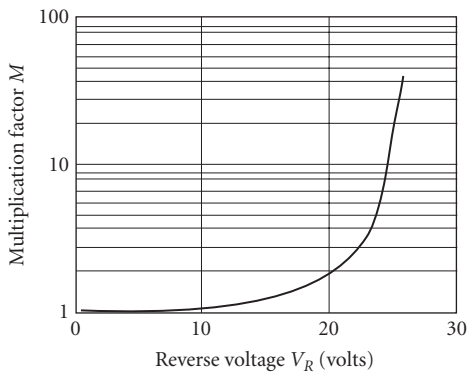


FIGURE 129 Gain as a function of reverse bias for germanium avalanche photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

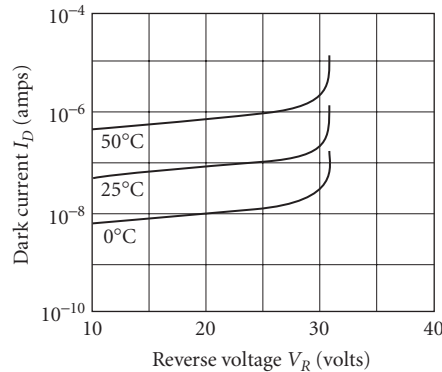


FIGURE 130 Dark current as a function of reverse bias for germanium avalanche photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

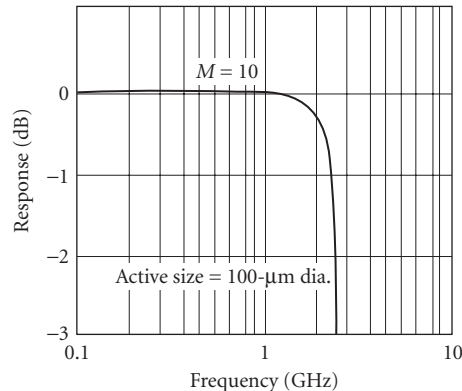


FIGURE 131 Frequency response of a germanium avalanche photodiode at two gain operating points. (EG&G Judson, *Infrared Detectors*, 1994.)

Operating temperature: Ambient or TE-cooled.

Time constant: 0.2 ns for 100- μm diameter; 0.3 ns for 300- μm diameter; both at 1.3 μm , $M = 10$ and with $R_L = 50 \Omega$. See Fig. 131.

Sensitive area: 100- and 300- μm diameter standard.

Recommended circuit: See circuit recommended for Si APD.

Manufacturers: Edmund Optics, Metrotech, North Coast Scientific, Teledyne Judson Technologies.

PbS Photoconductive lead sulfide was one of the earliest and most successful infrared detectors. Even today it is one of the most sensitive uncooled detectors in the 1.3- to 3- μm spectral region. With cooling, PbS sensitivity is competitive with other detectors out to about 4.2 μm ; however, its response time is slow.

Many PbS characteristics can be varied by adjusting the chemistry of the deposition process and/or the post-deposition heat treatment. These characteristics include spectral detectivity, time constant, resistance, and upper limit of operating temperature.²⁶ PbS is generally made by chemical reaction of Pb acetate and thiourea, except for high-temperature (373 K) applications, where evaporation is used. The material is deposited as a thin film (1 to 2 μm thick) on a variety of substrates, such as sapphire. With photolithographic processing, small sensitive areas can be made with comparatively high D^* values.

PbS may be tailored for ambient or room-temperature operation (ATO), intermediate or thermoelectrically cooled operation (ITO), and low-temperature or nitrogen-cooled operation (LTO). They are manufactured differently for particular temperature ranges, as shown in Table 4.

Sensitivity: $D^* 1.5 \times 10^{11}$ Jones at 295 K. See Figs. 132 to 137.

Responsivity: Depends on detector area, bias, resistance, and operating temperature (see Figs. 138 and 139).

Quantum efficiency: Generally limited by incomplete absorption in the thin film to 30 percent as estimated from blip D^* values.

Noise: Dominated by $1/f$ noise at low frequencies. See Figs. 135 and 140.

Time constant: Can be varied in manufacturing. Typical values are 0.2 ms at 295 K, 2 to 5 ms at 193 and 77 K. See Fig. 139.

Sensitive area: Typical sizes are square elements with dimensions of 0.5, 1, 2, and 5 mm on a side.

TABLE 4 PbS Performance Characteristics

	Typical operating temperature, K			
	350	273 (ATO)	193 (ITO)	77 (LTO)
Sensitivity	†	Figs. 132, 138	Figs. 133, 136, 138	Figs. 134, 138
$D^* (\lambda_{max})/D^* (500\text{ K})$		105	55	17
Noise, $V/Hz^{1/2}$		Fig. 140	Fig. 140	Fig. 140
Dark resistance, $M\Omega/sq$	<0.3	<2	<10	<20
Time constant, μs^\ddagger	50	100–500	5000	3000

†At 350 K, cutoff wavelength moves into $\sim 2.4\ \mu m$, with $D^* (\lambda_{max}) \approx 10^{10}$ Jones.

‡These are typical values; the time constant can be adjusted over two orders of magnitude in fabrication, but D^* is affected.

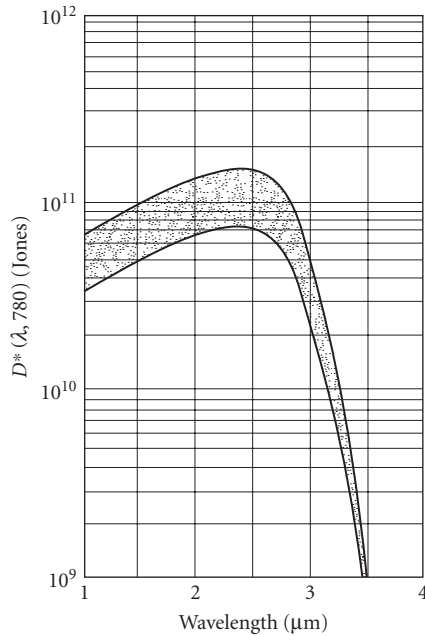


FIGURE 132 Range of spectral detectivities for PbS (ATO) at 295 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

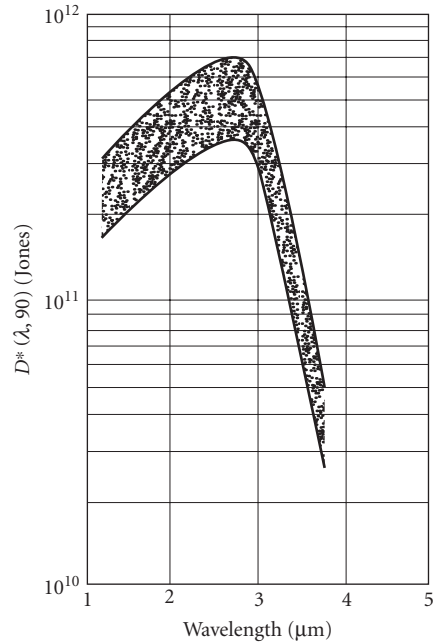


FIGURE 133 Range of spectral detectivities for PbS (ITO) at 193 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

Capacitance: <1 pF (limited by mounting configuration).

Recommended circuit: Standard photoconductor.

Stability: Exposure to visible and/or UV radiation can induce instability and drift. Stability will recover with storage in the dark, or by baking.

Sensitivity profile: Uniform within 10 percent.

Linearity: Excellent over broad range 10^{-8} to 10^{-3} W.

Manufacturers: Alpha Omega Instruments, Cal-Sensors, Edmund Optics, Electro-Optical Systems, Hamamatsu, New England Photoconductor, Teledyne Judson Technologies, OptoElectronics, Orielt.

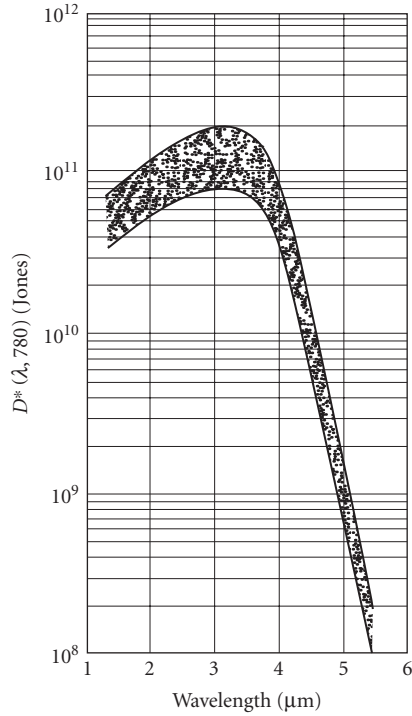


FIGURE 134 Range of spectral detectivities for PbS (LTO) at 77 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

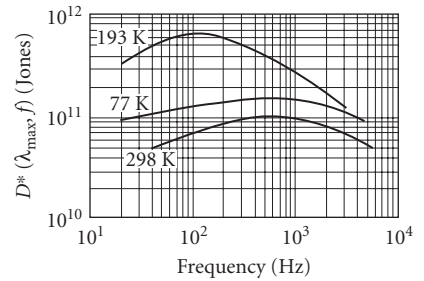


FIGURE 135 Example of detectivity vs. temperature for PbS detectors at various operating temperatures; 2π FOV, 295-K background. (Santa Barbara Research Center.)

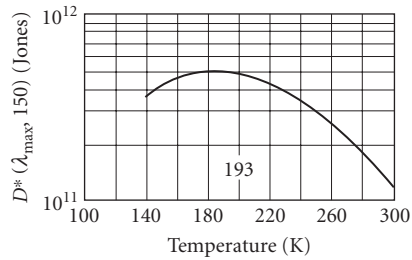


FIGURE 136 Example of detectivity versus temperature for PbS (ITO) detectors; 2π FOV, 295-K background. (Santa Barbara Research Center.)

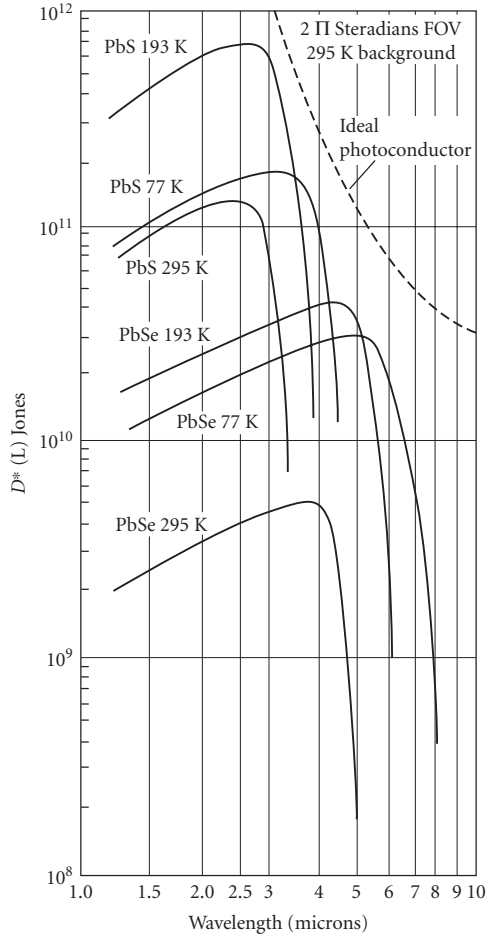


FIGURE 137 D^* versus wavelength for PbS and PbSe detectors operating at temperatures ranging between 77 K and 295 K. (CAL-SENSORS, *Infrared Detectors*.)

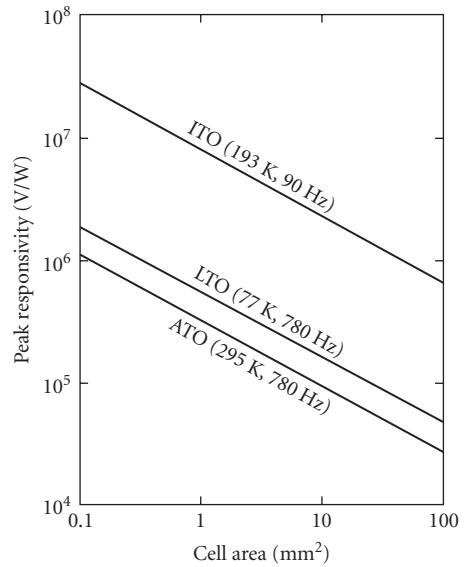


FIGURE 138 PbS typical peak responsivity versus cell area (actual values range within a factor of two of these shown.)

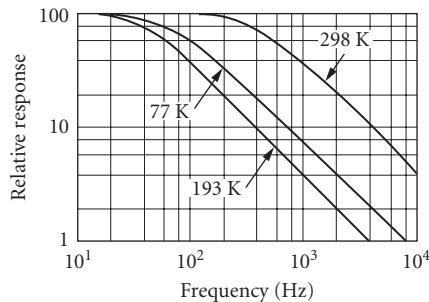


FIGURE 139 Example of signal versus frequency for PbS detectors. (Santa Barbara Research Center.)

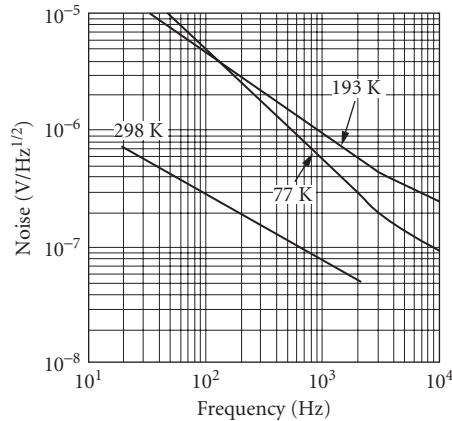


FIGURE 140 Example of noise versus frequency for PbS detectors. (Santa Barbara Research Center.)

InAs (Photovoltaic) InAs detectors are single-crystal, intrinsic, direct-bandgap photovoltaic devices for use in the 1- to 4- μm region (spectral cutoff varies with temperature). At room temperature, InAs provides good sensitivity and submicrosecond response times. At 195 K, InAs performance equals or better the sensitivity of any other detector in the 1 to 3.5- μm region. Devices with sapphire immersion lenses are available to increase signal responsivity for operation at higher temperatures where the detector is thermal-noise-limited. Compared to PbS and PbSe detectors, InAs has very little low frequency ($1/f$) noise if operated in the photovoltaic mode.

Sensitivity: D^* (peak) varies from 1.2×10^9 Jones at 295 K to 6×10^{11} Jones at 77 K. See Fig. 141.

Quantum efficiency: Maximum of about 75 percent without antireflection coating.

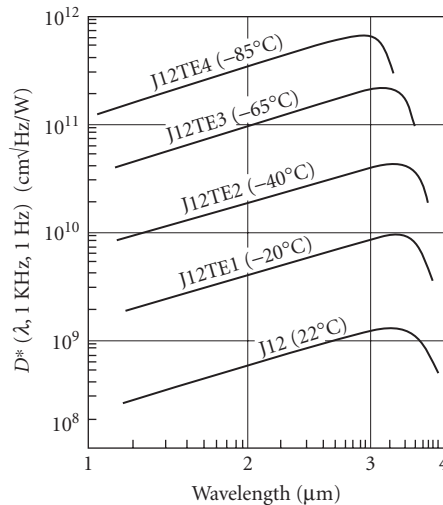


FIGURE 141 D^* versus wavelength for InAs detector operating at temperatures ranging between 77 K and 295 K. (Teledyne Judson Technologies, 2008.)

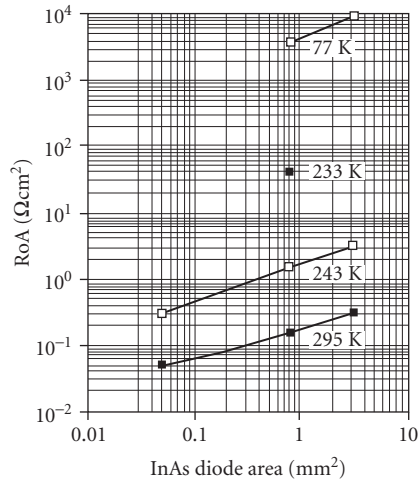


FIGURE 142 R_oA of InAs photodiodes shown as a function of diode area. Lower impedance per unit area for smaller devices indicates that these devices are surface-leakage-limited. (Santa Barbara Research Center.)

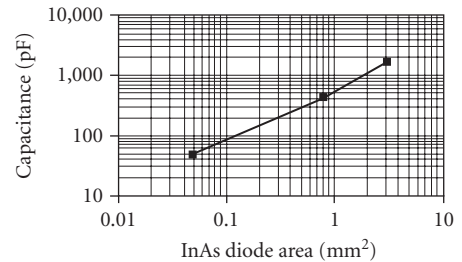


FIGURE 143 Capacitance of InAs photodiodes shown as a function of diode area. Capacitance will not change appreciably with temperature. (Teledyne Judson Technologies, 2008.)

Noise: Low impedance tends to make preamplifier noise dominate at room temperature; background limited (for 300 K background) at operating temperatures below 200 K.

Time constant: Less than 0.5 μ s at all temperatures when low values of load resistor are used to reduce the RC time constant.

Responsivity: 0.5 to 1.25 A/W at peak.

Dynamic resistance: See Fig. 142.

Diode capacitance: See Fig. 143.

Sensitive area: Standard sizes 0.25 to 2-mm diameters.

Operating temperatures: 77 to 300 K.

Linearity: Anticipated to be very good over many decades.

Sensitivity profile: ± 15 percent.

Recommended circuits:

Open circuit: PV InAs detectors with areas less than 2×10^{-2} cm require no bias when operated and can be connected directly into the input stage of amplifier (capacitor ensures elimination of dc bias from amplifier) (Fig. 144a).

Transformer: Useful when using InAs at zero bias, particularly at room temperature where diode impedance is low (Fig. 144b).

Reversed bias: At temperatures greater than 225 K considerable gain in impedance and responsivity is achieved by reverse-biasing (Fig. 144c).

Fast response: To utilize the short intrinsic time constant, it is sometimes necessary to load the detector to lower the RC of the overall circuit (reverse bias will also lower detector capacitance) (Fig. 144d).

Manufacturers: Electro-Optical Systems, Hamamatsu, Teledyne Judson Technologies.

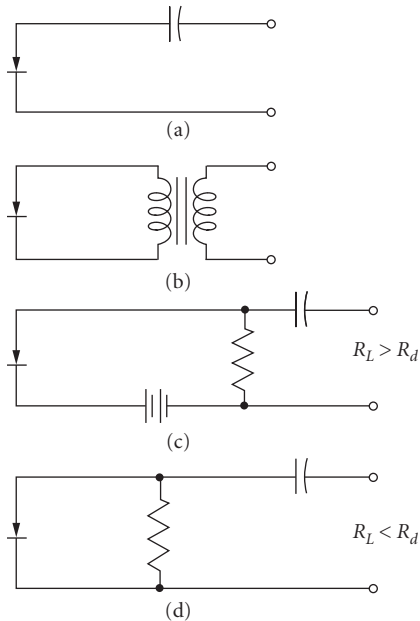


FIGURE 144 Recommended circuits for InAs detectors: (a) open circuit, (b) transformer, (c) reversed bias, (d) fast response.

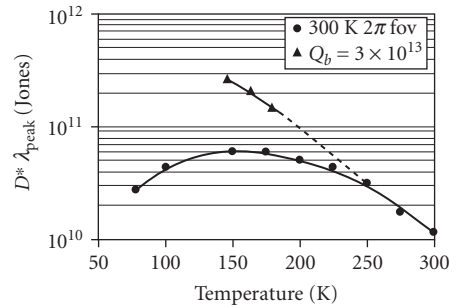


FIGURE 145 D^* of PbSe as a function of a function of temperature for two background flux conditions: high background of 2π field of view and reduced background of 3×10^{13} photons/cm²/s. D^* at the higher background flux reaches a maximum around 160 K because the background noise increases at lower temperatures due to the increase in long-wavelength spectral response of the detector. (Santa Barbara Research Center.)

PbSe Lead selenide is an intrinsic, thin-film photoconductor, whose long-wavelength spectral response and speed of response exceeds that of PbS. At room temperature, PbSe has peak D^* which can exceed 1×10^{10} Jones with a spectral cutoff out to 4.4 μm . At liquid-nitrogen temperature, InSb offers twice the D^* , largely because PbSe offers response out to 7 μm at 77 K, considerably longer than InSb. However, for intermediate temperatures, from 180 K to room temperature, PbSe offers competitive D^* combined with moderately fast response.²⁶ PbSe technology has made a significant advance in the past decade in some vendors being able to reproducibly make high-performance detectors.

Sensitivity: $D^* \approx 1 \times 10^{10}$ Jones at 300 K, increases with cooling (see Figs. 137 and 145). D^* is limited by $1/f$ noise at low frequencies (see Fig. 146).

Response: Figure 147 shows responsivity in amperes per watt for a high-quality detector with a length of 0.016 cm and width of 0.024 cm. Responsivity in volts per watt is obtained by multiplying A/W data by resistance (see Fig. 148). Responsivity will vary inversely with detector length (see Figs. 149 and 150).

Noise: Figure 151 shows the noise as a function of temperature for a detector with a length of 0.016 cm and width of 0.024 cm. Noise as a function of frequency is shown in Fig. 152.

Resistance: Figure 148 shows the resistance as a function of temperature for a detector with a length of 0.016 cm and width of 0.024 cm.

Capacitance: 1 pF (limited by mounting configuration).

Time constant: See Figs. 153 and 154. Time constant will be longer when detector is operated in reduced background flux condition.

Stability: Exposure to visible and/or UV radiation can induce instability and drift. Stability will recover with storage in the dark at room temperature.

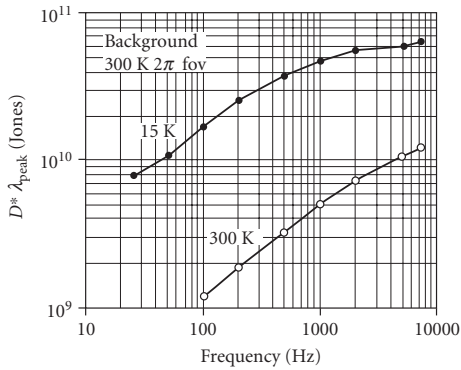


FIGURE 146 D^* of PbSe as a function of frequency for two temperatures. PbSe has considerable $1/f$ noise which reduces D^* at lower frequencies. (Santa Barbara Research Center.)

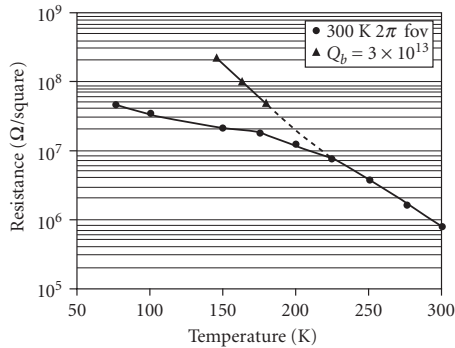


FIGURE 148 Resistance (Ω /square) of PbSe thin-films as a function of temperature for two background flux levels. At any temperature, the absolute value can be varied by altering the manufacturing process in chemical deposition and/or heat treatment. (Santa Barbara Research Center.)

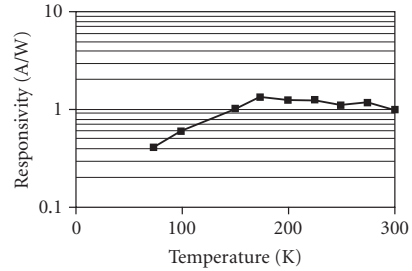


FIGURE 147 Responsivity in A/W of PbSe thin-film photoconductive detectors as a function of temperature for a high background flux level. Multiply by detector resistance to get V/W. (Santa Barbara Research Center.)

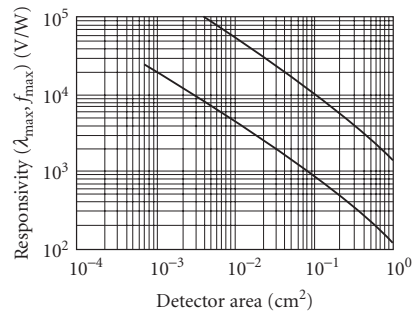


FIGURE 149 Expected range of peak responsivities versus detector size, typical PbSe (ATO) infrared detectors. (Santa Barbara Research Center.)

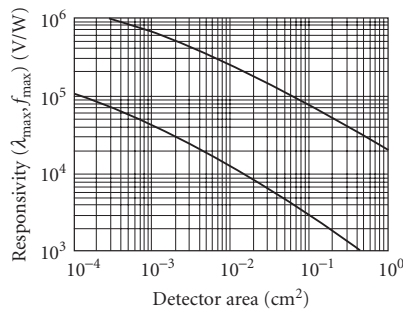


FIGURE 150 Expected range of peak responsivities versus detector size, typical PbSe (ITO and LTO) infrared detectors. (Santa Barbara Research Center.)

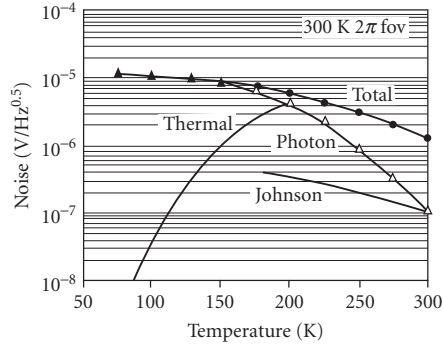


FIGURE 151 Noise voltage (per square root of bandwidth) of PbSe thin-film photoconductive detectors as a function of temperature for a high back ground flux level. Photon noise is dominant below 200 K. Thermal noise is dominant at higher temperatures. Total noise levels are well above typical preamplifier noise. (Santa Barbara Research Center.)

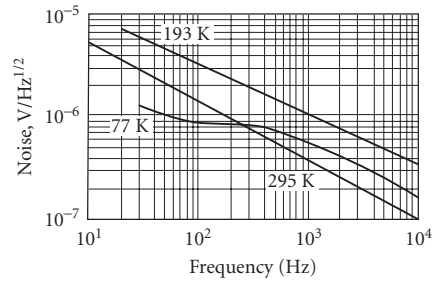


FIGURE 152 Example of noise versus frequency for PbSe detectors (ATO, ITO and LTO types) (1 × 1 mm). (Santa Barbara Research Center.)

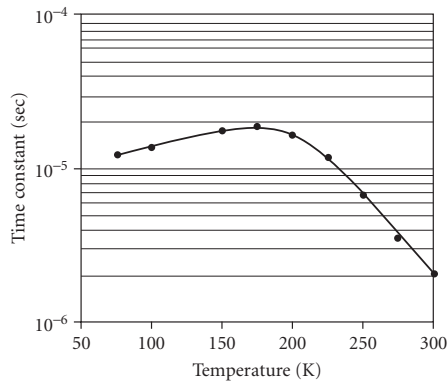


FIGURE 153 Time constant of PbSe thin-film photoconductive detectors as a function of temperature. (Santa Barbara Research Center.)

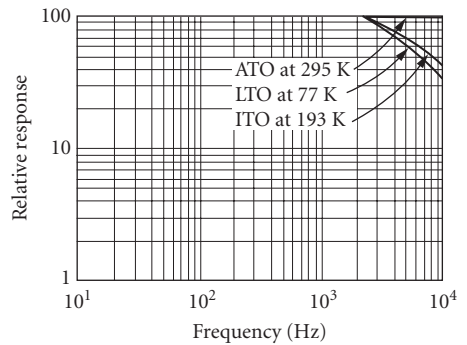


FIGURE 154 Example of signal versus frequency for PbSe detectors (ATO, ITO and LTO). (Santa Barbara Research Center.)

Recommended circuit: Standard photoconductor.

Operating temperature: 77 to 300 K.

Manufacturers: Cal-Sensors, Edmund Optics, Electro-Optical Systems, Hamamatsu

InSb Historically, indium antimonide material has been used for at least four different radiation detector types, two of which, the photoconductive and photoelectromagnetic types, are no longer widely used. We discuss here the intrinsic photovoltaic device. [The very far infrared bolometer (InSb bolometer) was previously discussed in Sec. 24.7.

At 77 K, InSb photodiodes offer background limited sensitivity at medium-to-high background flux conditions in the 1- to 5.5- μm spectral range. At lower temperatures, they provide sensitive detectors at low background flux levels such as in astronomy applications, but with a slightly shortened long-wavelength cutoff. Operation is possible up to as much as 145 K, but because the spectral

response increases with increasing temperature, the detector impedance drops rapidly leading to significant thermal noise.

Sensitivity: Spectral response out to 5.5 μm at 77 K (see Fig. 155). $D^* 1 \times 10^{11}$ Jones, increases with reduced background flux (narrow field of view and/or cold filtering) as illustrated in Fig. 156.

Quantum efficiency: ~ 60–70 percent without antireflection coating. >90 percent with antireflection coating.

Noise: Background current limited over wide range of background flux at 77 K (see Fig. 157).

Time constant: <1 μs .

Responsivity: 3 A/W at 5 μm without antireflection coating.

Noise equivalent power (NEP): Frequency dependence is shown in Fig. 158 for three detector sizes.

Capacitance: Typically 0.05 F/cm².

Impedance: Top-grade detectors have $1\text{--}5 \times 10^6 \Omega\text{cm}^2 R_0A$ product at 77 K, at zero bias, and without background flux.

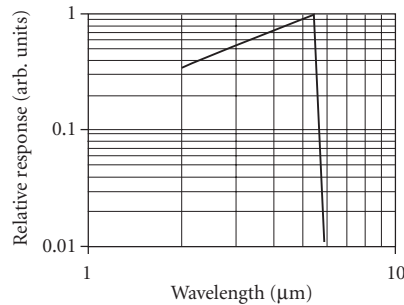


FIGURE 155 Relative spectral response per watt of an InSb photodiode without antireflection coating. the direct bandgap results in a sharp spectral cutoff. (Santa Barbara Research Center.)

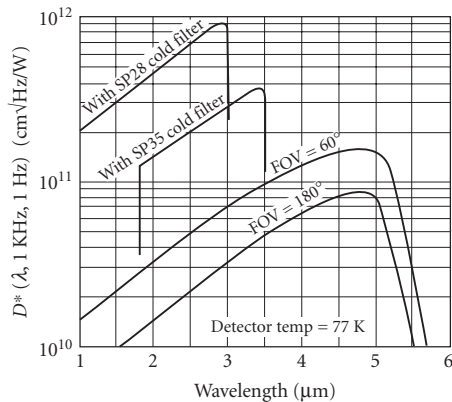


FIGURE 156 D^* as a function of wavelength for an InSb detector operating at 77 K. (Teledyne Judson Technologies, 2008.)

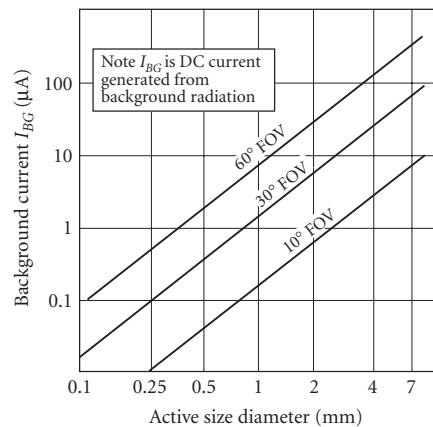


FIGURE 157 Background current as a function of photodiode area for InSb detectors operating at 77 K, shown at three values of the detector field of view. (Teledyne Judson Technologies, 2008.)

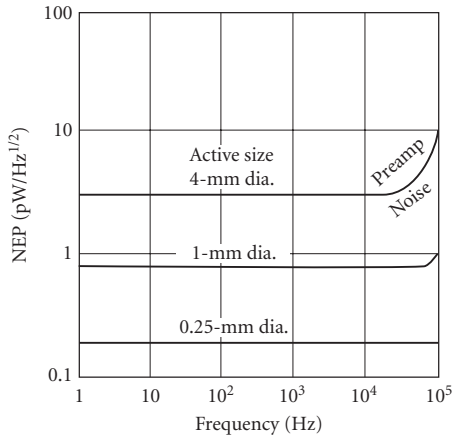


FIGURE 158 NEP as a function of frequency for three sizes of InSb photodiodes. (Teledyne Judson Technologies, 2008.)

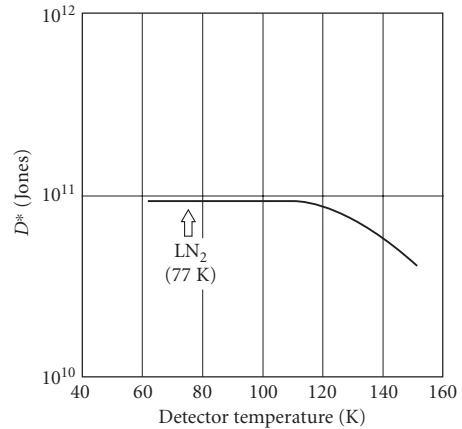


FIGURE 159 InSb photodiode D^* as a function of operating temperature between 77 K and 150 K, for a 2π (180°) FOV. (Teledyne Judson Technologies, 2008.)

Sensitive area: 0.04×0.04 -mm square to 1×5 -mm rectangle; 0.25 to 10-mm diameter.

Operating temperature: Normally 77 K; InSb can be used up to approximately 145 K (see Fig. 159).

Linearity: Linear to ~ 1 -mW/cm² flux.

Sensitivity profile: ± 15 percent or better.

Stability: Devices from some vendors are subject to “flashing,” where exposure to visible or UV flux causes a change in the insulating surface charge thereby causing a change in the diode impedance. The detector typically recovers at room temperature.

Recommended circuit: Same as for Si and Ge photodiodes; zero or reverse bias in combination with a load resistor and low-noise preamplifier. Low-impedance load resistor can be used for obtaining fast response, with consequences of reduced sensitivity.

Manufacturers: L3 Cincinnati Electronics, Edinburgh Instruments, New England Photoconductor, Teledyne Judson Technologies, Electro-Optical Systems, Hamamatsu, Infrared Associates.

Ge:Au Gold-doped germanium detectors are relatively fast single-crystal p-type impurity-doped photoconductors for the 2- to 9- μ m region. Although not the most sensitive detector anywhere in its range of spectral sensitivity, Ge:Au offers respectable sensitivity over a broad spectral region using liquid nitrogen cooling. Sensitivity can be improved by a factor of 2.5 by operating at $T < 65$ K (pumped liquid nitrogen or other cryogen). At these temperatures, Ge:Au becomes background limited.

Sensitivity: See Figs. 160 to 163.

Quantum efficiency: Dependent on wavelength, detector geometry (absorption thickness), antireflection coating, and enclosure (integration chamber can increase absorption). $D_{\lambda_{pk}}/D_{500\text{K}}^* = 2.7$ (see Fig. 164).

Noise: See Fig. 165.

Time constant: < 50 ns with full D^* [shorter response times (< 2 ns) can be tailored by heavy concentration of compensating (n -type) dopant and suitable bias circuit (see circuit discussion to follow). Heavy compensation increases resistance, and hence the incoherent signal-to-noise ratio becomes limited by the thermal noise of the load (typically a factor of 2 degradation in the signal-to-noise ratio). Quantum efficiency, however, is not significantly altered, so that a high compensation concentration does not hurt the coherent-detection signal-to-noise ratio].

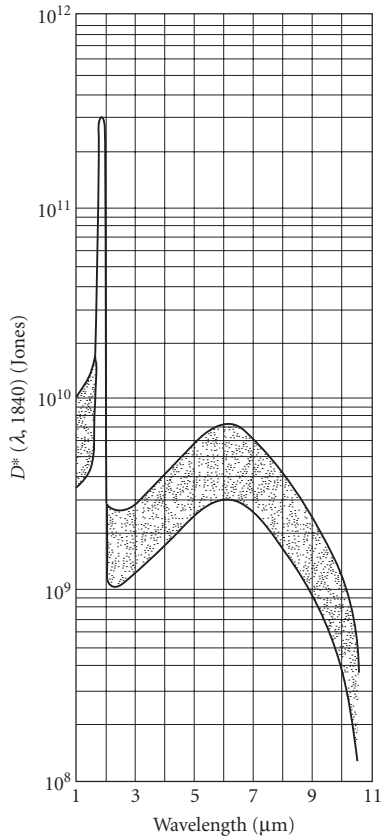


FIGURE 160 D^* versus λ for Ge:Au; $T = 77$ K, 2π FOV; 295-K background. (Santa Barbara Research Center.)

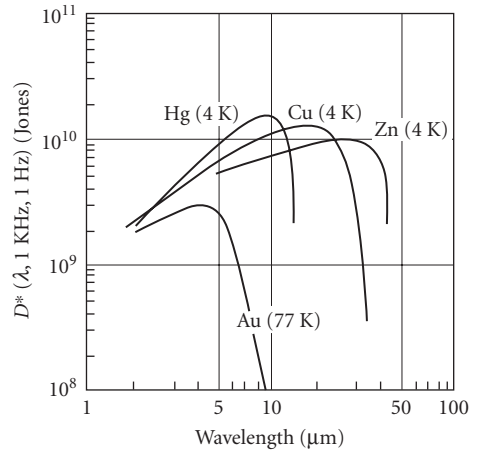


FIGURE 161 D^* as a function of wavelength for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn, for a 300-K 2π (180°) FOV background flux. (EG&G Judson, *Infrared Detectors*, 1994.)

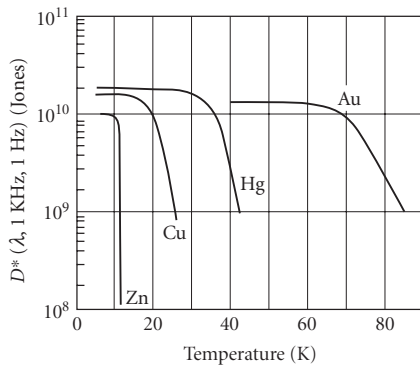


FIGURE 162 D^* as a function of operating temperature for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn for a 300-K 2π (180°) FOV background flux. (EG&G Judson, *Infrared Detectors*, 1994.)

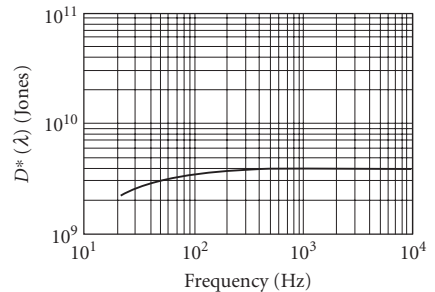


FIGURE 163 Typical D^* versus frequency ($T = 77$ K). (Santa Barbara Research Center.)

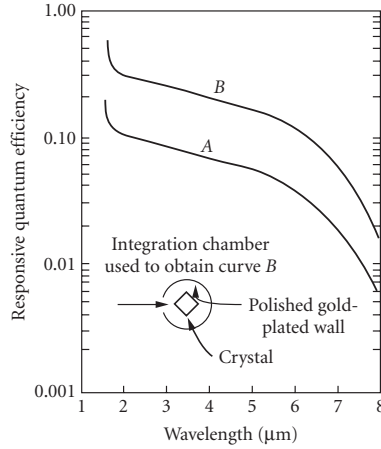


FIGURE 164 Quantum efficiency versus λ for Ge:Au ($T = 78$ K). (Santa Barbara Research Center, internal report.)

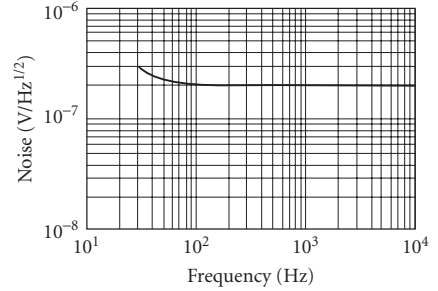


FIGURE 165 Typical noise spectrum for Ge:Au ($T = 77$ K, $A = 1 \times 1$ mm). (Santa Barbara Research Center.)

Responsivity: Dependent upon bias and geometry, typical values are 0.1 to 0.2 A/W at 77 K. Multiply by detector resistance to get V/W.

Dark resistance: Varies with background flux and effective quantum efficiency (see previous quantum efficiency discussion), range may be 20 k Ω to 5 M Ω , or much greater under very low background flux conditions if adequately cooled to limit thermally activated conductivity. (Also see previous time-constant discussion.)

Capacitance: Depends on device geometry and mounting, typically <1 pF.

Sensitive area: 1 to 5-mm diameter standard.

Operating temperature: < 85 K (normally 77 K, but see Fig. 162).

Recommended circuit: Standard photoconductive. See Fig. 166.

Manufacturers: No suppliers are presently known.

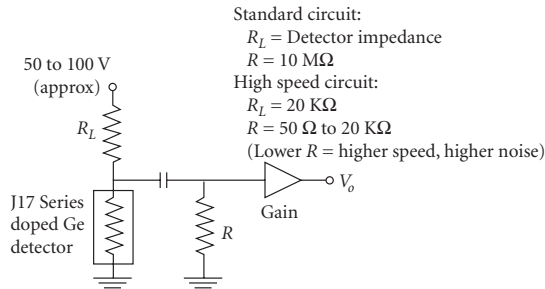


FIGURE 166 Basic operating circuit for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn, for a 300-K 2π (180°) FOV high background flux. If the detector is operated in very low background flux conditions, the detector impedance can become very high. Cooled JFET ($T > \approx 50$ K) or PMOS buffer amplifiers can be helpful in impedance matching under these conditions. (EG&G Judson, *Infrared Detectors*, 1994.)

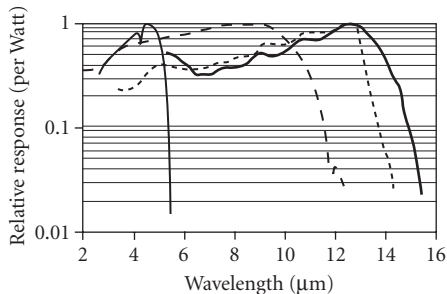


FIGURE 167 Relative spectral response per watt at 80 K for photoconductive HgCdTe detectors with antireflection coating. The curves are normalized to unity at peak value. The spectral cutoff can be adjusted by varying the ratio of HgTe to CdTe in the alloy. (Santa Barbara Research Center.)

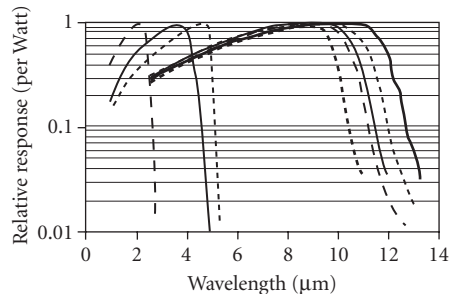


FIGURE 168 Relative spectral response per watt at 80 K for photovoltaic HgCdTe detectors without antireflection coating. The curves are normalized to unity at peak value. The spectral cutoff can be adjusted by varying the ratio of HgTe to CdTe in the alloy. (Santa Barbara Research Center.)

HgCdTe, HgZnTe, HgMnTe, etc. Mercury cadmium telluride is a direct-bandgap compound alloy semiconductor, made of chemical group II and VI elements, whose peak sensitivity at a particular temperature can be adjusted from 1 to 30 μm by varying the ratio of HgTe to CdTe (see Figs. 167 and 168). In addition to HgCdTe, other combinations of chemical groups II and VI elements can be used to produce similar variable spectral cutoff compound alloys, including HgZnTe, HgMnTe, HgCdZnTe, etc. For almost all purposes, HgCdTe will be as good as any other II-VI alloy detector, so we will speak here about it exclusively, but provide some data from the other alloys noted above. Both photoconductive (PC) and photovoltaic (PV) HgCdTe detectors are available for background-limited, high-speed, intrinsic photon detection in the SWIR, MWIR, and LWIR regions. Photoconductive devices with VLWIR response out to 25 μm are also available. HgCdTe detectors can be used at room temperature, with TE cooling, and at 77 K and lower temperatures. Sensitivity generally increases with cooling, depending upon the spectral cutoff and background flux. MWIR, LWIR, and VLWIR spectral range devices are generally operated at 77 K or lower temperature for maximum sensitivity, depending upon the background flux. The PC mode is advantageous when cooling is limited, since the thermal noise of a photoconductor increases less rapidly than a photodiode as the temperature is raised. Photoconductive devices with response out to 25 μm can be usefully operated at liquid nitrogen temperature and are popular for IR spectroscopy for this reason.

HgCdTe photoconductors are fabricated from thin ($\approx 10\text{-}\mu\text{m}$) single-crystal slices or epitaxial layers with metal contacts at each end of the element (see Fig. 8). They are low-impedance devices with 15 to 2000 Ω/square , depending upon the alloy composition, carrier concentration, operating temperature, background flux, and surface treatment. Photoconductor time constants at 77 K may be $\approx 2 \mu\text{s}$ for devices having a 12- μm cutoff, with longer time constants for shorter cutoffs, and shorter time constants for higher operating temperatures. In the case of small detector elements, the time constant may be reduced with increasing bias voltage because photoexcited carriers will be transported to the electrical contacts where they recombine. The spectral noise characteristics of PC HgCdTe typically exhibit $1/f$ noise out to a range of 50 Hz to 1 kHz or more, the value depending upon the detector quality, long-wavelength response, operating temperature, and background flux. White noise levels range from less than 1 nV/ $\sqrt{\text{Hz}}$ (where preamplifier noise may then dominate), up to 20 nV/ $\sqrt{\text{Hz}}$, depending upon detector quality, size, geometry, applied bias, temperature, and alloy composition. Photoconductive HgCdTe detectors are typically antireflection coated with a quarter-wave ZnS film, giving a peak quantum efficiency in the range of 85 to 90 percent, although this figure is only indirectly measured because the PC gain can be much greater than unity. Without antireflection coating, the quantum efficiency is typically 70 percent, limited by the optical index of ≈ 4 .

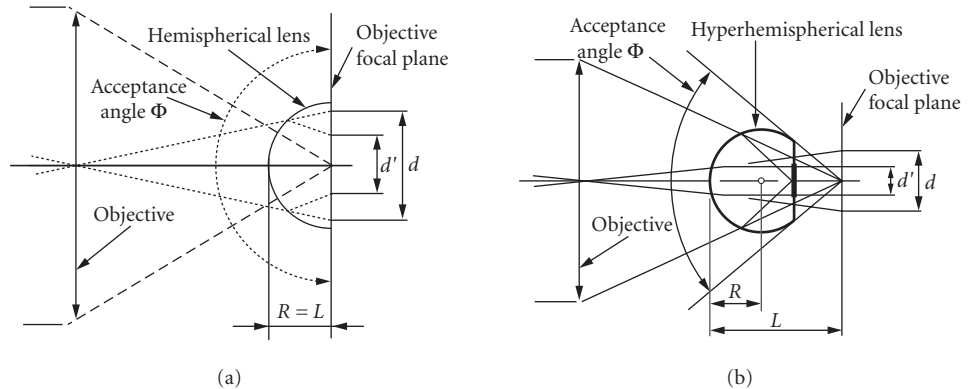


FIGURE 169 Schematic of an optically immersed HgCdZnTe detector with (a) hemispherical lens and (b) hyperhemispherical lens. Dimensions are summarized in Table 5. (*Vigo Systems*.)

TABLE 5 Summary of the Dimensions Depicted in Fig. 169, and Their Impact on D^* (*Vigo Systems*)

Parameter	Hemisphere	Hyperhemisphere
Distance, L	$L = R$	$L = R(n + 1)$
d/d'	n	n^2
$D^*_{\text{imm}}/D^*_{\text{non-imm}}$	n	n^2
Acceptance angle, Φ	$\Phi = 180^\circ$	$\Phi = 2 \arcsin(1/n)$
F/#	0.5	1.55

n = index of refraction (approx. 3.3 for GaAs and 2.7 for CdTe)

A class of HgCdTe detectors is offered for detection at TE-cooled and room temperature which are optically “immersed” with a hemispherical or hyperhemispherical lens of Ge, CdTe, GaAs, or other high-index material (see Fig. 169 and Table 5). The lens increases the effective area of the detector without increasing the detector noise, provided the noise is dominated by thermal rather than photon noise as is the case for minimal cooling. The lens must be in intimate contact with the detector surface ($\ll 1 \mu\text{m}$ spacing) to avoid total internal reflection of off-axis rays at the lens-detector interface. Immersed detectors offer up to a factor of n^2 (where n is the optical index) increased detector signal, which can mean an increase in D^* by the same factor for a thermal-noise-limited device. Operation of LWIR PC detectors at TE-cooled and room temperature is generally accompanied by increased $1/f$ noise which dominates out to higher frequencies.

Photovoltaic HgCdTe detectors ideally offer $\sqrt{2}$ higher D^* than detectors operating in the PC mode. Diodes are made in both n^+p and p^+n polarities, depending upon the manufacturer’s capabilities. The R_0A product of HgCdTe photodiodes varies significantly with temperature, spectral cutoff, and device quality. It also varies with the amount of background flux incident on the device. The R_0A product defines the maximum D^* in the limit of reduced background flux (see Fig. 170 and Fig. 19). In addition to theoretically higher D^* , high-quality PV HgCdTe detectors have lower $1/f$ noise than PC HgCdTe devices, with $1/f$ knee frequencies as low as 1 Hz or less. However, the noise of PV detectors increases more rapidly with increasing temperature than for PC detectors, making photodiodes less attractive for applications where cooling is limited. Photodiodes of high quality are more difficult to make than good photoconductors and can be expected to warrant a premium price. Antireflection coating is available from some diode producers, but is not routinely offered.

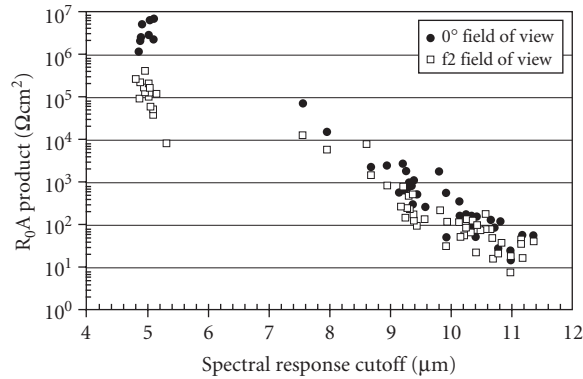


FIGURE 170 R_0A product “trendline” for small (25×25 to $100 \times 100 \mu\text{m}$) HgCdTe photodiodes at 77 K. Data is shown for devices with zero background flux (0° FOV) and with an F/2 field of view (29°) of a 300-K background. For 5- μm spectral cutoff material, the R_0A product is higher at 0° FOV by about an order of magnitude, compared with an F/2 background. At 10 μm , there is less of a difference between the two background conditions. Note that the R_0A product will generally be somewhat lower for larger area diodes.

Both PC and PV HgCdTe detectors are useful for infrared heterodyne detection. When sufficient local oscillator power is available, detector cooling becomes less important, since photon noise can dominate thermal noise at comparatively higher temperatures. Other things being equal, the photovoltaic detector has $\sqrt{2}$ sensitivity (signal-to-noise voltage) advantage over the photoconductor. For 10.6- μm heterodyne detection, 0.1×0.1 -mm HgCdTe *pin* photodiodes with sensitivity near the quantum limit of $\approx 2 \times 10^{-20}$ W/Hz are available with bandwidths up to several gigahertz. Ordinary photodiodes of the same area have bandwidths of several hundred megahertz. Photoconductors make better 10.6- μm heterodyne detectors when cooling is limited to TE-cooled temperatures of 180 K up to room temperature. At 180 K, TE-cooled photoconductors offer bandwidths of 50 to 100 MHz and heterodyne NEPs of 1 to 2×10^{-19} W/Hz. At room temperature, the NEP at 10.6 μm is limited to about 1×10^{-16} W/Hz. Immersion does not improve the performance of minimally cooled heterodyne detectors, since optical gain is already provided by the local oscillator.

Photoconductive HgCdTe

Sensitivity: Adjustable by varying alloy composition (see Figs. 167, 171 to 174).

Dark resistance: 15 to 2000 Ω/sq depending upon temperature, spectral cutoff, and surface passivation.

Responsivity: Varies with spectral cutoff, temperature, detector resistance, element length, and bias voltage or power. See Eq. (21) and Fig. 175 for detector elements with 50×50 - μm dimensions.

Noise: $1/f$ noise is dominant at frequencies below 50 to 1000 Hz for LWIR detectors at 77 K (greater for LWIR at room temperature or TE-cooled). Generation-recombination (thermal or photon) white noise is present beyond the $1/f$ region at a level of less than 10^{-9} V/ $\sqrt{\text{Hz}}$ to 2×10^{-8} V/ $\sqrt{\text{Hz}}$, depending upon spectral cutoff, background flux, responsivity, bias, and operating temperature. Noise and signal rolloff at high frequency is determined by the time constant. See Fig. 176 for an example of the noise spectrum of an LWIR detector at 77 K.

Operating temperature: 77 K and below to 300 K and above for short spectral cutoffs and/or with significant D^* reduction for operation at higher temperatures. Detector immersion can increase D^* at elevated temperatures where thermal noise is dominate.

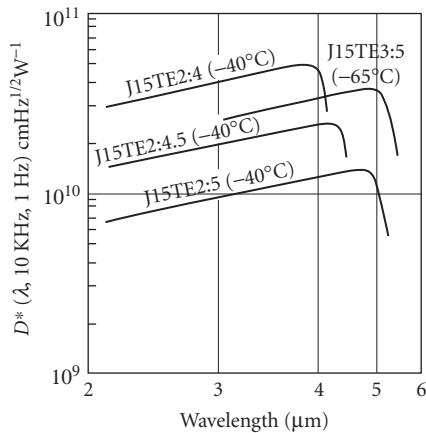


FIGURE 171 Typical D^* as a function of wavelength for a variety of MWIR HgCdTe photoconductors with thermoelectric cooling. (Teledyne Judson Technologies.)

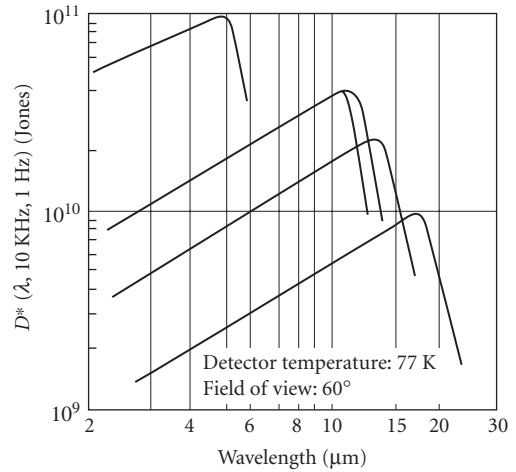


FIGURE 172 Typical D^* as a function of wavelength for a variety of LWIR and VLWIR HgCdTe photoconductors at 77 K. (Teledyne Judson Technologies.)

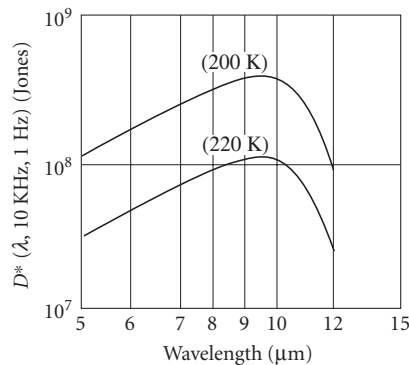


FIGURE 173 Typical D^* as a function of wavelength for LWIR HgCdTe photoconductors at 200 and 220 K. These units are cooled with three- or four-stage thermoelectric coolers. (Teledyne Judson Technologies.)

Linearity: At 77 K linearity begins to degrade at photon flux levels above $\sim 10^{-3}$ W/cm². At 200 K linearity begins to degrade at photon flux levels above ~ 1 W/cm².

Sensitive area: 0.025 to 4-mm linear dimensions.

Quantum efficiency: Typically >70 percent, 85 to 90 percent with antireflection coating.

Capacitance: Low, limited by mounting configuration.

Time constant: 1–2 μ s for LWIR at 77 K (see Fig. 176), depends on spectral cutoff, temperature, doping, and bias.

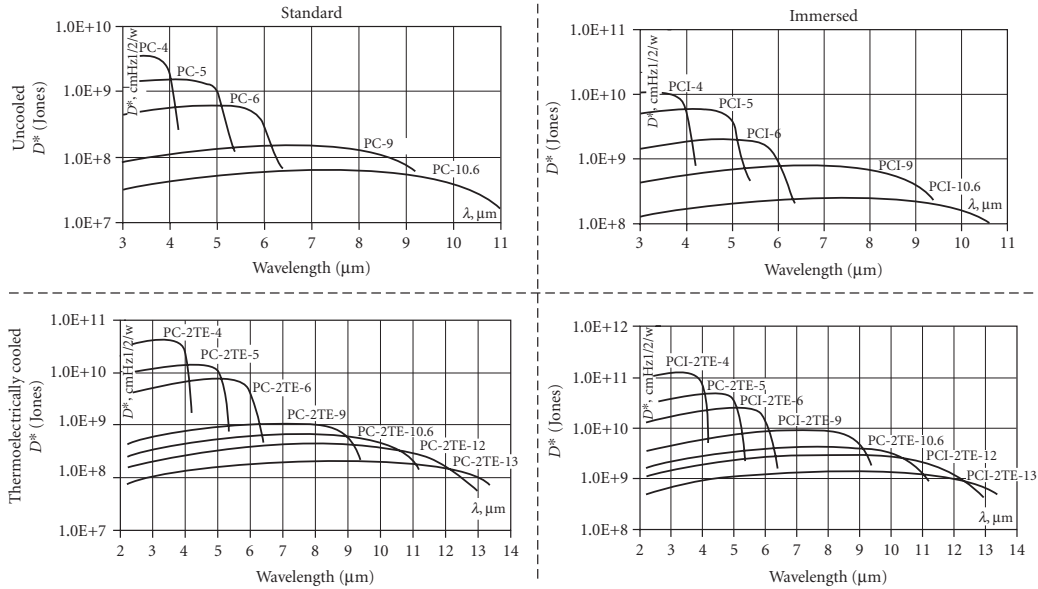


FIGURE 174 D^* for photoconductive HgCdTe detectors as a function of wavelength: *upper left*—ambient temperature operation; *upper right*—ambient temperature operation and immersed; *lower left*—thermoelectrically cooled; *lower right*—thermoelectrically cooled and immersed. (Vigo Systems.)

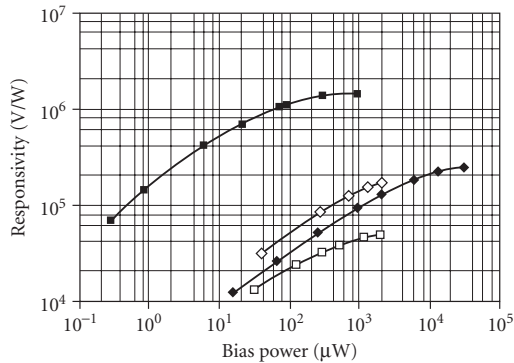


FIGURE 175 Range of peak responsivities for 12- μm cutoff HgCdTe photoconductors at 80 K. These devices have nominal dimensions of $50 \times 50 \mu\text{m}$, and resistance of 50 to $150 \Omega/\text{square}$. (Santa Barbara Research Center.)

Circuit: Standard photoconductive.

Manufacturers: Belov Technology, Boston Electronics, Hamamatsu, Infrared Associates, Kolmar Technologies, Oriel, Teledyne Judson Technologies, Vigo Systems.

Photovoltaic HgCdTe

Sensitivity: Adjustable by varying alloy composition (see Figs. 168, 170, 177, and 178). Also compare Fig. 170 with Fig. 19 for an estimate of the extent to which D^* may increase (up to the R_0A or shunt resistance limit) as the background flux is reduced.

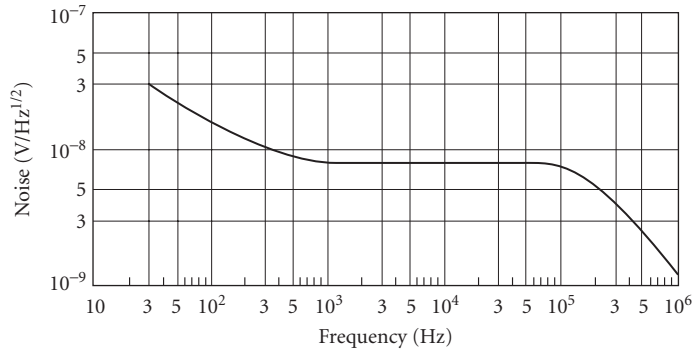


FIGURE 176 Variation of noise with frequency for photoconductive HgCdTe at 77 K. (GEC Marconi Infrared Ltd.)

Time constant: Depends on diode capacitance (area); 10 to 20 ns without bias; 0.5 to 3 ns with reverse bias (some trade-off of sensitivity). Low-capacitance *pin* devices with response out to several gigahertz (0.05 to 0.2-ns time constant) are also available for 10.6- μm CO₂ laser heterodyne detection.

Resistance: Refer to Fig. 170 for the R_0A product at zero bias corresponding to the detector cutoff wavelength (this figure shows very high quality diode impedances) and divide by the diode area. Large-area diodes will have somewhat lower R_0A product than shown in this figure. R_0A varies somewhat with background flux as can be noted from Fig. 170.

Operating temperature: Depends on spectral cutoff; 77 K and lower for LWIR and VLWIR detectors, up to room temperature for SWIR devices. Optical immersion and/or TE cooling will boost the performance for all spectral ranges compared with operation at ambient temperature—see Fig. 178.

Noise: High-quality devices may have flat noise response from 1 Hz out to the high-frequency limit of the time constant. $1/f$ noise may be present in lower quality devices and will increase with reverse bias.

Quantum efficiency: >50 percent (60–75 percent typical) without antireflection coating. Higher with antireflection coating.

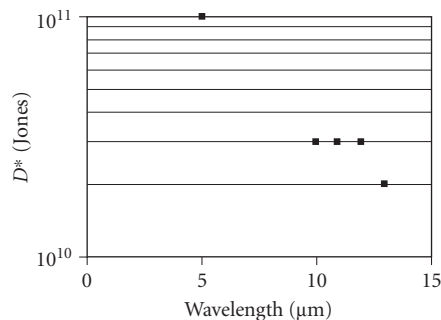


FIGURE 177 D^* specifications for small (50×50 to $250 \times 250 \mu\text{m}$) HgCdTe photodiodes at 77 K as a function of spectral cutoff. Data is shown for devices with 60° FOV background flux. (Fermionics, Mercury Cadmium Telluride MWIR and LWIR Detector Series.)

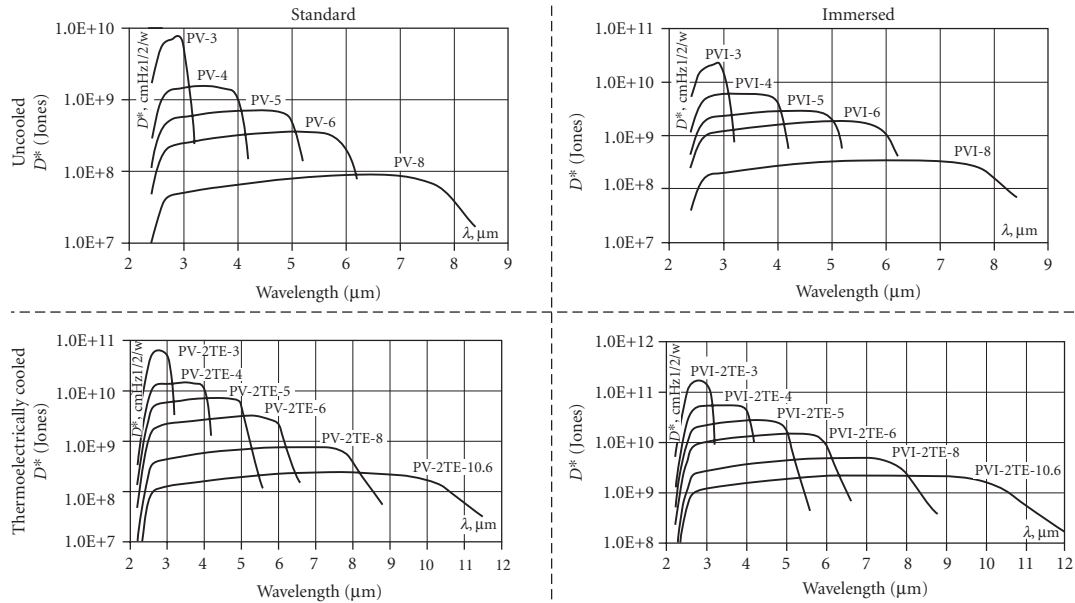


FIGURE 178 D^* for photovoltaic HgCdTe detectors as a function of wavelength: *upper left*—ambient temperature operation; *upper right*—ambient temperature operation and immersed; *lower left*—thermoelectrically cooled; *lower right*—thermoelectrically cooled and immersed. (Vigo Systems.)

Sensitive area: 0.0 to 0.25-mm square, 0.5- and 1-mm diameter.

Capacitance: Depends on junction doping, area, and applied bias (very slightly dependent on spectral cutoff). For standard *pn* junction devices at zero bias and 10^{15} cm^{-3} doping, capacitance is approximately $3 \times 10^4 \text{ pF/cm}$. Significantly lower for *pin* junction devices.

Circuit: Standard photovoltaic circuits, reverse-bias operation to enhance speed and zero bias to maximize D^* .

Manufacturers: Boston Electronics, Kolmar Technologies, Oriel, Raytheon Vision Systems, Vigo Systems.

PbSnTe PbSnTe offers an alternative semiconductor alloy system based upon the IV-VI chemical groups to the II-VI (HgCdMnTeSe) groups previously described for fabricating variable spectral cut-off detectors. Only photovoltaic detectors are available in PbSnTe. This technology has an advantage in the ease of material growth and in the fabrication of good quality photodiode junctions. It has a disadvantage in the very high dielectric constant of the material, combined with relatively high doping concentrations giving high-capacitance (comparatively slow) detectors. For low-frequency applications this is not a disadvantage.

Sensitivity: $D^*_{\text{peak}} > 10^{10}$ Jones (see Fig. 179).

Time constant: $> 50 \text{ ns}$.

Sensitive area: $1 \times 1 \text{ mm}$.

Operating temperature: 77 K.

Circuit: Standard photovoltaic.

Manufacturers: No suppliers are presently known.

Ge:Hg Mercury-doped germanium detectors are fast single-crystal impurity-doped photoconductors, sensitive out to 14 m. Ge:Hg is especially well suited for detection through the 8- to 13- μm

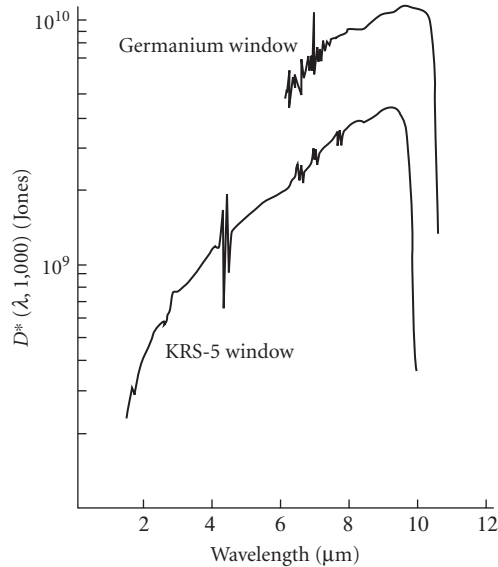


FIGURE 179 Photovoltaic PbSnTe D^* versus wavelength at 77 K and 60° FOV. (Barnes Engineering Co.)

atmospheric window and for detection of near-ambient sources. Unfortunately, its operating temperature must be kept less than 40 K, where it becomes 300-K background limited.

Sensitivity: See Figs. 161, 162, 180, and 181.

Quantum efficiency: 25 to 30 percent.

Noise: See Figs. 182 and 183.

Time constant: 100 ns with 50- Ω load for $T < 28$ K and electric fields < 30 V/cm. (Compensated material is available with ~ 5 -ns time constant with a 50- Ω load. Responsivity then is reduced by 5–10 \times and detectivity is reduced by 2.)

Responsivity: Depends on concentration of compensating impurities, bias, area, and background flux. See Figs. 183 to 185, $\sim 10^5$ V/W.

Dark resistance: Depends on area and FOV: ~ 100 k Ω for 180° FOV (see Fig. 186).

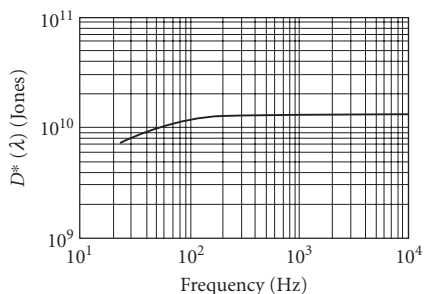


FIGURE 180 Typical D^* versus frequency at 30 K for Ge:Hg; 1×1 -mm area; essentially constant with temperature. (Santa Barbara Research Center.)

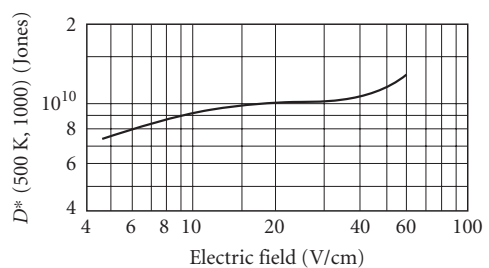


FIGURE 181 D^* versus electric field for Ge:Hg; $T = 5$ K, 90° FOV; 300-K background; 6×10^{-4} cm $^{-2}$ area; Irtran II window. (Reprinted by permission of Texas Instruments.)

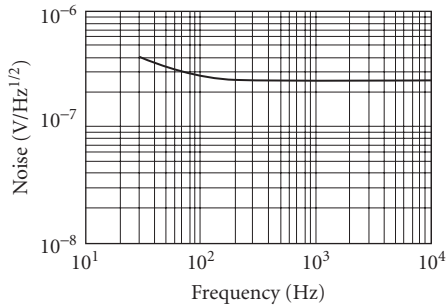


FIGURE 182 Typical noise versus frequency spectrum for Ge:Hg; 1×1 -mm area; essentially constant with temperature. (Santa Barbara Research Center.)

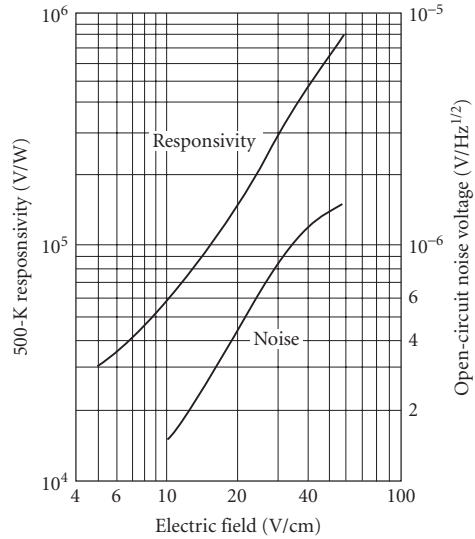


FIGURE 183 Open-circuit responsivity and noise voltage versus electric field for Ge:Hg; $T = 5$ K, 90° FOV; 300-K background; 6×10^{-4} cm² area; Irtran II window. (Reprinted by permission of Texas Instruments.)

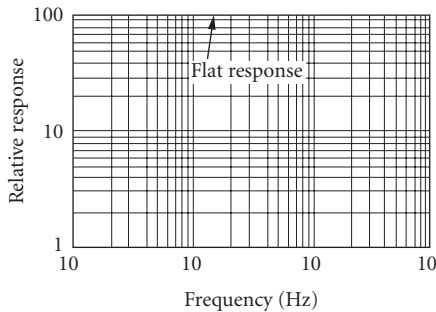


FIGURE 184 Ge:Hg typical relative response versus frequency; $T = 30$ K. (Santa Barbara Research Center.)

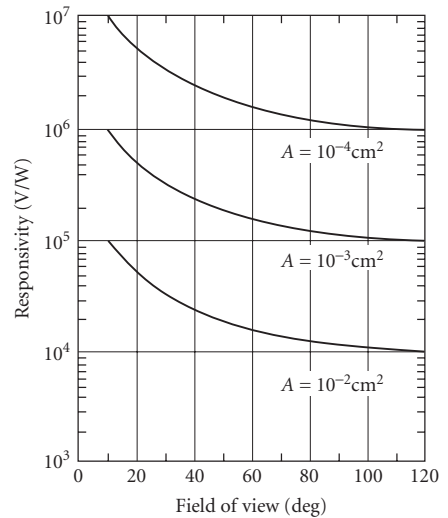


FIGURE 185 Open-circuit responsivity versus FOV for Ge:Hg at 5 K for various detector areas; 300-K background; 500-K blackbody source. (Reprinted by permission of Texas Instruments.)

Capacitance: < 1 pF.

Sensitive area: 1 to 5-mm diameter.

Operating temperature: See Fig. 162.

Linearity: 10^{-3} to 10^{-8} W (size dependent).

Sensitivity profile: ± 15 percent.

Recommended circuit: Standard photoconductive (see Ge: Au). See Fig. 187 for current-voltage characteristics.

Manufacturer: No suppliers are presently known.

Si:Ga Gallium in silicon forms an acceptor level with a binding energy of ~ 72 meV which is the basis of an infrared detector with spectral response out to approximately $17 \mu\text{m}$, as shown in Fig. 188. The exact spectral cutoff and quantum efficiency will vary slightly with the gallium doping concentration. Gallium-doped silicon requires cooling to 20 K or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 15 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in micrometers, η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 0.9 A/W.

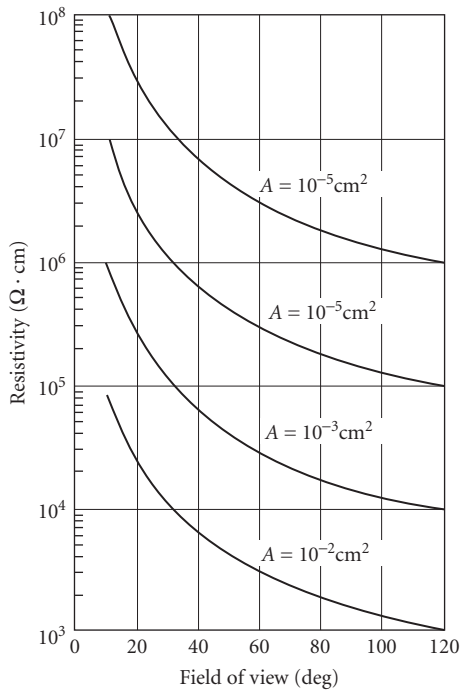


FIGURE 186 Open-circuit resistivity versus FOV for Ge:Hg at 5 K for various detector areas; 300-K background; 500-K blackbody source. (Reprinted by permission of Texas Instruments.)

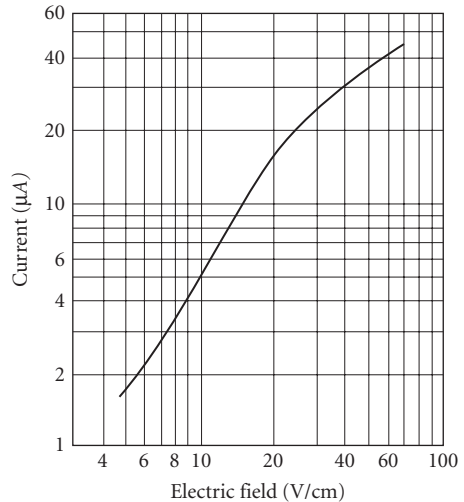


FIGURE 187 Ge:Hg bias current versus electric field; 90° FOV; $T = 5$ K, 300-K background; $6 \times 10^{-4} \text{ cm}^{-2}$ area; Irtran II window. (Reprinted by permission of Texas Instruments.)

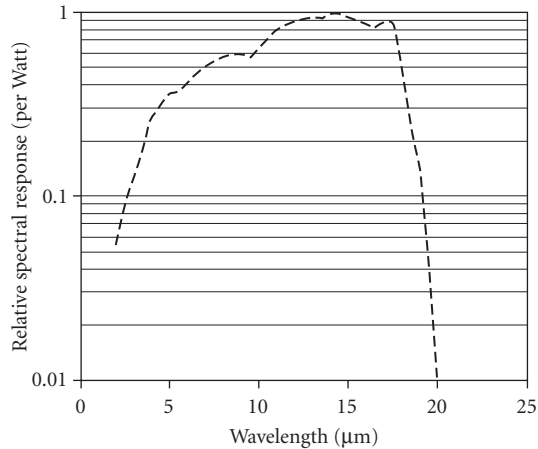


FIGURE 188 Relative spectral response per watt of Si:Ga as a function of wavelength. Data is normalized to unity at the peak response.

Time constant: <1 μs.

Resistance: Depends on background flux and detector bias (similar to Ge:Hg, see Fig. 186).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.2 to 2-mm square.

Operating temperature: <20 K.

Recommended circuit: Standard photoconductive.

Manufacturer: Infrared Laboratories.

Si:B Boron in silicon forms an acceptor level with a binding energy of ~45 meV which is the basis of this infrared detector with spectral response out to approximately 30 μm. The exact spectral cut-off and quantum efficiency will vary slightly with the boron-doping concentration. Boron-doped silicon requires cooling to about 15 K or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 10 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in micrometers, η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 2 A/W.

Time constant: <1 μs.

Resistance: Depends on background flux and detector bias (similar to Ge:Hg, see Fig. 186).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.2 to 2-mm square.

Operating temperature: <15 K.

Recommended circuit: Standard photoconductive (see Ge:Au).

Manufacturer: Infrared Laboratories.

Ge:Cu Copper-doped germanium detectors are fast, single-crystal, impurity-doped photoconductors, with high sensitivity in the broad region 2 to 30 μm. Operating temperature must be maintained below 20 K (ideally <14 K). Ge:Cu is then 300-K background-limited, and response time <50 ns.

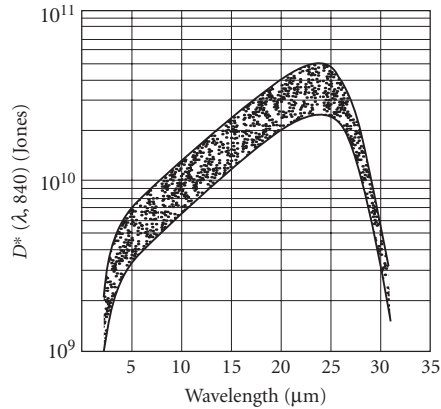


FIGURE 189 Range of spectral detectivities for Ge:Cu; 60° FOV. (Santa Barbara Research Center.)

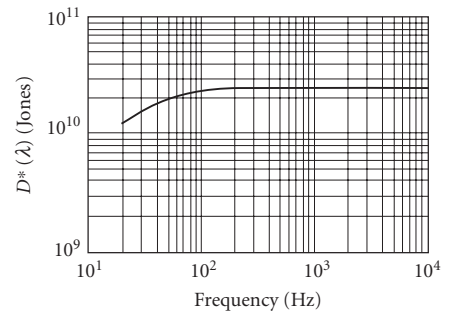


FIGURE 190 Typical D^* vs. frequency for Ge:Cu at 4.2 K. (Santa Barbara Research Center.)

Sensitivity: See Figs 161, 162, 189, and 190.

Noise: See Figs. 191 and 192.

Time constant: ~100 ns (can be doped to be faster, ~5 ns). (See discussion of time constant under Ge:Cu).

Responsivity: 10^5 V/W (see Figs. 192 to 194).

Dark resistance: Depends on FOV (~100 k Ω for 180° FOV).

Capacitance: <1 pF.

Sensitive area: 1 to 5-mm diameter.

Operating temperature: See Figs. 162 and 194.

Linearity: 10^{-3} – 10^{-8} W/cm² (depends on size).

Sensitivity profile: ± 15 percent.

Stability: Stable in all ambient storage environments tested.

Recommended circuit: See Figs. 166 and 195.

Manufacturer: No suppliers are presently known.

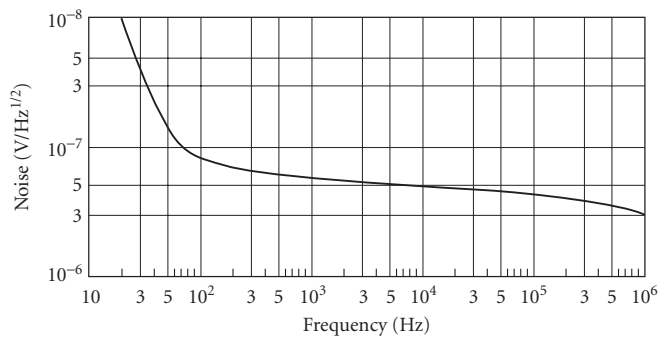


FIGURE 191 Typical noise frequency spectrum for Ge:Cu. (GEC Marconi Infra Red Ltd.)

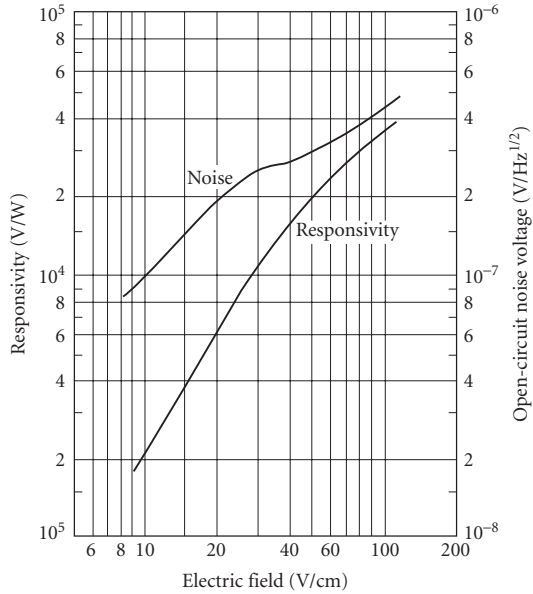


FIGURE 192 Typical noise and responsivity versus biasing field for Ge:Cu; 5 K, 60° FOV; 300-K background, $A = 10^{-2} \text{ cm}^2$, 500-K blackbody. (Reprinted by permission of Texas Instruments.)

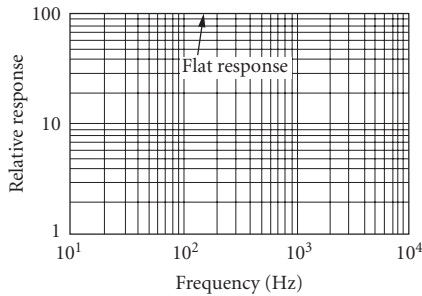


FIGURE 193 Relative response versus frequency for Ge:Cu at 4.2 K. (Santa Barbara Research Center.)

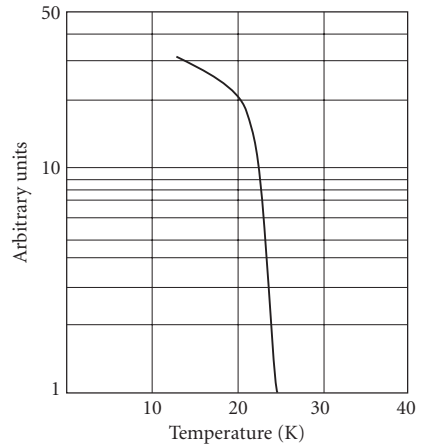


FIGURE 194 Relative responsivity versus temperature for Ge:Cu.

Ge:Zn Very similar to Ge:Cu except that cutoff wavelength moves out to 42 μm and operating temperature should be $< 10 \text{ K}$. A relatively low field breakdown limits the responsivity.

Sensitivity: See Figs 161, 162, and 196.

Noise: See Fig. 197 and 198.

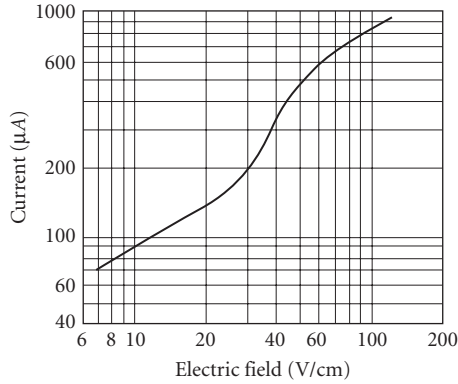


FIGURE 195 Typical current-voltage curve for Ge:Cu, 60° FOV; 300-K background; $A = 10^{-2} \text{ cm}^{-2}$. (Reprinted by permission of Texas Instruments.)

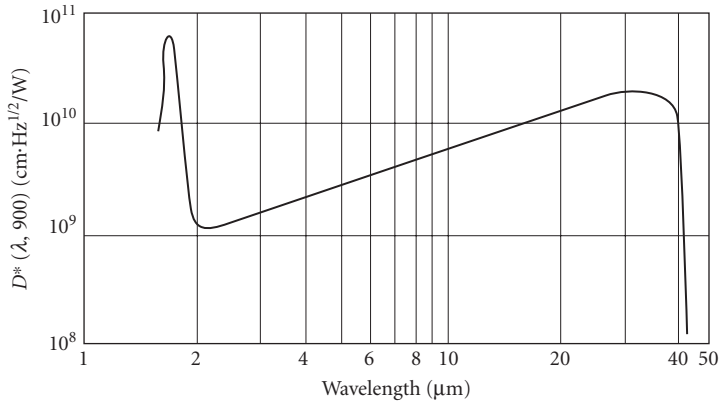


FIGURE 196 D^* versus wavelength for Ge:Zn.

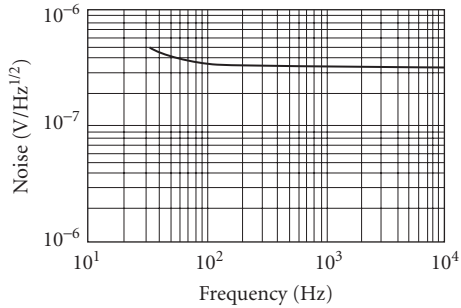


FIGURE 197 Typical noise-frequency spectrum for Ge:Zn at 4.2 K; $A = 1 \times 1 \text{ mm}$. (Santa Barbara Research Center.)

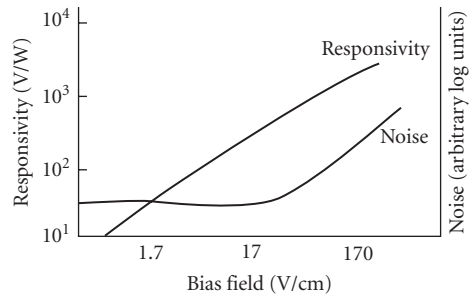


FIGURE 198 Signal and noise for Ge:Zn.

Responsivity: 10^3 V/W (see Fig. 198).

Time constant: <50 ns. (See discussion of time constant under Ge:Au.)

Dark resistance: 0.5 to 5 M Ω / sq (60°-FOV ambient background).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 1-, 2-, 3-, and 5-mm diameters.

Operating temperature: <10 K.

Recommended circuit: Standard photoconductive.

Manufacturer: No suppliers are presently known.

Ge:Ga The elements of B, Al, Ga, In, and Tl from chemical group III form shallow acceptor states (~10 meV) in germanium which are the basis of infrared detectors with spectral response out to approximately 120 μ m. Currently, gallium-doped germanium is commercially available, but germanium doped with other group III elements (Ge:B, Ge:Al, Ge:In, Ge:Tl) will give similar detector performance. The small binding energies associated with this detector require cooling to liquid helium temperatures (4.2 K) or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 7 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in μ m, η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 4 A/W.

Time constant: <1 μ s.

Resistance: Depends on background flux and detector bias.

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.5-, 1-, and 2-mm square.

Operating temperature: <4.2 K, best below 3 K.

Recommended circuit: Standard photoconductive.

Manufacturer: Infrared Laboratories.

Photographic In this paragraph we present only the spectral sensitivity of some typical photographic emulsions. See Chap. 29, "Photographic Films," for a more extensive coverage.

The term spectral sensitivity generally has a different meaning when applied to photographic detectors than it does when applied to the other detectors described in this chapter. It comes closer to responsivity than to minimum detectable power or energy.* In Fig. 199 sensitivity is the reciprocal of exposure, expressed in ergs per centimeter, required to produce

$$\text{Density} = \log \frac{1}{\text{transmittance}} = 0.3$$

above gross fog in the emulsion when processed as recommended.

Manufacturers: AGFA, Eastman Kodak, Fuji, Polaroid.

*Work is in progress to evaluate photographic materials in terms of minimum detectable energy, a concept involving the average number of photons necessary to produce a change in density (signal) equal to that of the fog-density fluctuations (noise); see Refs. 10, 26, 27, and 28.

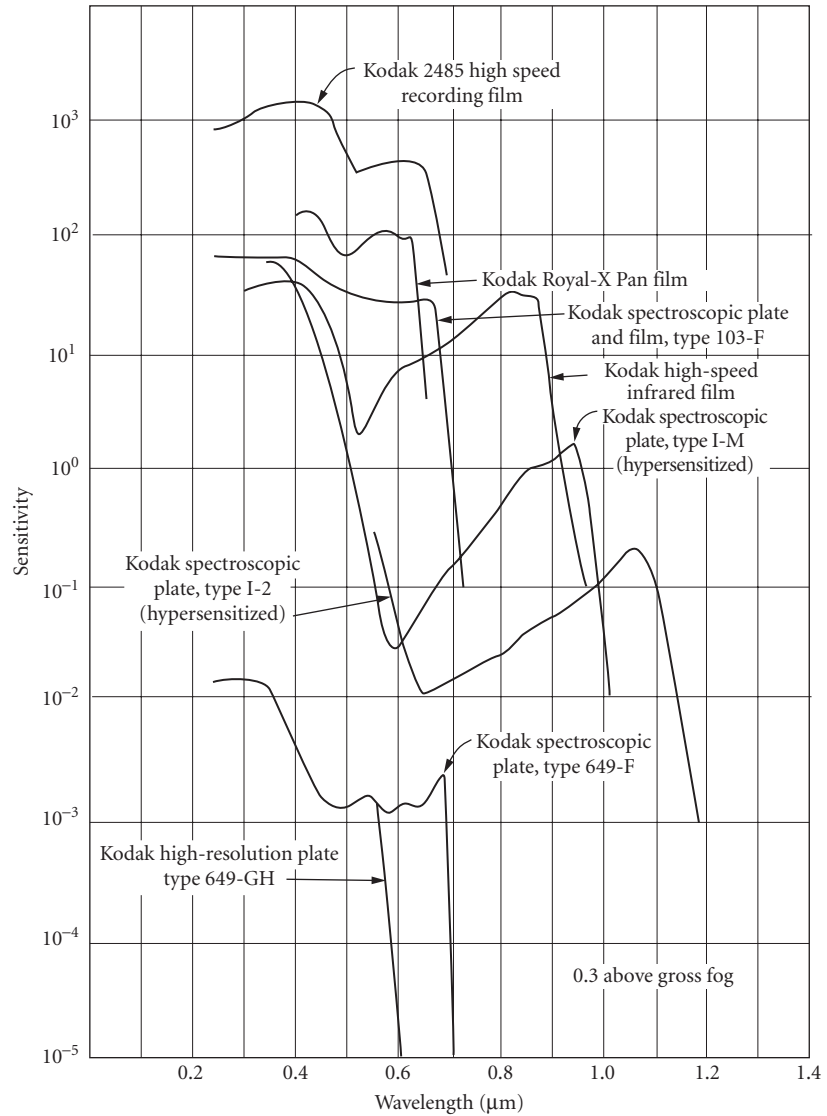


FIGURE 199 Sensitivity versus λ for typical photographic emulsions. (Eastman Kodak.)

24.8 REFERENCES

1. W. L. Wolfe, (ed.), *Handbook of Military and Infrared Technology*, Office of Naval Research, Washington, D.C., 1965.
2. Institute of Radio Engineers, "IRE Standards on Electron Tubes: Methods of Testing," *Proceedings of the IRE* 50(9): 1974-1975 (1962).
3. Radio Corporation of America, "Phototubes and Photocells," *Tech. Man. PT-60* (1963).

4. D. Vincent, *Fundamentals of Infrared Detector Operation and Testing*, Wiley, New York, 1989.
5. E. L. Dereniak and D. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984.
6. W. Rogatto, ed., *The Infrared and Electro-Optical Systems Handbook*, vol. 3, SPIE Press, Bellingham, Wash., 1993.
7. R. C. Jones, "Factors of Merit for Radiation Detectors," *J. Opt. Soc. Am.* **39**:344 (1949).
8. G. Bauer, "Ein halbleiter Hochohmbolometer mit Tafel, IV," *Phys. Z.* **44**:53 (1943).
9. J. A. R. Samson, *Techniques of Ultraviolet Spectroscopy*, Wiley, New York, 1967.
10. R. Jones, "On the Quantum Efficiency of Photographic Negatives: On the Minimum Energy Detectable by Photographic Materials," *Photogr. Sci. Eng.* **2**:57–65j, 191–204 (1958).
11. D. J. Fink, *Principles of Television Engineering*, McGraw-Hill, New York, 1940.
12. R. A. Smith, F. E. Jones, and R. P. Chasmar, *Detection and Measurement of Infrared Radiation*, Oxford University, London 1957.
13. S. F. Jacobs, and M. Sargent III, "Photon Noise Limited D^* for Low Temperature Backgrounds and Long Wavelengths," *Infrared Phys.* **10**(4):233–235 (1970).
14. P. W. Kruse, L. D. McLaughlin, and R. B. McQuistan, *Elements of Infrared Technology*, Wiley, New York, 1962.
15. P. B. Fellgett, "On the Ultimate Sensitivity and Practical Performance of Radiation Detectors," *J. Opt. Soc. Am.* **39**:970 (1949).
16. F. J. Low, and A. R. Hoffman, "The Detectivity of Cryogenic Bolometers," *Appl. Opt.* **2**:649 (1963).
17. F. J. Low, "Low-Temperature Germanium Bolometer," *J. Opt. Soc. Am.* **51**:1300 (1961).
18. E. H. Eberhardt, "Threshold Sensitivity and Noise Ratings of Multiplier Phototubes," *Appl. Opt.* **6**:251 (1967).
19. W. C. Livingston, "Enhancement of a Photocathode Sensitivity by Total Internal Rejection as Applied to an Image Tube," *Appl. Opt.* **5**:1335 (1966).
20. W. D. Gunter, G. R. Grant, and S. A. Shaw, "Optical Devices to Increase Photocathode Quantum Efficiency," *Appl. Opt.* **9**:251 (1970).
21. M. Cole and D. Ryer, "Cooling PM Tubes for Best Spectral Response," *Electro Optical Systems Design* **4**(6):16–19 (June 1972).
22. H. E. Bennett, "Accurate Method for Determining Photometric Linearity," *Appl. Opt.* **5**:1265 (1966).
23. R. E. Simon, A. H. Sommer, J. J. Tietjen, and B. F. Williams, "New High-Gain Dynode for Photomultipliers," *Appl. Phys. Lett.* **13**:355 (1968).
24. G. A. Morton, H. M. Smith, and H. R. Krall, "Pulse Height Resolution of High Gain First Dynode Photomultipliers," *Appl. Phys. Lett.* **13**:356 (1968).
25. S. M. Sze, *Physics of Semiconductor Devices*, 2d ed., Wiley, New York, p. 773, 1981.
26. T. H. Johnson, "Lead Salt Detectors and Arrays: PbS and PbSe," *Proc. SPIE* **443**: 60–94, (1984).
27. J. C. Marchant, "Exposure Criteria for the Photographic Detection of Threshold Signals," *J. Opt. Soc. Am.* **54**:79 (1964).
28. G. R. Bird, R. C. Jones, and A. E. Ames, "The Efficiency of Radiation Detection by Photographic Films: State-of-the-Art and Methods of Improvement," *Appl. Opt.* **8**:2389 (1969).

24.9 SUGGESTED READINGS

- Sommer, A. H., *Photoemissive Materials*, Wiley, New York, 1968.
- Sommers, H. S., Jr. and E. K. Gritchell, "Demodulation of Low-Level Broadband Optical Signals with Semiconductors," *Proc. IEEE* **54**:1553 (1966).
- Sommer, A. H., and W. B. Teusch, "Demodulation of Low-Level Broadband Optical Signals with Semiconductors, II: Analysis of the Photoconductive Detector," *Proc. IEEE* **52**:144 (1964).
- Sun, C. and T. E. Walsh, "Performance of Broadband Microwave-Biased Extrinsic Photoconductive Detectors at 10.6 μ m," *IEEE J. Quantum Electron.* **6**:450 (1970).

Abhay M. Joshi

*Discovery Semiconductors, Inc.
Cranbury, New Jersey*

Gregory H. Olsen

*Sensors Unlimited, Inc.
Princeton, New Jersey*

25.1 GLOSSARY

A	photodetector active area
A_0	incident photon flux
B	bandwidth of the photodetector
C	capacitance of the photodetector
D^*	detectivity
E	applied electric field
E_a	activation energy
E_g	bandgap of the semiconductor
E_i	impurity energy state
f	frequency
I_{diff}	diffusion current
I_{g-r}	generation-recombination current
IR_0	unity gain current
IR_1	reverse current generated by avalanche action
I_{tun}	tunneling current
k	Boltzmann's constant
L	distance traveled by a charge carrier
M	photocurrent gain
m	effective mass of a electron
N_A	acceptor impurity concentration on p side
N_D	donor impurity concentration on n side
n	refractive index of the AR coating
q	electron charge
R	sum of the detector series resistance and load resistance

R_o	detector shunt impedance
T	temperature in kelvin
t_n	transit time of electrons
t_p	transit time of holes
t_r	transit time of charge carriers (holes or electrons)
V	applied reverse bias in volts
V_B	breakdown voltage
V_{bi}	built-in potential of a p - n junction
W	depletion width of a p - n junction
α	absorption coefficient of the photodetector's absorption layer
ϵ_s	semiconductor permittivity
η	quantum efficiency of photodetector
θ	tunneling constant
λ	wavelength of incident photons (nm)
λ_{co}	detector cutoff wavelength (10 percent of peak response, nm)
μ	mobility of charge carriers (holes or electrons)
μ_n	mobility of electrons
μ_p	mobility of holes

25.2 INTRODUCTION

The approach of this chapter is descriptive and tutorial rather than encyclopedic. It is assumed that the reader is primarily interested in an overview of how things work. Among the many excellent references to be consulted for further details are Sze's book,¹ and the article by Forrest.² For the latest in photodetector developments, consult recent proceedings of the Society of Photo-optical and Instrumentation Engineers (SPIE) conference or the IEEE Optical Fiber Conference.

A photodetector is a solid-state sensor that converts light energy into electrical energy. According to Isaac Newton, light energy consists of small packets or bundles of particles called *photons*. Albert Einstein, who won a Nobel prize for the discovery of the photoelectric effect, showed that when these photons strike a metal they can excite electrons in it. The minimum photoenergy required to generate (excite) an electron is defined as the *work function* and the number of electrons generated is proportional to the intensity of the light. The semiconductor photodetectors are made from different semiconductor materials such as silicon, germanium, indium gallium arsenide, indium antimonide, and mercury cadmium telluride, to name a few. Each material has a characteristic bandgap energy E_g which determines its light-absorbing capabilities. Light is a form of electromagnetic radiation comprised of different wavelengths (λ). The range of light spectrum is split approximately as ultraviolet (0–400 nm), visible (400–1000 nm), near infrared (1000–3000 nm), medium infrared (3000–6000 nm), far infrared (6000–40,000 nm), extreme infrared (40,000–100,000 nm). The equation between bandgap energy E_g and cutoff wavelength (λ_c) is

$$\lambda_c = \frac{1.24 \times 10^3 \text{ nm}}{E_g \text{ (eV)}} \quad (1)$$

The smaller the bandgap (eV), the farther the photodetector “sees” into the infrared. Table 1 lists some prominent photodetector materials, their bandgaps, and cutoff wavelengths λ_c at room temperature (300 K).

Photodetectors find various applications in fiber-optic communications (800–1600 nm), spectroscopy (400–6000 nm), laser range finding (400–10,600 nm), photon counting (400–1800 nm),

TABLE 1 Important Photodetector Materials

Type	E_g (eV)	λ_c (nm)	Band
Silicon	1.12	1100	Visible
Gallium arsenide	1.42	875	Visible
Germanium	0.66	1800	Near-infrared
Indium gallium arsenide*	0.73–0.47	1700–2600	Near-infrared
Indium arsenide	0.36	3400	Near-infrared
Indium antimonide	0.17	5700	Medium-infrared
Mercury cadmium	0.7–0.1	1700–12500	Near-to-far-infrared

*The alloy composition of indium gallium arsenide and mercury cadmium telluride can be changed to alter the bandgap E_g .

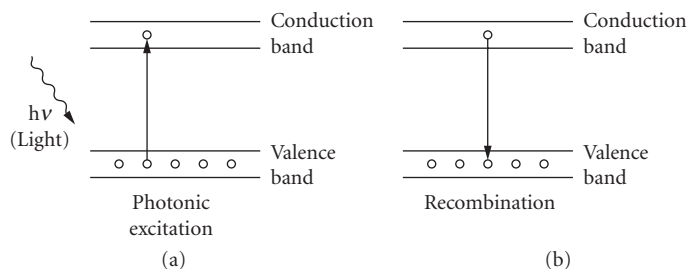
and satellite imaging (200–1200 nm), to name only a few topics. We will discuss three kinds of photodetectors: (1) Photoconductors, (2) *p-i-n* photodetectors (including avalanche diodes), and (3) photogates. Frequently, such detectors need to have high sensitivity, low noise, and high reliability. For fiber-optic applications, the frequency response and the cost can be a critical issue, whereas for infrared applications, many times it is the area of the photodetector. Large area (1 in diameter), high-sensitivity silicon avalanche photodetectors compete with conventional photomultiplier tubes for low light sensing applications in the visible. They offer the advantage of compact size and more rugged construction. We also discuss reliability issues concerning photodetectors and take notice of a few novel photodetector structures.

25.3 PRINCIPLE OF OPERATION

When an electron in the valence band receives external energy in the form of light, the electron may overcome the nuclear attraction and become a “free electron.” When light energy creates this transformation, it is termed *photonic excitation* (see Fig. 1a). The range of energies acquired by these free electrons is termed the *conduction band*. The energy difference between the bottom of this conduction band and the top of the valence band is termed the *energy bandgap* E_g and represents the minimum energy of light that the material can absorb.

However, under the influence of even a small external electric field, the free electrons can “drift” in a specific direction. This is the fundamental principle of a photodetector. Figure 2 shows the three kinds of photodetectors discussed in this chapter. A brief explanation of each kind follows.

All photodetectors can be characterized by their quantum efficiency, detectivity (sensitivity), and response time.¹ Quantum efficiency (QE) is perhaps the most fundamental property, as it determines just how efficiently the device converts incoming photons into conduction electrons.

**FIGURE 1** Photonic excitation and recombination in a semiconductor.

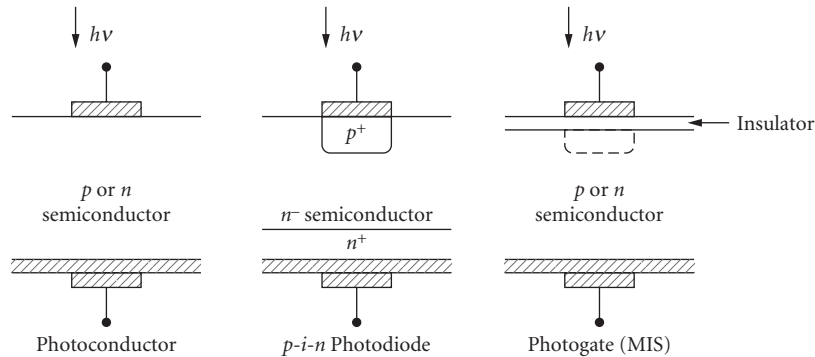


FIGURE 2 Types of photodetectors.

Usually expressed in percentage, quantum efficiency can range from under 1 percent for PtSi Schottky barrier infrared detectors to well over 90 percent for InGaAs *p-i-n* fiber-optic photodetectors. Responsivity (R) is a related term expressed in amps/watt, which determines how much photocurrent is produced by optical power of a given wavelength. Detectivity measures how *sensitive* a detector is; that is, not only its light conversion efficiency, but also its ability to detect low-level light signals. It is limited by the various noise currents (shot, $1/f$, etc.) introduced by the detector. Finally, response time describes how rapidly a detector can respond to a changing light signal. This ranges from milliseconds for certain types of PbS photoconductors to picoseconds for GaAs-like metal-insulator-semiconductor detectors.

These three parameters are frequently traded off. A large-area detector captures more light signal and thus might have greater detectivity. However, its larger capacitance would slow down the device. Similarly, response time in *p-i-n* detectors can be improved by thinning the absorbing region of the detector. However, this in turn cuts down quantum efficiency by reducing the total number of photons absorbed.

Photoconductor

A photoconductor, as the name implies, is a device whose conductivity increases with illumination. It acts as an open switch under dark (or no illumination) and as a closed switch under illumination. The simple equivalent electric circuit is shown in Fig. 3. This basic principle of a photoconductor finds numerous applications in relays and control circuits. An ideal switch should have low resistance in the closed position and, therefore, a pair of ohmic (not-rectifying) contacts are formed to the photoconductor. These ohmic contacts form the electrodes and usually have contact resistance of less than 10 ohms.

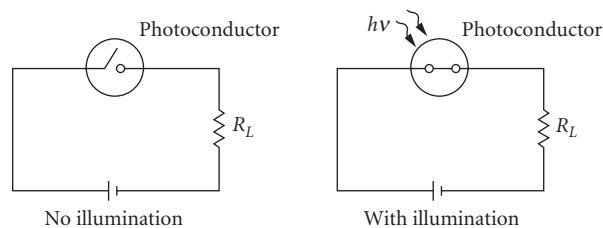


FIGURE 3 Equivalent circuit diagram of a photoconductor.

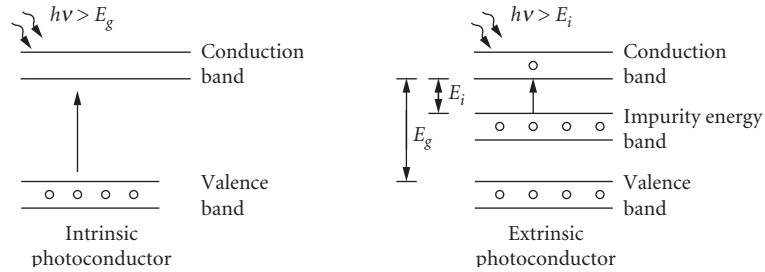


FIGURE 4 Functional diagram of an intrinsic and extrinsic photoconductor.

Types of Photoconductors The two principle photoconductors are (1) intrinsic and (2) extrinsic. In an intrinsic device there is no external impurity atom. However, when an external impurity (dopant) is added to a material, it is termed *extrinsic*. This impurity atom occupies an energy state between the valence band and the conduction band. The functional difference between the intrinsic and extrinsic photoconductor is seen in Fig. 4. For an intrinsic device, the photo excitation ($h\nu$) needs to have energy greater than the bandgap energy E_g and its cutoff wavelength λ_c is given by Eq. (1). But, for an extrinsic one, the photon excitation ($h\nu$) should exceed the impurity energy state E_i and its cutoff wavelength λ_{co} is

$$\lambda_{co} = \frac{1.24 \times 10^3 \text{ nm}}{E_i (\text{eV})} \quad (2)$$

For intrinsic photoconductors, it is extremely difficult to achieve bandgap energies E_g less than 0.1 eV (refer to Table 1). This limits its capability to see in the far infrared and extreme infrared (13,000–100,000 nm) and beyond. This is overcome by extrinsic devices whose E_i value is less than 0.1 eV and is normally done by doping germanium or silicon. However, the extrinsic photoconductor suffers from very low absorption coefficients and, hence, poor quantum efficiencies. Also, since ambient thermal energy can excite carriers, they have to be cooled to liquid nitrogen temperature (77 K) and below, whereas most intrinsic photoconductors can operate at room temperature (300 K).

Photo Gain The sensitivity of a photoconductor is determined by its gain. Photo gain is defined as the ratio of the output signal to the input optical signal. When photons impinge on a photoconductor, they generate electron-hole pairs and, under the influence of external fields, they are attracted toward the anode and cathode. A typical photoconductor is illustrated in Fig. 5 with L being the thickness of the active layer. The transit time (t_r) required for a charge carrier to travel a distance L is given by

$$t_r = \frac{L^2}{\mu V} \quad (3)$$

where V = applied voltage bias and μ = charge mobility.

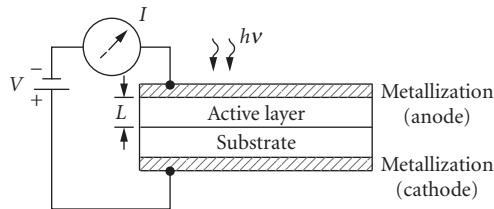


FIGURE 5 A typical photoconductor.

The mobility of electrons μ_n and that of holes μ_p is different, with μ_n being usually far higher than μ_p . This causes a difference in the transit time of electrons and holes. Hence, photon-generated electrons are swept away more quickly than holes which result in a positive charge in the active layer. To maintain the charge neutrality, new electrons are supplied by the external voltage source. Therefore, for one incident photon, more than one electron is circulated in the electric circuit. This results in an “effective gain.” Thus, photoconductor gain can be defined as the ratio of slower transit time t_p to faster transit time t_n .

$$M = \frac{t_p}{t_n} \quad \text{or} \quad \frac{\mu_n}{\mu_p} \quad (4)$$

The slower the transit time, the higher the gain; however, the bandwidth of the device is reduced. Hence, high-gain photoconductors will result in slow devices and vice versa. Such “high-gain, slow devices” can be best utilized for imaging applications.³ For high-speed optical communication applications in the 1000- to 1700-nm spectrum, InGaAs is the material of choice due to its high mobility. Several reports have been published on high-speed InGaAs photoconductors that find practical applications in optical receivers.⁴⁻⁹ (Also see Chap. 26, “High-Speed Photodetectors,” by John Bowers and Yih G. Wey.)

***p-i-n* Photodiode**

Unlike photoconductors, a photodiode has a *p-n* junction, usually formed by diffusion or epitaxy. In a photoconductor, metal contacts are made to either *n*- or *p*-type material. However, a photodiode consists of both *n*- and *p*-type materials across which a natural electric field is generated. This field is known as the *built-in potential* V_{bi} , and its value depends on the bandgap of its material. A silicon *p-n* junction has V_{bi} of 0.7 V whereas in germanium it is 0.3 V. The higher the bandgap E_g , the larger the built-in potential V_{bi} . An important physical phenomenon called *depletion* occurs when a *p*-type semiconductor is merged with an *n*-type semiconductor. After an initial exchange of charge, a potential is built up to prevent further flow of charge. This built-in field creates the depletion width W , which is a region free of any charge carriers and is given by¹⁰

$$W = \sqrt{\frac{2\epsilon_s(N_A + N_D)}{q(N_A N_D)}(V_{bi} - V)} \quad (5)$$

where N_A and N_D are impurity concentrations of *p* side and *n* side, respectively, q is the electron charge, V_{bi} is the built-in potential, and ϵ_s is the semiconductor permittivity, V is the applied bias and is negative for reverse-bias operation. As seen from Eq. (5), under reverse bias the depletion width W increases causing a decrease in the capacitance of the photodiode. A *p-i-n* photodiode is similar to a parallel plate capacitor with the anode-cathode being the two plates and the depletion width W being the separating medium. A typical InGaAs *p-i-n* photodiode is shown in Fig. 6 and its capacitance is given by¹

$$C = \frac{\epsilon_s A}{W} \quad (6)$$

where ϵ_s is the semiconductor permittivity and A is the active area of the photodiode. From first principles, a decrease in capacitance improves the bandwidth B of the photodetector according to

$$B = \frac{0.35}{2.2RC} \quad (7)$$

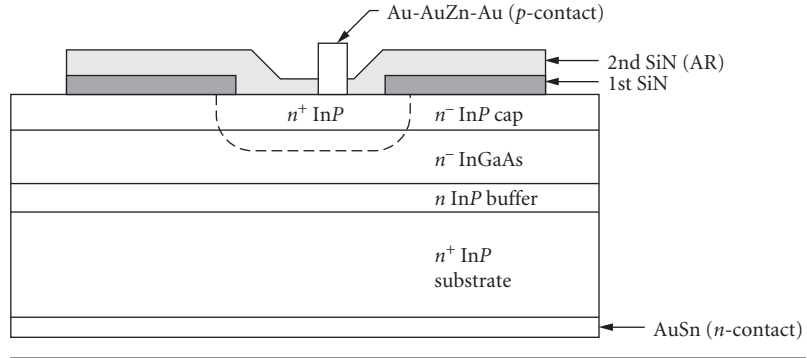


FIGURE 6 A typical InGaAs p - i - n photodiode.

where, R is the sum of the detector series resistance and load resistance. For a detailed analysis on high-speed photodetectors, see Chap. 26 by Bowers and Wey.

Dark Current For applications ranging from optical communications (III-V compound semiconductors) to infrared sensing (Si, Ge, III-IV, and II-IV compound semiconductors), a p - i - n photodiode must have high sensitivity and low noise. These are largely determined by the dark currents originating in the device. Several authors have published papers on dark currents in InGaAs^{11–13} and HgCdTe.^{14,15} The three major components of dark current are (1) diffusion current, (2) generation-recombination current, and (3) tunneling current.

Diffusion current In the nondepleted region of the photodiode, electron-hole pairs are formed by the ambient temperature. These thermally generated carriers diffuse toward the depletion region and produce the diffusion current.

$$I_{\text{diff}} \propto e^{-E_g/kT} \quad (8)$$

where E_g is the bandgap of the photodiode material, k is Boltzmann's constant, and T is the ambient temperature in kelvin. It is clear from Eq. (8) that the diffusion current is higher in a low-bandgap material. Therefore, InSb ($E_g = 0.17$ eV) has far higher diffusion current than silicon ($E_g = 1.12$ eV) and in fact this makes InSb almost useless at room temperature. To overcome this excessive diffusion current, InSb photodiodes are cooled to liquid nitrogen temperature (77 K).

Generation-recombination current The current generated in the depletion region of the photodiode is called the *generation-recombination current*. When impurity trap levels are present within the forbidden gap E_g , trapped carriers can be elevated to the conduction band with less energy than for diffusion current. This “trap-assisted” current is given by

$$I_{g-r} \propto \sqrt{(V_{\text{bi}} - V)} e^{-E_g/2kT} \quad (9)$$

From Eq. (5), the depletion width W is proportional to $V_{\text{bi}} - V$. Hence,

$$I_{g-r} \propto W e^{-E_g/2kT} \quad (10)$$

Generation-recombination current is proportional to the volume of the depletion width and, hence, is reverse-bias-dependent, whereas the diffusion current in Eq. (8) is bias-independent. For high-bandgap semiconductors with bandgaps above 1.0 eV (e.g., silicon), the generation current

usually dominates over the diffusion current at room temperature. However, for low-bandgap material such as indium antimonide, the diffusion current is dominant over generation current at room temperature.

Tunneling current When the electric field in a reverse-biased p - n junction exceeds 10^5 V/cm, a valence band electron can jump to the conduction band due to the quantum mechanical effect¹⁰ called *tunneling* which occurs at high field and with geometrically narrow energy barriers. The tunneling current is given by

$$I_{\text{tun}} \propto EV \exp\left(\frac{-\theta\sqrt{m}}{E} E_g^{3/2}\right) \quad (11)$$

where E is the applied electric field, m is the effective mass of an electron, and θ is a dimensionless constant whose value depends on the tunneling barrier height. Higher doping levels at the p - n junction lead to a narrower depletion width which causes higher electric fields, thus increasing the amount of tunneling current. Low-bandgap photodiodes exhibit much more tunneling than do higher-bandgap diodes. Tunneling shows a weak dependence on temperature, the only minor change being caused by the temperature dependence of the bandgap E_g . This leads to a *decreasing* breakdown voltage with an increasing temperature as opposed to an *increasing* breakdown voltage exhibited by the avalanche effect.

Quantum Efficiency, Responsivity, and Absorption Coefficient Quantum efficiency η is defined as the ratio of electron-hole pairs generated for each incident photon. In a nonavalanche p - i - n photodiode, quantum efficiency is less than unity. Responsivity R is a measured quantity in amps/watt or volts/watt and is related to quantum efficiency by

$$\eta = \frac{(1240)R}{\lambda} \quad (12)$$

where λ is the wavelength in nm of incident photons and R is the responsivity in amps/watt. The value of η is determined by the absorption coefficient α of the semiconductor material and the penetration distance x in the absorbing layer. The light flux A at a distance x with the absorption layer is

$$A = A_0 e^{-\alpha x} \quad (13)$$

where, A_0 is the incident photon flux and α is a strong function of wavelength λ . Figure 7 shows its typical values for a 1- μm -thick undoped $\text{In}_x\text{Ga}_{1-x}\text{As}$, $0 < x < 0.25$.¹⁶ For optimized η , the reflectivity at the semiconductor surface has to be minimized. Hence, an antireflection (AR) coating of proper thickness is deposited on the photodiode surface. For single-layer AR coatings, the proper “quarter-wave” thickness (L) of the AR coating is

$$L = \frac{\lambda}{4n} \quad (14)$$

where n is its refractive index. With good AR coatings, InGaAs photodiodes can achieve quantum efficiencies above 95 percent at 1300 to 1500 nm. For the visible region, silicon photodiodes show high η (90 percent) in the 800-nm range, and the mid-infrared InSb has a typical η of 80 percent at 5000 nm.

Avalanche Photodiodes Avalanche photodiodes (APDs) will be briefly discussed here. For a detailed treatment, see Chap. 26, “High-Speed Photodetectors,” by Bowers and Wey. An avalanche photodiode is a p - i - n diode with a net efficiency or gain greater than unity. This is obtained through the process of “impact ionization” by operating the photodiode at a sufficiently high reverse bias. The typical operating voltage for an InGaAs APD is 75 V, while that for silicon can be as high as 400 V. The

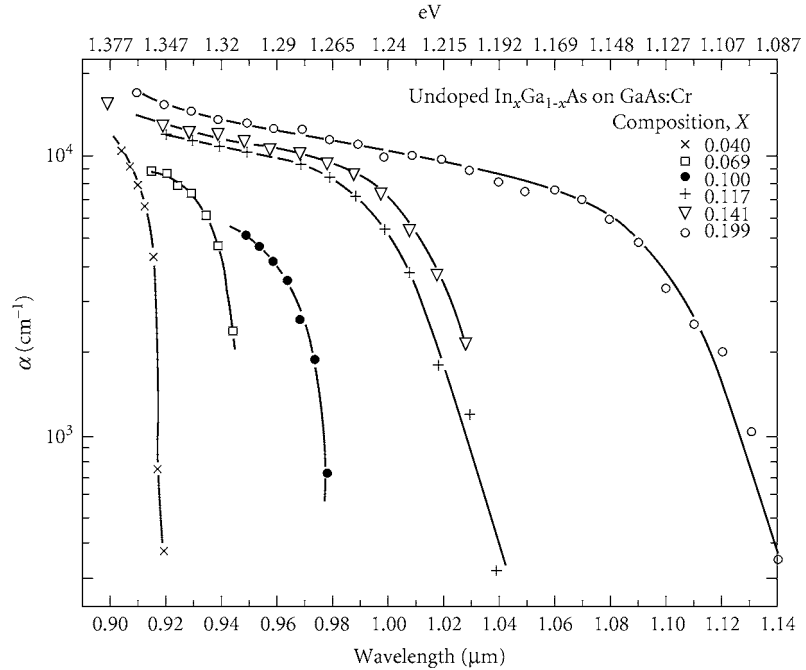


FIGURE 7 Absorption coefficients for 1- μm -thick undoped $\text{In}_x\text{Ga}_{1-x}\text{As}$, $0 < x < 0.25$.¹⁶

impact ionization process is described in Fig. 8. Under the influence of a high electric field ($>5 \times 10^5$ V/cm), electron A gains sufficient kinetic energy to strike atom B with a tremendous force and knock out an electron hole pair $B'-B''$. A is called the “parent” and $B'-B''$ the “child” charge carriers. The child electron B'' moves through a critical distance S and acquires enough kinetic energy to create its own child particles $D'-D''$. The sum effect of the impact ionization of a number of electrons is termed *avalanche multiplication*. Because of this avalanche action, the gain in an APD exceeds unity, reaching useful values above 10 for InGaAs and several hundred for silicon before the multiplied noise begins to exceed the multiplied signal. A solid-state APD is a fast device with gain-bandwidth products that can exceed 20 GHz.^{17,18} In spite of high operating bias, an APD can be designed for low noise operation¹⁹ and used for numerous applications such as photon-counting, laser pulse detection,²⁰ and fiber-optic communication. The gain or avalanche multiplication M of an APD is given by

$$M = \frac{IR_1}{IR_0} \quad (15)$$

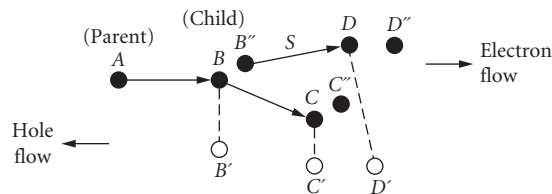


FIGURE 8 Impact ionization process.

where IR_1 is the reverse current generated by avalanche action and IR_0 is the unity gain current. At voltage breakdown of the APD, the multiplication factor M tends to infinity. An empirical relation between the multiplication factor M and reverse bias V is given by^{21,22}

$$M = \frac{1}{1 - (V/V_B)^n} \quad (16)$$

where V is the applied reverse bias and V_B is the breakdown voltage. The factor n varies between 3 and 6, depending on the semiconductor material and its substrate type. Typical gains are on the order of 10 to 20 for germanium and InGaAs APDs, and above 100 for silicon APDs. Due to their lower noise, InGaAs and silicon APDs have better sensitivity than their germanium counterparts.

Extended Wavelength (1000–3000 nm) Photodetectors

Detector materials used for the 1000 to 3000 nm spectrum include InSb, InAs, PbS, HgCdTe, and recently InGaAs. PbS is an inexpensive, reasonably sensitive detector that can operate at relatively high temperatures, even at room temperature. Its major drawback is its slow (typically milliseconds) response time. InAs has higher sensitivity over the 1000 to 3500 nm spectrum and fast response time, but must be cooled thermoelectrically (to 230 K) or cryogenically (to 77 K). InSb has similar properties out to 5500 nm but must definitely be cryogenically cooled. HgCdTe has high sensitivity and speed and it can be operated at room temperature. Indium gallium arsenide was originally developed for fiber-optic applications out to 1.7 μm (using $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$) but it can be used out to 2500 nm by increasing its indium content to $\text{In}_{0.8}\text{Ga}_{0.2}\text{As}$. InGaAs appears to be the *best* detector material for high-temperature operation in the 1000 to 3000 nm spectrum. It has a 10 to 100 times advantage in shunt resistance at room temperature compared to HgCdTe—the previously used material for this wavelength.

Table 2 contains a summary of the data. It is difficult to find data at exactly the same cutoff wavelengths and temperature with the same area device. R_0A was determined (in cases where it was given as such) by simply multiplying two numbers, where, R_0 is the shunt impedance of the detector, and A is the active area of the photodetector.

Photogate (Metal-Insulator-Semiconductor Detector)

The advent of silicon charge-coupled devices (CCDs) has revolutionized the television industry and introduced one of the most popular consumer items to millions of people around the world—the CCD camcorder. From the sandy shores of Hawaii to the ski slopes of Colorado, people have captured

TABLE 2 Comparison of R_0A Values in HgCdTe and InGaAs ($\Omega\text{-cm}^2$)

λ_{co} (nm)	$R_0A(T)$	
	HgCdTe	InGaAs
1400	4×10^4 (292 K)	2.5×10^5 (300 K)
	7×10^6 (230 K)	1.3×10^8 (220 K)
1700	2×10^2 (300 K)	2.5×10^5 (300 K)
	2×10^5 (220 K)	1.3×10^8 (220 K)
2100	7×10^1 (300 K)	2.5×10^3 (300 K)
	7×10^3 (220 K)	6.5×10^5 (220 K)
2500	1×10^1 (300 K)	1.3×10^2 (300 K)
	1×10^3 (210 K)	1.0×10^5 (210 K)

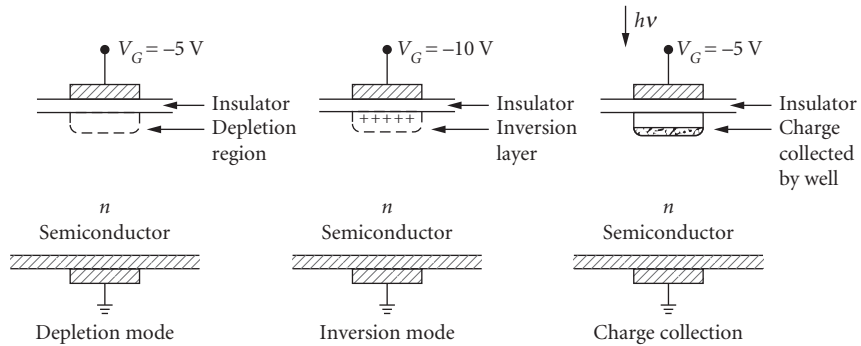


FIGURE 9 Schematic diagram of a photogate.

life's best moments with a CCD camcorder. With its superior imaging quality and noise performance of a few electrons per pixel (<20), a CCD has diverse uses from space imaging to chemical analysis spectroscopy. A CCD is a matrix of metal oxide semiconductor (MOS) devices operated in the "depletion" mode. Each individual MOS device is called a *photogate* and its schematic diagram is shown in Fig. 9. Consider an n -type semiconductor with a negative potential applied to its gate. It repels the negatively charged electrons and create a depletion layer. As the negative potential on the gate is increased, the volume of the depletion region increases further into the bulk. However, the surface potential at the semiconductor-insulator interface also becomes more negative. Finally, with increased gate bias, the surface potential becomes sufficiently high to attract minority carriers (holes). This creates a positive charge at the semiconductor-insulator interface and is termed the *inversion layer*. In a MOS transistor, the inversion layer forms a conducting channel between the source and the drain, and the gate bias needed to achieve inversion is termed the *threshold voltage*. Usually a photogate is operated in depletion at a gate bias lower than the threshold voltage. When incident photons create hole-electron pairs, the minority carriers drift away to the depletion region and the volume of the depletion region shrinks. The total amount of charge that a photogate can collect is defined as its *well capacity*. The total well capacity is decided by the gate bias, the insulator thickness, the area of the electrodes, and the background doping of the semiconductor. Numerous such photogates with proper clocking sequence form a CCD imaging array. For in-depth understanding of CCDs, we refer to Chap. 32 "Visible Array Detectors," by Timothy J. Tredwell.

25.4 APPLICATIONS

The main commercial uses of photodetectors include optical communications and infrared sensing. Although these applications often overlap, optical communication typically involves transmitting data over an optical fiber at higher rates. The format is increasingly digital (telecommunications and data links) at rates from 1 Mbit/s to over 2 Gbit/s.

However, one growing application is cable TV where analog data rates from 1 to over 1000 MHz are most often found. Infrared sensing mostly involves nonfiber applications at sub-MHz analog rates. The property to be detected is usually the amplitude (in watts) and wavelength of the incoming radiation. In digital applications, the wavelength and individual pulse amplitude are relatively fixed, and successful communication occurs simply by distinguishing when the pulse is "on" or "off." Although very weak pulses must sometimes be detected, the actual amplitude of the pulse is irrelevant. The ultimate "resolving power" of the detector is when a weak pulse can no longer be distinguished from background noise, i.e., when the incoming signal strength S equals the background

noise strength N or when the signal-to-noise ratio $S/N = 1$. Thus, the strengths of individual pulses are unimportant as long as the presence of a pulse can be detected. Information is conveyed by the timing sequence of the pulses rather than by the amplitude of the individual pulse. Analog applications, on the other hand, depend critically on the frequency content and amplitude of the transmitted signal. In an AM cable TV transmission system, the detector must be able to linearly reproduce the incoming optical signal as an electrical current of the same frequency content and amplitude, and to minimize intermodulation and harmonic distortion that is invariably produced in the detection of an AM signal.

Infrared applications often involve spectroscopy whereby the detected electrical signal depends on both the optical wavelength and strength of the incoming infrared signal. Thus, the detector must be carefully calibrated in terms of “responsivity” (electrical amps/optical watt) versus wavelength in order to accurately identify the nature of the incoming signal. Identification of gases (e.g., methane, which absorbs light near 1650 nm) depends on these properties. Other “infrared” applications include spectroscopy, remote sensing from satellite, and general laboratory detection. Not *all* infrared applications are analog, however. One notable digital application is LIDAR (Light detection and ranging), which essentially is a form of laser radar. High-intensity light pulses are emitted into the atmosphere (or a gas) which absorbs, scatters, and reemits the laser pulse. The character of the light pulses detected near the source can be used to determine the nature of the gas particles that interact with the light: the absorbing wavelength, the gas density (velocity), and the amount present. Applications include remote pollution monitoring and “windshear detection,” whereby the presence of abrupt changes in wind velocity can be instantly detected at distances of several miles. This application²³ has been demonstrated with laser wavelengths of 2060 and 10,600 nm for use on an aircraft. The 2060-nm system works better in severe storms and does not require a cryogenically cooled detector (as does the 10,600-nm system).

One important noise source in infrared applications is the so-called $1/f$ noise which becomes noticeable at frequencies below 10 MHz. Although poorly understood, this noise is thought to originate at heterointerfaces such as semiconductor-metal contacts and heteroepitaxial interfaces. Photodiode arrays are often used to detect low-light-level signals of a few hundred photons, and they must integrate the signal for 1 second or more. However, with longer integration times, $1/f$ noise may become noticeable and can degrade the S/N ratio and thus, impose an upper limit on the effectiveness of longer integration times. Limiting $1/f$ noise becomes critical for numerous infrared sensing applications and research indicates that surface depletion width at the semiconductor-insulator interface to be a major source of $1/f$ noise in the InGaAs photodiodes.²⁴

One important area for detectors is the array configuration used both for spectroscopy (linear) and imaging (two-dimensional). Linear arrays are used in the so-called multichannel analyzers whereby the detector is placed behind a fixed grating and the instrument functions as “motionless” or “instant” spectrometer with each pixel corresponding to a narrow band of wavelengths. The resolution of the instrument is determined by the number and spacing of pixels, so *narrow* pixel geometries are needed along one direction whereas *tall* pixel geometries are needed along the perpendicular direction to enhance the light collection.

One “mixed” infrared/fiber-optic application is the use of large-area (typically 3 mm diameter) detectors for optical power meters: the optical equivalent of a voltmeter which accurately measures the amount of optical power in watts or dBm (number of decibels above or below 1 mW contained in an incoming beam). The large area ensures large collection efficiency. The most important parameter here is the responsivity and the uniformity of response across the detector.

A “figure-of-merit” for infrared detectors is D^* (D-star), whereby detectors of differing area can be compared. It is related to the noise equivalent power (NEP) in watts, the lowest power a detector can detect at a signal-to-noise ratio of 1 as

$$D^*(\lambda, f, B) = (AB)^{1/2} / \text{NEP} \quad (17)$$

where A is the detector area. The optical bandwidth B (often taken to be 1 Hz), frequency of signal modulation f , and operating wavelength λ must be stated.

25.5 RELIABILITY

In today's global economy of severe competitiveness, new product development and innovation are incomplete without quality assurance and reliability. Reliability is the assurance that a device will perform its stated functions for a certain period of time under stated conditions, and considerable research has been done to improve the reliability of photodetectors.^{25–28} The two major industrial standards for testing semiconductor device reliability are (1) test methods and procedures for microelectronics (MIL. STD. 883C), and (2) Bellcore technical advisory (TA-TSY-00468). The former standard is generic to the semiconductor industry, while the latter is specifically developed for fiber-optic optoelectronic devices.

The tests performed under MIL. STD. 883C comprise the following major groups: (1) *environmental tests*, e.g., moisture resistance, burn-in, seal, dew point, thermal shock, (2) *mechanical tests*, e.g., constant acceleration, mechanical shock, vibration, solderability, and bond strength, and (3) *electrical tests*, e.g., breakdown voltage, transition time measurements, input currents, terminal capacitance, and electrostatic discharge (ESD) sensitivity classification. Under the Bellcore Technical Advisory, each photodetector lot undergoes visual inspection, optical and electrical characterization, and screening. Visual inspection removes any photodiodes with faulty wire bonds or cracks in the glass window or in the insulating films. Table 3 lists the electrical and optical testing performed on every photodiode. After testing, all the devices are sent for screening (burn-in), e.g., some $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ photodiodes are burned-in at 200°C for 20 hours at –20 V reverse bias to weed out any infant mortality.

Photodetector Life Test (Accelerated Aging)

To predict the lifetime or *mean-time-to-failure* (MTTF), accelerated aging tests are carried out on groups of diodes at several elevated temperatures. For example, the MTTF for 300 μm diameter InGaAs photodiodes was determined on groups of 20 screened devices at elevated temperatures of 200, 230, and 250°C. The failure criterion was a 25 percent increase in the room temperature dark current value.^{28,29} The total lifetest extended over a time period of several years, and every week the samples were cooled to room temperature to check their dark current. Failed devices were removed from the sample population and the remaining good ones put back at the elevated temperature.

From the temperature-dependence of the data, it was observed that the failure mechanism is thermally activated. The Arrhenius relationship calculates the activation energy E_a for thermally activated failure²⁹ as

$$\text{MTTF}(T) = Ce^{(E_a/kT)} \quad (18)$$

where, C is a constant, k is Boltzmann's constant (8.63×10^5 eV/K), and T is the temperature in kelvin. Figure 10 shows the MTTF for three batches of 300 μm diameter $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ photodiodes. Using least-squares fit to the data, the calculated activation energy E_a is 1.31 eV with a correlation coefficient r^2 of 0.99. From Eq. (18) and an experimentally determined activation energy

TABLE 3 Electrical and Optical Testing of Photodetectors

Tests or Measurement	Parameter	Symbol	References
Optical response	Responsivity	R	Bellcore Technology Advisory
	Gain	G	
Electrical performance	Dark current	I_d	TA-TSY-00468 Issue 2, July 1988
	Breakdown voltage	V_{br}	

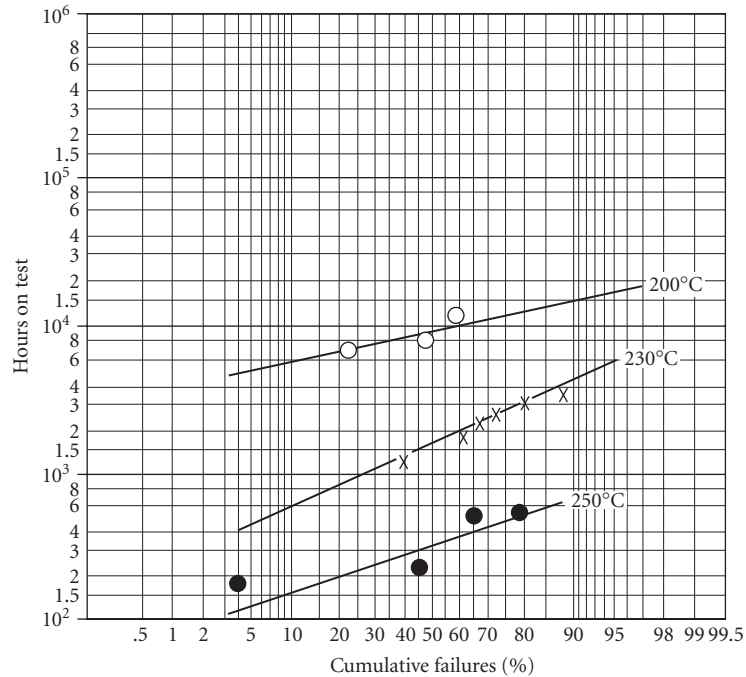


FIGURE 10 Time versus cumulative percent failure of InGaAs diodes at 200, 230, and 250°C lifetest.

of 1.31 eV, the MTTF at 25°C is calculated to be 1.34×10^{14} hours. Such a “geological” lifetime may seem to be an overkill, even for the electronics industry. However, when thousands of these devices are working together in a single system (e.g., telephones), the net MTTF of all these devices chained together may be on the order of only a few years. Thus, continuing improvements in reliability must be an ongoing process. Reliability in most photodetectors is determined by a number of factors including: (1) material quality, (2) processing procedures, (3) planar technology versus mesa technology, and (4) amount of leakage current. Poor material quality can introduce crystal defects such as vacancies and dislocations which can increase the dark current. Higher dark current has been directly linked to lower MTTF.²⁸ Device processing is probably the most crucial item in photodetector reliability. The dielectric (typically silicon nitride) used in planar detector processing serves as a diffusion mask in *p-n*-junction formation and a passivant (termination) for the junction so produced. Any surface states or impurities introduced here can directly increase leakage current and degrade reliability.

An important milestone in detector reliability was the changeover mesa to planar structures.^{2,30} Just as the transistors in the 1950s were first made in mesa form, so were the optical photodetectors of the 1980s, due to their simplicity and ease of fabrication. However, in both cases, reliability issues forced the introduction of the more complex planar structure. A sketch of a mesa and planar photodiode is illustrated in Fig. 11. A mesa photodiode typically is formed by wet chemical etching of an epitaxially grown *p-n* crystal structure, while in a planar process, a *p-n* junction is formed by diffusing a suitable *p* or *n* dopant in an *n*- or *p*-type crystal. It has been shown a planar structure to be more reliable than a mesa one³⁰ because a *p-n* junction is never exposed to ambient conditions in a well-designed planar process. Exposure of the *p-n* junction can cause surface corrosion leading to increased leakage current and, in effect, poorer reliability.³¹

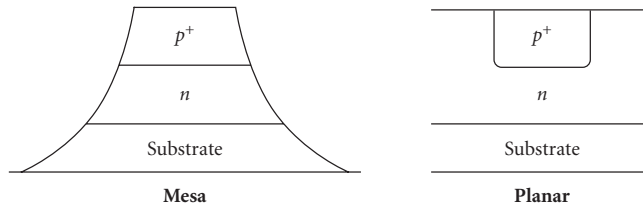


FIGURE 11 Sketch of a mesa and planar photodetector.

25.6 FUTURE PHOTODETECTORS

A lateral $p-i-n$ photodiode and a quantum well infrared photodetector (QWIP) have been developed with better characteristics compared to the photodetector structures. A long wavelength QWIP in the 8000- to 12000-nm band³²⁻³⁵ has posed a severe challenge to the present favorite mercury cadmium telluride photodetectors, while a medium wavelength QWIP in the 3000- to 5000-nm band³⁶ may compete with indium antimonide and platinum silicide photodetectors. A QWIP made from GaAs/AlGaAs heterosystems promises to have higher detectivity (D^*), higher yield due to well-established 3-in wafer GaAs technology, and easier monolithic integration with circuit electronics. A lateral $p-i-n$ diode, as the name implies, has charge carrier flow in a lateral direction compared to the vertical direction in a conventional (vertical) photodiode structure. Because of its process compatibility and simple fabrication, a lateral $p-i-n$ photodiode can be suitably integrated on an optoelectronic integrated circuit (OEIC) chip^{37,38} having numerous field-effect transistors. An OEIC has a lower noise floor due to the reduced stray capacitances and inductances compared to that of hybrid detector-amplifier packages and finds applications in high-speed digital data communication.

Lateral $p-i-n$ Photodetector

The vertical $p-i-n$ structure in Fig. 6 has high sensitivity, low noise, low capacitance, better reliability, and an easy manufacturing process. However, such a vertical structure is nonplanar and therefore harder to integrate on an OEIC. The nonplanarity is also an issue with lasers and LEDs, and optical integration demands surface-emitting LEDs and lasers (SLEDs and SLASERs) over the conventional edge-emitting sources (ELED and ELASER). The cross section of an AlGaAs/GaAs lateral $p-i-n$ photodiode is shown in Fig. 12.

The higher-bandgap AlGaAs layer acts as a surface barrier, reducing the leakage currents. The low-bandgap GaAs layer absorbs the incoming light, and the generated carriers flow to the W-Zn

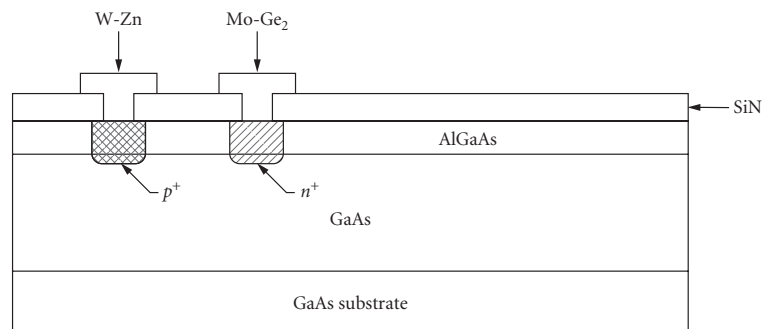


FIGURE 12 Cross section of an AlGaAs/GaAs lateral $p-i-n$ photodiode.³⁸

and Mo-Ge₂ ohmic contacts which act as the *p* and *n* regions, respectively. The diffusion of zinc and germanium in the ohmic contacts forms a compositionally graded barrier at the AlGaAs/GaAs interface, rather than an abrupt interface. This smooth barrier helps the lateral *p-i-n* photodiode to have better speed than a lateral metal-semiconductor-metal photodetector. A comfortable spacing of 3 to 5 μm between the *p* and *n* regions gives high quantum efficiency and low capacitance, thus providing all the desirable properties of a vertical *p-i-n* structure and yet being easier to integrate.

Quantum Well Infrared Photodetector

Quantum well infrared photodetectors (QWIPs) offer *long* wavelength (5000–10,000 nm) infrared detection by using materials whose bandgap normally allows them only to absorb light in the *short* wavelength (~10,000 nm) region, e.g., GaAs/AlGaAs. The use of thin (<500 Å) layers allows the absorbing wavelength to be controlled by material *geometry* rather than its *chemistry*.³⁹

Before discussing the QWIPs, we take the liberty of explaining a few basic terms and concepts of quantum physics. Superlattices or quantum well structures consist of a stack of ultrathin semiconductor layers normally 50 to 500 Å in thickness. Molecular beam epitaxy (MBE) techniques are frequently employed to grow these structures because their characteristically slow growth rate of a few Angstroms per seconds which helps achieve abrupt heterointerfaces. Two semiconductors of different compositions, when stacked together, form a heterointerface. Type III-V compound semiconductors such as AlGaAs/GaAs and InAlAs/InGaAs are the best candidates for growing quantum well structures, as they can be easily doped and their alloy composition readily changed to form semiconductor layers of different bandgaps. Tailoring the bandgap can alter the heterobarriers, creating exciting device results. When quantum well layers have thicknesses less than the electron mean free path (typically 50 to 100 Å), electron and holes cannot have their normal three-dimensional motion. This restricts carriers to move in two dimensions in the plane of the layer.^{2,39} Because of this quantized motion, a new band of discrete energy levels is generated. Carriers no longer obey Boltzmann's statistics¹ and optical absorption becomes more complicated than the conventional band-to-band absorption given by Eq. (1). The absorption of light energy by a quantum well structure can cause an electron to jump from "multiple valence subbands" to "multiple conduction subbands," thereby enabling it to absorb light wavelengths not decided by the *material* properties (bandgap) of the semiconductor layers alone, but by its *geometrical* properties as well.

In QWIPs, the light energy transfers an electron in a bound state to an excited state in the continuum.⁴⁰ Figure 13 shows an AlGaAs/GaAs quantum well structure with *L* being the width of the

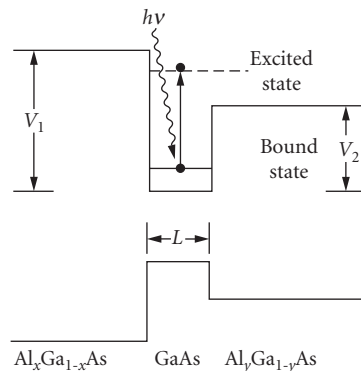


FIGURE 13 Infrared detection with an AlGaAs/GaAs quantum well.⁴⁰

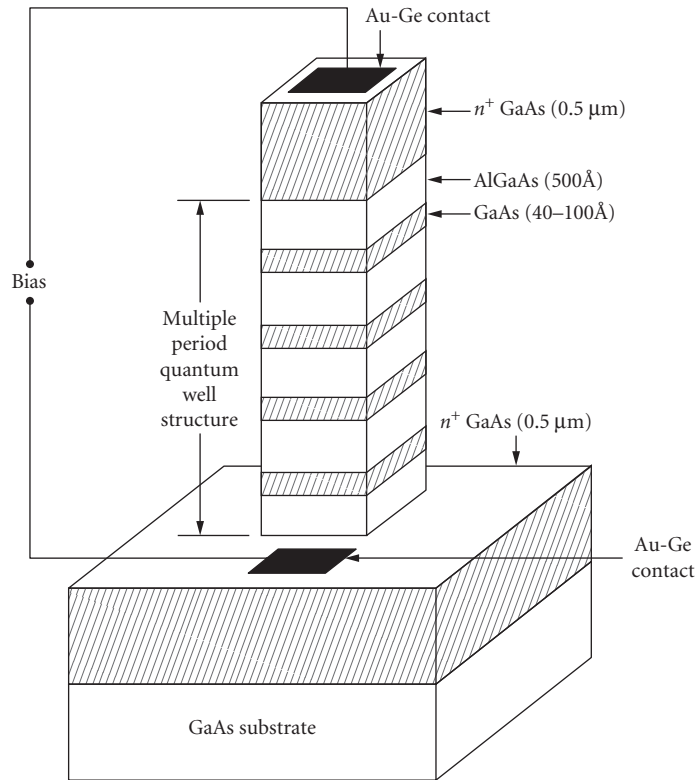


FIGURE 14 Multiple period AlGaAs/GaAs quantum well infrared photodetectors.

well and V_1, V_2 being the barrier heights. The electron excited by the IR radiation is swept out of the doped GaAs well by applying an external electrical field. By controlling the barrier heights V_1, V_2 , and quantum well width L , the spectral response of a QWIP can be changed for the desired IR window of 3000 to 5000 or 8000 to 12,000 nm.⁴⁰ A multiple period quantum well infrared photodetector is illustrated in Fig. 14.³⁴ The n^+ -doped ($2 \times 10^{18} \text{ cm}^{-3}$) GaAs quantum wells are 40 to 100 Å and the undoped AlGaAs barriers are of 500 Å thickness. The multiple period stack is sandwiched between two n^+ GaAs-doped contacts. This photodetector has exhibited a blackbody D^* of $1 \times 10^{10} \text{ cm}/(\text{Hz/W})^{1/2}$ at 68 K for a cutoff wavelength of 10,700 nm. InGaAs/AlInAs superlattices have exhibited blackbody D^* of $2 \times 10^{10} \text{ cm}/(\text{Hz/W})^{1/2}$ at 120 K with peak responsivity at 4000 nm.³⁶

In summary, QWIPs have high detectivity, good uniformity, high yield, multiple spectral windows, and intrinsic radiation hardness for numerous imaging and spectroscopy applications.⁴¹

25.7 ACKNOWLEDGMENTS

We sincerely acknowledge the support of EPITAXX, Inc., Amy Vasger, and Jennifer Romano (Sensors Unlimited) for preparing the manuscript, and Jim Rue for technical discussions.

25.8 REFERENCES

1. S. M. Sze, *Physics of Semiconductor Devices*, Wiley, 1981.
2. S. R. Forrest, "Optical Detectors for Lightwave Communication," *Optical Fiber Telecommunications II*, 1988, pp. 569–599.
3. Z. S. Huang and T. Ando, "A Novel Amplified Image Sensor with a Si:H Photoconductor and MOS Transistors," *IEEE Trans. on Elect. Dev.*, vol. 37, no. 6, 1990, pp. 1432–1438.
4. J. C. Gammel, G. M. Metzger, and J. M. Ballantyne, "A Photoconductor Detector for High Speed Fiber Communication," *IEEE Trans. on Elect. Dev.* vol. ED-28, no. 7, 1981, pp. 841–849.
5. M. V. Rao, P. K. Bhattacharya, and C. Y. Chen, "Low Noise $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$: Fe Photoconductive Detectors for Optical Communications," *IEEE Trans. on Elect. Dev.* vol. ED-33, no. 1, 1986, pp. 67–71.
6. J. C. Gammel, H. Ohno, and J. M. Ballantyne, "High Speed Photoconductive Detectors Using GaInAs," *IEEE J. of Quantum Elect.* vol. QE-17, no. 2, 1981, pp. 269–272.
7. C. Y. Chen, Y. M. Pang, K. Alavi, A. Y. Cho, and P. A. Garbinski, "Interdigitated $\text{Al}_{0.48}\text{In}_{0.52}\text{AsGa}_{0.47}\text{In}_{0.53}\text{As}$ Photoconductive Detectors," *App. Phys. Lett.* 44, 1983, pp. 99–101.
8. M. V. Rao, G. K. Chang, and W. P. Hong, "High Sensitivity, High Speed InGaAs Photoconductive Detector," *Elect. Lett.* vol. 26, no. 11, 1990, pp. 756–757.
9. J. Degani, R. F. Leheny, R. E. Nahory, M. A. Pollack, J. P. Heritage, and J. C. DeWinter, "Fast Photoconductive Detector Using p- $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ with Response to 1.7 μm ," *Appl. Phys. Lett.* 38, 1981, pp. 27–29.
10. S. M. Sze, *Semiconductor Devices—Physics and Technology*, Wiley, 1985.
11. G. H. Olsen, "Low Leakage, High Efficiency, Reliable VPE InGaAs 1.0–1.7 μm Photodiodes," *IEEE Elect. Dev. Lett.* vol. EDL-2, no. 9, 1981, pp. 217–219.
12. S. R. Forrest, "Performance of $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}$ Photodiodes with Dark Current Limited by Diffusion, Generation Recombination and Tunneling," *IEEE J. of Quantum Elect.* vol. QE-17, no. 2, 1981, pp. 217–226.
13. O. K. Kim, B. V. Dutt, R. J. McCoy, and J. R. Zuber, "A Low Dark-Current, Planar InGaAs *p-i-n* Photodiode with a Quaternary InGaAsP Cap Layer," *IEEE J. of Quantum Elect.* vol. QE-21, no. 2, 1985, pp. 138–143.
14. Y. Yoshida, Y. Hisa, T. Takiguchi, and Y. Komine, "Reduction of Surface Leakage Current in $\text{Cd}_{0.2}\text{Hg}_{0.8}\text{Te}$ Photodiode," *Proc. SPIE* vol. 972, 1988, pp. 39–43.
15. J. C. Flachet, M. Royer, Y. Carpentier, and G. Pichard, "Emission and Detection in the 1 to 3 μm Spectral Range with $\text{Hg}_{1-x}\text{Cd}_x\text{Te}$ Diodes," *Proc. SPIE* vol. 587, 1985, pp. 149–155.
16. R. E. Enstrom, P. J. Zanucchi, and J. R. Apert, "Optical Properties of Vapor-Grown $\text{In}_x\text{Ga}_{1-x}\text{As}$ Epitaxial Films on GaAs and $\text{In}_x\text{Ga}_{1-x}\text{P}$ Substrates," *J. of Appl. Phys.* vol. 45, no. 1, 1974, pp. 300–306.
17. H. W. Ruegg, "An Optimized Avalanche Photodiode," *IEEE Trans. on Elect. Dev.* vol. ED-14, no. 5, 1967, pp. 239–251.
18. K. Taguchi, T. Torikai, Y. Sugimoto, K. Makita, and H. Ishihara, "Planar Structure InP/InGaAsP/InGaAs Avalanche Photodiodes with Preferential Lateral Extended Guard Ring for 1.0–1.6 μm Wavelength Optical Communication Use," *J. of Lightwave Tech.* vol. 6, no. 11, pp. 1643–1655.
19. S. R. Forrest, R. G. Smith, and O. K. Kim, "Performance of $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ Avalanche Photodiodes," *IEEE J. of Quantum Elect.* vol. QE-18, no. 12, 1982, pp. 2040–2047.
20. R. J. McIntyre, "The Distribution of Gains in Uniformly Multiplying Avalanche Photodiodes: Theory," *IEEE Trans. on Elect. Dev.* vol. ED-19, no. 6, 1972, pp. 703–713.
21. A. S. Grove, *Physics and Technology of Semiconductor Devices*, Wiley, New York, 1967.
22. S. L. Miller, "Avalanche Breakdown in Germanium," *Phys. Rev.* vol. 99, 1955, p. 1234.
23. M. E. Storm, "Coherent 2 μm Sources Burst into Windshear Detection," *Laser Focus* vol. 21, 1991, pp. 117–122.
24. A. M. Joshi, G. H. Olsen, V. S. Ban, E. Mykietyn, M. J. Lange, and D. R. Mohr, "Reduction of 1/f Noise in Multiplexed Linear $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ Detector Arrays via Epitaxial Doping," *IEEE Trans. on Elect. Dev.* vol. 40, no. 2, 1993, pp. 303–308.
25. R. U. Martinelli and R. E. Enstrom, "Reliability of Planar InGaAs/InP Photodiodes Passivated with BoroPhosho-Silicate Glass," *J. of Appl. Phys.* vol. 63, no. 1, 1988, pp. 250–252.
26. A. K. Chin, F. S. Chen, and F. Ermanis, "Failure Mode Analysis of Planar Zinc-Diffused $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ *p-i-n* Photodiodes," *J. of Appl. Phys.* vol. 55, no. 6, 1984, pp. 1596–1606.

27. Y. Kuhara, H. Tercuchi, and H. Nishizawa, "Reliability of InGaAs/InP Long Wavelength *p-i-n* Photodiodes Passivated with Polyimide Thin Film," *IEEE J. of Lightwave Tech.* vol. LT-4, no. 7, 1986, pp. 933–937.
28. A. M. Joshi, G. H. Olsen, and S. R. Patil, "Reliability of InGaAs Detectors and Arrays," *Proc. SPIE* vol. 1580, 1991, pp. 34–40.
29. S. R. Forrest, V. S. Ban, G. Gasparian, D. Gay, and G. H. Olsen, "Reliability of Vapor Grown In_{0.53}Ga_{0.47}As/InP *p-i-n* Photodiodes With Very High Failure Activation Energy," *IEEE Elect. Dev. Lett.* vol. 9, no. 5, 1988, pp. 217–219.
30. C. P. Skrimshire, J. R. Farr, D. F. Sloan, M. J. Robertson, P. A. Putland, J. C. D. Stokoe, and R. R. Sutherland, "Reliability of Mesa and Planar InGaAs *p-i-n* Photodiodes," *IEEE Proc.* vol. 137, part J, no. 1, 1990, p. 7478.
31. R. R. Sutherland, J. C. D. Stokoe, C. P. Skrimshire, B. M. Macdonald, and D. F. Sloan, "The Reliability of Planar InGaAs/InP *p-i-n* Photodiodes with Organic Coatings for Use in Low Cost Receiver," *Proc. SPIE* vol. 1174, 1989, pp. 226–232.
32. B. F. Levine, C. G. Bethea, G. Hasnain, V. O. Shen, E. Pelve, R. R. Abbott, and S. J. Hsieh, "High Sensitivity Low Dark Current 10 μm GaAs Quantum Well Infrared Photodetectors," *Appl. Physics Lett.* vol. 56, no. 9, 1990, pp. 851–853.
33. B. F. Levine, C. G. Bethea, V. O. Shen, and R. J. Malik, "Tunable Long-Wavelength Detectors Using Graded Quantum Wells Grown by Electron Beam Source Molecular Beam Epitaxy," *Appl. Phys. Lett.* vol. 57, no. 4, 1990, pp. 383–385.
34. G. Hasnain, B. F. Levine, S. Gunapala, and N. Chand, "Large Photoconductive Gain In Quantum Well Infrared Photodetectors," *Appl. Phys. Lett.* vol. 57, no. 6, 1990, pp. 608–610.
35. E. Pelve, F. Beltram, C. G. Bethea, B. F. Levine, V. O. Shen, S. J. Hsieh, and R. R. Abbott, "Analysis of the Dark Current in Doped Well Multiple Quantum Well AlGaAs Infrared Photodetectors," *J. of Appl. Phys.* vol. 66, no. 11, 1989, pp. 5656–5658.
36. G. Hasnain, B. F. Levine, D. L. Sivco, and A. Y. Cho, "Mid-Infrared Detectors in the 3–5 μm Band Using Bound to Continuum State Absorption in InGaAs/InAlAs Multi-quantum Well Structures," *Appl. Phys. Lett.* vol. 56, no. 8, 1990, pp. 770–772.
37. S. N. Subbarao, D. W. Bechtel, R. J. Menna, J. C. Connolly, R. L. Camisa, and S. Y. Narayan, "2–4 GHz Monolithic Lateral *p-i-n* Photodetector and MESFET Amplifier on GaAs-on-Si," *IEEE Trans. on Microwave Theory and Tech.* vol. 38, no. 9, 1990, pp. 1199–1202.
38. S. Tiwari, J. Burroughs, M. S. Milshtein, M. A. Tischler, and S. L. Wright, "Lateral *p-i-n* Photodetectors with 18 GHz Bandwidth at 1.3 μm Wavelength and Small Bias Voltages," *Tech. Dig. of IEEE Int. Elect. Dev. Mtg.* 1991, pp. 421–425.
39. D. S. Chemla, "Quantum Wells for Photonics," *Physics Today*, May 1995, pp. 56–64.
40. D. D. Coon and R. P. G. Karunasiri, "New Mode of IR Detection Using Quantum Wells," *App Phys. Lett.* vol. 45, no. 6, 1984, pp. 649–651.
41. B. F. Levine, C. G. Bethea, J. W. Stayt, K. G. Glogovski, R. E. Leibenguth, S. D. Gunapala, S. S. Pei, and J. M. Kuo, "Long Wavelength GaAs/Al_xGa_{1-x}As Quantum Well Infrared Photodetectors (QWIPs)," *Proc. SPIE* vol. 1540, 1991, pp. 232–238.

25.9 ADDITIONAL READING

- Dereniak, E. L. and D. G. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984.
- Olsen, G. H., "Reliable Operation of Lattice Mismatched InGaAs Detectors on Silicon," *Tech. Dig. of IEEE Int. Elect. Dev. Mtg.*, 1990, pp. 145–147.

This page intentionally left blank.

HIGH-SPEED PHOTODETECTORS

John E. Bowers and Yih G. Wey

*Department of Electrical and Computer Engineering
University of California
Santa Barbara, California*

26.1 GLOSSARY

A	area
A^{**}	modified effective Richardson constant
$A_e(A_h)$	electron (hole) ionization parameters
B	bit rate
C_j	junction capacitance
C_p	pad capacitance
$D_e(D_h)$	diffusion coefficient for electrons (holes)
E	electric field
$e_n(e_p)$	emission functions for electrons (holes)
F	frequency response
f	frequency
f_{3dB}	3-dB bandwidth
G	photoconductor gain
H	transfer function
h	Planck's constant
I_d	dark current
I_{dm}	multiplied dark current
I_{du}	unmultiplied dark current
I_{ph}	photocurrent
i	current
$\langle i_{na}^2 \rangle$	amplifier noise power
J	current density
J_{DIFF}	diffusion component of current density
J_{DRIFT}	drift component of current density

$J_e(J_h)$	electron (hole) component of current density
k	ratio of electron to hole ionization coefficient
k_b	Boltzmann constant
L	absorption layer thickness
$L_e(L_h)$	diffusion length for electrons (holes)
L_s	series inductance
M	multiplication factor
$M_n(M_p)$	electron (hole) initiated multiplication factor
m	electron mass
$n(p)$	electron (hole) density
P	input optical flux
q	electron charge
R	reflectivity
R_L	load resistance
R_s	series resistance
$R_1(R_2)$	reflectivity of the surface (substrate) mirror in a resonant detector
T	temperature
t	time
$t_e(t_h)$	transit time for electrons (holes)
V_B	breakdown voltage
V_j	junction voltage
v_e, v_h	electron and hole velocities
W	thickness of the depleted region
x	position
α	absorption coefficient
α_{FC}	free carrier absorption inside the absorption layer
α_{FCx}	free carrier absorption outside the absorption layer
α_{IB}	interband absorption
α_i	electron ionization rate
α_s	scattering loss
β	propagation constant in a waveguide photodetector
β_i	hole ionization rate
Γ	confinement factor
ϵ	permittivity
η	quantum efficiency
κ	coupling coefficient to the waveguide of a waveguide detector
λ	wavelength
$\mu_e(\mu_h)$	mobility for electrons (holes)
ν	optical frequency
σ	charge density
σ_d	noise current spectral density
τ_e, τ_h	trapping time at a heterojunction for electrons (holes)
τ_{tr}	transit time
$\phi_{bc}(\phi_{bv})$	barrier height for the conduction band (valence band)
ω	angular frequency

26.2 INTRODUCTION

High-speed photodetectors are required for telecommunications systems, for high-capacity local area networks, and for instrumentation. Many different detector structures and materials are required to cover this range of applications. Silicon is one of the most commonly used detector materials for wavelengths from 0.4 to 1.0 μm , while Ge photodetectors are used at longer wavelengths up to 1.8 μm . Silicon and germanium have indirect bandgaps at these wavelengths, which result in relatively small bandwidth-efficiency products. Consequently, for high-speed applications, direct bandgap semiconductors such as III-V materials are more important and are the focus of this chapter. $\text{Ga}_{0.47}\text{In}_{0.53}\text{As}$ with a cutoff wavelength of 1.65 μm is especially useful for telecommunication photodetectors at 1.3 and 1.55 μm . GaAs has a cutoff wavelength around 0.9 μm and is ideal for visible and near-infrared applications.

This chapter will focus on the physics and technology of high-speed photodetectors. The next section discusses the different structures that are possible. Later sections discuss some specific results and motivations for particular structures. The primary limitations to detector speed are discussed, followed by a description of specific photodetector systems. To supplement this chapter, the reader should refer to excellent chapters and articles written specifically about photodetectors,¹ photoconductors,² pin detectors,³ avalanche photodetectors,⁴ phototransistors,⁵ and receivers.^{6,7}

26.3 PHOTODETECTOR STRUCTURES

Many photodetector structures have been demonstrated and many more structures are possible. In this section, we classify the different possible structures and identify a few of the trade-offs. The optimum structure for a given application depends on the required bandwidth, efficiency, saturation power, linearity, ease of integration, and leakage current.

There are four common types of photodetectors: (1) photovoltaic detectors, (2) photoconductive detectors, (3) avalanche photodetectors (APD), and (4) phototransistors. Photovoltaic detectors have blocking contacts and operate under reverse bias. The blocking contact can be a reverse-biased p - n junction or a Schottky contact. The photoconductive detector has identical, nonblocking contacts such as two $n+$ regions in an undoped sample. The avalanche photodetector has a similar configuration to a photovoltaic detector except that it has a high-field region that causes avalanching and results in gain in the detector. Improvements to the basic APD design include separate avalanche and gain regions (SAM APDs), and staircase APDs to increase the ratio of electron to hole (or hole to electron) multiplication rate. Phototransistors are three-terminal devices which have an integrated electronic gain region.

The second criterion is the contact type and configuration. The photogenerated carriers may be collected by means of (1) a vertical current collector, often a p - n or Schottky junction, (2) an interdigitated metal-semiconductor-metal (MSM) structure, or (3) a laterally grown or etched structure. These options are illustrated in Fig. 1a. PIN junctions are usually formed during the growth steps and tend to have low leakage current and high reliability.^{8,9} Schottky junctions are simple to fabricate, but tend to have a large leakage current on narrow-gap semiconductors, such as InGaAs. MSM structures have the advantage of lower capacitance for a given cross-sectional area, but often have longer transit times, limited by the lithography capabilities possible in production. Experimental demonstrations with very fine lines (50 nm) have yielded high-speed devices with good quantum efficiencies. MSM detectors tend to be photoconductive detectors, but one could lower the capacitance for a given area of a PIN detector by using an interdigitated MSM structure with p and n regions under alternating metal fingers.

The third important aspect of photodetector design is the orientation of the light with respect to the wafer and the current collection region (Fig. 1b). Most commercial photodetectors are vertically illuminated and the device area is 10 μm in diameter or larger, which allows simple, high-yield packaging with single-mode optical fibers, or easy alignment to external bulk optics. The problem with this configuration is that the absorbing layer must be thin for a high-speed detector to keep

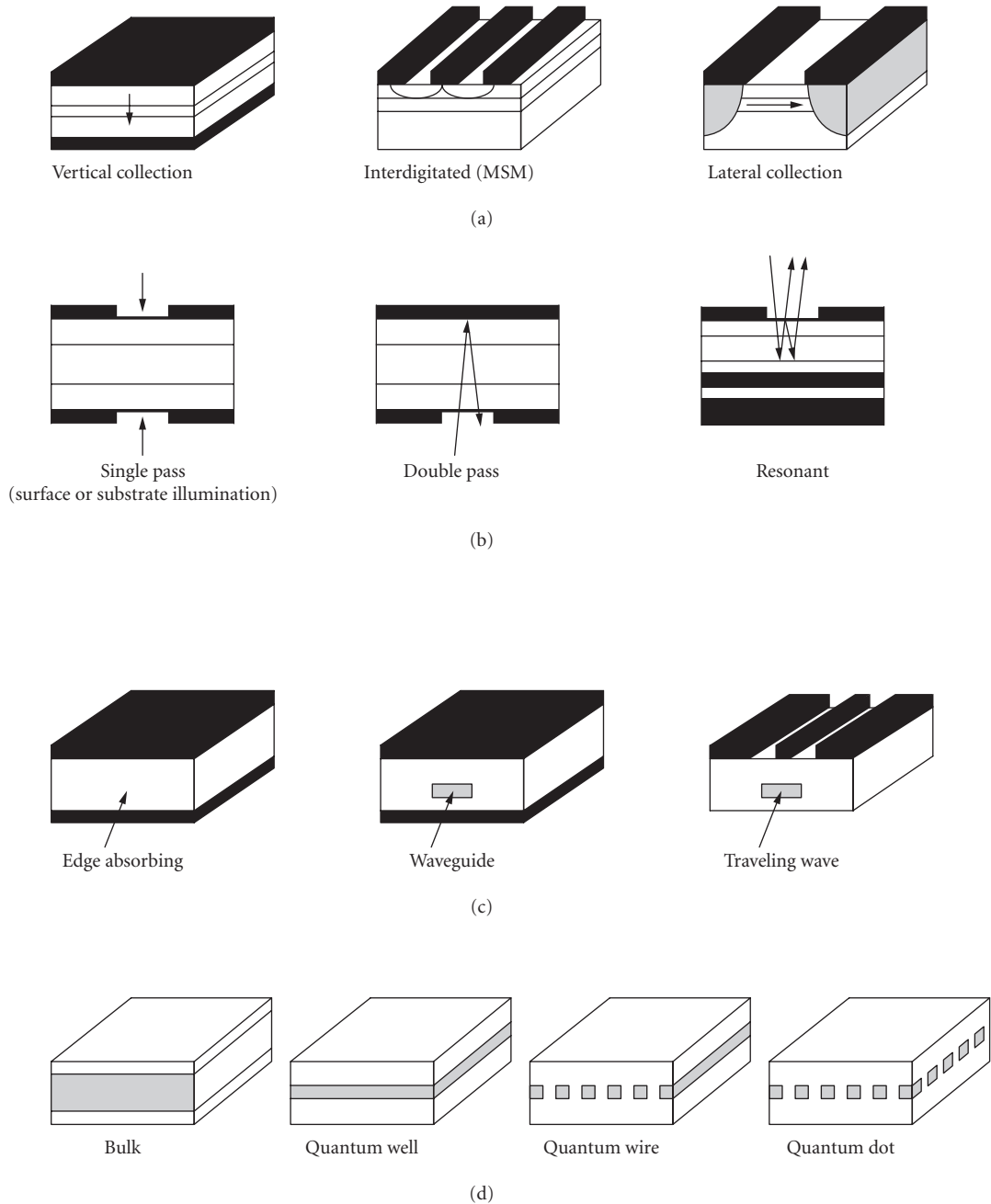


FIGURE 1 Schematic drawings of different types of photodetectors: (a) electrical configuration; (b) optical configuration—vertical illumination; (c) optical configuration—horizontal illumination; and (d) absorbing material.

the transit time of photogenerated carriers short. Consequently, the quantum efficiency is low, and single-pass vertically illuminated photodetectors tend to have bandwidth efficiency products around 30 GHz.³ Bandwidth efficiency products are discussed in greater detail in Sec. 26.5. The bandwidth efficiency product can be increased by allowing two passes by reflecting the light off a metal layer or dielectric mirror.¹⁰ Making a resonant cavity with multiple reflections at particular wavelengths should allow bandwidth efficiency products in excess of 100 GHz.^{11–13} As will be seen below, essentially 100 percent quantum efficiency is possible with bandwidths up to 20 GHz, so there is no need for resonant detectors unless the required bandwidth is above 20 GHz or wavelength selectivity is needed as in a wavelength division multiplexed (WDM) system.

The other class of optical inputs are horizontally illuminated photodetectors (Fig. 1c). The simplest configuration is an edge-illuminated detector. The primary problem with an edge-illuminated detector is that the light is not guided. Diffraction of the incident light causes absorption to occur outside of the high-field region, and slow diffusion tails in the impulse response occur. A solution to this problem is the waveguide detector, where an optical waveguide confines the light to the high-field absorption region.^{14–16} The waveguide efficiency product of this structure can be 100 to 200 GHz. However, it is limited by the capacitance of the structure, particularly if thin intrinsic layers are used for ultrahigh-speed devices. A solution to the capacitance limitation is a traveling wave photodetector where the incoming optical beam is velocity matched with the generated microwave signal.^{17,18} The bandwidth efficiency product is then limited only by loss on the electrical transmission lines, and bandwidth efficiency products of hundreds of GHz are possible. Traveling wave detectors and, to a lesser extent, waveguide detectors have the important advantage that the volume of the light absorption can be quite large and, consequently, these detectors have much higher saturation powers.¹⁹ Velocity matching in these structures requires quite narrow waveguides. Wu and Itoh²⁰ have suggested separating the parts of the optical waveguide with microwave delay lines to achieve velocity matching.

The fourth issue is the type of absorbing material (Fig. 1d) (1) bulk; (2) quantum well; (3) quantum wire; (4) quantum dot; (5 to 7) strained quantum well, wire, or dot; (8) *n-p-i-i* structure. The vast majority of commercial and experimental detectors use bulk material. However, quantum well detectors²¹ are becoming increasingly important in photonic integrated circuits (PICs) because the absorbing quantum well material is also used in other parts of the PIC such as the laser. Quantum wire photodetectors²² have potential advantages in terms of higher-bandwidth-efficiency products, but uniform quantum wires are rather difficult to fabricate. Quantum dot detectors have even higher peak absorption coefficients and more wavelength selectivity, but will probably have problems with slow impulse responses due to trapping of the carriers by the heterojunction. In other quantum-confined detectors, the carriers can be extracted along the quantum wire or well, and this problem can be avoided.²²

The final classification is by means of the lifetime of the material. Conventional detectors have material lifetimes of typically 1 ns and achieve speed by using high field for rapid carrier collection. A second approach is to use low temperature (LT) grown material which has a very short lifetime, perhaps as low as 1 ps. A third approach is to damage the material by means of ion implantation. The final approach is to grow or diffuse in traps into the material such as iron²³ or gold.

If we combine these classifications, we find that 2600 types of photodetectors are possible, and additional subgroups such as superlattice APDs or SAGM APDs increase the total even further. In reality, about 100 types of detectors have been demonstrated. One of the points of this section is that improvements in one type of detector, such as adding a resonant cavity to a PIN detector, can be applied to other types of detectors, such as adding a resonant cavity to an APD. In the following section, we discuss in more detail some of the real limitations to the speed of a detector, and then apply this knowledge to a few important types of detectors.

26.4 SPEED LIMITATIONS

Generally speaking, the bandwidths of most photodetectors are limited by the following factors: (1) carrier transit time, (2) RC time constant, (3) diffusion current, (4) carrier trapping at heterojunctions, and (5) packaging. These limiting factors will be discussed in turn with specific application to *p-i-n* photodiodes.

Carrier Transit Time

In response to light absorbed in a material, the photogenerated carriers in the active region will travel across the high-field region and then be collected by the electrodes. As an example, Fig. 2a shows the p - i - n structure with a photogenerated electron-hole charge sheet of density σ . In response to the electric field, the electron will travel to the right and the hole to the left. This induces a

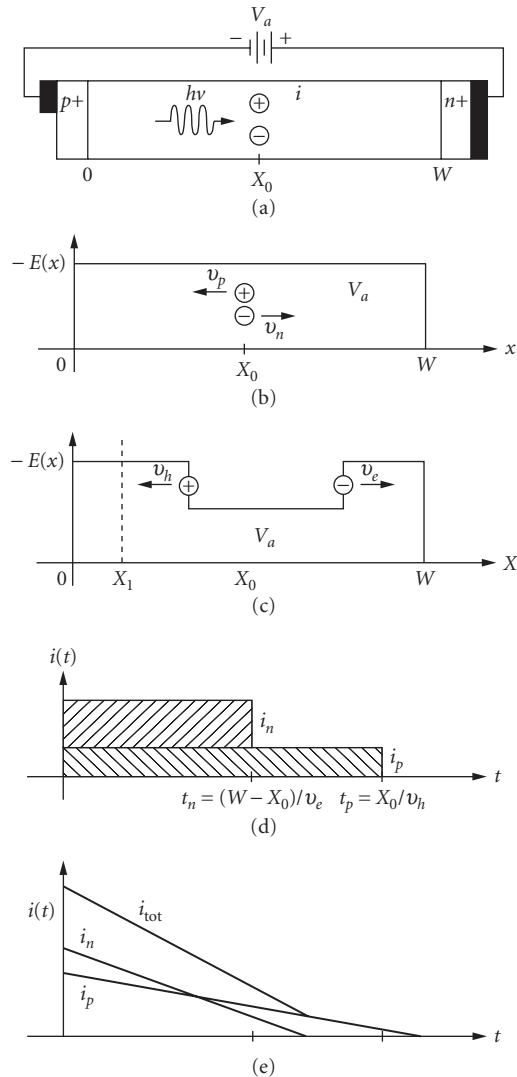


FIGURE 2 (a) Biased p - i - n structure. (b) Electrical field at the time when the electron-hole pairs are generated. (c) The perturbed electrical field due to the separated electron-hole pair. (d) The photocurrent due to single electron-hole pair. (e) The total photocurrent due to uniform illumination across the photodiode.

displacement current and reduces the internal electric field (Fig. 2b and c), which is the cause of saturation in photodetectors. From Gauss's law the difference in the electric field at the position of electron or hole is

$$\Delta E = \frac{q\sigma}{\epsilon} \quad (1)$$

where q is the electron charge, and ϵ is the permittivity. Due to the constant total voltage across the depletion region, the reduced electric field between the electron and hole will be compensated by the increased electric field outside. The rate of change of the electric field at the position $X = X_1$ is

$$\frac{\partial E}{\partial t} = -\frac{(v_e + v_h)\Delta E}{W} \quad (2)$$

where v_e and v_h are the saturation velocities for electrons and holes, respectively. The assumption of saturation velocities is valid at high fields. The displacement current is hence given by

$$i(t) = -\epsilon A \frac{\partial E}{\partial t} = \frac{qAv_e\sigma}{W} + \frac{qv_h\sigma A}{W} \quad (3)$$

The current consists of two components due to the electron and hole currents. The electron current lasts for a time duration of $(W - X_0)/v_e$ and hole current of X_0/v_h . This is shown in Fig. 2d. Here, we note that if the fast carrier (i.e., electron) travels a longer distance, then we have a shorter pulse. The total electron and hole currents are given by

$$i_e(t) = \frac{qv_e A}{W} \int_0^w n(x, t) dx \quad (4)$$

$$i_h(t) = \frac{qv_h A}{W} \int_0^w p(x, t) dx \quad (5)$$

where $n(x, t)$ and $p(x, t)$ represent the electron and hole densities in the depletion region. The total current is the sum of Eqs. (4) and (5).

RC Time Constant

The RC time constant is determined by the equivalent circuit parameters of photodiode. For example, the intrinsic response of the p - i - n diode can be modeled as a current source in parallel with a junction capacitor. The diode series resistance, parasitic capacitance, and load impedance form the external circuit. Figure 3 shows the equivalent circuit of the p - i - n photodiode. The junction capacitance is defined by the edge of the depletion region (or space charge region). The series resistance is due to the ohmic contacts and bulk resistances. In addition, the parasitic capacitance depends on the metallization

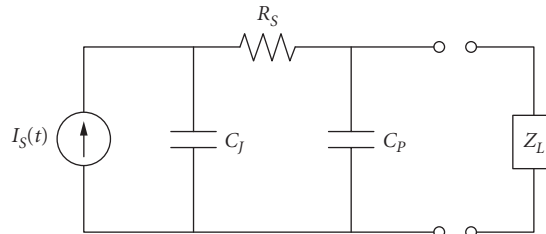


FIGURE 3 Equivalent circuit of a photodiode.

geometry. If the diode series resistance is R_s and a load resistance R_L is used to terminate the device, then the electrical 3-dB bandwidth can be approximated as

$$f_{3\text{dB}} = \frac{1}{2\pi(C_J + C_p)(R_L + R_s)} \quad (6)$$

If the photodiode is bonded by a section of gold wire, additional series inductance will be included in the load impedance. The 3-dB bandwidth due to parasitics in this case is then given in Ref. 3.

Diffusion Current

Diffusion current is important in detectors in which significant absorption occurs in regions outside the high-field region. This effect is reduced to some extent by recombination in these highly doped contact layers. Those carriers within about one diffusion length of the depletion region will have a chance to diffuse into the active region. This diffusion current will contribute a slow tail to the detector impulse response (Fig. 4). The electron diffusion current at the edge of the depletion region is given by

$$J_e = qD_e \frac{\partial n}{\partial x} = qD_e \frac{\Delta n}{L_e} \quad (7a)$$

and

$$J_h = -qD_h \frac{\partial p}{\partial x} = qD_h \frac{\Delta p}{L_h} \quad (7b)$$

where D_e (D_h) and L_e (L_h) are the diffusion coefficient and diffusion length, respectively, for electrons (holes). The diffusion process is a relatively slow process compared with the drift process. Assuming the photocarrier density is n , with the Einstein relation and Eq. (7a), the electron diffusion current can be written as

$$J_{\text{DIFF}} = qn\mu_e \left(\frac{kT/q}{L_e} \right) \quad (8)$$

and the drift current term for electron can be written as

$$J_{\text{DRIF}} = qn\mu_e \mathbf{E} \quad (9)$$

where μ_e is the electron mobility and \mathbf{E} is the electric field. For most devices the electric field inside the depletion region is an order of magnitude larger than $(kT/q)L_e$. For example, the hole diffusion

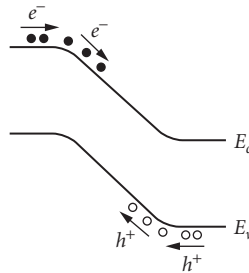


FIGURE 4 The origins of diffusion current.

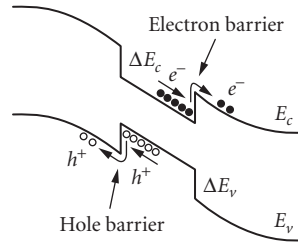


FIGURE 5 Heterostructure carrier trapping effect.

length for GaAs is typically around 10 μm and electric field is very often over 10 kV/cm. However, the diffusion-current terms could last as long as the carrier lifetime and the charge content in the tail can be as large as the drift component due to the slow diffusion times. For high-speed detectors, the diffusion-current problem can be eliminated with a double-heterostructure design that limits the absorbing regions to the high-field intrinsic regions.²⁴

Carrier Trapping

Heterojunctions in photodetectors cause carrier trapping of electrons at conduction band discontinuities and trapping of holes at valence band discontinuities (Fig. 5). Hole trapping is a significant problem in long-wavelength photodetectors because of the large valence band discontinuity at the InP/InGaAs heterojunction. Usually, the emission rate is approximated by thermionic emission. If the interface deep-level recombination rate is significant, the total emission rate will be the sum of the two emission rates. The emission functions for electrons and holes are given by

$$e_n(t) = (1/\tau_e) \exp(-t/\tau_e) u(t) \quad (10a)$$

$$e_p(t) = (1/\tau_h) \exp(-t/\tau_h) u(t) \quad (10b)$$

where τ_e (τ_h) represents the emission time constant for electron (hole) and $u(t)$ is the step function. The rates of thermionic emission of trapped carrier are related to the Schottky barrier height due to the bandgap discontinuity:

$$1/\tau_e = B \exp(-\phi_{bc}/kT) \quad (11)$$

where B is a constant and ϕ_{bc} is the barrier height for the conduction band. The response of the carrier-trap current in time domain is often obtained by convolving an intrinsic current source with the emission function. Since the applied bias will reduce the barrier height, sufficient device bias therefore will increase the emission rate. In order to reduce the barrier height, superlattice or compositional grading is often added at the heterointerface.²⁴

Packaging

The external connections to the photodetector often limit the detector performance. Another problem is that the photodiode is a high impedance load and the device has a reflection coefficient close to unity. One solution to this problem is to integrate a matching resistor with the device.²⁴ This can usually be added using the lower contact layer without adding any additional mask or process steps. Figure 6 shows a Smith chart plot of the impedance of a typical photodetector along with the impedance of a device with an integrated matching resistor. A good match up to 40 GHz

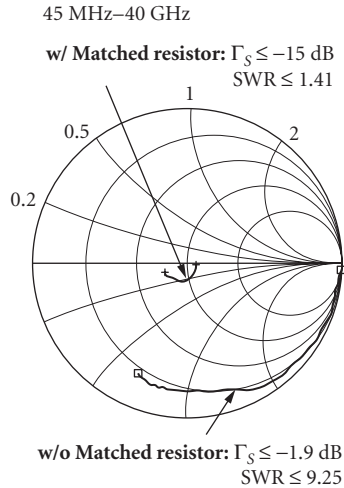


FIGURE 6 Smith chart of the impedance of typical photodiodes and photodiodes with integrated matching resistors.

is achieved. The disadvantage of a load resistor is the reduction in effective quantum efficiency by a factor of two since half of the photocurrent goes through the matching resistor. However, since the load resistance is now one-half, the RC time constant is also cut in half. Bandwidths in excess of 100 GHz have been achieved with quite large devices ($7 \times 7 \mu\text{m}^2$) in this way. A second problem with very high speed devices is the difficulty in building external bias circuits without resonances in the millimeter range. The necessary bias capacitor and load resistor can be integrated with *p-i-n* photodetector without adding any additional mask or process steps by using a large-area *p-i-n* region as the capacitor and using the lower contact layer as the series resistor.²³ A photograph of the device is shown in Fig. 7 along with the device performance. In this case, bandwidths in excess of 100 GHz were achieved.

Optical fiber alignment and packaging are now quite standard. Simplified alignment by means of holes etched in the substrate of back-illuminated photodiodes may allow passive alignment of optical fibers. The photodetectors must be antireflection coated to reduce the reflection to air or optical epoxy. Single dielectric layers are typically used to minimize the reflection at one wavelength. Braun et al.²⁵ have achieved minimum reflectivity at multiple wavelengths with one dielectric layer by using one of the semiconductor layers in a multiple antireflectivity design.

26.5 P-I-N PHOTODETECTORS

Vertically Illuminated *p-i-n* Photodiode

In order to increase the frequency response of the vertically illuminated *p-i-n* photodiode, the efficiency is always sacrificed. As the active layer thickness is reduced, the transit time decreases, and the optical absorption decreases, and there is a trade-off between the efficiency and speed. The external quantum efficiency for a surface-illuminated *p-i-n* diode is given by

$$\eta = (1 - R) \times (1 - e^{-\alpha L}) \quad (12)$$

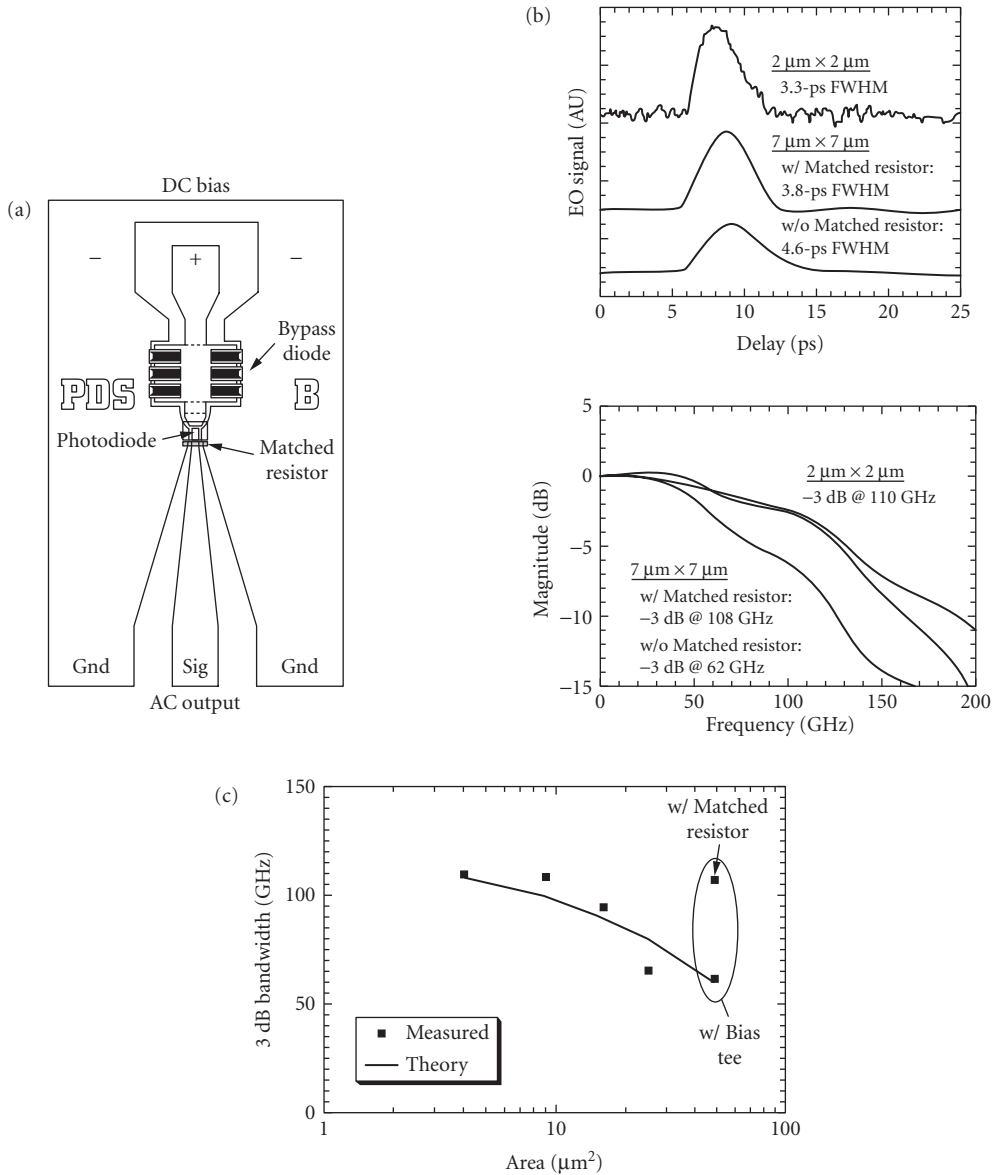


FIGURE 7 (a) Schematic diagram of a $p-i-n$ photodiode with integrated matching resistor and bias circuit. (b) Impulse response of a $2 \times 2 \mu\text{m}$ pin detector compared to $7 \times 7 \mu\text{m}$ detectors with and without matching resistors. (After Ref. 24.) (c) Dependence of measured bandwidth on detector area and comparison to the theory presented in the text.

where R is the surface reflection, α is the absorption coefficient, and L is the active layer thickness. Since the absorption coefficient is a function of wavelength $\alpha = \alpha(\lambda)$, usually α decreases as λ increases. Thus, the diode intrinsic response is wavelength dependent. We can easily see the effect of light absorption on the transit-time-limited bandwidth by comparing the transit-time response^{3,25} for two limiting cases for $\alpha L \rightarrow 0$ and $\alpha L \rightarrow \infty$ when $t_e = t_h = \tau_r$. The transit-time frequency response for a uniformly illuminated detector is

$$|F(\omega)_{\alpha L \rightarrow 0}| = \frac{2}{\omega \tau_{tr}} \left[1 + \frac{\sin^2\left(\frac{\omega \tau_{tr}}{2}\right)}{\left(\frac{\omega \tau_{tr}}{2}\right)^2} - 2 \frac{\sin(\omega \tau_{tr})}{(\omega \tau_{tr})} \right]^{1/2} \quad (13)$$

For electron-hole pairs generated near the p side of the intrinsic region, electrons travel across the i region, and the frequency response is given by

$$|F(\omega)_{\alpha L \rightarrow \infty}| = \left| \frac{\sin\left(\frac{\omega \tau_{tr}}{2}\right)}{\left(\frac{\omega \tau_{tr}}{2}\right)} \right| \quad (14)$$

For these two limits, the transit-time-limited bandwidths are $f_{3dB(\alpha L=0)} = 0.45/\tau_{tr}$ and $f_{3dB(\alpha L \rightarrow \infty)} = 0.55/\tau_{tr}$, respectively. For long-wavelength high-speed $p-i-n$ diodes, the absorption layer is often very thin so that $1 - \exp(-\alpha L) \approx \alpha L$. The bandwidth efficiency product for transit-time-limited $p-i-n$ diode is given by³

$$\eta f_{3dB} = 0.45 \alpha v_s \quad (15)$$

Figure 8 shows the calculated 3-dB bandwidth for GaInAs/InP $p-i-n$ diodes on the device area versus thickness plane for wavelength $\lambda = 1.3 \mu\text{m}$. The horizontal axis is the active layer thickness (which

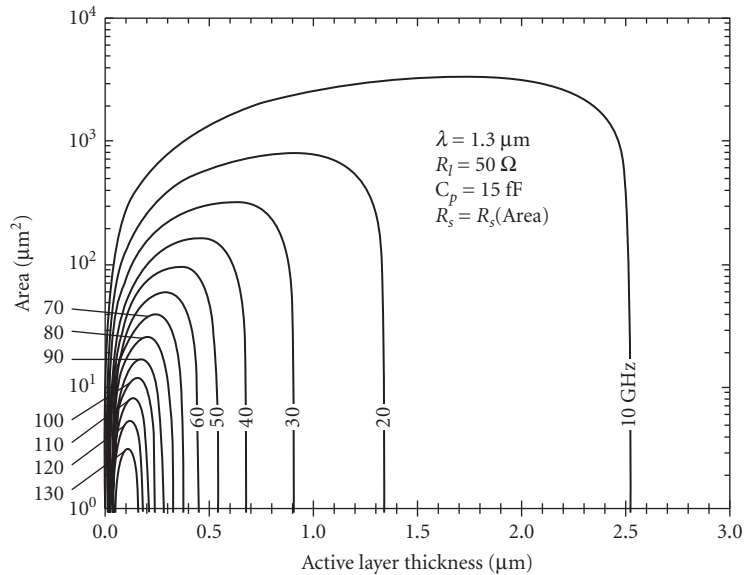


FIGURE 8 Calculated 3-dB bandwidth contours for a GaInAs pin vertically illuminated photodiode.

corresponds to the quantum efficiency) and the vertical axis is the device area. As we can see, when the device active layer thickness decreases, the quantum efficiency of p - i - n diode also decreases due to the insufficient light absorption in the active layer. The capacitance decreases as the device area decreases. Optimization of the device bandwidth is reached when transit-time-limited bandwidth approximately equals the RC limited bandwidth.

To minimize the bonding-pad capacitance, a semi-insulating substrate and thick polyimide layer are often used. Sometimes the series inductance of the bond wire is used to resonate the parasitic capacitance, and this results in a slightly peaked response with an increased 3-dB corner frequency. The electrical transfer function with series inductance is given by

$$H(\omega) = \frac{R_L}{[1 - \omega^2(R_S R_L C_J C_P + L_S(C_J + C_P))] - j[\omega(R_L(C_J + C_P) + R_S C_J) - \omega^3 R_S C_J C_P L_S]} \quad (16)$$

where L_S is the series inductance.

To achieve high detector bandwidth, *double heterostructure* InP/GaInAs/InP p - i - n photodiodes have been fabricated to reduce the diffusion-current problem. However, carrier trapping can limit the impulse response. This effect can be characterized by the emission function $e_{e,h}(t) = (1/\tau_{e,h}) \exp(-t/\tau_{e,h})$ where $\tau_{e,h}$ is emission time for electron (hole). The current-source response due to the electron and hole trapping at the heterointerfaces for p -side illumination is given by

$$\begin{aligned} \frac{i_s(\omega)}{i_s(0)} = & \frac{1}{(1 - e^{-\alpha L})} \left\{ \left(\frac{1 - e^{-j\omega\tau_e}}{j\omega\tau_e} - e^{-\alpha L} \frac{1 - e^{-\alpha L} e^{-j\omega\tau_e}}{j\omega\tau_e - \alpha L} \right) \left(\frac{1}{1 + j\omega\tau_e} \right) \right. \\ & \left. + \left(\frac{1 - e^{-\alpha L} e^{-j\omega\tau_h}}{j\omega\tau_e + \alpha L} - e^{-\alpha L} \frac{1 - e^{-j\omega\tau_h}}{j\omega\tau_h} \right) \left(\frac{1}{1 + j\omega\tau_h} \right) \right\} \quad (17) \end{aligned}$$

and for n -side illumination is given by

$$\begin{aligned} \frac{i_s(\omega)}{i_s(0)} = & \frac{1}{(1 - e^{-\alpha L})} \left\{ \left(\frac{1 - e^{-\alpha L} e^{-j\omega\tau_e}}{j\omega\tau_e + \alpha L} - e^{-\alpha L} \frac{1 - e^{-j\omega\tau_e}}{j\omega\tau_e} \right) \left(\frac{1}{1 + j\omega\tau_e} \right) \right. \\ & \left. + \left(\frac{1 - e^{-j\omega\tau_h}}{j\omega\tau_h} - e^{-\alpha L} \frac{1 - e^{-\alpha L} e^{-j\omega\tau_h}}{j\omega\tau_h - \alpha L} \right) \left(\frac{1}{1 + j\omega\tau_h} \right) \right\} \quad (18) \end{aligned}$$

where $\tau_{e,h}$ is the electron (hole) transit time. Other than the original p - i - n diode response, the extra terms $1/(1 + j\omega\tau_{e,h})$ are due to the trapping effect. For InGaAs/InP heterostructure p - i - n diodes, the valence band offset is larger than the conduction offset and the hole effective mass is much larger than the electron effective mass. Thus, hole trapping is worse than the electron trapping in an InGaAs/InP p - i - n diode.

Waveguide p - i - n Photodiode

The main advantages of waveguide detectors are the very thin depletion region resulting in a very short transit time and the long absorption region resulting in a high bandwidth photodetector with a high saturation power (Fig. 1c). Due to the thin intrinsic layer, it can often operate at zero bias.²⁷ The absorption length of a waveguide detector is usually *designed to be long enough* ($>5 \mu\text{m}$) to ensure full absorption. The waveguide structure design (Fig. 9) is often required to have low coupling loss due to modal mismatch and reasonable effective absorption coefficient. The external quantum efficiency of a waveguide p - i - n detector is^{28,29}

$$\eta = \kappa(1 - R) \frac{\Gamma \alpha_{\text{IB}}}{\alpha} (1 - e^{-\alpha L}) \quad (19)$$

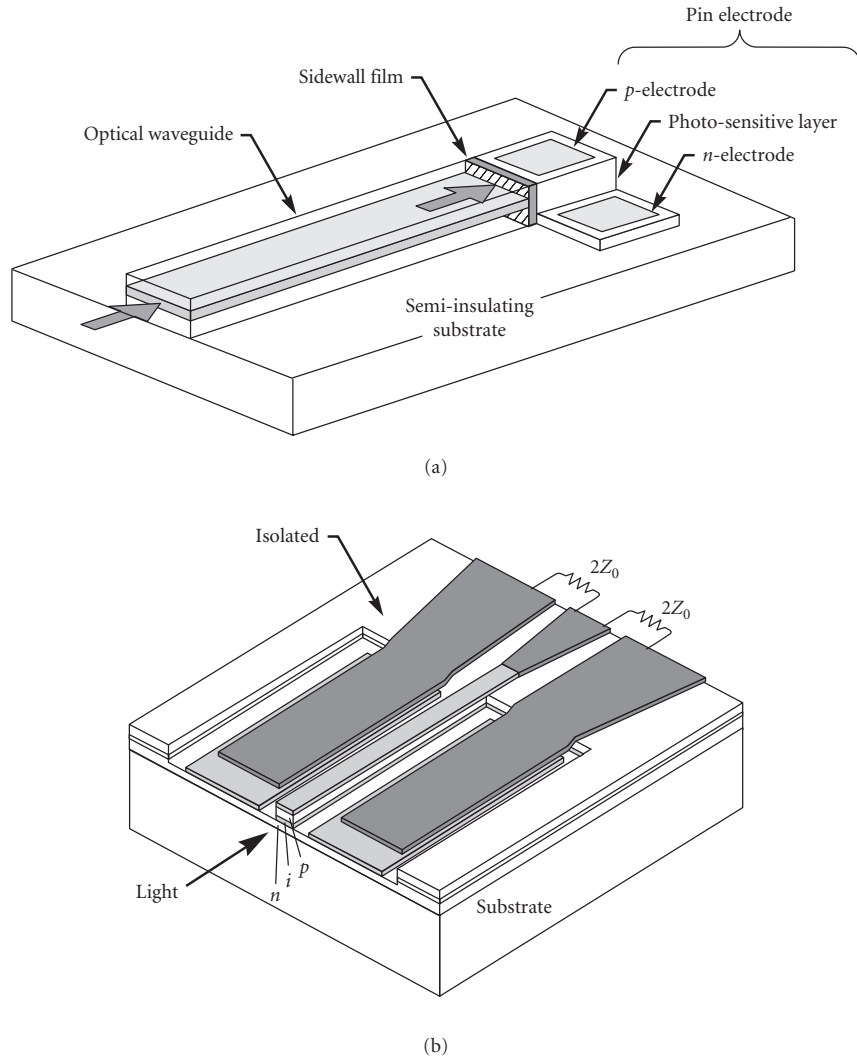


FIGURE 9 Schematic diagram of (a) a waveguide photodetector (after Ref. 14) and (b) a traveling wave photodetector (after Ref. 17).

where κ is the coupling efficiency due to the modal mismatch, Γ is the mode confinement factor, α_{IB} is the interband absorption. The loss coefficient α is given by

$$\alpha = \Gamma \alpha_{IB} + \Gamma \alpha_{FC} + (1 - \Gamma) \alpha_{FCx} + \alpha_s \quad (20)$$

where α_{FC} , α_{FCx} are the free carrier absorption loss inside and outside the absorption layer, α_s is the scattering loss. Kato et al.¹⁴ reported an InGaAs waveguide *p-i-n* diode with bandwidth of 40 GHz. The detector quantum efficiency is 44 percent at 1.55- μm wavelength. The coupling loss estimated by an overlap integral was 2.1 dB. To reduce the coupling loss, it is important to have a good design of the layer structure to reduce modal mismatch and to coat the facet with an antireflecting (AR) film.

Resonant p - i - n Photodiode

A resonant detector utilizes the multiple passes in a Fabry-Perot resonator to achieve high quantum efficiency with thin absorbing layers (Fig. 1b). Since the speed of light is about three orders of magnitude faster than the carrier velocities, the quantum efficiency can be increased without significant pulse broadening due to the effective optical transit time.

The schematic diagram of a resonant cavity-enhanced photodetector is shown in Fig. 10. The efficiency of the resonant detector is given by

$$\eta = \left[\frac{(1 + R_2 e^{-\alpha d})}{1 - 2\sqrt{R_1 R_2} e^{-\alpha d} \cos(2\beta L + \phi_1 + \phi_2) + R_1 R_2 e^{-2\alpha d}} \right] \times (1 - R_1) \times (1 - e^{-\alpha d}) \quad (21)$$

where R_1, R_2 are mirror reflectivities, ϕ_1, ϕ_2 are mirror phase shifts, $\beta = 2\pi/n\lambda$ is the propagation constant, and d is the thickness of active region. The quantum efficiency has its maximum when $2\beta L + \phi_1 + \phi_2 = 2m\pi$ ($m=1, 2, 3, \dots$) and the quantum efficiency is then

$$\eta = \left[\frac{(1 + R_2 e^{-\alpha d})}{(1 - \sqrt{R_1 R_2} e^{-\alpha d})^2} \right] \times (1 - R_1) \times (1 - e^{-\alpha d}) \quad (22)$$

Figure 11 shows the calculated resonant quantum efficiency versus normalized absorption coefficient αd .¹¹ High quantum efficiency is possible even from thin absorption layers.

However, in terms of the fabrication, the material growth, and the structure design, building a high-speed resonant photodetector is not a simple task. The required low resistance and low capacitance with incorporated mirror structure is difficult to achieve due to the significant resistance of the multiple heterojunction mirror stack.

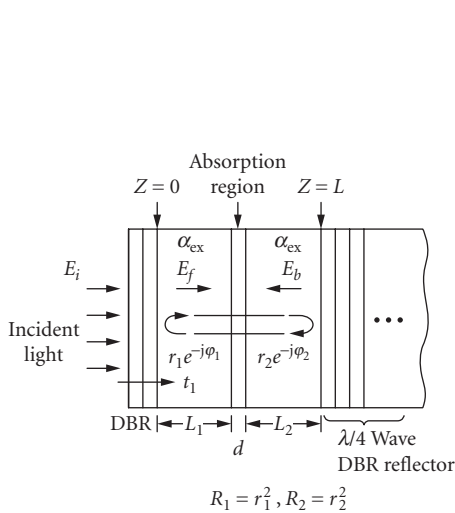


FIGURE 10 Schematic diagram of a resonant photodetector. (After Ref. 13.)

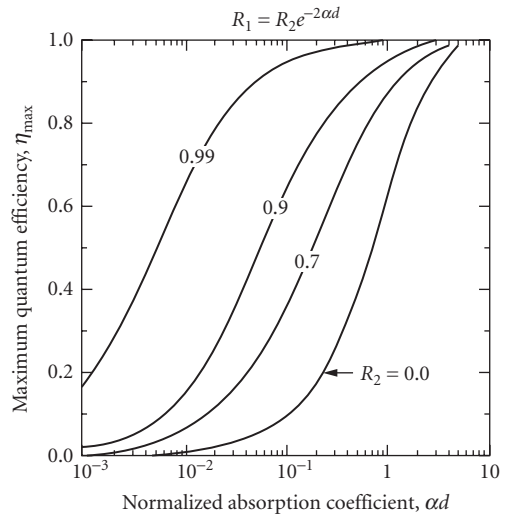


FIGURE 11 Dependence of quantum efficiency on mirror design in a resonant photodetector. (After Ref. 13.)

26.6 SCHOTTKY PHOTODIODE

Schottky photodiodes^{30–32} are especially attractive for integration with FETs and III-V integrated circuits because of their simple material structure and easy fabrication. Figure 12 shows the Schottky barrier structure. For front-illuminated devices, the metal is very thin so that the light can penetrate the metal with very little loss. The J - V characteristic of a Schottky diode is given by³³

$$J = J_0 \left[\exp\left(\frac{qV}{k_B T}\right) - 1 \right] \quad (23)$$

where

$$J_0 = A^{**} T^2 \exp\left(-\frac{q\Phi_b}{k_B T}\right) \quad (24)$$

and ϕ_b is the barrier height and A^{**} is the modified effective Richardson constant.³⁴

The dynamics of photogenerated carriers in a Schottky diode are similar to those of a p - i - n diode (Fig. 2). The dynamics of both electrons and hole must be included in the analysis of a Schottky photodiode, resulting in expressions for the Schottky diode response given in Eqs. (13) to (18) with the exception that there is no diffusion current from the metal layer. The equivalent circuit of a Schottky diode is the same as that of a p - i - n diode.

In high-speed applications, GaAs Schottky diodes in the short-wavelength region with bandwidths over 200 GHz have been reported.³⁶ These devices can be combined with FETs or sampling diodes.^{35,36} Figure 13 shows an integrated Schottky photodiode with a diode sampling circuit. A pair of short voltage pulses are generated by the nonlinear transmission line and a differentiator. The short voltage pulses are used to control the sampling capacitors to measure the photodiode signal. The sampled signal is then passed through a low-pass filter to extract the equivalent time domain waveform. Impulse responses of under 2 ps have been demonstrated in this way.³⁶

In the long-wavelength region, GaInAs Schottky diodes experience high dark current problems due to the relatively low Schottky barrier height at the metal/GaInAs interface.³⁷ An InP or quaternary layer is usually added at the interface in order to increase the Schottky barrier height.³⁸ A graded bandgap layer (e.g., GaInAsP) is then needed at the GaInAs/InP interface to reduce hole trapping.

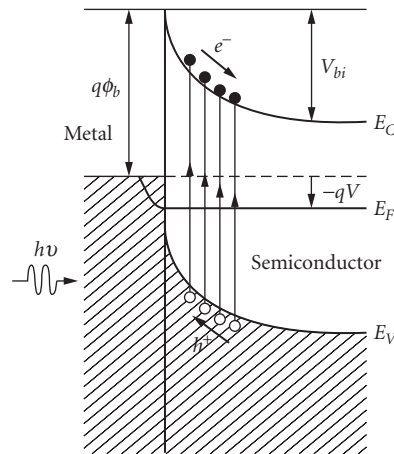


FIGURE 12 Schematic diagram of a Schottky barrier photodiode.

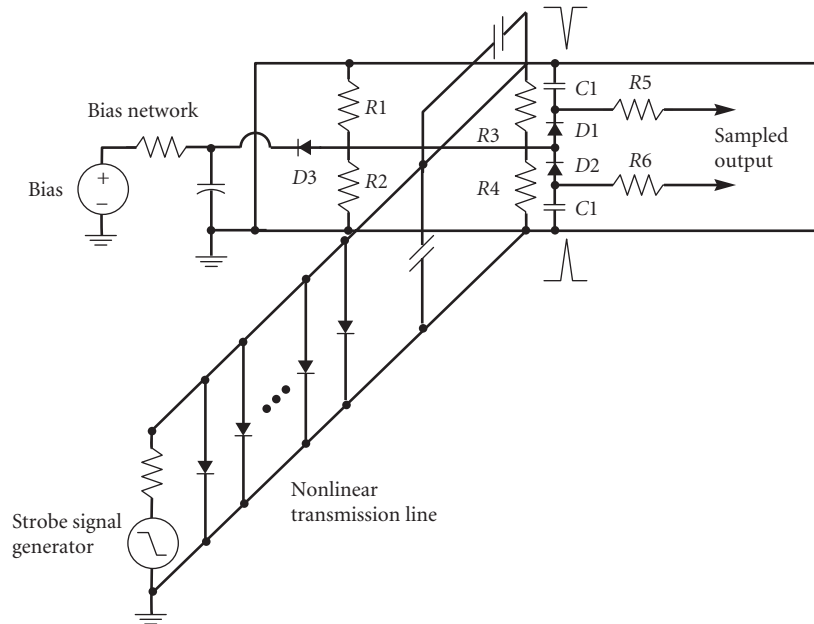


FIGURE 13 Integrated Schottky photodiode and sampling circuit. (After Ref. 35.)

26.7 AVALANCHE PHOTODETECTORS

High-speed avalanche photodiodes (APDs) are widely used in fiber communication. APDs with gain-bandwidth (GB) products in excess of 100 GHz have been reported.^{39–43} In long-wavelength applications, InGaAs/InP APDs are better than Ge APDs due to their lower dark current and lower multiplication noise. The Ge APD also has a limited spectral response at 1.55- μm wavelength. The maximum achievable GB product of InGaAs/InP APDs is predicted to be around 140 GHz,⁴⁴ while the gain-bandwidth product of Si APDs in the near-infrared region can have GB products of over 200 GHz.⁴⁵ InGaAs/InAlAs superlattice avalanche photodiodes have a lower ionization ratio ($k = 0.2$)⁴⁶ than bulk avalanche photodiodes, and lower noise and higher gain-bandwidth product can be achieved.

High-speed GaInAs/InP APDs make use of separated absorption and multiplication layers (SAM APD). Figure 14 shows the simplified one-dimensional APD structure. The narrow bandgap n -GaInAs layer absorbs the incident light. The layer is usually thick ($>1 \mu\text{m}$) to ensure high quantum efficiency. The electric field in the absorption layer is high enough for carriers to travel at saturated velocities, yet is below the field where significant avalanching occurs and the tunneling current is negligible. The wide bandgap InP multiplication layer is thin (a few tenths of a micron) to have shorter multiplication buildup time.^{47,48} The bias is applied to the fully depleted absorption layer in order to obtain effective carrier collection efficiency and, at the same time, electric field in the multiplication region must be high enough to achieve avalanche gain. A guard ring is usually added to prevent premature avalanche breakdown (or microplasma) at the corner of the diffusion edge. To reduce the hole pileup effect, a graded bandgap layer (e.g., superlattice or compositional grading) is often added at the heterointerface between the absorption layer and multiplication layer. This is the so-called separated absorption, grading, multiplication avalanche photodiode (SAGM APD).

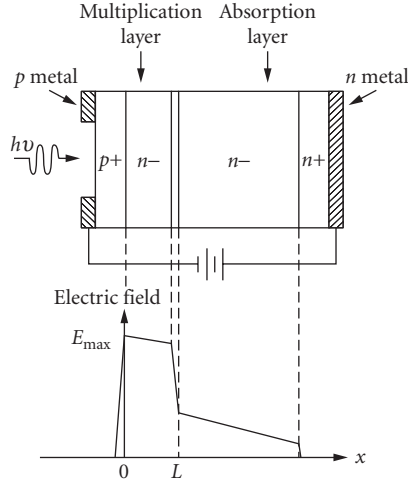


FIGURE 14 Schematic diagram of a SAGM APD.

The multiplication process in APDs can be described by the electron and hole ionization coefficients α_i and β_i . The field dependencies of ionization coefficients for electrons and holes are given by

$$\alpha_i(x) = A_e \exp(-B_e/E(x)) \quad (25a)$$

$$\beta_i(x) = A_h \exp(-B_h/E(x)) \quad (25b)$$

where $A_{e,h}$ and $B_{e,h}$ are constant parameters.⁴⁹ Since the electric field is generally position-dependent, the ionization coefficients are also position-dependent. With Eq. (25a, b) and the electric field distribution, the position-dependence of ionization coefficients can be derived. The multiplied photocurrent in the avalanche region ($0 \leq x \leq W$) including injected electron current density $J_n(0)$, injected hole current density $J_p(0)$, and photo-generation of electron-hole pairs $g(x)$ was derived by Lee et al.⁵⁰ The total photocurrent density is given by

$$J = \frac{J_p(w) \exp\left[-\int_0^w (\alpha_i - \beta_i) dx\right] + J_n(0) + q \int_0^w g(x) \exp\left[-\int_0^x (\alpha_i - \beta) dx'\right] dx}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (26)$$

The electron-initiated and hole-initiated multiplication factors, M_n and M_p , can be obtained by putting $J_p(w) = g(x) = 0$ and $J_n(0) = g(x) = 0$, respectively in Eq. (26):

$$M_n = \frac{J}{J_n(0)} = \frac{1}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (27a)$$

$$M_p = \frac{J}{J_p(w)} = \frac{\exp\left[-\int_0^w (\alpha_i - \beta_i) dx\right]}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (27b)$$

The bandwidth of an APD is limited by the device RC time constant when the multiplication gain M is low (i.e., $M < \alpha_i/\beta_i$). As the multiplication gain increases above the ratio of the electron and hole ionization coefficients (i.e., $M > \alpha_i/\beta_i$), the avalanche buildup time becomes the dominant limitation on 3-dB bandwidth and the product of the multiplication gain and 3-dB bandwidth reaches a constant. The multiplication factor M as a function of frequency was derived by Emmons⁵¹ and is given by

$$M(\omega) \approx \frac{M_o}{\{1 + \omega^2 M_o^2 \tau_1^2\}^{1/2}}, \quad M_o > \alpha_i/\beta_i \quad (28a)$$

$$\tau_1 \approx N(\alpha_i/\beta_i)\tau \quad (28b)$$

where τ_1 is the effective transit time, τ is the multiplication-region transit time, and $N(\beta_i/\alpha_i)$ is a number changing between 1/3 and 2 as β_i/α_i varies from 1 to 10^{-3} . The dc multiplication factor M_o is given by Miller.⁵²

$$M_o = \frac{1}{1 - (V_j/V_B)^n} \quad (29)$$

where V_B is the breakdown voltage, V_j is the junction voltage, and n is an empirical factor ($n < 1$).

The total APD dark current consists of two components. I_{du} is the unmultiplied current which is mainly due to the surface leakage current. I_{dm} is the bulk dark current experiencing the multiplication process. The total dark current is expressed by

$$I_d = I_{du} + MI_{dm} \quad (30)$$

where M is the avalanche gain. The noise current spectral density due to the dark current is given by

$$\sigma_d^2 = 2qI_{du} + 2qI_{dm}M^2F(M) \quad (31)$$

where $F(M)$ is the avalanche excess noise factor derived by McIntyre.⁵³ Excess noise factors for electron-initiated or hole-initiated multiplication are given by

$$F(M) = F_e(M) = [kM + (1-k)(2-1/M)] \quad (32a)$$

$$F(M) = F_h(M) = \left[\frac{1}{k}M + \left(1 - \frac{1}{k}\right)(2-1/M) \right] \quad (32b)$$

where k is the ratio of the ionization coefficient of holes to electrons ($k = \beta_i/\alpha_i$), and k is assumed to be a constant independent of the position. Figure 15 shows the sensitivity of an APD receiver as a function of k_{eff} which is obtained by weighting the ionization rates over the electric field profile. From Fig. 15 we can see that the smaller the k factor is, the smaller the noise factor is and the better the receiver sensitivity is. Ge APDs have k values close to unity (0.7 to 1.0). GaInAs/InP APDs using an InP multiplication region have $1/k$ values from 0.3 to 0.5. Silicon is an excellent APD material since its k value is 0.02. Therefore, a Si APD has an excellent low dark current noise density and is predominantly used at short wavelengths compared with Ge APDs and GaInAs/InP APDs which are used at longer wavelengths.

In optical receiver applications,^{54,55} the photodetector is used with a low-noise amplifier. The dark current noise power is given by

$$\langle i_{nd}^2 \rangle = 2qI_{du}BI_2 + 2qI_{dm}M^2F(M)BI_2 \quad (33)$$

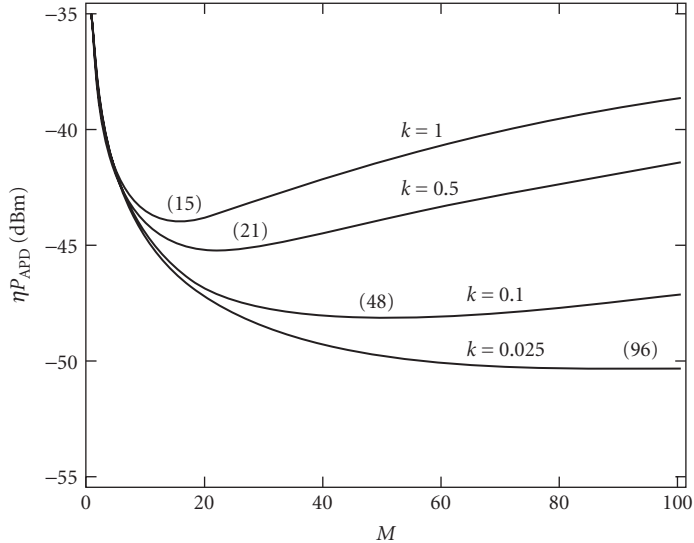


FIGURE 15 Dependence of APD receiver sensitivity on β/α in SAGM APDs. (After Ref. 7.)

where B is the receiver bit rate and I_2 is a parameter depending on the input optical pulse shape. The receiver sensitivity penalty⁵⁶ is given in terms of parameter ε_N ,

$$\bar{\eta P} = (1 + \varepsilon_N) \bar{\eta P}_o \quad (34)$$

where $\bar{\eta P}_o$ is the sensitivity with zero dark current. For example, $\varepsilon_N = 0.023$ for a 0.1-dB penalty. The maximum allowable dark current for a given sensitivity for a p - i - n FET receiver is

$$I_{du} = \frac{\varepsilon_N (2 + \varepsilon_N)}{2qBI_2} \langle i_{na}^2 \rangle \quad (35)$$

where $\langle i_{na}^2 \rangle$ is the amplifier noise power and is proportional to B^3 above 100 MBits/s. Therefore, the maximum allowable dark current is proportional to B^2 . For APD receivers, the maximum allowable dark current I_{dm} as a function of bit rate can be approximated by assuming that the sensitivity penalty is within 1 or 2 dB and optimum gain is constant. In Fig. 16, as we can see, the dark current is proportional to B at lower bit rates and $B^{1.25}$ at higher bit rates. So, as the bit rate increases, the maximum allowable dark current increases.

26.8 PHOTOCONDUCTORS

High-speed photoconductors⁵⁷⁻⁶² have become more important not only because of their simplicity in fabrication and ease of integration with MESFET amplifiers but also because of their useful applications for photodetector and photoconductor sampling gates. Usually the photoconductive film has a high density of defects with the trap energy levels deep within the bandgap to shorten the material lifetime and the detector impulse response. The characteristics of the photoconductive films include (1) high resistivity due to the fact that Fermi level is pinned at the midgap, (2) enhanced optical absorption for photon energy below the bandgap due to the introduction of new bandgap states, and (3) easy fabrication of ohmic contacts possibly due to the enhancement of tunneling through the narrow Schottky barrier with a pinned Fermi level.

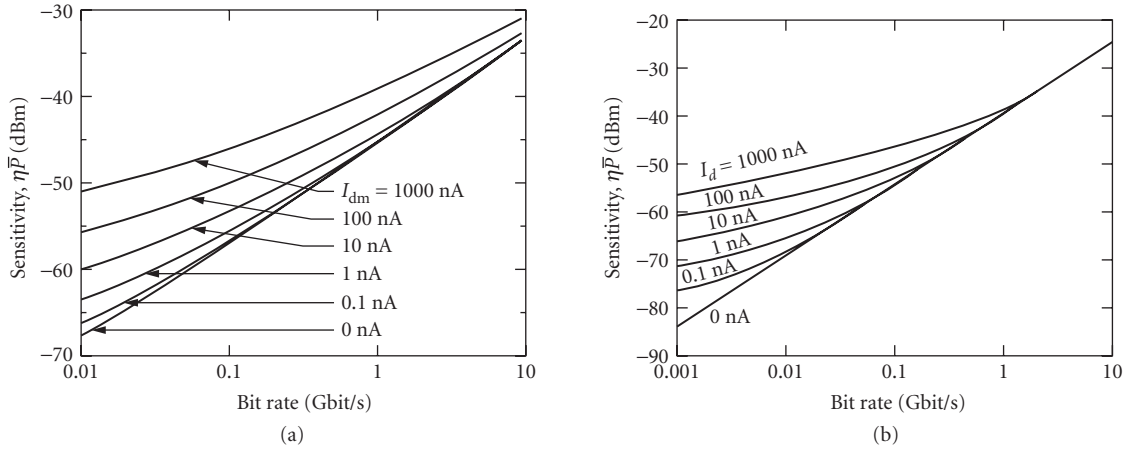


FIGURE 16 Dependence of receiver sensitivity on dark current for an (a) SAGM APD and (b) PIN detector, each with a GaAs FET preamplifier. (After Ref. 6.)

Figure 17 shows a typical photoconductor on a microstrip line structure. The photoconductive film is formed on top of a semi-insulating substrate. A microstrip transmission line consists of microstrip electrodes on top and ground plane on bottom. Under a steady-state illumination, the photogenerated carrier will experience high electrical field and travel to the electrodes. The photocurrent is

$$I_{\text{ph}} = \frac{q\eta GP}{h\nu} \quad (36)$$

where η is the external quantum efficiency, G is the photoconductor gain, and P is optical input flux. The photoconductor gain G is given by

$$G = \frac{\tau}{\tau_{\text{tr}}} \quad (37)$$

which is the ratio of carrier lifetime τ to the carrier transit time τ_{tr} . The frequency response of a photoconductive detector is plotted in Fig. 18 for different material lifetimes. In a detector without damage sites, the gain can be quite large at low frequencies. We can see from this figure how the

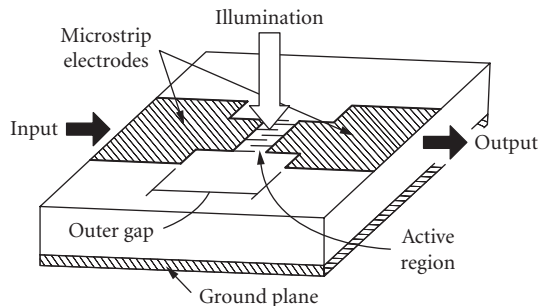


FIGURE 17 Schematic diagram of a high-speed photoconductor.

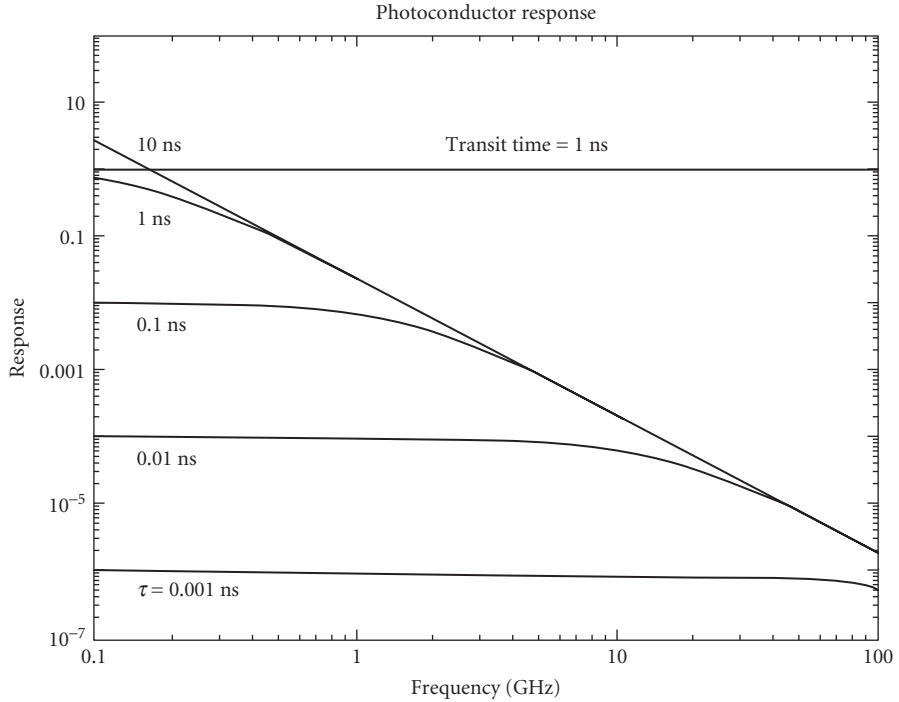


FIGURE 18 Frequency response of a photoconductor.

increased bandwidth is achieved at the expense of quantum efficiency. Using smaller finger separations, higher quantum efficiency can be achieved for a particular bandwidth.

The standard microstrip line configuration has reflection problems in the thickness direction of the substrate and the dispersion characteristics of a microstrip line is worse than that of a coplanar stripline.⁵⁷ Coplanar striplines with “sliding-contact” excitation can have zero capacitance to first order.⁵⁸ The photoconductor using coplanar stripline has been very successful in generating short electrical pulses. To measure the short electric pulse, several techniques can be used such as photoconductor sampling or electro-optic sampling. Both of the above techniques can provide subpicosecond resolution. The coplanar strip line configuration with sliding contact and sampling gate is shown in Fig. 19a. The equivalent circuit is shown in Fig. 19b.⁵⁸ The infinite capacitances represent that the line extends without end in both directions. The generated electrical signal due to the time-varying resistance $R_s(t)$ is

$$V_{\text{out}}(t) = V_b \frac{Z_o}{Z_o + R_s(t) + R_c} \quad (38)$$

where R_c is the contact resistance. If the excitation intensity is sufficiently low to keep $R_s(t) \gg Z_o$, then

$$V_{\text{out}}(t) = V_b \frac{Z_o}{R_s(t) + R_c} \quad (39)$$

The photoconductor resistance $R_s(t)$ can be related to the photoexcited electron-hole pair density $n(t)$:⁵⁹

$$R_s(t) = \frac{L}{qn(t)(\mu_e + \mu_h)wd_e} \quad (40)$$

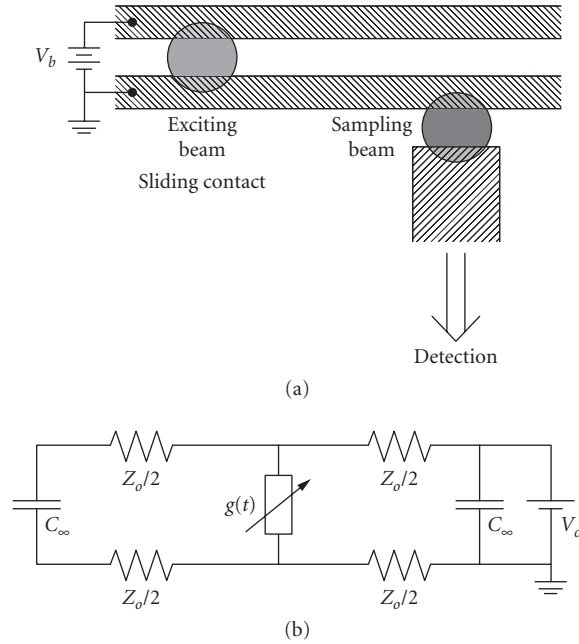


FIGURE 19 (a) Coplanar circuit layout of a photoconductor with sliding contact and (b) Equivalent circuit.

where L is the gap width, w is the width of photoconductive volume and d_e is the effective absorption length. When the pulse width is of the same order of magnitude as carrier lifetime and much shorter than the carrier transit time across the switch gap, the electron-hole pair density is given by

$$n(t) = e^{-t/\tau} \int_0^t e^{t'/\tau} \frac{\eta P_o(t')(1-R)}{h\nu w d_e} dt' \quad (41)$$

Here, we notice that the carrier density is an exponential decay function, so $G(t) = 1/R_s(t)$ is also an exponential decaying function with a time constant τ . This can be explained as a result of the convolution of the laser pulse with an exponential function with carrier life time τ .

The low temperature (LT) grown GaAs⁶⁰ can have both high carrier mobility and subpicosecond carrier lifetime when being compared with that of the ion-implanted photoconductor.⁶¹ The dislocation density in LT GaAs is about the same as that in GaAs epitaxial layer grown at normal substrate temperatures such that the LT GaAs has a mobility as high as that of the bulk material. The resistivity of the LT GaAs is greater than that of semi-insulating GaAs ($>10^{17} \Omega\text{-cm}$) due to its high deep-level concentration. The LT GaAs photoconductive-gap switch in a coplanar strip transmission line configuration has obtained a 1.6-ps (FWHM) response with a 3-dB bandwidth of 220 GHz. Chen et al.⁶² reported a high-speed photodetector utilizing LT GaAs MSM photoconductor. To achieve reasonable quantum efficiency and high-speed response, the optimum design requires carrier transit time approximately equal to carrier lifetime. With this requirement, the carriers not collected fast enough by the electrodes will be consumed by recombination. The response of a 0.2- μm finger and space MSM photodetector was measured by electro-optic sampling system. A 1.2 ps (FWHM) response with a 3-dB bandwidth of 350 GHz is obtained.

26.9 SUMMARY

Photodetector performance has steadily improved over the past decade. High-speed detectors are now available at a variety of wavelengths from 1.65 to 0.4 μm . MSM photoconductors have demonstrated the shortest impulse responses of under a picosecond. For applications that require high speed and high efficiency, the best results have been obtained using two passes through a p - i - n photodetector (30 percent quantum efficiency with 110-GHz bandwidth). Many applications require a high saturation power, and waveguide photodetectors have achieved the best results (20-GHz bandwidth with 0.5-A/W responsivity and 10-mW saturation power). Traveling wave photodetectors appear to offer the ultimate results in high-speed, high-responsivity, high-saturation power detectors. The combination of high-speed photodetectors with optical amplifiers is resulting in superb sensitivity of all bit rates, but requires the fabrication of high-speed photodetectors with at least 10-dBm saturation power.

An increasing amount of attention is being paid to integrating high-speed photodetectors with electronic and photonic circuits. Integration with electronic circuits increases the performance by eliminating the parasitics and limited bandwidth of bonding pads, wires, and connectors. Integration with optical waveguides decreases the optical loss associated with coupling from one device to another and reduces the packaging cost. Integration of photodetectors with optical amplifiers and wavelength tuning elements is a particularly important research direction.

26.10 REFERENCES

1. T. P. Lee and T. Li, "Photodetectors," *Optical Fiber Telecommunications*, S. E. Miller and A. G. Chynoweth (eds.), Academic Press, New York, 1979.
2. D. H. Auston, *Picosecond Optoelectronic Devices*, chap. 4, C. H. Lee (ed.), Academic Press, New York, 1984, pp. 73–117.
3. J. E. Bowers and C. A. Burrus, Jr., "Ultrawide-Band Long-Wavelength p - i - n Photodetectors," *J. Lightwave Tech.*, vol. LT-5, no. 10, October 1987.
4. F. Capasso, "Physics of Avalanche Photodiodes," *Semiconductors and Semimetals*, 22D, W. T. Tsang (ed.), Academic Press, New York, 1985.
5. J. Campbell, *Semiconductors and Semimetals*, 22D, W. T. Tsang (ed.), Academic Press, New York, 1985.
6. B. Kasper, "Receiver Design," *Optical Fiber Telecommunications II*, S. E. Miller and I. P. Kaminow (eds.), Academic Press, Boston, 1988.
7. S. R. Forrest, "Avalanche Photodetector Receiver Sensitivity," *Semiconductors and Semimetals*, vol. 22D, *Lightwave Communications Technology: Photodetectors*, W. Tsang (ed.), Academic Press, New York, 1985.
8. S. R. Sloan, "Processing and Passivation Techniques for Fabrication of High Speed InP/InGaAs/InP Mesa Photodetectors," *Hewlett Packard Journal*, Oct. 1989, p. 69.
9. K. Carey, S. Y. Wang, J. S. C. Chang, and K. Nauka, "Leakage Current in GaInAs/InP Photodiodes Grown by OMVPE," *J. Crystal Growth*, vol. 98, 1989, p. 90.
10. Y. G. Wey, D. L. Crawford, K. Giboney, J. E. Bowers, M. J. Rodwell, P. Silvestre, M. J. Hafich, and G. Y. Robinson, "Ultrafast Graded Double-Heterostructure GaInAs/InP Photodiode," *Appl. Phys. Lett.*, vol. 58, no. 19, 1991, p. 2156.
11. M. S. Unlu, K. Kishino, J. Chyi, L. Aresenault, J. Reed, and S. N. Mohammad, "Resonant Cavity Enhanced AlGaAs/GaAs Heterojunction Phototransistors with an Intermediate InGaAs Layer in the Collector," *Appl. Phys. Lett.*, vol. 57, no. 8, 20 Aug. 1990, p. 750.
12. A. Chin and T. Y. Chang, "Enhancement of Quantum Efficiency in Thin Photodiodes through Absorptive Resonance," *J. Lightwave Tech.*, vol. 9, no. 3, March 1991, p. 321.
13. K. Kishino, M. S. Unlu, J.-I. Chyi, J. Reed, L. Aresenault, and H. Morkoc, "Resonant Cavity-Enhanced (RCE) Photodetectors," *IEEE J. Quantum Electron.*, vol. 27, no. 8, Aug. 1991, p. 2025.
14. K. Kato, S. Hata, A. Kozen, J. Yoshida, and K. Kawano, "High-Efficiency Waveguide InGaAs Pin Photodiode with Bandwidth of over 40 GHz," *IEEE Photon. Tech. Lett.*, vol. 3, no. 5, May 1991, p. 473.
15. D. Wake, S. N. Judge, T. P. Spooner, M. J. Harlow, W. J. Duncan, I. D. Henning, and M. J. O'Mahony, "Monolithic Integration of 1.5 μm Optical Pre-amplifier and PIN Photodetector with a Gain of 20 dB and a Bandwidth of 35 GHz," *Electron. Lett.*, vol. 26, no. 15, July 19, 1990, pp. 1166–1168.

16. R. J. Deri, N. Yasuoka, M. Makiuchi, H. Hamaguchi, O. Wada, A. Kuramata, and R. J. Hawkins, "Integrated Waveguide/Photodiodes with Large Bandwidth and High External Quantum Efficiency," *IEEE Photon. Technol. Lett.*, vol. 2, 1990, pp. 496–498.
17. K. S. Giboney, M. J. W. Rodwell, and J. E. Bowers, "Travelling-Wave Photodetectors," *Photon. Tech. Lett.*, vol. 4, no. 12, Dec. 1992, pp. 1363–1365.
18. H. Taylor, O. Eknoyan, C. S. Park, K. N. Choi, and K. Chang, "Traveling Wave Photodetectors," *SPIE Proc. on Optoelectronic Signal Processing for Phased Array Antennas II*, 1990, p. 59.
19. A. R. Williams, A. L. Kellner, X. S. Jiang, and P. K. L. Yu, "InGaAs/InP Waveguide Photodetector with High Saturation Intensity," *Electron. Lett.*, vol. 28, 1992, p. 2258.
20. M. Wu and T. Itoh, "Ultrafast Photonic to Microwave Transformer (PMT)," *LEOS Topical Meeting on Optical Microwave Interactions*, Paper W1.2, 1993.
21. A. Larsson et al., *J. Quantum Electron.*, vol. 24, 1988, p. 787.
22. D. L. Crawford, R. Nagarajan, and J. E. Bowers, "Comparison of Bulk and Quantum Wire Photodetectors," *Appl. Phys. Lett.*, vol. 58, no. 15, April 1991, pp. 1629–1631.
23. D. Kuhl, F. Hieronymi, E. H. Bottcher, and D. Bimberg, "High-Speed Metal-Semiconductor-Metal Photodetectors on InP: Fe," *IEEE Photon. Tech. Lett.*, vol. 2, no. 8, August 1990, p. 574.
24. Y. G. Wey, K. S. Giboney, J. E. Bowers, M. J. W. Rodwell, P. Silvestre, P. Thiagarajan, and G. Y. Robinson, "110 GHz GaInAs/InP *p-i-n* Photodiodes with Integrated Bias Tees and Matched Resistors," *IEEE Photonic Tech. Lett.*, August 1993.
25. D. M. Braun, "Design of Single Layer Antireflection Coatings for InP/InGaAs/InP Photodetectors for the 1200–1600 nm Wavelength Range," *Appl. Opt.*, vol. 27, 1988, pp. 2006–2011.
26. G. Lucovsky, R. F. Schwarz, and R. B. Emmons, "Transit-Time Considerations in *p-i-n* Diodes," *J. Appl. Phys.*, vol. 35, no. 3, March 1961.
27. J. E. Bowers and C. A. Burrus, "High Speed Zero Bias Waveguide Photodetectors," *Electron. Lett.*, vol. 22, 1986, p. 905.
28. A. Alping, R. Tell, and S. T. Eng, "Photodetection Properties of Semiconductor Laser Diode Detectors," *J. Lightwave Tech.*, vol. LT-4, 1986, pp. 1662–1668.
29. A. Alping, "Waveguide *p-i-n* Photodetectors: Theoretical Analysis and Design Criteria," *IEEE Proceedings*, vol. 136, part J, no. 3, June 1989.
30. S. Y. Wang and D. Bloom, "100 GHz Bandwidth Planar GaAs Schottky Photodiode," *Electron. Lett.*, vol. 19, no. 14, 7 July, 1983, p. 554.
31. D. G. Parker and P. G. Say, "Indium Tin Oxide/GaAs Photodiodes for Millimetric-Wave Applications," *Electron. Lett.*, vol. 22, no. 23, 6 Nov. 1986, p. 1266.
32. H. Kamiyama, Y. Kobayashi, T. Nagatsuma, and T. Kamiya, "Very Short Electrical Pulse Generation by a Composite Planar GaAs Photodetectors," *Jpn. J. Appl. Phys.*, vol. 29, Sept. 1990, p. 1717.
33. S. M. Sze, *Physics of Semiconductor Devices*, 1981, pp. 255–263.
34. M. Missouf and E. H. Rhoderick, "On the Richardson Constant for Aluminum/Gallium Arsenide Schottky Diodes," *J. Appl. Phys.*, vol. 69, no. 10, 15 May 1991, p. 7142.
35. M. Kamegawa, K. Giboney, J. Karin, S. Allen, M. Case, R. Yu, M. J. W. Rodwell, and J. E. Bowers, "Picosecond GaAs Monolithic Optoelectronic Sampling Circuit," *Photonics Technology Lett.*, vol. 3, no. 6, June 1991, pp. 567–569.
36. E. Ozbay, K. D. Li, and D. M. Bloom, "2.0 psec GaAs Monolithic Photodetector," *IEEE Photon. Tech. Lett.*, vol. 3, no. 6, June 1991, p. 570.
37. N. Emeis, H. Schumacher, and H. Beneking, "High-Speed GaInAs Schottky Photodetector," *Electron. Lett.*, vol. 21, no. 5, 28 Feb. 1985, p. 181.
38. L. Yang, A. S. Sudbo, R. A. Logan, T. Tanbun-Ek, and W. T. Tsang, "High Performance Fe:InP/GaAs Metal/Semiconductor/Metal Photodetectors Grown by Metalorganic Vapor Phase Epitaxy," *IEEE Photon. Tech. Lett.*, vol. 2, no. 1, January 1990, p. 56.
39. T. Mikawa, H. Kuwatsuka, Y. Kito, T. Kumai, M. Makiuchi, S. Yamazaki, O. Wada, and T. Shirai, "Flip-Chip InGaAs Avalanche Photodiode with Ultra Low Capacitance and Large Gain-Bandwidth Product," *Tech. Digest, ThO₂, OFC 1991*, p. 186.
40. H. Kuwatsuka, T. Mikawa, S. Miura, N. Yasuoka, T. Tanahashi, and O. Wada, "An Al_xGa_{1-x}Sb Avalanche Photodiode with Gain Bandwidth Product of 90 GHz," *Photon. Tech. Lett.*, vol. 2, no. 1, Jan 1990, p. 54.

41. F. Capusso, H. M. Cox, A. L. Hutchinson, N. A. Olsson, and S. G. Hummel, "Pseudo-Quaternary GaInAsP Semiconductors: A New $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ Graded Gap Superlattice and Its Applications to Avalanche Photodiodes," *Appl. Phys. Lett.*, vol. 45, no. 11, 1 December 1984, pp. 1193–1195.
42. L. E. Tarof, "Planar InP/InGaAs Avalanche Photodiodes with a Gain-Bandwidth Product Exceeding 100 GHz," *Tech. Digest*, ThO3, OFC 1991, p. 187.
43. K. Taguchi, T. Torikai, Y. Sugimoto, K. Makito, and H. Ishihara, "Planar-Structure InP/InAsP/InGaAs Avalanche Photodiodes with Preferential Lateral Extended Guard Ring for 1.0–1.6 μm Wavelength Optical Communication Use," *J. Lightwave Tech.*, vol. 6, no. 11, Nov. 1988, p. 1643.
44. T. Shiba, E. Ishimura, K. Takahashi, H. Namizaki, and W. Susaki, "New Approach to the Frequency Response Analysis of an InGaAs Avalanche Photodiode," *J. Lightwave Tech.*, vol. 6, no. 10, Oct. 1988, p. 1502.
45. K. Berchtold, O. Krumpolz, and J. Suri, "Avalanche Photodiodes with a Gain-Bandwidth Product of More Than 200 GHz," *Appl. Phys. Lett.*, vol. 26, no. 10, 15 May 1975, p. 585.
46. T. Kagawa, H. Asai, and Y. Kawamura, "An InGaAs/InAlAs Superlattice Avalanche Photodiode with a Gain Bandwidth Product of 90 GHz," *IEEE Photon. Tech. Lett.* vol. 3, no. 9, September 91, pp. 815–817.
47. H. Imai and T. Kaneda, "High-Speed Distributed Feedback Lasers and InGaAs Avalanche Photodiodes," *J. Lightwave Tech.*, vol. 6, no. 11, Nov. 1988, p. 1643.
48. H. C. Hsieh and W. Sargeant, "Avalanche Buildup Time of an InP/InGaAsP/InGaAs APD at High Gain," *J. Quantum Electron.*, vol. 25, no. 9, Sept. 1989, p. 2027.
49. F. Osaka, T. Mikawa, and T. Kaneda, "Impact Ionization of Electrons and Holes in (100)-Oriented Ga $_{1-x}$ In $_x$ As $_y$ P $_{1-y}$," *IEEE J. Quantum Electron.*, vol. QE-21, no. 9, September 1985, pp. 1326–1338.
50. C. A. Lee, R. A. Logan, R. L. Batdorf, J. J. Kleimack, and W. Wiegmann, "Ionization Rates of Holes and Electrons in Silicon," *Phys. Rev.*, vol. 134, 1964, pp. A761–A773.
51. R. B. Emmons, "Avalanche-Photodiode Frequency Response," *J. Appl. Phys.*, vol. 38, no. 9, August 1967, p. 3705.
52. S. L. Miller, "Avalanche Breakdown in Germanium," *Phys. Rev.*, vol. 99, Aug. 1955, pp. 1234–1241.
53. R. J. McIntyre, "Multiplication Noise in Uniform Avalanche Junctions," *IEEE Trans. Electron. Devices*, vol. ED-13, Jan. 1966, pp. 164–168.
54. B. L. Kasper and J. C. Campell, "Multigigabit-per-Second Avalanche Photodiode Lightwave Receivers," *J. Lightwave Tech.*, vol. LT-5, no. 10, October 1987, p. 1351.
55. M. Brain and T. P. Lee, "Optical Receiver for Lightwave Communication Systems," *J. Lightwave Tech.*, vol. LT-3, no. 6, December 1985, p. 1281.
56. T. V. Muoi, "Receiver Design for High-Speed Optical-Fiber Systems," *J. Lightwave Tech.*, vol. LT-2, no. 3, June 1984, p. 243.
57. J. A. Caldmain and G. Mourou, "Subpicosecond Electrooptic Sampling: Principles and Applications," *IEEE J. Quantum Electron.*, vol. QE-22, 1986, pp. 69–78.
58. D. R. Grischkowsky, M. B. Ketchen, C.-C. Chi, I. N. Duling, III, N. J. Halas, J.-M. Halbout, and P. G. May, "Capacitance Free Generation and Detection of Subpicosecond Electrical Pulses on Coplanar Transmission Lines," *IEEE J. Quantum Electron.*, vol. 24, no. 2, February 1988, pp. 221–225.
59. W. C. Nunnally and R. B. Hammond, "Optoelectronic Switch for Pulsed Powers," *Picosecond Optoelectronic Devices*, C. H. Lee (ed.), Academic Press, Orlando, Fla., 1984, pp. 373–398.
60. F. W. Smith, H. Q. Le, V. Diadiuk, M. A. Hollis, A. R. Calawa, S. Gupta, M. Frankel, D. R. Dykaar, G. A. Mourou, and T. Y. Hsiang, "Picosecond GaAs-Based Photoconductive Optoelectronic Detectors," *Appl. Phys. Lett.*, vol. 54, no. 10, 6 March 1989, p. 890.
61. N. G. Paulter, A. J. Gibbs, and D. N. Sinha, "Fabrication of High-Speed GaAs Photoconductive Pulse Generators and Sampling Gates by Ion Implantation," *IEEE Trans. Electron Device*, vol. 35, no. 12, December 1988, pp. 2343–2348.
62. Y. Chen, S. Williamson, and T. Brock, "1.2 ps High Sensitivity Photodetector/Switch Based on Low-Temperature-Grown GaAs," Postdeadline papers, CPDP 10/591, CLEO 1991.

SIGNAL DETECTION AND ANALYSIS

John R. Willison

*Stanford Research Systems, Inc.
Sunnyvale, California*

27.1 GLOSSARY

A	dimensionless material constant for $1/f$ noise
C	capacitance (farads)
I	current (amps)
$I_{\text{shot noise}}$	shot noise current (amps)
k	Boltzmann's constant
q	electron charge (Coulombs)
R	resistance (ohms)
S/N	signal-to-noise ratio
T	temperature (Kelvin)
$V_{\text{Johnson, rms}}$	RMS Johnson noise voltage (V)
Δf	bandwidth (Hz)

27.2 INTRODUCTION

Many optical systems require a quantitative measurement of light. Applications range from the very simple, such as a light meter using a photocell and a d'Arsenval movement, to the complex, such as the measurement of a fluorescence lifetime using time-resolved photon counting.

Often, the signal of interest is obscured by noise. The noise may be fundamental to the process: photons are discrete quanta governed by Poisson statistics which gives rise to shot noise. Or, the noise may be from more mundane sources, such as microphonics, thermal emf's, or inductive pickup.

This chapter describes methods for making useful measurements of weak optical signals, even in the presence of large interfering sources. The chapter will emphasize the electronic aspects of the problem. Important details of optical systems and detectors used in signal recovery are covered in Chap. 24, "Photodetectors," by Paul R. Norton; Chap. 25, "Photodetection," by Abhay M.

Joshi and Gregory H. Olsen; and Chap. 26, “High-Speed Photodetectors,” by John E. Bowers and Yih G. Wey in this volume.

27.3 PROTOTYPE EXPERIMENT

Figure 1 details the elements of a typical measurement situation. We wish to measure light from the source of interest. This light may be obscured by light from background sources. The intensity of the source of interest, and the relative intensity of the interfering background, will determine whether some or all of the techniques shown in Fig. 1 should be used.

Optics

An optical system is designed to pass photons from the source of interest and reject photons from background sources. The optical system may use spatial focusing, wavelength, or polarization selection to preferentially deliver photons from the source of interest to the detector.

There are many trade-offs to consider when designing an optical system. For example, if the source is nearly monochromatic and the background is broadband, then a monochromator may be used to improve the signal-to-background ratio of the light reaching the detector. However, if the source of interest is an extended isotropic emitter, then a monochromator with narrow slits and high f -number will dramatically reduce the number of signal photons from the source which can be passed to the detector. In this case, the noise of the detector and amplifiers which follow the optical system may dominate the overall signal-to-noise ratio (S/N).

Photodetectors

There are many types of nonimaging photodetectors. Key criteria to select a photodetector for a particular application include sensitivity for the wavelength of interest, gain, noise, and speed. Important details of many detector types are given in other chapters in the *Handbook*. Operational details (such as bias circuits) of photomultipliers which are specific to boxcar integration and photon counting are discussed in Sec. 27.5.

Amplifiers

In many applications, the output of the detector must be amplified or converted from a current to a voltage before the signal may be analyzed. Selection criteria for amplifiers include type (voltage or transconductance), gain, bandwidth, and noise.

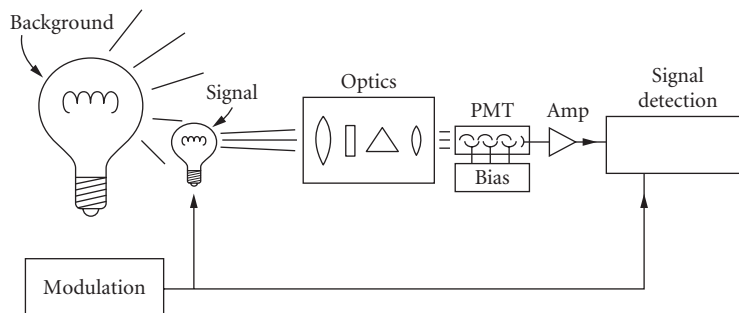


FIGURE 1 Prototypical optical measurement.

Signal Analysis

There are two broad categories of signal analysis, depending on whether or not the source is modulated. Modulating the source allows the signal to be distinguished from the background. Often, source modulation is inherent to the measurement. For example, when a pulsed laser is used to induce a fluorescence, the signal of interest is present only after the laser fires. Other times, the modulation is “arranged,” as when a cw source is chopped. Sometimes the source cannot be modulated or the source is so dominant over the background that modulation is unnecessary.

27.4 NOISE SOURCES

An understanding of noise sources in a measurement is critical to achieving signal-to-noise performance near theoretical limits. The quality of a measurement may be substantially degraded by a trivial error. For example, a poor choice of termination resistance for a photodetector may increase current noise by several orders of magnitude.¹

Shot Noise

Light and electrical charge are quantized, and so the number of photons or electrons which pass a point during a period of time are subject to statistical fluctuations. If the signal mean is M photons, its standard deviation (noise) will be \sqrt{M} , hence the $S/N = M/\sqrt{M} = \sqrt{M}$. The mean M may be increased if the rate is higher or the integration time is longer. Short integration times or small signal levels will yield poor S/N values. Figure 2 shows the S/N , which may be expected as a function of current level and integration time for a shot-noise limited signal.

“Integration time” is a convenient parameter when using time-domain signal recovery techniques. “Bandwidth” is a better choice when using frequency-domain techniques. The rms noise current in the bandwidth Δf Hz due to a “constant” current, I amps, is given by

$$I_{\text{shot noise}} = \sqrt{(2qI\Delta f)} \quad (1)$$

where $q = 1.6 \times 10^{-19}$ C

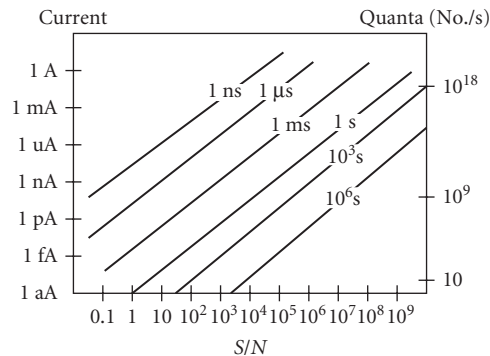


FIGURE 2 Signal-to-noise versus flux and measurement time.

Johnson Noise

The electrons which allow current conduction in a resistor are subject to random motion which increases with temperature. This fluctuation of electron density will generate a noise voltage at the terminals of the resistor. The rms value of this noise voltage for a resistor of R ohms, at a temperature of T Kelvin, in a bandwidth of Δf Hz is given by

$$V_{\text{Johnson,rms}} = \sqrt{(4kTR\Delta f)} \quad (2)$$

where k is Boltzmann's constant. The noise voltage in a 1-Hz bandwidth is given by

$$V_{\text{Johnson,rms}}(\text{per } \sqrt{\text{Hz}}) = 0.13 \text{ nV} \times \sqrt{(R(\text{ohms}))} \quad (3)$$

Since the Johnson noise voltage increases with resistance, large-value series resistors should be avoided in voltage amplifiers. For example, a 1-k Ω resistor has a Johnson voltage of about 4.1 nV/ $\sqrt{\text{Hz}}$. If detected with a 100-MHz bandwidth, the resistor will show a noise of 41- μV rms, which has a peak-to-peak value of about 200 μV .

When a resistor is used to terminate a current source, or as a feedback element in a current-to-voltage converter, it will contribute a noise current equal to the Johnson noise voltage divided by the resistance. Here, the noise current in a 1-Hz bandwidth is given by

$$I_{\text{Johnson,rms}}(\text{per } \sqrt{\text{Hz}}) = 130 \text{ pA} / \sqrt{R(\text{ohms})} \quad (4)$$

As the Johnson noise current increases as R decreases, small-value resistors should be avoided when terminating current sources. Unfortunately, small terminating resistors are required to maintain a wide frequency response. If a 1-k Ω resistor is used to terminate a current source, the resistor will contribute a noise current of about 4.1 pA/ $\sqrt{\text{Hz}}$, which is about 1000 \times worse than the noise current of an ordinary FET input operational amplifier.

1/f Noise

The voltage across a resistor carrying a constant current will fluctuate because the resistance of the material used in the resistor varies. The magnitude of the resistance fluctuation depends on the material used: carbon composition resistors are the worst, metal film resistors are better, and wire wound resistors provide the lowest 1/f noise. The rms value of this noise source for a resistance of R ohms, at a frequency of f Hz, in a bandwidth of Δf Hz is given by

$$V_{1/f,\text{rms}} = IR \times \sqrt{(A\Delta f/f)} \quad (5)$$

where the dimensionless constant A has a value of about 10^{-11} for carbon. In a measurement in which the signal is the voltage across the resistor (IR), the $S/N = 3 \times 10^5 \sqrt{(f/\Delta f)}$. Often, this noise is a troublesome source of low-frequency noise in voltage amplifiers.

Nonessential Noise Sources

There are many discrete noise sources which must be avoided in order to make reliable low-level light measurements. Figure 3 shows a simplified noise spectrum on log-log scales. The key features in this noise spectra are frequencies worth avoiding: diurnal drifts (often seen via input offset drifts with temperature), low frequency (1/f) noise, power line frequencies and their harmonics, switching power supply and crt display frequencies, commercial broadcast stations

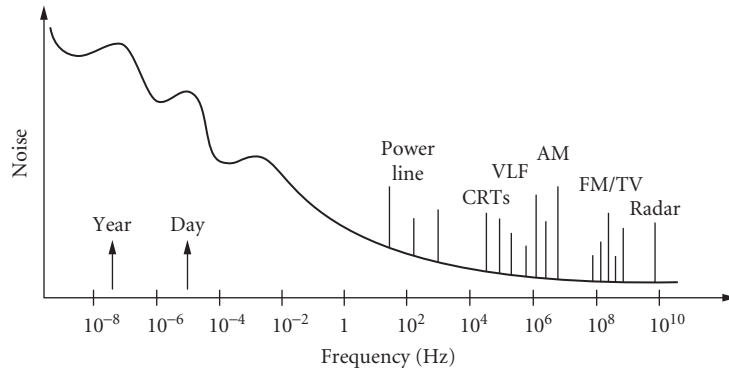


FIGURE 3 Simplified noise spectrum.

(AM, FM, VHF, and UHF TV), special services (cellular telephones, pagers, etc.), microwave ovens and communications, to RADAR and beyond.

Your best alternatives for avoiding these noise sources are

1. Shield to reduce pickup.
2. Use differential inputs to reject common mode noise.
3. Bandwidth limit the amplifier to match expected signal.
4. Choose a quiet frequency for signal modulation when using a frequency-domain detection technique.
5. Trigger synchronously with interfering source when using a time-domain detection technique.

Common ways for extraneous signals to interfere with a measurement are illustrated in Fig. 4a to f.

Noise may be injected via a stray capacitance as in Fig. 4a. The stray capacitance has an impedance of $1/j\omega C$. Substantial currents may be injected into low-impedance systems (such as transconductance inputs), or large voltages may appear at the input to high-impedance systems.

Inductive pickup is illustrated in Fig. 4b. The current circulating in the loop on the left will produce a magnetic field which in turn induces an emf in the loop on the right. Inductive noise pickup may be reduced by reducing the areas of the two loops (by using twisted pairs, for example), by increasing the distance between the two loops, or by shielding. Small skin depths at high frequencies allow nonmagnetic metals to be effective shields; however, high- μ materials must be used to shield from low frequency magnetic fields.

Resistive coupling, or a “ground loop,” is shown in Fig. 4c. Here, the detector senses the output of the experiment plus the IR voltage drop from another circuit which passes current through the same ground plane. Cures for ground-loop pickup include grounding everything to the same point, using a heavier ground plane, providing separate ground return paths for large interfering currents, and using a differential connection between the signal source and amplifier.

Mechanical vibrations can create electrical signals (microphonics) as shown in Fig. 4d. Here, a coaxial cable is charged by a battery through a large resistance. The voltage on the cable is $V = Q/C$. Any deformation of the cable will modulate the cable’s capacitance. If the period of vibration which causes the deformation is short compared to the RC time constant then the stored charge on the cable, Q , will remain constant. In this case, a 1-ppm modulation of the cable capacitance will generate an ac signal with an amplitude of 1 ppm of the dc bias on the cable, which may be larger than the signal of interest.

The case of magnetic microphonics is illustrated in Fig. 4e. Here, a dc magnetic field (the earth’s field or the field from a permanent magnet in a latching relay, for example) induces an emf in the signal path when the magnetic flux through the detection loop is modulated by mechanical motion.

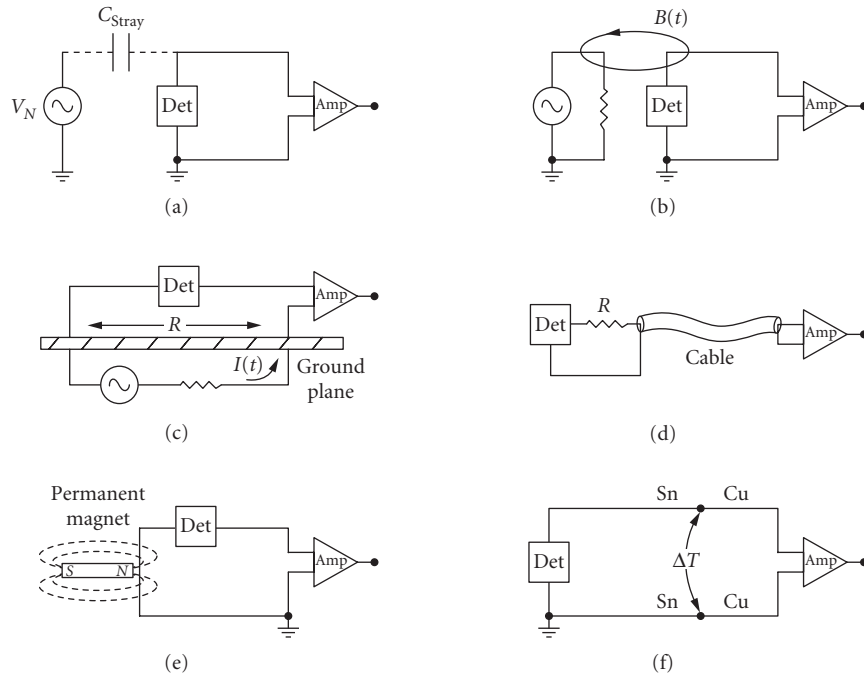


FIGURE 4 Coupling of noise sources.

Unwanted thermocouple junctions are an important source of offset and drift. As shown in Fig. 4f, two thermocouple junctions are formed when a signal is connected to an amplifier. For typical interconnect materials (copper, tin) one sees about $10 \mu\text{V}/^\circ\text{C}$ of offset. These extraneous junctions occur throughout instruments and systems: their impact may be eliminated by making ac measurements.

27.5 APPLICATIONS USING PHOTOMULTIPLIERS

Photomultiplier tubes (PMTs) are used for detection of light from about 200 to 900 nm. Windowless PMTs can be used from the near UV through the x-ray region, and may also be used as particle detectors. Their low noise, high gain, wide bandwidth, and large dynamic range have placed them in many applications. They are the only detectors which may be recommended for low-noise photon counting applications.^{2,3}

In this chapter, we are primarily concerned with the electrical characteristics of PMTs. Understanding these characteristics is important if we are to realize the many desirable features of these devices.

A schematic representation of a PMT, together with a typical bias circuit, is shown in Fig. 5. While the concepts depicted here are common to all PMTs, the particulars of biasing and termination will change between PMT types and applications. PMTs have a photocathode, several dynodes (6 to 14), and an anode. They are usually operated from a negative high voltage, with the cathode at the most negative potential, each successive dynode at a less negative potential, and the anode near ground. An incident photon may eject a single photoelectron from the photocathode which will strike the first dynode with an energy of a few hundred volts. A few (2– to 5) electrons will be ejected from the first dynode by the impact of the photoelectron: these electrons will in turn strike

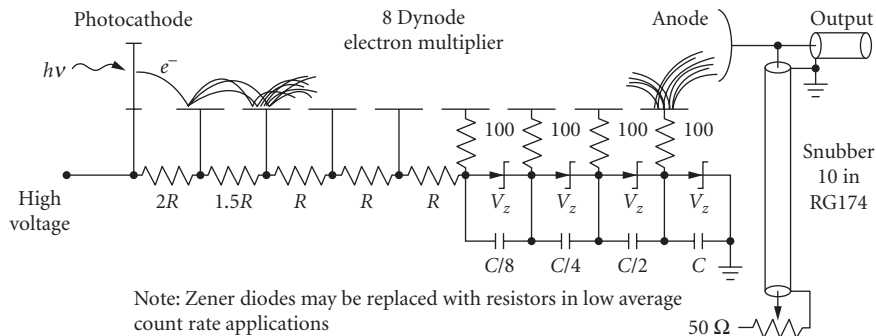


FIGURE 5 PMT base for photon counting or fast integration.

the second dynode, ejecting more electrons. The process continues at each dynode until all of the electrons are collected by the anode.

Quantum Efficiency

The quantum efficiency (QE) of a PMT is a measure of the probability that a photon will eject a photoelectron at the photocathode. The QE depends on the type of material used in the cathode and the wavelength of light. QEs may be as high as 10 to 30 percent at their peak wavelength. The cathode material will also affect the dark count rate from the PMT: a cathode with good red sensitivity may have a high dark count rate.

Gain

A PMT's gain depends on the number of dynodes, the dynode material, and voltage between the dynodes. PMT gains range from 10^3 to 10^7 . The anode output from the PMT will typically go to an electronic amplifier. To avoid having the system noise be dominated by the amplifier's noise, the PMT should be operated with enough gain so that the dark current times the gain is larger than the amplifier's input current noise.

Bandwidth

The frequency response, speed, rise time, and pulse-pair resolution of PMTs depend on the structure of the dynode multiplier chain. The leading edges of the anode output have transition times from 2 to 20 ns. Trailing edges are usually about three times slower. Much faster PMTs, with rise times on order 100 ps, use microchannel plate multipliers.

When using gated integrators to measure PMT outputs, the pulse width of the anode signal should be less than the gate width so that timing information is not lost. For photon counting, the pulse width should be smaller than the pulse-pair resolution of the counter/discriminator to avoid saturation effects. When using lock-in amplifiers, pulse width is usually not important, since the slowest PMTs will have bandwidths well above the modulation frequency.

Pulse Height

In pulsed experiments, the criterion for a detectable signal often depends on the electrical noise environment of the laboratory and the noise of the preamplifier. In laboratories with Q-switched lasers or pulsed discharges, it is difficult to reduce the noise on any coaxial cable below a few millivolts.

A good, wide bandwidth preamplifier will have about $1.5 \text{ nV}/\sqrt{\text{Hz}}$, or about $25\text{-}\mu\text{V}$ rms over a 300-MHz bandwidth. Peak noise will be about 2.5 times the rms noise, so it is important that the PMT provide pulses of greater than 1-mV amplitude.

Use manufacturer's specifications for the current gain and rise time to estimate the pulse amplitude from the PMT:

$$\text{Amplitude (mV)} = 4 \times \text{gain (millions)}/\text{rise time (ns)} \quad (6)$$

This formula assumes that the electrons will enter a $50\text{-}\Omega$ load in a square pulse whose duration is twice the rise time. (Since the rise time will be limited by the bandwidth of the preamplifier, use the larger of the amplifier or PMT rise times in this formula.)

If the PMT anode is connected via a $50\text{-}\Omega$ cable to a large load resistance, then the pulse shape may be modeled by the lumped parameters of the cable capacitance (about $100 \text{ pF}/\text{meter}$ for RG-58) and the termination resistance. All of the charge in the pulse is deposited on the cable capacitance in a few nanoseconds. The voltage on the load will be $V = Q/C$ where C = cable capacitance. This voltage will decay exponentially with a time constant of RC where R is the load resistance in ohms. In this case, the pulse height will be

$$\text{Amplitude (mV)} = 160 \times \text{gain (millions)}/\text{cable } C \text{ (pF)} \quad (7)$$

The current gain of a PMT is a strong function of the high voltage applied to the PMT. Very often, PMTs will be operated well above the high voltage recommended by the manufacturer, and thus substantially higher current gains ($10\times$ to $100\times$ above specs). There are usually no detrimental effects to the PMT as long as the anode current is kept well below the rated value.

Dark Counts

PMTs are the quietest detectors available. The primary noise source is thermionic emission of electrons from the photocathode and from the first few dynodes of the electron multiplier. PMT housings which cool the PMT to about -20°C can dramatically reduce the dark counts (from a few kHz to a few Hz). The residual counts arise from radioactive decays of materials inside the PMT and from cosmic rays.

PMTs which are specifically designed for photon counting will specify their noise in terms of the rate of output pulses whose amplitudes exceed some fraction of a pulse from a single photon. More often, the noise is specified as an anode dark current. Assuming the primary source of dark current is thermionic emission from the photocathode, the dark count rate is given by

$$\text{Dark count (kHz)} = 6 \times \text{dark current (nA)}/\text{gain (millions)} \quad (8)$$

PMT Base Design

PMT bases which are designed for general-purpose applications are not appropriate for photon counting or fast-gated integrator applications (gates < 10 to 20 ns). General-purpose bases will not allow high count rates, and often cause problems such as double counting and poor plateau characteristics. A PMT base with the proper high-voltage taper, bypassing, snubbing, and shielding is required for good time resolution and best photon counting performance.

Dynode Biasing A PMT base provides bias voltages to the PMTs photocathode and dynodes from a single, negative, high-voltage power supply. The simplest design consists of a resistive voltage divider. In this configuration the voltage between each dynode, and thus the current gain at each dynode, is the same. Typical current gains are three to five, so there will typically be four electrons leaving the first dynode, with a variance of about two electrons. This large relative variance (due to the small

number of ejected electrons) gives rise to large variations in the pulse height of the detected signal. Since statistical fluctuations in pulse height are dominated by the low gain of the first few stages of the multiplier chain, increasing the gain of these stages will reduce pulse-height variations and so improve the pulse-height distribution. This is important for both photon counting and analog detection. To increase the gain of the first few stages, the resistor values in the bias chain are increased to increase the voltage in the front end of the multiplier chain. The resistor values are tapered slowly so that the electrostatic focusing of electrons in the multiplier chain is not adversely affected.⁴

Current for the electron multiplier is provided by the bias network. Current drawn from the bias network will cause the dynode potentials to change, thus changing the tube gain. This problem is of special concern in lifetime measurements. The shape of exponential decay curves will be changed if the tube gain varies with count rate. To be certain that this is not a problem, lifetime measurements should be repeated at reduced intensity. The problem of gain variation with count rate is avoided if the current in the bias network is about 20 times the output current from the PMT's anode.

There are a few other methods to avoid this problem which do not require high bias currents. These methods depend on the fact that the majority of the output current is drawn from the last few dynodes of the multiplier:

1. Replace the last few resistors in the bias chain with Zener diodes. As long as there is some reverse current through a Zener, the voltage across the diodes is nearly constant. This will prevent the voltage on these stages from dropping as the output current is increased.
2. Use external power supplies for the last few dynodes in the multiplier chain. This approach dissipates the least amount of electrical power since the majority of the output current comes from lower-voltage power supplies. However, it is the most difficult to implement.
3. If the average count rate is low, but the peak count rate is high, then bypass capacitors on the last few stages may be used to prevent the dynode voltage from dropping (use $20\times$ the average output current for the chain current). For a voltage drop of less than 1 percent, the stored charge on the last bypass capacitor should be $100\times$ the charge output during the peak count rate. For example, the charge output during a 1-ms burst of a 100-MHz count rate, each with an amplitude of 10 mV into $50\ \Omega$ and a pulse width of 5 ns, is $0.1\ \mu\text{C}$. If the voltage on the last dynode is 200 Vdc, then the bypass capacitor for the last dynode should have a value given by

$$C = 100Q/V = 100 \times 0.15\text{C}/200\text{V} = 0.05\ \mu\text{F} \quad (9)$$

The current from higher dynodes is smaller so the capacitors bypassing these stages may be smaller. Only the final four or five dynodes need to be bypassed, usually with a capacitor which has half the capacitance of the following stage. To reduce the voltage requirement for these capacitors, they are usually connected in series.

Bypassing the dynodes of a PMT may cause high-frequency ringing of the anode output signal. This can cause multiple counts for a single photon or poor time resolution in a gated integrator. The problem is significantly reduced by using small resistors between the dynodes and the bypass capacitors.

Snubbing Snubbing refers to the practice of adding a network to the anode of the PMT to improve the shape of the output pulse for photon counting or fast-gated integrator applications. This "network" is usually a short piece of $50\text{-}\Omega$ coax cable which is terminated into a resistor of less than $50\ \Omega$. The snubber will delay, invert, and sum a small portion of the anode signal to itself. Snubbing should not be used when using a lock-in amplifier since the current conversion gain of a $50\text{-}\Omega$ resistor is very small.

There are four important reasons for using a snubber network:

1. Without some dc resistive path between the anode and ground, anode dark current will charge the signal cable to a few hundred volts (last dynode potential). When the signal cable is connected to an amplifier, the stored charge on the cable may damage the front end of the instrument. PMT bases without a snubber network should include a $100\text{-M}\Omega$ resistor between the anode and ground to protect the instruments.

2. The leading edge of the output current pulse is often much faster than the trailing edge. A snubber network may be used to sharply increase the speed of the trailing edge, greatly improving the pulse pair resolution of the PMT. This is especially important in photon counting applications.
3. Ringing (with a few nanoseconds period) is very common on PMT outputs. A snubber network may be used to cancel these rings which can cause multiple counts from a single photon.
4. The snubber network will help to reverse terminate reflections from the input to the preamplifier.

The round-trip time in the snubber cable may be adjusted so that the reflected signal cancels anode signal ringing. This is done by using a cable length with a round-trip time equal to the period of the anode ringing.

Cathode Shielding Head-on PMTs have a semitransparent photocathode which is operated at negative high voltage. Use care so that no objects near ground potential contact the PMT near the photocathode.

Magnetic Shielding Electron trajectories inside the PMT will be affected by magnetic fields. A field strength of a few gauss can dramatically reduce the gain of a PMT. A magnetic shield made of a high permeability material should be used to shield the PMT.

PMT Base Summary

1. Taper voltage divider for higher gain in the first stages.
2. Bypass last few dynodes in pulsed applications.
3. Use a snubber circuit to shape the output pulse for photon counting or fast-gated integration.
4. Shield the tube from electrostatic and magnetic fields.

27.6 AMPLIFIERS

Several considerations are involved in choosing the correct amplifier for a particular application. Often, these considerations are not independent, and compromises will be necessary. The best choice for an amplifier depends on the electrical characteristics of the detector, and on the desired gain, bandwidth, and noise performance of the system.

Voltage Amplifiers

High Bandwidth Photon counting and fast-gated integration require amplifiers with wide bandwidth. A 350-MHz bandwidth is required to preserve a 1-ns rise time. The input impedance to these amplifiers is usually 50 Ω in order to terminate coaxial cables into their characteristic impedance. When PMTs (which are current sources) are connected to those amplifiers, the 50- Ω input impedance serves as the current-to-voltage converter for the PMT anode signal. Unfortunately, the small termination resistance and wide bandwidth yield a lot of current noise.⁵

High Input Impedance It is important to choose an amplifier with a very high input impedance and low-input bias current when amplifying a signal from a source with a large equivalent resistance. Commercial amplifiers designed for such applications typically have a 100-M Ω input impedance. This large input impedance will minimize attenuation of the input signal and reduce the Johnson noise current drawn through the source resistance, which can be an important noise source. Field effect transistors (FETs) are used in these amplifiers to reduce the input bias current to the amplifiers. Shot noise on the input bias current can be an important noise component, and temperature drift of the input bias current is a source of drift in dc measurements.⁶

The bandwidth of a high-input impedance amplifier is often determined by the RC time constant of the source, cable, and termination resistance. For example, a PMT with 1 meter of RG-58 coax (about 100 pF) terminated into a 1-M Ω resistor will have a bandwidth of about 1600 Hz. A smaller resistance would improve the bandwidth, but increase the Johnson noise current.

Moderate Input Impedance Bipolar transistors offer an input noise voltage which may be several times smaller than the FET inputs of high-input impedance amplifiers, as low as $1 \text{ nV}/\sqrt{\text{Hz}}$. Bipolar transistors have larger input bias currents, hence larger shot noise current, and so should be used only with low-impedance ($<1 \text{ k}\Omega$) sources.

Transformer Inputs When ac signals from very low source impedances are to be measured, transformer coupling offers very quiet inputs. The transformer is used to step up the input voltage by its turns-ratio. The transformer's secondary is connected to the input of a bipolar transistor amplifier.

Low Offset Drift Conventional bipolar and FET input amplifiers exhibit input offset drifts on the order of $5 \mu\text{V}/\text{C}$. In the case where the detector signal is a small dc voltage, such as from a bolometer, this offset drift may be the dominant noise source. A different amplifier configuration, chopper-stabilized amplifiers, essentially measure their input offsets and subtract the measured offset from the signal. A similar approach is used to "autozero" the offset on the input to sensitive voltmeters. Chopper-stabilized amplifiers exhibit very low input offsets with virtually no input offset drift.

Differential The use of "true-differential" or "instrumentation" amplifiers is advised to provide common mode rejection to interfering noise, or to overcome the difference in grounds between the voltage source and the amplifier. This amplifier configuration amplifies the difference between two inputs, unlike a single-ended amplifier, which amplifies the difference between the signal input and the amplifier ground. In high-frequency applications, where good differential amplifiers are not available or are difficult to use, common mode choke may be used to isolate disparate grounds.

Transconductance Amplifiers

When the detector is a current source (or has a large equivalent resistance) then a transconductance amplifier should be considered. These amplifiers (current-to-voltage converters) offer the potential of lower noise and wider bandwidth than a termination resistor and a voltage amplifier; however, some care is required in their application.⁷

A typical transconductance amplifier configuration is shown in Fig. 6. An FET input op amp would be used for its low-input bias current. (Op amps with input bias currents as low as 50 fA are readily available.) The detector is a current source, I_o . Assuming an ideal op amp, the transconductance gain is $A = V_{\text{out}}/I_{\text{in}} = R_f$, and the input impedance of the circuit is R_{in} to the op amp's virtual null. (R_{in} allows negative feedback, which would have been phase shifted and attenuated by the

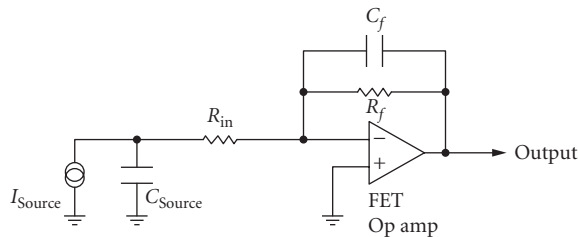


FIGURE 6 Typical transconductance amplifier.

source capacitance at high frequencies, to assure stability.) Commercial transconductance amplifiers use R_s 's as large as $10\text{ M}\Omega$, with R_{in} 's, which are typically $R_f/1000$. A low-input impedance will ensure that current from the source will not accumulate on the input capacitance.

This widely used configuration has several important limitations which will degrade its gain, bandwidth, and noise performance. The overall performance of the circuit depends critically on the source capacitance, including that of the cable connecting the source to the amplifier input. Limitations include:

1. The "virtual null" at the inverting input to the op amp is approximately R_f/A_v , where A_v is the op amp's open loop gain at the frequency of interest. While op amps have very high gain at frequencies below 10 Hz (typically a few million), these devices have gains of only a few hundred at 1 kHz. With an R_f of $1\text{ G}\Omega$, the virtual null has an impedance of $5\text{ M}\Omega$ at 1 kHz, hardly a virtual null. If the impedance of the source capacitance is less than the input impedance, then most of the ac input current will go to charging this capacitance, thereby reducing the gain.
2. The configuration provides high gain for the voltage noise at the noninverting input of the op amp. At high frequencies, where the impedance of the source capacitance is small compared to R_{in} , the voltage gain for noise at the noninverting input is R_f/R_{in} , typically about 1000. As FET input op amps with very low bias currents tend to have high-input-voltage noise, this term can dominate the noise performance of the design.
3. Large R_f 's are desired to reduce the Johnson noise current; however, large R_f 's degrade the bandwidth. If low values of R_f are used, the Johnson noise current can dominate the noise performance of the design.
4. To maintain a flat frequency response, the size of the feedback capacitance must be adjusted to compensate for different source capacitances.

As many undesirable characteristics of the transconductance amplifier can be traced to the source capacitance, a system may benefit from integrating the amplifier into the detector, thereby eliminating interconnect capacitance. This approach is followed in many applications, from microphones to CCD imagers.

27.7 SIGNAL ANALYSIS

Unmodulated Sources

For unmodulated sources, a strip-chart recorder, voltmeter, A/D converter, or oscilloscope may be used to measure the output of the amplifier or detector. In the case of low-light-level measurement, continuous photon counting would be the method of choice.

A variety of problems are avoided by modulating the signal source. When making dc measurements, the signal must compete with large low-frequency noise sources. However, when the source is modulated, the signal may be measured at the modulation frequency, away from these large noise sources.

Modulated Sources

When the source is modulated, one may choose from gated integration, boxcar averaging, transient digitizers, lock-in amplifiers, spectrum analyzers, gated photon counters, or multichannel scalars.

Gated Integration A measurement of the integral of a signal during a period of time can be made with a gated integrator. Commercial devices allow gates from about 100 ps to several milliseconds. A gated integrator is typically used in a pulsed laser measurement. The device can provide shot-by-shot data which is often recorded by a computer via an A/D converter. The gated integrator is

recommended in situations where the signal has a very low duty cycle, low pulse repetition rate, and high instantaneous count rates.⁸

The noise bandwidth of the gated integrator depends on the gate width: short gates will have wide bandwidths, and so will be noisy. This would suggest that longer gates would be preferred; however, the signal of interest may be very short-lived, and using a gate which is much wider than the signal will not improve the S/N .

The gated integrator also behaves as a filter: the output of the gated integrator is proportional to the average of the input signal during the gate, so frequency components of the input signal which have an integral number of cycles during the gate will average to zero. This characteristic may be used to “notch out” specific interfering signals.

It is often desirable to make gated integration measurements synchronously with an interfering source. (This is the case with time-domain signal detection techniques, and not the case with frequency-domain techniques such as lock-in detection.) For example, by locking the pulse repetition rate to the power-line frequency (or to any submultiple of this frequency) the integral of the line interference during the short gate will be the same from shot to shot, which will appear as a fixed offset at the output of the gated integrator.

Boxcar Averaging Shot-by-shot data from a gated integrator may be averaged to improve the S/N . Commercial boxcar averagers provide linear or exponential averaging. The averaged output from the boxcar may be recorded by a computer or used to drive a strip-chart recorder. Figure 7 shows a gated integrator with an exponential averaging circuit.

Lock-In Amplifiers Phase-sensitive synchronous detection is a powerful technique for the recovery of small signals which may be obscured by interference that is much larger than the signal of interest. In a typical application, a cw laser which induces the signal of interest will be modulated by an optical chopper. The lock-in amplifier is used to measure the amplitude and phase of the signal of interest relative to a reference output from the chopper.⁹

Figure 8 shows a simplified block diagram for a lock-in amplifier. The input signal is ac-coupled to an amplifier whose output is mixed (multiplied by) the output of a phase-locked loop which is locked to the reference input. The operation of the mixer may be understood through the trigonometric identity

$$\cos(\omega_1 t + \Phi) * \cos(\omega_2 t) = \frac{1}{2} \{ \cos[(\omega_1 + \omega_2)t + \Phi] + \cos[(\omega_1 - \omega_2)t + \Phi] \} \quad (10)$$

When $\omega_1 = \omega_2$ there is a dc component of the mixer output, $\cos \Phi$. The output of the mixer is passed through a low-pass filter to remove the sum frequency component. The time constant of the filter is selected to reduce the equivalent noise bandwidth: selecting longer time constants will improve the S/N at the expense of longer response times.

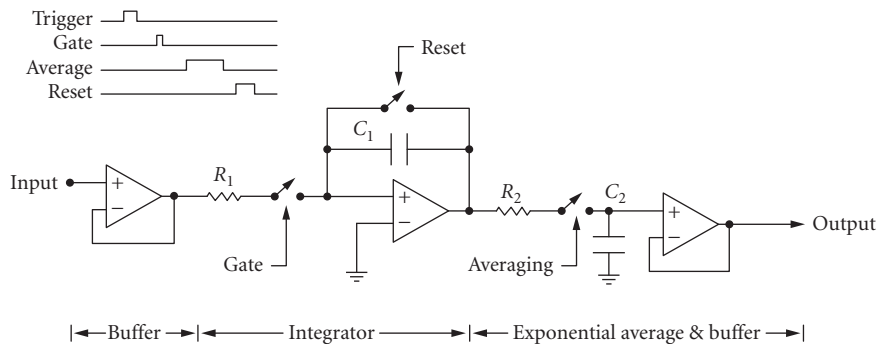


FIGURE 7 Gated integrator and exponential averager.

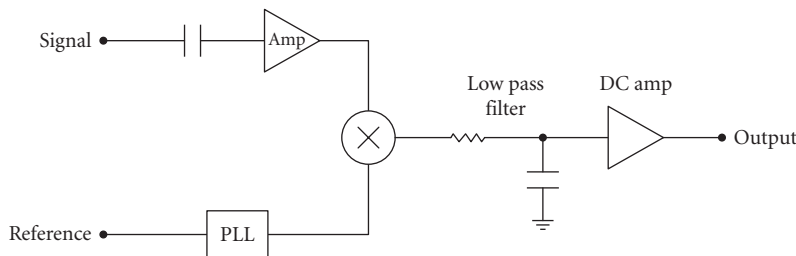


FIGURE 8 Lock-in amplifier block diagram.

The simplified block diagram shown in Fig. 8 is for a “single-phase” lock-in amplifier, which measures the component of the signal at one set phase with respect to the reference. A dual-phase lock-in has another channel which measures the component of the signal at 90° relative to the first channel, which allows simultaneous measurement of the amplitude and phase of the signal.

Digital signal processing (DSP) techniques are rapidly replacing the older analog techniques for the synchronous detection of the signal. In these instruments, the input signal is digitized by a fast, high-resolution A/D converter, and the signal’s amplitude and phase are determined by high-speed computations in a digital signal processor. To maintain the 100-kHz bandwidth of the analog designs, the DSP designs must complete a quarter million 16-bit A/D conversions and 20 million multiply-and-accumulate operations each second. Many artifacts of the analog designs are eliminated by the DSP approach; for example, the output drift and dynamic range of the instruments are dramatically improved.¹⁰

Photon Counting Photon counting techniques offer several advantages in the measurement of light: very high sensitivity (count rates as low as 1 per minute can be a usable signal level), large dynamic range (signal levels as high as 100 MHz can be counted, allowing a 195-dB dynamic range), discrimination against low-level noise (analog noise below the discriminator thresholds will not be counted), and ability to operate over widely varying duty cycles.¹¹

Key elements of a photon counting system include a high-gain PMT operated with sufficiently high voltage so that a single photoelectron will generate an anode pulse of several millivolts into a 50- Ω load, a fast discriminator to generate logic pulses from anode signals which exceed a set threshold, and fast-gated counters to integrate the counts.

Transient Photon Counting In situations where the time evolution of a light signal must be measured (LIDAR, lifetime measurements, chemical kinetics, etc.) transient photon counters allow the entire signal to be recorded for each event. In these instruments, the discriminated photon pulses are summed into different bins depending on their timing with respect to a trigger pulse. Commercial instruments offer 5-ns resolution with zero dead-time between bins. The time records from many events may be summed together in order to improve the S/N .¹²

Choosing the “Best” Technique Which instrument is best suited for detecting signals from a photomultiplier tube? The answer is based on many factors, including the signal intensity, the signal’s time and frequency distribution, the various noise sources and their time-dependence and frequency distribution.

In general, the choice between boxcar averaging (gated integration) and lock-in detection (phase-sensitive detection) is based on the time behavior of the signal. If the signal is fixed in frequency and has a 50 percent duty cycle, lock-in detection is best suited. This type of experiment commonly uses an optical chopper to modulate the signal at some low frequency. Signal photons occur at random times during the “open” phase of the chopper. The lock-in detects the average difference between the signal during the “open” phase and the background during the “closed” phase.

To use a boxcar averager in the same experiment would require the use of very long, 50 percent duty cycle gates since the photons can arrive anywhere during the “open” phase. Since the gated integrator is collecting noise during this entire gate, the signal is easily swamped by the noise. To correct for this, baseline subtraction can be used where an equal gate is used to measure the background during the “closed” phase of the chopper and subtracted from the “open” signal. This is then identical to lock-in detection. However, lock-in amplifiers are much better suited to this, especially at low frequencies (long gates) and low signal intensities.

If the signal is confined to a very short amount of time, then gated integration is usually the best choice for signal recovery. A typical experiment might be a pulsed laser excitation where the signal lasts for only a short time (100 ps to 1 μ s) at a repetition rate of up to 10 kHz. The duty cycle of the signal is much less than 50 percent. By using a narrow gate to detect signal only when it is present, noise which occurs at all other times is rejected. If a longer gate is used, no more signal is measured but the detected noise will increase. Thus, a 50 percent duty cycle gate would not recover the signal well and lock-in detection is not suitable.

Photon counting can be used in either the lock-in or the gated mode. Using a photon counter is usually required at very low signal intensities or when the use of a pulse height discriminator to reject noise results in an improved S/N . If the evolution of a weak light signal is to be measured, a transient photon counter or multichannel scaler can greatly reduce the time required to make a measurement.

27.8 REFERENCES

1. P. Horowitz and W. Hill, *The Art of Electronics*, Cambridge, New York, 1989, p. 428–447.
2. Photomultiplier Tubes, Hamamatsu Company catalog, 1988.
3. Photomultipliers, Thorn EMI Company catalog, 1990.
4. G. A. Morton and H. M. Smith, “Pulse Height Resolution of High Gain First Dynode Photomultipliers,” *Appl. Phys. Lett.* vol. 13, 1968, p. 356.
5. Model SR445 Fast Preamplifier, *Operation and Service Manual*, Stanford Research Systems, 1990.
6. Model SR560 Low Noise Preamplifier, *Operation and Service Manual*, Stanford Research Systems, 1990.
7. Model SR570 Low Noise Current Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1992.
8. Fast Gated Integrators and Boxcar Averagers, *Operation and Service Manual*, Stanford Research Systems, 1990.
9. Model SR510 Lock-in Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1987.
10. Model SR850 DSP Lock-in Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1992.
11. Model SR400 Gated Photon Counter, *Operation and Service Manual*, Stanford Research Systems, 1988.
12. Model SR430 Multichannel Scaler/Averager, *Operation and Service Manual*, Stanford Research Systems, 1989.

This page intentionally left blank.

William L. Wolfe

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

Paul W. Kruse

*Consultant
Edina, Minnesota*

28.1 GLOSSARY

DTGS	deuterated triglycine sulfate
p	pyroelectric coefficient
R_e	electrical resistance
R_{th}	thermal resistance
\mathfrak{R}	responsivity
S	Seebeck coefficient
TGS	triglycine sulfate
Z	figure of merit
τ_e	electrical time constant
τ_{th}	thermal time constant

28.2 THERMAL DETECTOR ELEMENTS¹

Introduction

Thermal detectors (transducers) of optical radiation are generally considered to be those devices that absorb the radiation, increase their own temperature, and provide a resultant electrical signal. There are several types, divided according to the physical mechanism that converts the temperature change to a resultant electrical one. The oldest are bolometers and thermocouples. The bolometer changes its electrical resistance as a result of the temperature increase; the thermocouple changes its contact potential difference. There are several different types of bolometers, including thermistor, semiconducting, superconducting, carbon, and metallic. They may also be subdivided according to whether they operate at room or cryogenic temperature. Thermocouples vary according to the materials that are joined, and are sometimes connected in series to generate thermopiles.

Pyroelectric detectors make use of the property of a change in the internal polarization as a function of the change in temperature, the pyroelectric effect. Golay cells and certain variations make use of the expansion of a gas with temperature. All of these detectors are governed by the fundamental equation of heat absorption in the material. Many reviews and two books of collected reprints² provide additional information.

Thermal Circuit Theory

In the absence of joulean heating of the detector element, the spectrum of the temperature difference $d\Delta T$ is given in terms of the spectrum of the absorbed power \tilde{P} (the power is P)

$$d\Delta\tilde{T} = \frac{\epsilon\tilde{P}}{G(1+i\omega\tau)} \quad (1)$$

where G is the thermal conductance, given by the product of the thermal conductivity times the cross-sectional area of the path to the heat sink and divided by the length of the path to that heat sink. The time constant τ is the product of the thermal resistance and the heat capacitance. The thermal resistance is the reciprocal of the thermal conductance, while the thermal capacitance is the thermal capacity times the mass of the detector. In the absence of joulean heating, this is a simple, single time constant thermal circuit, for which the change in temperature is given by

$$d\tilde{T} = \frac{\epsilon\tilde{P}}{G(1+i\omega\tau)} \quad (2)$$

The absorbed power is equal to the incident power times the absorptance α of the material:

$$\tilde{P} = \epsilon\tilde{P}_i \quad (3)$$

The absorptance α is usually written as ϵ (which is legitimate according to Kirchhoff's law) since α is also used for the relative temperature coefficient of resistance (some writers use η).

$$\alpha = \frac{1}{R} \frac{dR}{dT} \quad (4)$$

As radiation is absorbed, part of the heat is conducted to the sink. Some of it gives rise to an increase in temperature. Some is reradiated, but this is usually quite small and is ignored here. The dc responsivity of a thermal detector is proportional to the emissivity and to the thermal resistance. The greater proportion of radiation that is absorbed, the greater will be the responsivity. The less heat that is conducted to the sink, the greater will be the temperature rise. The time constant is a true thermal time constant, the product of thermal resistance and capacitance. The greater the heat capacitance, the more heat necessary for a given temperature increase, and the less heat conducted to the sink, the more available for temperature increase. A high absorptance is accomplished by the use of a black coating, and a sufficient amount of it. Thus, there is a direct conflict between high speed and high responsivity.

The Ideal Thermal Detector³⁻⁶

The ideal thermal detector has a noise that is associated only with the thermal fluctuations of the heat loss to the heat sink, and this coupling is purely radiative. Then the noise equivalent power (NEP) is given by

$$\text{NEP} = \sqrt{16A\sigma kT^5/\epsilon} \quad (5)$$

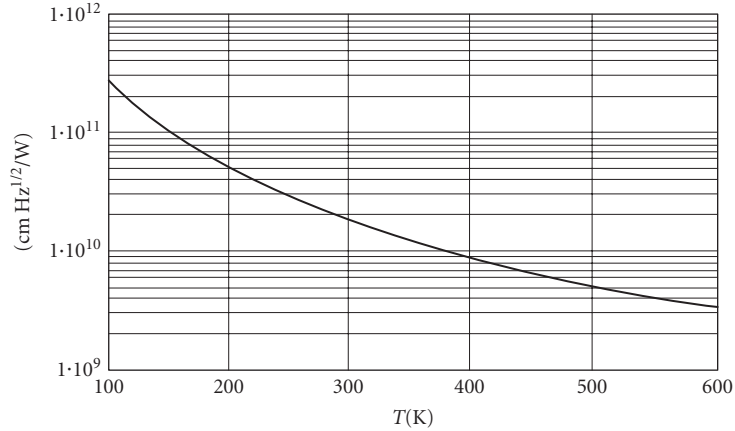


FIGURE 1 Theoretical specific detectivity for ideal thermal detectors.

where it is assumed that the detector is irradiated by a hemisphere of blackbody radiation at the same temperature T as the detector. The corresponding specific detectivity, assuming that the signal varies as the area and the noise as its square root, is

$$D^* = \frac{\epsilon^{1/2}}{4\sqrt{\sigma k T^5}} \quad (6)$$

where the detector and background are at the same temperature.

For circumstances in which the detector is in a cooled chamber, the total radiation from the sources at various temperatures must be calculated. Figure 1 shows the specific detectivity of a background limited ideal thermal detector as a function of the temperature of the surround.

No detector is ideal, and every one will be limited by the signal loss due to incomplete absorption at the surface and any transmission losses by the optical system that puts the radiation on the detector. The detector will also have noise that arises from its conductive coupling to the heat sink, and probably Johnson noise as well. The conductive mean square power fluctuation is given by

$$\langle P^2 \rangle = 4kT^2G \quad (7)$$

The Johnson noise power density is $4kT$. Therefore, the total mean square power fluctuation is given by

$$\langle P^2 \rangle = 4kT[GT + 4\epsilon A \sigma T^4 + 1] \quad (8)$$

Bolometers

Most single-element bolometers are connected in a voltage divider network, as shown in Fig. 2. A stable voltage supply is used to develop a current and consequent voltage drop across the two resistors. One is the detector, while the other should be a matching element to eliminate signals arising from a change in the ambient temperature. It should match the detector in both resistance and in the

temperature coefficient of resistance. Usually another, but blinded, detector is used. The expression for power conservation is

$$\begin{aligned}
 C \frac{d\Delta T}{dt} + G\Delta T &= \frac{d(i^2 R)}{dt} \Delta T + P \\
 C \frac{d\Delta T}{dt} + G\Delta T &= \frac{V^2(R_1 - R)}{(R_1 + R)^3} \frac{dR}{dT} \Delta T + P \\
 C \frac{d\Delta T}{dt} + \left[G - \frac{V^2 R \alpha}{(R_1 + R)^2} \frac{(R_1 - R)}{(R_1 + R)} \right] \Delta T &= P
 \end{aligned} \tag{9}$$

The solution to this is a transient that has an RC time constant, where R is the reciprocal of the bracketed term, and C is the thermal capacitance, and the same steady-state term given above. The transient decays as long as G is greater than the rest of the bracket, but the detector burns up if not. This is still another reason for matching the resistances. The dc responsivity is a function of the construction parameters, including the path to the sink, the bias voltage, and the relative change of resistance with temperature. The different bolometers are divided according to how their resistances change with temperature. (R_1 and R represent slightly different values of R_D .)

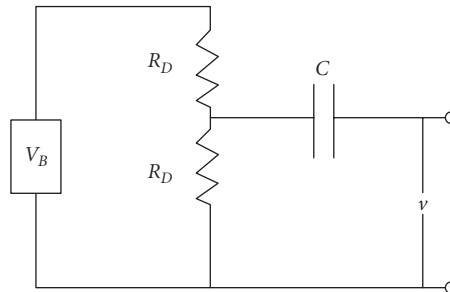


FIGURE 2 Balanced voltage divider circuit for a thermal detector.

Metal Bolometers These have a linear change in resistance with temperature that may be expressed as

$$R = R_0 [1 + \gamma(T - T_0)] \tag{10}$$

Therefore the thermal coefficient is

$$\alpha = \frac{\gamma}{1 + \gamma(T - T_0)} \tag{11}$$

This coefficient always decreases with temperature, and burnout does not occur. The coefficient is approximately equal to the inverse of the temperature, and is therefore never very high.

Semiconductor Bolometers These have an exponential change of resistance with temperature, given by

$$R = R_0 e^{\beta/T} \tag{12}$$

so that

$$\alpha = -\beta/T^2 \tag{13}$$

The value of β depends upon the particular material. These detectors can burn out. Two basic types exist: (1) those that are used at low temperatures and (2) those that are used at about room temperature.

The most used low-temperature bolometer⁷ is germanium in a bath of liquid helium. Pure germanium is transparent in the infrared, but with enough compensated doping it becomes a good conductor with a high-temperature coefficient of resistance.⁸ Typical concentrations are about 10^{16}cm^{-3} of gallium and 10^{15} of indium. Even these are not sufficient at wavelengths shorter than $10\ \mu\text{m}$ since the free-carrier absorption is proportional to wavelength. In such a case a black coating is sometimes used. Improvements have been made since Low's first work.⁹⁻¹¹

Superconducting Bolometers These make use of the extremely large thermal coefficient of resistance at the transition temperature.¹²⁻¹⁴ Originally they needed to be controlled very carefully, or a small change in ambient conditions (on the order of 0.01 K) could cause an apparent signal of appreciable magnitude. A more recent version¹⁵ incorporates an evaporated thin film on an anodized aluminum block that is coupled to a helium bath by a brass rod. The detector has a time constant of about $3\ \mu\text{s}$ due to this high thermal conductance and a good NEP of about $10^{-13}\ \text{WHz}^{-1/2}$. It still must be controlled to about $10^{-5}\ \text{K}$, and this is accomplished with a heater current and control circuit.

Recently developed materials not only have high-temperature transition points but also have more gradual transitions, and provide a better compromise between good responsivity and the requirement for exquisite control.¹⁶

Carbon Bolometers These are a form of semiconductor bolometers that have been largely superseded by germanium bolometers. They are made of small slabs of carbon resistor material, connected to a metal heat sink by way of a thin mylar film. Although their responsivities are comparable to germanium bolometers, their noise is several orders of magnitude higher.¹⁷

Thermocouples and Thermopiles^{18,19}

A *thermocouple* is made by simply joining two dissimilar conductors. A good pair has a large relative Seebeck coefficient and gives rise to a potential difference. The materials also have large electrical conductivities and small thermal ones, so there is little voltage drop across the length and a small thermal gradient. Although there are many different couples (many are not even used for radiation detection), those most often used for this application are bismuth telluride, copper, and constantan. The expression for the responsivity is given in terms of the relative Seebeck coefficient S_{12} (the difference in the voltage change with temperature between the two materials) and the expression derived above for the thermal circuit

$$\mathfrak{R} = \frac{S\epsilon}{G(1+i\omega\tau)} \quad (14)$$

Good materials are those that have a large Seebeck coefficient, a high electrical conductivity, and a small thermal conductivity, and the figure of merit is often defined as

$$Z_{12} = \frac{S_{12}^2}{\left[\sqrt{G_1/\sigma_1} + \sqrt{G_2/\sigma_2}\right]^2} \quad (15)$$

Thermopiles are arrays of thermocouples connected in series. They are manufactured in two ways. Some are carefully wound wires with junctions aligned in the desired pattern, while others are evaporated with the pattern determined by masking operations. Most of the "bulk" thermopiles are wrapped on appropriate mandrels to obtain rigidity. Both kinds are obtainable in a variety of sizes

and patterns that correspond to such things as spectrometer slits, centering annuli, and staggered arrays for moderate-sensitivity thermal imaging.

The Golay Cell²⁰

This detector is used mostly for laboratory operations, as it is slow and fragile, although it has high sensitivity. It is a gas-filled chamber that has a thin membrane at one end and a blackened detector area at the other. Light on the blackened surface causes the increase in temperature; this is transferred to the gas which therefore expands. The membrane bulges, and the amount of the bulge is sensed by some sort of optical lever²¹ or even change in capacity of an electrical element.²² Other versions do not use a blackened surface, but allow the radiation to interact with the gas directly, in which case they are spectral detectors that are “tuned” to the absorption spectrum of the gas.²³

Pyroelectric Detectors²⁴

Some crystals which do not have a center of symmetry experience an electric field along a crystal axis. This internal electric field results from an alignment of electric dipoles (known as polarization), and is related to the crystal temperature. In these ferroelectric crystals, this results in a charge being generated and stored on plates connected to the crystal. Polarization disappears above the so-called Curie temperature that is characteristic of each material. Thus, below the Curie temperature, a change in temperature results in a current. The equation for the response of a pyroelectric detector is

$$\mathfrak{R} = \frac{\omega p A_d \epsilon R_e R_{th}}{(1 + i\omega\tau_{th})(1 + i\omega\tau_e)} \quad (16)$$

where ω is the radian frequency, p is the pyroelectric coefficient, A_d is the detector area, R_e is the electrical resistance, R_{th} is the thermal resistance, τ_{th} is the thermal time constant, and τ_e is the electrical time constant. The relation is shown in Fig. 3, where the responsivity is plotted as a function of frequency. In the low-frequency region the responsivity rises directly as the frequency. This is a result of the ac operation of a pyroelectric. At the (radian) frequency that is the reciprocal of the slower (usually the thermal) time constant, the response levels off. This is the product of the ac rise and the thermal rolloff. Then, when the frequency corresponding to the shorter time constant is reached, the response rolls off.

Type II pyroelectric detectors work on a slightly different mechanism, which is still not fully understood. The electrodes are on the sensing surface of the detector and parallel to the polar axis. In these crystals, the temperature change is not uniform at the onset of radiation and the primary

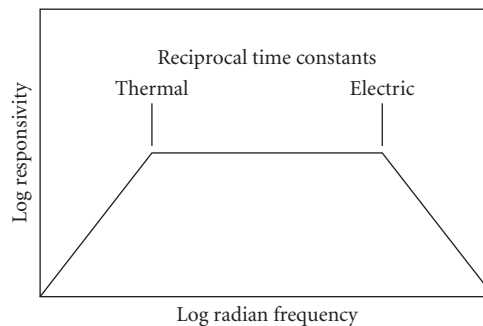


FIGURE 3 Responsivity asymptotes versus frequency.

TABLE 1 General Properties of Thermal Detectors

Type	Operating Temperature (K)	$D^* \times 10^8$ (cmHz ^{1/2} W ⁻¹)	NEP $\times 10^{-10}$ (WHz ^{-1/2})	Time Constant (m)	Size (mm ²)
Silicon bolometer	1.6		3×10^{-5}	8	0.25–0.70
Metal bolometer	2–4	1		10	
Thermistor bolometer	300	1–6		1–8	0.01–10
Germanium bolometer	2–4		0.005	0.4	1.5
Carbon bolometer	2–4		0.03	10	20
Superconducting bolometer (NbN)	15		0.2	0.5	5×0.25
Thermocouples	300		2–10	10–40	0.1 \times 1 to 0.3 \times 3
Thermopiles	300			3.3–10	1–100
Pyroelectrics	300	2–5		10–100 [†]	2 \times 2
Golay cell	300	10	0.6	10–30	10

[†]Shorter values can be obtained at the expense of NEP (for laser detection).

and secondary pyroelectric effects take place, thereby generating a body electric charge distribution.²⁵ Materials most often used for these detectors are TGS (triglycine sulfate), DTGS (deuterated TGS), Li₂SO₄, LiNbO₃, LiTaO₃, and PLZT (lead lanthanum zirconate titanate). TGS is the most used for specialized sensor systems, but has a relatively low Curie point. For higher-temperature operation, usually LiTaO₃ or PLZT is used in the general laboratory environment.

The two advantages of the pyroelectric detector over the other thermal detectors, bolometers, and thermopiles, are its responsivity and its capability of rapid response. The response time and responsivity are traded by choice of the load resistor in the circuit. For instance, with a 100-M Ω load the time constant can be 1 ms and the responsivity 100 V/W, but with a 1-M Ω load the values would be 10 μ s and 1 V/W.

Summary of Elemental Thermal Detector Properties

Although the user should contact suppliers for detailed information, this section provides overall property information about thermal detectors. There are several cautions about summary data. Most detectors can be tailored to have somewhat different properties. Improvements have often been made since the publication of these results. Not all parameters are available in all combinations. Table 1 does, however, give the general flavor of the performance of different thermal detectors.

28.3 ARRAYS

Introduction

As pointed out earlier, thermal detector response is governed by the thermal response time, which is the ratio of the pixel heat capacity C to the thermal conductance G of the heat leakage mechanism. High pixel responsivity is associated with high thermal isolation, i.e., low thermal conductance. Thermal detector design is driven by the thermal isolation structure. It is the structure which determines the extent to which the pixel performance can approach the temperature fluctuation noise limit and, ultimately, the background fluctuation noise limit. Given the value of G associated

with the heat loss mechanism, the pixel heat capacity must be designed appropriately to attain the required thermal response time. Response times in the millisecond range are compatible with high thermal isolation; response times in the microsecond range are not. Thus, two-dimensional arrays of thermal detectors which operate at TV frame rates (30 Hz in the United States) are under development for applications in thermal imagers.

Noise Equivalent Temperature Difference

Whereas elemental detectors are usually described by such figures of merit as NEP and D^* , arrays have been described by a noise equivalent temperature difference (NETD) associated with their use in a camera under certain specific conditions. It is defined as the change in temperature of a blackbody which fills the field of view of a pixel of an infrared imaging system that gives rise to a change of unity in the signal-to-noise ratio at the output of the system. The measurement of the NETD should, however, be with the flooding of several pixels to avoid fringing effects and with an SNR (signal-to-noise ratio) well above 1 to obtain good accuracy. The pixel is defined as the subtense of a single element of the array. The NETD can be written in several different forms. Perhaps the simplest is

$$\text{NETD} = \frac{\sqrt{A_d B}}{D^* (dP_d / dT)} \quad (17)$$

where D^* is the specific detectivity, A_d is the area of a single pixel, B is the system bandwidth and (dP_d / dT) represents the change in power on the detector element per unit change in temperature in the spectral band under consideration. This form does not include the system noise, which is often included by the manufacturers in their calculations. In Eq. (17) the change in power with respect to temperature is

$$\frac{dP_d}{dT} = \frac{A_d \tau_o}{4FN^2} \int_{\lambda_1}^{\lambda_2} \frac{dM}{dT} d\lambda \quad (18)$$

where τ_o is the optics transmission, FN is the focal ratio (defined as the effective focal length divided by the entrance pupil diameter), and M is the radiant emittance of the source. This is almost the definition of the specific detectivity. The NETD can also be written in terms of the responsivity \mathfrak{R} , since the detectivity and responsivity are related in the following way:

$$D^* = \frac{\sqrt{A_d B}}{P} \frac{V_s}{V_N} = \sqrt{A_d B} \frac{\mathfrak{R}}{V_N} \quad (19)$$

where V_s is the signal voltage at the sensor and V_N is the rms noise voltage of a pixel in the bandwidth B . Therefore

$$\text{NETD} = \frac{V_N}{\mathfrak{R} (\partial P_d / \partial T)_{\lambda_1 - \lambda_2}} \quad (20)$$

The power on the detector is related to the power on the aperture by the optical transmission τ_o . The expression can also be formulated in terms of the source radiance, L

$$\text{NETD} = \frac{4FN^2 \sqrt{B}}{D^* \tau_a \tau_o \pi D \Delta \theta (\partial L / \partial T)_{\lambda_2 - \lambda_1}} \quad (21)$$

where D is the diameter of the aperture, $\Delta\theta$ is the angular subtense of a pixel, L is the source radiance, and τ_a is the atmospheric transmission. One last form can be generated by recognizing that, for an isotropic radiator, the radiance is the radiant emittance divided by π :

$$\text{NETD} = \frac{4FN^2V_N}{A_D \tau_a \tau_o \mathfrak{R}(T_s) (\partial M / \partial T)} \quad (22)$$

In this form of the expression for NETD, it is not necessary that the noise be white, nor is it necessary that the noise not include system noise. Whether or not system noise is included should be clearly stated.

Theoretical Limits

Figure 4 illustrates the theoretical limits of thermal arrays having the parameters shown and operating at 300 and 85 K as a function of thermal conductance. The performance of real thermal arrays with those parameters lies on or above the sloping line. As the conductance G is reduced (better thermal isolation), the noise equivalent temperature difference NETD is reduced (improves) until the background limit is reached, when radiant power exchange between the array and the background becomes the dominant heat transfer mechanism. Reducing the detector temperature to 85 K appropriate to a bolometer operating at the transition edge of the high-temperature superconductor $\text{YBa}_2\text{Cu}_3\text{O}_{7-x}$ (YBCO) reduces the NETD by $\sqrt{2}$ and allows the limit to be reached with less thermal isolation (higher G value).

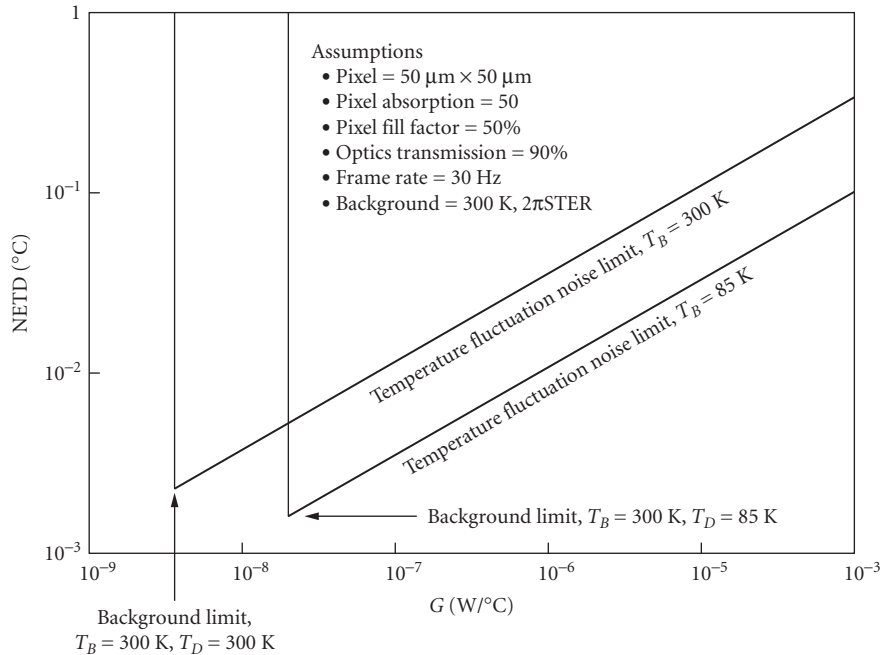


FIGURE 4 Temperature fluctuation noise limit and background fluctuation noise limit of uncooled and cryogenic thermal detector arrays.

Arrays fall into two categories: monolithic and hybrid. Monolithic arrays are prepared on a single substrate, e.g., silicon, upon which the detecting material is deposited in the form of a thin film which is subsequently processed into an array. Hybrid arrays are prepared in two parts: (1) the read-out electronics arrays, usually in silicon, and (2) the detecting material array, usually in wafer form which is thinned by lapping, etching, and polishing. These two arrays are mated by a technique such as flip-chip bonding. Here the interconnection at each pixel must have a sufficiently high electrical conductivity, yet a sufficiently low thermal conductivity—a difficult requirement. If array cost considerations are important, then the monolithic approach, especially in silicon, is the more desirable.

Resistive Bolometer Arrays

The development of resistive bolometric arrays has proceeded along two paths: uncooled arrays and cryogenic arrays. Large, uncooled bolometric arrays have been developed at Honeywell by a team lead by R. A. Wood.^{26,27} Silicon microstructure technology is employed to produce the arrays, a process resembling the fabrication of integrated circuits. Twelve arrays are prepared on a 4-in-diameter silicon wafer. Each monolithic array consists of 240×336 pixels; each pixel is $50 \times 50 \mu\text{m}$. The detecting material is a thin film of vanadium oxide. A Si_3N_4 membrane having a thermal conductance of $1 \times 10^{-7} \text{ WC}^{-1}$ supports the vanadium oxide at each pixel, as shown in Fig. 5. Bipolar transistors implanted in the silicon substrate act as pixel switches for the matrix-addressed array. The response is optimized for the 8- to $14\text{-}\mu\text{m}$ spectral interval. The thermal response time is adjusted for a 30-Hz frame rate. Each pixel is addressed once per frame by a $5\text{-}\mu\text{s}$ pulse. A thermoelectric stabilizer maintains the array at ambient temperature. Other than a one-shot shutter, the camera has no moving parts.

The measured NETD of the camera with F/1 optics at 300 K is 0.04 K. Given the G value of $1 \times 10^{-7} \text{ WK}^{-1}$ it can be seen from Fig. 4 that the array is within a factor of 4 of the temperature fluctuation noise limit. Furthermore, the pixel thermal isolation is so complete that there is no measurable thermal spreading among the pixels.

Linear resistive bolometric arrays of the high-temperature superconductor YBCO on silicon microstructures have been prepared by Johnson et al.,²⁸ also of Honeywell. A two-dimensional array is under development.²⁹ The monolithic arrays operate at the transition edge from 70 to 90 K. As was true for the uncooled arrays, the superconducting ones employed a silicon nitride membrane to support the thin film and provide thermal isolation. Excess noise at the contacts limited the performance

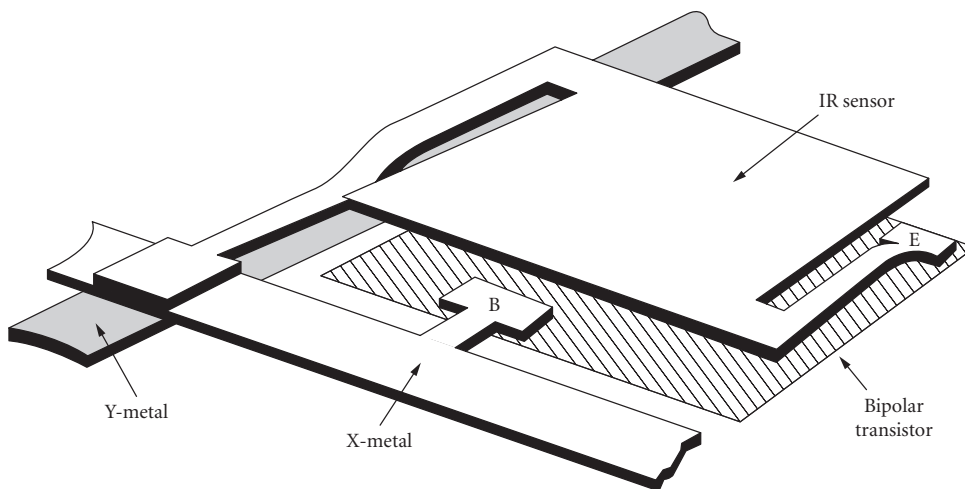


FIGURE 5 Monolithic silicon microbolometer.^{26,27} (© 1992 IEEE.)

of the 12-element linear array. With no excess noise, the calculated NETD²⁹ of a 240×336 array with $50 \mu\text{m}$ pixels and F/1 optics would be 0.002 K , which is near the 300 K background limit, as shown in Fig. 4.

Pyroelectric Hybrid Arrays

Linear and two-dimensional pyroelectric uncooled arrays have been under development for more than two decades.^{20–32} The arrays employ hybrid construction, in which a bulk pyroelectric ceramic material such as lithium tantalate or lead zirconate is mechanically thinned, etched, and polished, then bump-bonded to a silicon substrate containing readout electronics,³³ as shown in Fig. 6. Reticulation is usually employed to prevent lateral heat conduction through the pyroelectric material. The theoretical system NETD of a two-dimensional uncooled array with F/1 optics is estimated to be 0.1 K .³³ Two-dimensional uncooled arrays operating in the $8\text{--}14\text{-}\mu\text{m}$ region having 100×100 pixels, each $100 \times 100 \mu\text{m}$, are available commercially.³⁴ Their NETD (with F/1 speed) is 0.35 K . A two-dimensional pyroelectric monolithic array employing a thin film of lead titanate on a silicon microstructure is under development.³⁵

Ferroelectric bolometer arrays, also known as field-enhanced pyroelectric arrays, have been developed by Texas Instruments.^{36,37} Operation depends upon the temperature dependence of the spontaneous polarization and dielectric permittivity in a ferroelectric ceramic near the Curie temperature. Barium strontium titanate (BST), the selected material, has its composition (barium-to-strontium ratio) adjusted during preparation so the Curie point is 22°C . A thermoelectric stabilizer is employed to hold the BST near 22°C such that the absorbed infrared radiation changes the temperature and thus the dielectric properties. The effect is similar to the pyroelectric effect; however, a voltage is applied to enhance the signal. Construction of this array is naturally similar to that of the pyroelectric array, described above, as shown in Fig. 7. Reticulation of the ceramic is frequently applied to these arrays as well. A radiation chopper is required as both the pyroelectric and ferroelectric effects depend upon the change in temperature. The Texas Instruments BST array,

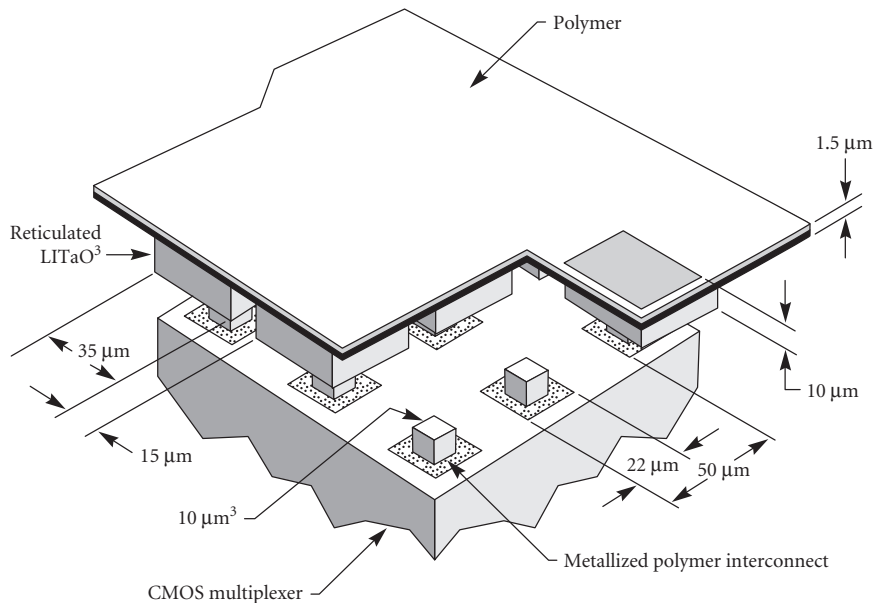


FIGURE 6 Hybrid pyroelectric array structure.³³

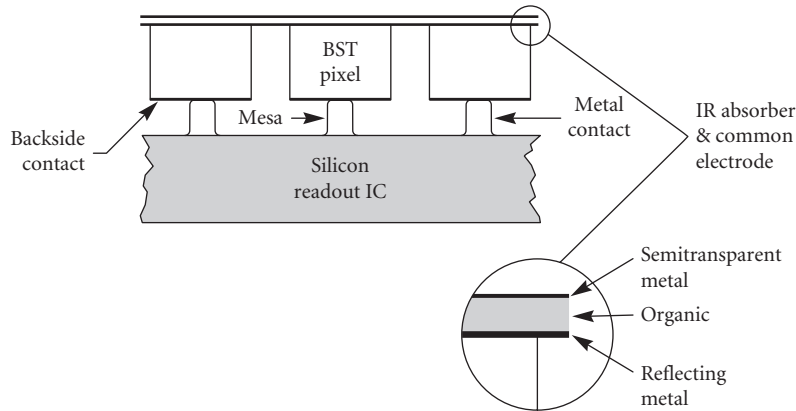


FIGURE 7 Ferroelectric bolometer hybrid array.^{36,37}

incorporating 80,000 pixels, each about $50 \times 50 \mu\text{m}$, which are matrix addressed, has an NETD of less than 0.1°C (with $F/1$ optics).

Thermoelectric Arrays

Thermoelectric arrays prepared by silicon micromachining have been described by Choi and Wise.³⁸ Series-connected, thin-film thermocouples, i.e., a thermopile, are prepared on a silicon microstructure, the “hot” junctions (receiving the thermal radiation) on a silicon nitride/silicon dioxide membrane and the “cold” shield junctions on the surrounding silicon substrate. Both 64- and 96-pixel microthermopile linear arrays in silicon microstructures have been prepared by Honeywell,³⁹ each microthermopile consisting of several nickel iron/chromium thermocouples connected in series, as shown in Fig. 8. The “hot” junctions are deposited on silicon nitride membranes, whereas the “cold” junctions are on the silicon substrates. A camera incorporating the linear array has been employed to image moving targets such as automobiles. With an $F/0.73$ lens, the measured NETD is 0.10°C .

Since the first publication of this *Handbook*, many advances have been made in these arrays. The suppliers have improved sensitivity somewhat but have increased the number of pixels and decreased their size. The reader should check with the manufacturers for the latest information.

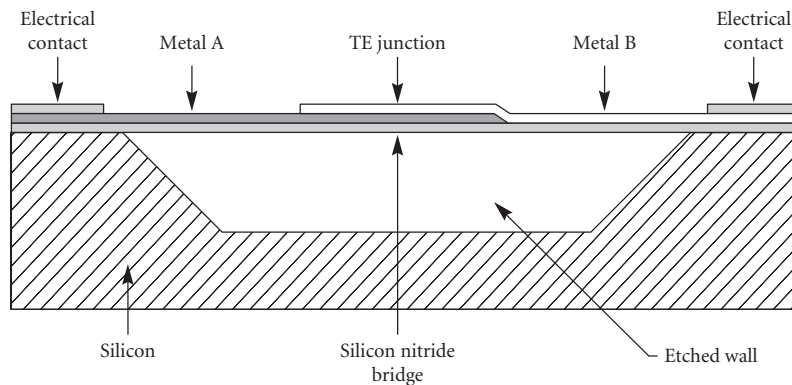


FIGURE 8 Monolithic thermoelectric array.³⁹ (© 1991 Instrument Society of America. Reprinted with permission from the Symposium for Innovation in Measurement Science.)

28.4 REFERENCES

1. P. W. Kruse, L. D. McGlauchlin, and R. B. McQuistan, *Elements of Infrared Technology*, John Wiley and Sons, New York, 1962.
2. R. D. Hudson and J. W. Hudson, *Infrared Detectors*, Halsted Press, New York, 1975; and F. R. Arams, *Infrared to Millimeter Wave Detectors*, Artech House, Dedham, Mass., 1973.
3. E. H. Putley, *Optical and Infrared Detectors*, R. J. Keyes (ed.), Springer-Verlag, Berlin, 1980, chap. 3.
4. R. A. Smith, F. E. Jones, and R. P. Chasmar, *The Detection and Measurement of Infrared Radiation*, Oxford University Press, 1968.
5. R. C. Jones, "The Ultimate Sensitivity of Radiation Detectors," *J. Opt. Soc. Am.* **37**:879 (1974).
6. S. Nudelman, "The Detectivity of Infrared Detectors," *Appl. Opt.* **1**:627 (1962).
7. F. J. Low, "Low Temperature Germanium Bolometer," *J. Opt. Soc. Am.* **51**:1300 (1961).
8. S. Zwerdling, R. A. Smith, and J. P. Thierault, "A Fast High Responsivity Bolometer for the Very Far Infrared," *Infrared Physics* **8**:271 (1968).
9. N. Coron, "Infrared Helium Cooled Bolometers in the Presence of Background Radiation: Optimal Parameters and Ultimate Performance," *Infrared Physics* **16**:411 (1976).
10. N. Coron, G. Dambier, and J. Le Blanc, *Infrared Detector Techniques for Space Research*, V. Manno and J. Ring (eds.), Reidel, Dordrecht, 1972.
11. N. Coron, G. Dambier, J. Le Blanc, and J. P. Moliac, "High Performance, Far Infrared Bolometer Working Directly in a Helium Bath," *Rev. Sci. Instr.* **46**:492 (1975).
12. W. H. Andrews, R. M. Milton, and W. De Sorbo, "A Fast Superconducting Bolometer," *J. Opt. Soc. Am.* **36**:518 (1946).
13. R. M. Milton, "A Superconducting Bolometer for Infrared Measurements," *Chem. Rev.* **39**:419 (1946).
14. N. J. Fuson, "The Infrared Sensitivity of Superconducting Bolometers," *J. Opt. Soc. Am.* **38**:845 (1948).
15. G. Gallinaro and R. Varone, "Construction and Calibration of a Fast Superconducting Bolometer," *Cryogenics* **15**:292 (1975).
16. K. B. Bhasin and V. O. Heinen (eds.), "Superconductivity Applications for Infrared and Microwave Devices," *Proc. SPIE* **1292** (1990). (Includes many other references.)
17. W. S. Boyle and K. F. Rodgers, *J. Opt. Soc. Am.* **49**:66 (1959).
18. D. F. Hornig and B. J. O'Keefe, "Design of Fast Thermopiles and the Ultimate Sensitivity of Thermal Detectors," *Rev. Sci. Instr.* **18**:7 (1947).
19. P. B. Felgett, "Dynamic Impedance and the Sensitivity of Radiation Thermocouples," *Proc. Phys. Soc.* **B62**:351 (1949).
20. M. J. E. Golay, "Theoretical Considerations in Heat and Infrared Detection with Particular Reference to the Pneumatic Detector," *Rev. Sci.* **18**:347 (1947); "Pneumatic Infrared Detector," *ibid.*, **18**:357 (1947); "Theoretical and Practical Sensitivity of the Pneumatic Detector," *ibid.*, **20**:816 (1949).
21. J. R. Hickey and D. B. Daniels, "Modified Optical System for the Golay Detector," *Rev. Sci. Instr.* **40**:732 (1969).
22. M. Chatanier and G. Gauffre, *IEEE Transactions Instr. and Meas.* **IMEE** **179** (1973).
23. A detector once made by Patterson Moos and cited by R. DeWaard and E. Wormser, "Description and Properties of Various Thermal Detectors," *Proc. IRE* **47**:1508 (1959).
24. E. H. Putley, *Semiconductors and Semimetals*, vol. 5, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1970, chap. 6, "The Pyroelectric Detector;" vol. 12, 1977, chap. 7.
25. Zu-Sheng Wang and Jian-Qi Zhang, "The Mechanism of Type II Pyroelectric Detectors," *Infrared Phys.* **33**(6):481–486 II (1993).
26. R. A. Wood, C. J. Han, and P. W. Kruse, "Integrated Uncooled Infrared Detector Imaging Array," *Proc. of the 1992 IEEE Solid State Sensor and Actuator Workshop*, Hilton Head Island, S.C., pp. 132–135.
27. R. A. Wood, "Uncooled Thermal Imaging with Monolithic Silicon Focal Plane Arrays," *Proc. SPIE* **2020**: Infrared Tech. XIX (1993).
28. B. R. Johnson, T. Ohnstein, C. J. Han, R. Higashi, P. W. Kruse, R. A. Wood, H. Marsh, and S. B. Dunham, "High- T_c Superconductor Microbolometer Arrays Fabricated by Silicon Micromachining," *IEEE Trans. Appl. Superconductivity* **3**:2856 (1993).

29. B. R. Johnson and P. W. Kruse, "Silicon Microstructure Superconducting Microbolometer Infrared Arrays," *Proc SPIE* **2020**:Infrared Technology XIX (1993).
30. E. H. Putley, "The Pyroelectric Detector," *Semiconductors and Semimetals*, vol. 5, *Infrared Detectors*, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1970.
31. P. A. Manning, D. E. Burgess, and R. Watton, "A Linear Pyroelectric Array IR Sensor," *Proc SPIE* **590**:2 (1985).
32. R. Watton and M. V. Mansi, "Performance of a Thermal Imager Employing a Hybrid Pyroelectric Detector Array with MOSFET Readout," *Proc. SPIE* **865**:79 (1987).
33. N. Butler and S. Iwasa, "Solid State Pyroelectric Imager," *Proc SPIE* **1685**:146 (1992).
34. GEC-Marconi Materials Technology Ltd., 9360 Ridgehaven Court, San Diego, CA 92123.
35. B. E. Cole, R. D. Horning, and P. W. Kruse, "PbTiO₃ Deposited by an Alternating Dual-Target Ion-Beam Sputtering Technique," *Ferroelectric Thin Films II*, A. I. Kingon, E. R. Myers, and B. Tuttle (eds.), *Materials Research Society Symposium Proc.*, **243**:185 (1992).
36. C. Hanson, H. Beratan, R. Owen, M. Corbin, and S. McKenney "Uncooled Thermal Imaging at Texas Instruments," *Infrared Detectors: State of the Art, Proc. of SPIE* **1735**:17 (1992).
37. C. M. Hanson, "Uncooled Ferroelectric Thermal Imaging," *Proc. SPIE* **2020**: Infrared Technology XIX (1993).
38. I. H. Choi and K. D. Wise, "A Silicon-Thermopile-Based Infrared Sensing Array for Use in Automated Manufacturing," *IEEE Trans, on Electron Devices* **ED-33**:72 (1986).
39. M. Listvan, M. Rhodes, and M. L. Wilson, "On-Line Thermal Profiling for Industrial Process Control," *Proc. of the Instrument Society of America, Symposium for Innovation in Measurement Science*, Geneva, NY, August 1991.

PART

6

**IMAGING
DETECTORS**

This page intentionally left blank.

Joseph H. Altman

*Institute of Optics
University of Rochester
Rochester, New York*

29.1 GLOSSARY

A	area of microdensitometer sampling aperture
a	projective grain area
D	optical transmission density
D_R	reflection density
DQE	detective quantum efficiency
$d(\mu)$	diameter of microdensitometer sampling aperture stated in micrometers
E	irradiance/illuminance (depending on context)
\mathcal{G}	Selwyn granularity coefficient
g	absorbance
H	exposure
IC	information capacity
M	modulation
M_e	angular magnification
m	lateral magnification
NEQ	noise equivalent quanta
$P(\lambda)$	spectral power in densitometer beam
Q'	effective Callier coefficient
q	exposure stated in quanta/unit area
R	reflectance
S	photographic speed
$S(\lambda)$	spectral sensitivity
S/N	signal-to-noise ratio of the image
T	transmittance
$T(\nu)$	modulation transfer factor at spatial frequency ν

t	duration of exposure
$WS(\nu)$	value of Wiener (or power) spectrum for spatial frequency ν
γ	slope of D-log H curve
ν	spatial frequency
$\rho(\lambda)$	spectral response of densitometer
$\sigma(D)$	standard deviation of density values observed when density is measured with a suitable sampling aperture at many places on the surface
$\sigma(T)$	standard deviation of transmittance
$\phi(\tau)$	Autocorrelation function of granular structure

29.2 STRUCTURE OF SILVER HALIDE PHOTOGRAPHIC LAYERS

The purpose of this chapter is to review the operating characteristics of silver halide photographic layers. Descriptions of the properties of other light-sensitive materials, such as photoresists, can be found in Ref. 4.

Silver-halide-based photographic layers consist of a suspension of individual crystals of silver halide, called *grains*, dispersed in gelatin and coated on a suitable “support” or “base.” The suspension is termed an *emulsion* in the field. The grains are composed of AgCl, AgClBr, AgBr, or AgBrI, the listing being in order of increasing sensitivity. Grain size ranges from less than 0.1 μm (“Lippmann” emulsions) to 2 to 3 μm , depending on the intended use of the coating. The number of grains per square centimeter of coating surface is usually very large, of the order of 10^6 to 10^8 grains/ cm^2 . The weights of silver and gelatin coated per unit area of support vary depending on intended use; usually both fall in the range 1 to 10 g/m^2 . The silver-to-gel ratio may also vary depending on intended use. Typically, the emulsion may be about 30 to 40 percent silver by weight, but some special-purpose materials, such as films to record Schumann-wavelength-region radiation, contain very little gelatin.

For modern materials, both the emulsion and the coating structure can be very complex. The emulsion layer is much more than silver halide in gelatin, containing additional agents such as hardeners, antifoggants, fungicides, surfactants, static control agents, etc. Likewise, the coating structure may be very complex. Even some black-and-white materials consist of layers of two different emulsions coated one over the other and a thin, clear layer of gelatin is often coated over the emulsion(s) to provide some mechanical protection. In the case of color films, as many as 15 layers may be superimposed, some of them of the order of 1 μm thick. The thickness of the complete coating varies from about 3 μm to about 25 μm in normal films.

Commercially available emulsions are coated on a variety of glass, plastic (film), and paper supports (or “bases”). Glass is used for mechanical rigidity, spatial stability, or surface flatness.

Two different types of plastic are available commercially as film supports: cellulose acetate and polyethylene terephthalate (trade names “Cronar,” “Estar,” and “Mylar”). Of the two types, Mylar is superior in strength, flexibility, and spatial stability. However, the material is birefringent and its physical properties may be different in orthogonal directions. Also, these directions may not be aligned with the length or width of the sample. The anisotropic properties arise from the method of manufacture. Although not as tough as Mylar, cellulose acetate is, of course, fully adequate for most purposes. Also, this material is isotropic and easier to slit and perforate. Typical supports for roll films are around 4 mils (102 μm) thick, and for sheet films, 7 mils (178 μm), and other thicknesses are available. Most films are also coated on the back side of the support. The “backing” may be a layer of clear gelatin applied for anticurl protection, or of gelatin dyed with a dye that bleaches during processing, and provides both anticurl and antihalation protection. Lubricants and antistatic agents may also be coated, either on the front or back of the film. Properties of supports are discussed in Ref. 1.

29.3 GRAINS

The grain is the radiation-sensing element of the film or plate. It is a face-centered cubic crystal, with imperfections in the structure. For the most part, the grains act as independent receptors. In general, the larger grains are faster. The properties of the individual grains are controlled by the precipitation conditions and the after-precipitation treatment. Details of these matters are proprietary, but some discussion is given in Refs. 2 and 3. From the user's standpoint, the important fact is that when the grain is exposed to sufficient radiation it forms a "latent-image speck" and becomes *developable* by a solid-state process called the "Gurney-Mott mechanism." An excellent review of grains and their properties is given by Sturmer and Marchetti in chap. 3 of Ref. 4.

29.4 PROCESSING

The exposed halide layer is converted to a usable image by the chemical processes of development and fixation.

Development consists of reducing exposed crystals from silver halide to metallic silver, and a developing agent is an alkaline solution of mild reducer that reduces the grains having latent image specks, while not attacking the unexposed grains. Generally, once development starts the entire grain is reduced if the material is allowed to remain in the developer solution. Also, in most cases adjacent grains will not be affected, although developers can be formulated that will cause adjacent grains to be reduced ("infectious developers"). The number of quanta that must be absorbed by a grain to become developable is relatively small, of the order 4 to 40, while the developed grain contains on the order of 10^6 atoms. The quantum yield of the process is thus very high, accounting for the speed of silver-based materials.

The remainder of the process consists essentially of removing the undeveloped halide crystals which are still light-sensitive. The "fixer" is usually an acid solution of sodium thiosulfate $\text{Na}_2\text{S}_2\text{O}_3$, called "hypo" by photographers. The fixing bath usually serves as a gelatin hardener also. The thio-sulfate reacts with the halide of the undeveloped grains to form soluble silver complexes, which can then be washed out of the emulsion layer. It is worth noting that proper washing is essential for permanent images. Additional treatments to promote permanence are available. Processing is discussed in detail in Refs. 2, 3, and 5, and image permanence in Ref. 6.

The exposed and processed silver halide layer thus consists of an array of grains of metallic silver, dispersed in a gelatin matrix. In color films, the silver is removed, and the "grains" are tiny spheres of dyed gelatin (color materials will be discussed below). Either type of grain acts as an absorber; in addition, the metallic silver grains act as scatterers. The transmittance or reflectance of the layer is thus reduced and, from the user's standpoint, this change constitutes the response of the layer.

29.5 EXPOSURE

From fundamental considerations it is apparent that the dimensions of exposure must be energy per unit area. Exposure H is defined by

$$H = Et \quad (1)$$

where t is the time for which the radiation is allowed to act on the photosensitive layer, and therefore E must be the irradiance on the layer. The symbol H is used here for exposure in accordance with international standards, but it should be noted that in many publications, especially older ones, E is used for exposure and I for irradiance, so that the defining expression for exposure becomes $E = It$.

Strictly speaking, H and E in Eq. (1) should be in radiometric units. However, photographic exposures are customarily stated not in radiometric but in photometric units. This is done mostly

for historical reasons; the English scientists Hurter and Driffield, who pioneered photographic sensitometry in 1891, measured the incident flux in their experiments in lumen per square meter, or lux. Their unit of exposure was thus the lux-second (old term, meter-candle-second). Strictly speaking, of course, weighting the incident flux by the relative visibility function is wrong or at least unnecessary, but in practice it works well enough because in most cases the photographer wishes to record what he or she sees, i.e., the visible spectrum. Conversion between radiometric and photometric units is discussed by Altman, Grum, and Nelson.⁷

Also, it should be noted that equal values of the exposure product (Et) may produce different outputs on the developed film because of a number of *exposure effects* which are described in the literature.⁸ A complete review of radiometry and photometry is given in Chap. 34 in this volume of the *Handbook*.

29.6 OPTICAL DENSITY

As noted above, the result of exposure and processing is a change in the transmittance or reflectance of the layer. However, in photography, the response is usually measured in terms of the *optical density*, hereafter called the “density” in this chapter. For films (transmitting samples), density is defined by

$$D = -\log T = \log 1/T \quad (2)$$

where T is the transmittance. (Note: throughout this chapter, “log” indicates the base-10 logarithm.) For either silver grains or color grains, the *random dot* model of density predicts that

$$D = 0.434nag \quad (3)$$

where n = the number of grains per unit area of surface
 a = the average projective grain area
 g = the absorbance of the grain

Absorbance in turn is defined as $g = 1 - (T + R)$ where T and R are the transmittance and reflectance of the grain. For silver grains, the absorbance is taken as unity. The above expression, sometimes known as “Nutting’s law,” is based on a geometric approach, and does not take into account any scattering by the grains. However, of course, opaque silver particles on a clear background will act as scatterers, and in fact multiple scattering usually occurs in developed silver layers. This produces an increase in the density of such layers by a factor of 2 to 3 times from that predicted by Nutting’s law. For color films, the refractive index of the gelatin in the dyed spheres is only negligibly different from that of the surround so that such layers are not scatterers. Even in the case of silver films, however, Nutting’s law provides a useful model. Since for a given population of grains a and g will remain effectively constant, the law states that density should vary as the number of grains per unit area of surface. This fact is easily verified with a microscope.

Transmission Density

Transmission density is measured in a densitometer. It is worth noting that the device actually measures the transmittance of the sample and then displays the negative log of the result. In a normal or *macro* densitometer the sampling aperture area A is typically 1 mm² or more in size. When A is small, say, 0.1 mm² or less, the device becomes a *micro* densitometer. Microdensitometers present special problems and will be discussed below.

Because the scattered light may not reach the sensor of the densitometer when silver layers are measured, it is necessary to specify the angular subtenses of both the incident (influx) and emergent (efflux) beams at the sample. Clearly, if scattered light is lost to the sensor, the indicated density of the sample will *increase*. Four types of transmission density are described in an ISO standard,⁹ of

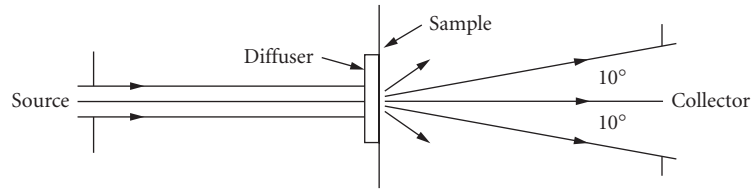


FIGURE 1 Optical system for measuring ISO/ANSI diffuse density with a 20° collection angle. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

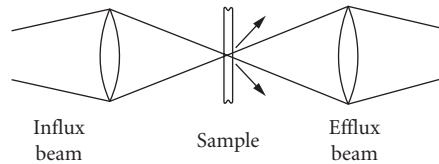


FIGURE 2 Optical conditions for projection density measurement.

which two are principally important to the user. The first of these is *diffuse density*, which is diagrammed in Fig. 1. As can be seen, a collimated incident beam illuminates an opal glass diffuser. The emulsion side of the sample is placed in contact with this diffuser, and the flux contained within a cone angle of $\pm 10^\circ$ is collected and evaluated by the sensor. The reverse of this arrangement yields the same density values and is also permitted by the standard. This is the type of density normally measured in practice. Physically, it corresponds to the conditions of contact printing.

The other case that is important in practice is projection density, which is diagrammed in Fig. 2. This case simulates the use of the layer in an optical system. As the figure shows, light is lost to the efflux system in projection density, the exact amount depending on the numerical aperture of the optics involved and the scattering characteristics of the sample. Thus, the projection density of a silver film is usually greater than the diffuse density. The effective Callier coefficient Q' may be defined by

$$Q' = \frac{\text{project density}}{\text{diffuse density}} \quad (4)$$

This factor can be measured experimentally. For silver films and $f/2$ optics, $Q' \approx 1.3$; for color films $Q' \approx 1.0$.

Nonneutral (Color) Density

In many cases, silver densities can be treated as neutrals. For dye densities, i.e., color films, measured density depends on the spectral characteristics of both the dye and the densitometer. The spectral response of the instrument is given by

$$\rho_\lambda = P_\lambda S_\lambda F_\lambda \quad (5)$$

where ρ_λ = the response at wavelength λ

P_λ = the power in the densitometer beam at λ

S_λ = the spectral sensitivity of the sensor at λ

F_λ = the transmittance of the densitometer optics at λ , specifically including any filters placed in the densitometer beam

The measured density of a nonneutral layer is then

$$D = \log \left[\frac{\int_{\lambda_1}^{\lambda_2} \rho_\lambda d\lambda}{\int_{\lambda_1}^{\lambda_2} \rho_\lambda T_\lambda d\lambda} \right] \quad (6)$$

where T_λ is the transmittance of the layer at wavelength λ , and the wavelength limits are set by the distributions. The response ρ_λ of the system is adjusted to be equal to that of the readout device with which the film is to be used.

Thus, for example, if the sample is to be viewed by an observer, ρ_λ is made equal to the visibility function, and the resulting measurement is called visual density, etc. Instrument responses have been standardized for sensitometry of color films.¹⁰

Reflection Density

When the emulsion is coated on paper the density is measured by reflection. Reflection density is then defined by

$$D_R = -\log R \quad (7)$$

where R is reflectance, measured under suitable geometric conditions. The measurement of reflection density is also described in the standards literature.¹¹

29.7 THE D-LOG H CURVE

In routine sensitometry, samples receive a series of exposures varying by some constant factor, such as $\times 2$ or $\times \sqrt{2}$. After processing, the measured densities are plotted against the common logarithm of the exposures that produced them. The resulting curve is known as the “D-log H curve,” or the “H & D” curve (after Hurter and Driffled, the previously mentioned pioneers in the field). A typical D-log H curve is shown in Fig. 3.

As shown, the curve is divided into three regions, known as the “toe,” “straight-line portion,” and “shoulder,” respectively. The fact that an appreciable straight-line portion is found in many cases is

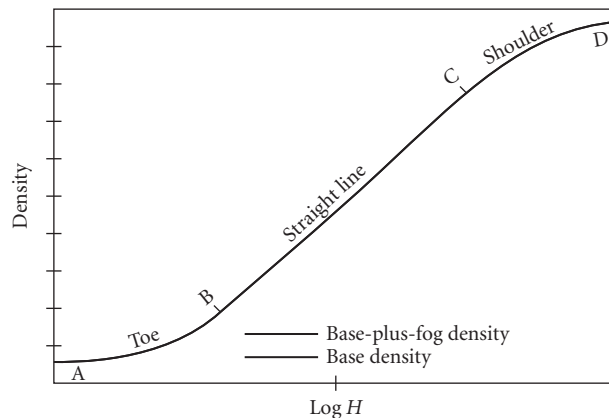


FIGURE 3 Typical D-log H curve for a negative photographic material. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

not an indication that the film is a linear responder in this region because, of course, both axes of the plot are logarithmic. It is worth noting here that the equation of the straight-line portion of this curve can be written

$$D = \gamma(\log H - \log C) \quad (8)$$

where γ is the slope of the straight-line portion and C is the exposure at the point where the extrapolated straight line cuts the exposure axis. Taking antilogarithms, Eq. (8) becomes

$$T = \left(\frac{H}{C}\right)^{-\gamma} \quad (8a)$$

If $\gamma = -1$, $T = (1/C)H$, and for this special case, the system becomes linear over the exposure range corresponding to the straight-line portion. A negative value of γ indicates a *positive* image.

A number of useful performance parameters for films are taken from their D-log H curves, as follows.

1. *Fog*: For most films, a certain number of grains will be reduced even though they have received no exposure at all, or insufficient exposure to form a latent image speck. The resulting density is called *fog*. Since it is not exposure-related, and since it tends to veil any information recorded in the toe, excessive fog is very undesirable. For many purposes, the fog density plus the density of the support are subtracted from the gross density to give the value of the net density resulting from the exposure. More complicated formulas for correcting the film's response for fog grains have been proposed, but are rarely used.
2. *Gamma*: Traditionally, the slope of the straight-line portion of the D-log H curve is called the "gamma." Gamma is a crude measure of the contrast with which the original object is reproduced; it would be a good measure of this contrast if the object were in fact recorded entirely on the straight-line portion of the response curve. However, for many purposes, notably pictorial photography, an appreciable part of the toe is used. This fact led Niederpruem, Nelson, and Yule to propose the use of an average gradient that included part of the toe as a measure of the contrast of the reproduction.¹² This quantity is called the *contrast index*. Since it includes part of the toe, contrast index is less than gamma. Since information is often recorded in the toe, it is convenient to generalize the meaning of γ to refer to the gradient anywhere along the D-log H curve, and this is done in this chapter. Note that in this case, the traditional gamma is the maximum value the gradient attains.
3. *Latitude*: Latitude can be defined as the log exposure range between the point in the toe and the point in the shoulder between which the gradient is equal to or greater than the minimum value required for acceptable recording. Clearly, the latitude of the film must be at least equal to the log exposure range of the object for proper recording. In many practical cases the film's latitude easily exceeds the required range. Note that the latitude is determined in part by the maximum density that the film can produce.
4. *Speed*: Speed is defined by the general expression

$$S = \frac{K}{H_{\text{ref}}} \quad (9)$$

where K is a proportionality factor and H_{ref} is the exposure required to produce some desired effect. Since the desired effect varies depending on the type of film and the application, H_{ref} also varies. Also, H_{ref} can be stated in either radiometric or photometric units. If H_{ref} is given in radiometric units, the proportionality factor K is set to unity, and the resulting values are termed "radiometric speeds." Although radiometric speeds are the fundamental speed values, they are rarely used in practice because, as previously noted, exposures are usually given in photometric units. In this case, the factor K takes on different values that depend not only on H_{ref} , but also on the characteristics of the exposure meter, which is standardized.¹³ Varying the factor K allows a single meter to be used with

all kinds of films and applications, which is a practical necessity. Thus, for example, the photometric speed (usually simply the “speed”) of black-and-white pictorial films is evaluated from

$$S = \frac{0.8}{H_{0.1}} \quad (9a)$$

where $H_{0.1}$ is the exposure in lux-seconds required to produce a density of 0.1 above the densities of the base plus fog in a specified process. This density level has been shown empirically to be predictive of excellent tone reproduction quality in the print. Similarly, for color slide films

$$S = \frac{10}{H_m} \quad (10)$$

where H_m is the exposure to reach a specified position on the film’s D-log H curve. Again, H_m was established empirically. The above two examples show how the two factors involved in determining a photometric or practical speed can vary. A number of other speeds have been defined and are described in the literature.¹⁴

Variation of Sensitometry with Processing

The rate of reduction of the exposed photographic grain depends on the characteristics of the grain itself, the formulation of the developer, and its temperature. In general, the reaction is allowed to continue until substantially all exposed grains have been reduced, and ended before fog becomes excessive. Many modern films are hardened and able to withstand processing temperatures up to, say, 40°C. Development times are often chosen on the basis of convenience and usually run in the order of 5 to 10 minutes in nonmachine processing. As development time and/or temperature are increased, the amount of density generated naturally increases. Thus for a given film a whole family of response curves can be produced, as shown in Fig. 4. As development time is lengthened, gamma and contrast index also increase. Typical behavior of these parameters is shown in the inset of Fig. 4.

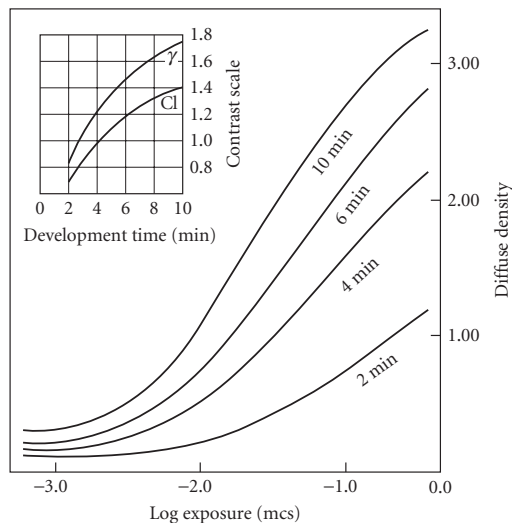


FIGURE 4 A family of characteristic curves for development times as shown, with corresponding plots of contrast index and gamma, mcs or the meter-candle second, as stated on p. 29.6, is the old term for lux-second or lxs. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

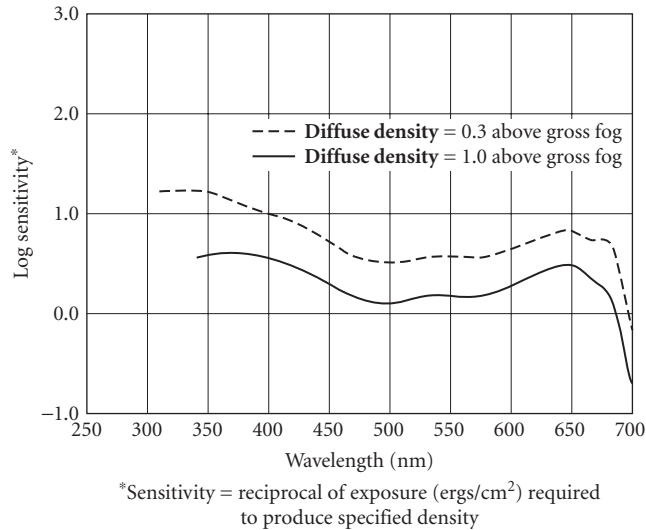


FIGURE 5 Spectral sensitivity curves for a modern negative material sensitized to about 690 nm. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

29.8 SPECTRAL SENSITIVITY

The spectral absorption of the silver halide grain extends only to about 500 nm, and thus the inherent sensitivity of the grain is limited to regions shorter than that limit. However, it was discovered by Vogel in 1873 that sensitivity could be extended to longer wavelengths by dyeing the grains, and over the years the effective sensitivity range was extended into the infrared. Presently, materials are available usefully sensitized out to about 1.2 μm , but it should be noted that IR materials tend to have poor shelf life, and may require special handling. For our purposes, we may say that the sensitizing dye absorbs the incident energy in the required spectral region and transfers the energy to the grain in a manner that produces the required latent-image speck. The mechanisms are discussed in Ref. 3, chap. 10.

The spectral sensitivity of a photographic layer is usually specified by a family of curves showing $\log(1/H_D)$ vs λ , where H_D is the exposure in ergs/cm² of wavelength λ required to produce some stated density. Spectral sensitivity values are thus radiometric speeds. Typical curves are shown in Fig. 5. In practical work, three broad classes of sensitization are recognized, which are called “color-blind” (or blue-sensitive), orthochromatic (additionally sensitized to green), and panchromatic (additionally sensitized to green and red). Most modern general-purpose materials are panchromatic.

In general, the shape of the spectral sensitivity curve follows that of the spectral absorption of the layer. It should also be noted that the gradient of the D-log H curve will be affected by the absorption of the layer. Gamma may therefore vary as a function of wavelength, and is generally somewhat lower in the blue and UV regions of the spectrum. This means that if the material is being used as a radiometer, it must be calibrated at the wavelength(s) of interest.

29.9 RECIPROCITY FAILURE

It was noted above that the exact response of a photographic layer may change due to exposure effects. Of these, the phenomenon of *reciprocity failure* is the most important in practical photography.

By definition [Eq. (1)] exposure is simply the product of the irradiance and time, and nothing in this definition specifies the magnitude of either factor. However, the developed *density* resulting from a given calculated exposure is often found to depend on the *rate* at which the radiation is supplied, all other factors being held constant. Broadly speaking, exposure times of about 0.01 to 1.0 second are most efficient in producing density, the exact values depending on the film involved. Times much outside the above range tend to produce lower density for the same calculated exposure. The emulsion-maker has some ways of minimizing the effect, and usually attempts to optimize the response for the exposure times expected to be used with the material. Reciprocity failure is discussed in detail in Ref. 3, chap. 4, sec. II.

The loss of efficiency for short-time and correspondingly high-irradiance exposures is termed “high-intensity reciprocity failure,” and that for long exposure times (low irradiances) is termed “low-intensity reciprocity failure.” The Gurney-Mott mechanism explains both types of failure well. Note also that the names are misnomers; the terms should, of course, be high and low *irradiance*.

Only limited data are available, but the gradient of the D-log H curve tends to decrease as exposure times are shortened. An example is shown in data published by Hercher and Ruff.¹⁵ Limited data also indicate that the speed loss due to high-intensity failure stabilizes for times shorter than about 10^{-5} second.¹⁶ Essentially, the amount of reciprocity failure is independent of exposing wavelength.¹⁷ This is to be expected.

Reciprocity failure may be a considerable problem for photographers working in specialized time domains, such as oscilloscope photography or astronomy. Astronomers have been able to devise user treatments for minimizing low-intensity failure.¹⁸ For the practical photographer, reciprocity failure sometimes appears as a problem in color photography. If the RF of the three sensor layers differs, the resulting picture may be “out of balance,” i.e., grays may reproduce as slightly tinted. This is very undesirable, and correction filter recommendations are published for some films for various exposure times.

29.10 DEVELOPMENT EFFECTS

Besides exposure effects the final density distribution in the developed image may be affected by “development effects” arising from chemical phenomena during development. Various names such as “border effect,” “fringe effect,” “Eberhard effect,” etc., are applied to these phenomena; what we shall here term “edge effect” may be important in practice.

Consider a sheet of black-and-white film developing in a tray, and assume for purposes of discussion that there is no motion of the developer. Since developing agent must be oxidized as halide is reduced, and since the by-products of this reaction may themselves be development inhibitors, it can be seen that local variations of developer activity will be produced in the tray, with the activity decreasing as density increases. Agitation of the solution in the tray reduces the local variations, but usually does not eliminate them entirely, because it is the developer that has diffused into the gelatin matrix that is actually reacting. Now an “edge” is a boundary between high- and low-density areas, as shown in Fig. 6. Because of the local exhaustion and the diffusion phenomena, the variation of developer activity within the layer will be as shown by the dotted line in the figure. The result is that the developed density near the edge on the low-density side tends to decrease, and on the high-density side tends to increase, as also shown in the figure. In other words, the density distribution at the edge is changed; this actually occurs to some degree in much practical work and has interesting consequences, as will be discussed below.

The local exhaustion of the developer may also be important in color films where development in one layer (see below) may affect the response of an adjacent layer. In color photography, the phenomenon is called “interimage effect.”

29.11 COLOR PHOTOGRAPHY

Color photography has been extensively reviewed by Kapecki and Rodgers.¹⁹ With one exception at the time of writing, all commercially available color films employ subtractive color reproduction. The exception is an instant film for color slides marketed by Polaroid Corporation, which employs

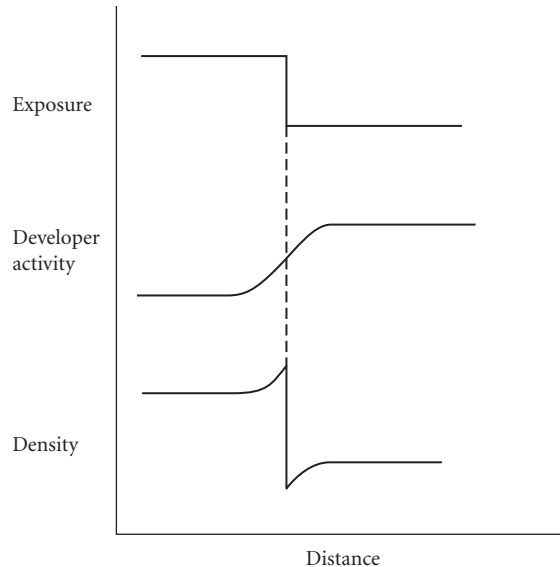


FIGURE 6 Distributions of developer activity and density resulting from a step-function input in the presence of edge effect.

additive color. The mechanism involves a reseau of very fine red, green, and blue stripes, which provide the color separation during the taking of the photograph, and also the color when the (reversed) image is projected.²⁰

The basic structure of most other color films is similar to that shown in Fig. 7. The incoming light first encounters a nonspectrally-sensitized emulsion layer, which records the blue-light elements of the scene. The next layer is a yellow, or minus-blue filter, the purpose of which is to prevent any blue light from reaching the other two emulsion layers. This yellow filter is usually composed of “Carey Lea,” or colloidal silver dispersed in gelatin. Such sols are yellow. The reason for using Carey Lea silver is that all metallic silver is removed from the film during processing anyway, and the necessary removal of the filter layer is thus accomplished automatically.

Moving downward in the stack, the next layer is an ortho-sensitized emulsion. Since any blue light has been blocked by the yellow filter, this layer records the green-light elements of the scene. The final layer in the stack is sensitized to red light but not to green light and this layer serves to record the red elements of the scene. Modern films usually contain many more than the four layers indicated here, but the operating principles of the film can be discussed in terms of such a “tripack.” In accordance with the principles of subtractive color reproduction, the images in the three separation

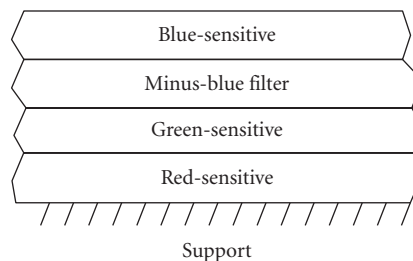


FIGURE 7 Schematic tripack structure of a color film.

records are then converted to yellow, magenta, and cyan dyes, respectively, and the final image is composed of these dyes, without the developed silver, which is either removed chemically or left behind, depending on the exact material.

Within this broad framework, films can be separated into two classes: chromogenic and nonchromogenic. In the former class, the dyes are not coated in the film, but are formed during the processing by a reaction called “coupling.” In coupling, the by-products of the halide reduction reaction serve as components for a second reaction in which dye is formed; the amount of dye thus increases as density increases. The components of the dye-forming reaction (i.e., other than the development by-products) may be present in the developing solution (Eastman Kodak Co., “Kodachrome”) or, more generally, coated within the various layers (Kodacolor, Polaroid “One Film,” Agfachrome, Fujicolor, etc.). After the required dyes are formed, the developed silver is removed by a chemical process termed “bleaching.”

The advantage of incorporating the couplers in the various layers of the tripack is simpler processing, but because of the additional material in the layers, such films tend to be not as sharp as the nonincorporated-coupler types. Chromogenic color films are available both as slide materials, in which the film undergoes a reversal process,²¹ and as color negative—color print materials in which a dye negative is formed and then printed onto a color paper whose structure is fundamentally similar to that of the films.

The principal example of the nonchromogenic film is the Polaroid Instant Color Film. In this film the yellow, magenta, and cyan dyes are actually coated in the structure, along with the blue, green, and red-sensitive emulsions. When development takes place in a given layer, the corresponding dye is immobilized. The dye that has *not* been immobilized in the three layers migrates to a “receiver” layer, where it is mordanted. Since the amount of dye that migrates *decreases* as the original density *increases*, the result is a positive color image formed in the receiver. The material has been described in more detail in a paper by Land,²² and also in chap. 6 of Ref. 4.

The image in a color film thus essentially consists of three superimposed dye images. Typical spectrophotometric curves for dyes formed by coupling reactions are shown in Fig. 8. The density of any one of these dye layers taken by itself is known as an *analytical density*. Note, however, that all the dyes show some “unwanted” absorption—that is, absorption in spectral regions other than the specific region that the dye is supposed to control. Thus the total density of the layer at any wavelength is the sum of the contributions of all three dyes; this type of density is known as “integral density.” The integral density curve is also shown in Fig. 8.

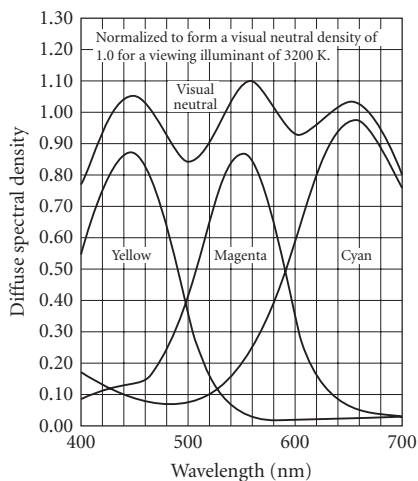


FIGURE 8 Spectral dye density curves for a color film. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

The reproduction of color by photographic systems has been discussed by Hunt²³ and many others.²⁴ In general, exact colorimetric reproduction is not achieved, but for most purposes the reproduction of the hues and luminances in the original scene is satisfactory. Evans,²⁵ in fact, observes that under the right conditions the “magnitude of the reproduction errors that can be tolerated is astonishing.” One aspect that is critical for many workers, especially expert photographers, is the ability of the system to produce good “balance,” i.e., to reproduce a neutral as a neutral. This requirement is so important that one of the types of density that is measured for dye layers is the *equivalent neutral density* (END) which is defined as the visual neutral density that a dye layer would produce if combined with the correct amounts of the other two dyes (whatever those correct amounts may be). When the ENDs of the three layers are equal, the system is in balance, and color-film sensitometry is therefore often done in terms of ENDs. Further discussion of color sensitometry and densitometry may be found in Ref. 3, chap. 18.

29.12 MICRODENSITOMETERS

As indicated above, a microdensitometer is a densitometer designed to measure the density of a small area. The sampling apertures are typically slits which may be as narrow as 1 to 2 μm in nominal width. The sample is scanned over the aperture, creating a record of density as a function of position on the sample surface, i.e., distance.

In practical instruments, the small sampling aperture dimensions are achieved by projecting an enlarged image of the film onto a physical aperture. The optical system produces some effects not encountered in macrodensitometers, as follows:

1. In general, microdensitometers measure projection, or semispecular, density. As already noted, projection density is higher than diffuse density for silver layers, and the exact value of the effective Callier coefficient Q' depends in part on the numerical aperture of the optical system. Thus two microdensitometers fitted with optics of different NAs may give different density values for the same sample. Furthermore, macrodensity data are usually in terms of diffuse density, so that data from the microdensitometer should be corrected if intercomparisons are to be made. The effective Callier coefficient for the specific optics-sample combination is easily determined by measuring suitable areas of sample both in the microdensitometer and a macrodensitometer, and taking the ratios of the values.
2. The presence of stray light in the system tends to lower the measured density. This problem is especially troublesome in microdensitometry because of the types of images that are often encountered. Thus, for example, when an interface between clear and dense areas—that is, an edge—is scanned, stray light will distort the record in the manner shown in Fig. 9. If the image is that of a star or spectroscopic line, this behavior results in an artificially low density reading. It is very important to control stray light as completely as possible in microdensitometry.

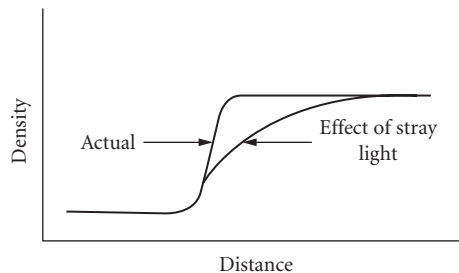


FIGURE 9 Effect of stray light in the microdensitometer on the apparent density distribution at an edge.

3. One feature of the usual microdensitometer optical system, installed for control of stray light, is the “preaperture,” a field stop that limits the area of the sample that is illuminated. Because of this preaperture, and the optical system, a normal microdensitometer is a partially coherent system. This means that the instrument may respond to path-length differences in the sample as well as to density differences. This is undesirable since in practice it is the density differences that carry the information. Partial coherence in microdensitometers has been studied by Thompson and by Swing among others, and the results are summarized by Dainty and Shaw.²⁶ It has been shown that the coherence effects can be minimized by satisfying two conditions. In the first condition, the width W of the preslit

$$W \geq \frac{4\lambda}{NA_{\text{in}}} \quad (11a)$$

The second condition is

$$\frac{NA_{\text{in}}}{NA_{\text{eff}}} = 1 + \frac{\nu_s}{\nu_o} \quad (11b)$$

where λ = the wavelength of the light
 NA_{in} = numerical aperture of the influx optics
 NA_{eff} = numerical aperture of the efflux optics
 ν_s = maximum spatial frequency in the sample
 ν_o = spatial frequency cutoff of the scanning objective

If, for example $\nu_o = 3\nu_s$, then $NA_{\text{in}}/NA_{\text{eff}} \geq 1.3$. A microdensitometer arranged to minimize coherence problems is called a *linear* microdensitometer. Note that the two conditions above conflict with conditions commonly adopted to control stray light.

29.13 PERFORMANCE OF PHOTOGRAPHIC SYSTEMS

The following discussion of performance is limited to those aspects which are properties of the system, and excludes such aspects as the skill of the photographer, etc. Of those aspects, which we may term “technical” quality parameters, the most important is *tone reproduction*. The subject is divided into two areas: subjective tone reproduction and objective tone reproduction. Subjective tone reproduction is concerned with the relation between the brightness sensations produced in the observer’s mind when the scene is viewed and when the picture is viewed. Since the sensation of brightness depends markedly on the viewing conditions and the observer’s state of adaptation, the subjective tone reproduction of a given picture is not constant, and this aspect of the general subject is not often measured in the photographic laboratory. It is discussed by Kowalski.²⁷ Objective tone reproduction is concerned with the reproduction of the luminance and luminance differences of the scene as luminances in the final output. Tone reproduction studies apply equally well to projected images, prints, transparencies, and video images, but a negative-print system is usually assumed for discussion. It is easy to show that the log luminance of a print area, $\log L_p = C - D_p$, where C is a constant determined by the illuminance on the print and D_p is the density of the print area. Thus tone reproduction curves are usually plotted in the form of the print density versus the log luminance of the corresponding scene element (Fig. 10). Although both the scene luminances and the print densities are fixed quantities, the viewer’s reaction still depends on the illumination level at which the picture is seen.

It has been shown empirically that for paper prints viewed under typical “room lighting conditions,” the preferred tone reproduction curve is the solid line in Fig. 10. Perfect objective tone reproduction, defined as the case where $\Delta D_p = -\Delta \log L_{sc}$ for all scene luminance levels, would be

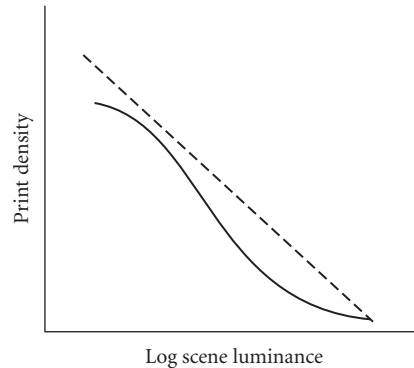


FIGURE 10 Preferred tone reproduction for viewing under typical room lighting conditions.

the dotted line in the figure. Thus under typical room lighting conditions viewers prefer a reproduction that has somewhat more contrast and less density than the “perfect” result. This preference will, however, change with illumination and stray light levels.

The exact shape of the tone reproduction curve obtained with a given system depends on the shape of the negative and positive D-log H curves, and also on the stray-light characteristic of the camera system. (Excessive stray light in the camera can be very deleterious to tone reproduction quality.) Using a graphic method devised by Jones and described by Nelson and others,²⁸ the effect of the D-log H curves and the stray light on the tone reproduction can be studied.

29.14 IMAGE STRUCTURE

The other two technical quality parameters of a photographic system are its sharpness and graininess, to use the most familiar terms for these properties. These properties are often lumped together under the general term “image structure.” Actually, in photoscience the term *sharpness* and *graininess* are reserved for the subjective aspects of the phenomena, and measurement of these properties requires psychometric testing. Since such testing is expensive and time-consuming, objective correlates of both properties have been defined, and methods for measurement have been established. The objective correlate of sharpness is termed *acutance*, and of graininess is termed *granularity*. Image structure data for various kinds of materials are given in Table 1, along with speed and contrast values.

29.15 ACUTANCE

The original proposal for measuring acutance was made by Higgins and Jones²⁹ in 1952. It involved calculating a value from a microdensitometer trace of a test edge. However, the visual processes that occur when an observer views an edge are complex, and the straightforward calculation proposed by Higgins and Jones fails to predict the sensation of sharpness produced by some edge distributions. At about the same time, the optical transfer function (OTF) and related concepts began to be widely used in optics, and these were soon applied to photographic materials also.

The concepts of the point and line spread functions are essentially identical in optics and photoscience; in the photographic layer the smearing of the point (or line) input is caused by diffraction around the grains, refraction through them, and reflection from them. These phenomena are usually

TABLE 1 Performance Data for Types of Materials^a

Type	SPD ^b	γ	Gran. ^c	ν_{50} ^d
B&W microfilm	80	3.0	6	>200
B&W very slow camera neg.	25	0.5–3.5	5–7	80–145
B&W slow camera neg.	100	0.5–1.1	8–9	65–120
B&W fast camera neg.	400	0.5–1.0	10–14	50–100
B&W very fast camera neg.	1600–3200	0.5–1.0	18	70
Color neg. very slow	25–50	0.65	4–5	40–60
Color neg. slow	100–160	0.60–0.80	4–6	30–70
Color neg. fast	400	0.65–0.80	5–7	25–40
Color neg. very fast	1000–1600	0.80	8–11	25–35
Color rev. slide very slow	25–50	1.8–2.3	9–10	30–40
Color rev. slide slow	100	2.0–3.0	10–13	25–30
Color rev. slide fast	400	2.0–2.4	15–20	20
Color rev. slide very fast	800–1600	2.2–2.8	22–28	16–20
Instant print films ^e		–1.7 ^f	NA	3–4
—Black and white	3000	–1.6	NA	2.5
—Color	100			

^aData are as of early 1993 and were obtained from publications of the manufacturers listed in Sec. 29.21. They are presented as published. Note that products are frequently changed or improved.

^bSpeeds are calculated in different ways for various classes of product. The values given are suitable for use with standard exposure meters.

^cValues represent 1000 \times the standard deviation of the diffuse density, measured at an average density of 1.0 using a 48- μ m circular aperture. The exact granularity of a print depends on the characteristics of the print material and the printer as well as the granularity of the negative.

^dValues show the spatial frequency at which the modulation transfer function is 50 percent.

^eMTF values apply to the final print.

^fNegative sign arises from the definition of gamma for the case of a positive image.

lumped together and termed “scattering,” and have been treated by Gasper and dePalma.³⁰ Likewise, the concept of the optical transfer function, or the Fourier transform of the LSF, is basically the same in optics and photography. However, three important differences should be noted for the photographic case. (1) The emulsion is isotropic, so that the PSF and LSF are always symmetrical, and the complex OTF reduces to the modulation transfer function (MTF) only. (2) Unlike lenses, photographic layers are stationary, but are generally nonlinear. Therefore, all data and calculations must be in terms of *exposure* or allied quantities. When the calculations are complete, the results are converted to density via the D-log H curve. (3) The presence of edge effects tends to raise the MTF curve, so that for low frequencies the measured response values are often found to be greater than 100 percent. The subject is treated in detail by Dainty and Shaw.³¹ Typical MTF curves for a film, showing the overshoot, are given in Fig. 11. Data for MTF curves of various types of films are also given in Table 1. The value given shows the spatial frequency for which the transfer factor drops to 50 percent.

The chain relating the MTF curve to image sharpness is the same as in optics: a high MTF curve transforms to a narrow spread function, and this in turn indicates an abrupt transition of exposure—and therefore density—across the edge. Thus MTF response values greater than unity, although mathematically anomalous, indicate improved image sharpness, and this is found to be the case in practice. As a matter of fact, edge effects are often introduced deliberately to improve sharpness. This is done either by adding suitable compounds to the coating itself, or by adjusting the developer formulation. The process is similar to the electronic “crispening” often used in television.

An index of sharpness can be computed from the MTF data by a procedure first suggested by Crane and later modified by Gendron.³² These workers were interested mainly in films, but they recognized that the film is one component of a system; e.g., a color slide system involves a camera lens, the film (and process), a projector lens, the screen, and the observer’s eye. The MTF of the system is then the cascaded MTFs of these components. Gendron suggested that the area under the cascaded MTF be taken as the stimulus that produced the sensation of sharpness. A formula was proposed

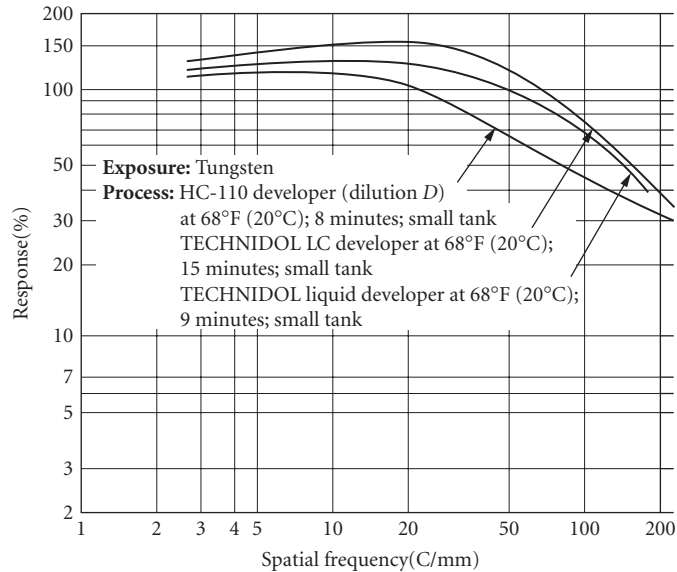


FIGURE 11 Modulation transfer functions of Kodak Technical Pan Film for three conditions of development. Note that the response at low frequencies exceeds 100 percent. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

that produced a sharpness index scaled to 100. This index, now called CMT-Acutance (to distinguish it from the original Higgins-Jones acutance) has been found to correlate well with subjective data, and is used in the industry. Note that the treatment above has been simplified for the sake of brevity; details are available in Gendron's paper.

29.16 GRAININESS

The granular nature of the photographic image is one of its most significant characteristics. It may appear to the observer as an unpleasant roughness in what should be uniform areas, but it may also interfere with extracting information from the image. In the former case it is an aesthetic problem; in the latter case the structure is equivalent to noise in a communications channel. In either case it is desirable to measure the phenomenon objectively and in engineering terms.

The procedure now used for making these objective measurements was proposed by Selwyn in 1935.³³ He postulated that if the density of a uniformly exposed and processed layer were measured at many places using a suitable sampling aperture, the population of density values so obtained would be approximately gaussianly distributed around the mean. This being so, the variability for a given mean density is completely described by the standard deviation of the values. This quantity is termed the *rms-granularity*, $\sigma(D)$, and has indeed been shown to correlate with the subjective graininess.³⁴ It might be noted that calculating the standard deviation of the density is mathematically improper, since the underlying transmittance values are being multiplied. To avoid this problem the rms-granularity is formally defined by

$$\sigma(D) = \frac{0.434\sigma(T)}{\bar{T}} \quad (12)$$

where T is transmittance. However, it can be shown that when $\sigma(T)$ is small compared to \bar{T} , the error involved in calculating directly in density is small, and this is often done in practice.

Selwyn also showed that the measured value of $\sigma(D)$ depended upon A , the area of the sampling aperture. The product $\sigma(D)A^{1/2}$ may be termed the *Selwyn coefficient* \mathcal{G} ; for black-and-white films exposed to light, Selwyn showed that it should be constant, and this relation is called "Selwyn's law." Thus $\sigma(D)A^{1/2}$ is a measure of sample graininess no matter what the size of the sampling aperture. Unfortunately, the Selwyn coefficient does not remain constant with changes of aperture size for very important classes of samples. Selwyn's law may fail for prints and enlargements, black-and-white or color, for many color materials even if not enlarged, and also for radiographs, especially screened radiographs. For such materials $\sigma(D)$ still increases as A decreases, but not at a rate sufficient to keep the Selwyn coefficient constant, and it (the coefficient) is therefore not useful as an objective measure of graininess.

Stultz and Zweig³⁵ found that they obtained good correlation between $\sigma(D)$ and the sensation of graininess when the sampling aperture was selected in accordance with the rule

$$d(\mu) \cdot M_\theta \approx 515 \quad (13)$$

where $d(\mu)$ is the diameter of the sampling aperture in μm , and M_θ is the angular magnification³⁶ at which the photograph is seen by the viewer. M_θ is readily calculated from the relation

$$M_\theta = \frac{m}{4V} \quad (13a)$$

where m is the ordinary lateral magnification between the film image and the image presented to the viewer, and V is the viewing distance in meters.

An American standard³⁷ exists for the measurement of rms-granularity. This standard specifies that samples be scanned with a $48\text{-}\mu\text{m}$ -diameter aperture; for such an aperture, rms-granularity values for commercial films range from about 0.003 to 0.050 at an average density of 1.0. In practice, these values are often multiplied $\times 1000$ to eliminate the decimals. It is worth noting that, experimentally, the measurement of rms-granularity is subject to many sources of error, such as sample artifacts. The standard discusses sources of error and procedures for minimizing them, and is recommended reading for those who must measure granularity.

While in practice rms-granularity serves well as an objective correlate of graininess, the situation is complicated by the fact that there are two broadly different types of granular pattern. Silver grains are small, opaque, and in nearly all cases are situated randomly and independently in the coating. The granular structure in an enlargement, however, is composed of clusters of print-stock grains that reproduce the exposure pattern coming from the enlarged negative grains. This type of granular pattern tends to be large and soft-edged compared to the pattern arising from the primary grains. The patterns found in such samples as color films and screened radiographs are generally similar to those found in enlargements. Microdensitometer traces of these two structures are illustrated in Fig. 12. A little thought will show that the two patterns shown in the figure might have the same mean and standard deviation, and yet the two patterns look entirely different. When different types of patterns are involved, the rms-granularity above is not a sufficient descriptor. The work by Bartleson which showed the correlation between graininess and rms-granularity was done with color negative films having similar granular structures.

As discussed by Dainty and Shaw,³⁸ further objective analysis of granular patterns may be carried out in terms of the *autocorrelation function*:

$$\phi(\tau) = \lim_{x \rightarrow \infty} \frac{1}{2x} \int_{-x}^{+x} \delta(x) \delta(x + \tau) dx \quad (14)$$

where $\delta(x) = D_x - \bar{D}$
 D_x = the density reading at point x on the sample
 \bar{D} = mean density
 $\delta(x + \tau) = D_{x+\tau} - \bar{D}$
 τ = a small increment of distance

Note carefully that for the sake of simplicity it has been assumed that the sample is scanned by a very long, narrow slit, so that the autocorrelation function reduces to a one-dimensional function.

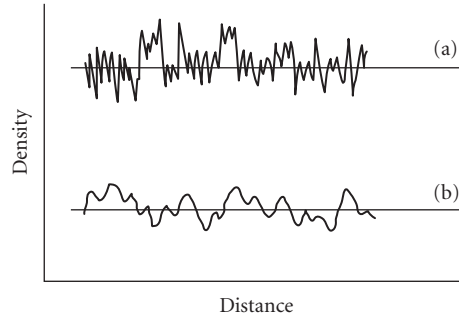


FIGURE 12 Microdensitometer traces of (a) primary silver grains and (b) granular structure of a print.

If the sample is scanned by a small circular aperture it will be two-dimensional. The use of a slit is common in practice. For the case where $\tau = 0$, we have

$$\phi(0) = \lim_{x \rightarrow \infty} \frac{1}{2x} \int_{-x}^{+x} [\delta(x)]^2 dx = \sigma^2 \quad (15)$$

Let $\phi(\tau)$ be calculated for several different values of τ , including values smaller than the slit width. Since when $\tau < w$ the slit will contain some of the same grains at points x and $x + \tau$, correlation is observed. If there is no spatial correlation in the actual sample, $\phi(\tau) \rightarrow 0$ for $\tau > w$. Positive values of $\phi(\tau)$ for $\tau > w$ are an indication of spatial correlation in the sample; that is, large-scale grain.

In practice, it is convenient to carry out the analysis of large-scale patterns in the spatial frequency domain. The Wiener-Khintchin theorem states that what is called the “Wiener spectrum” or *power spectrum of the granularity* distribution is the Fourier transform of the autocorrelation function; that is,

$$\text{WS}(\nu) = \int_{-\infty}^{\infty} \phi(\tau) e^{i2\pi\nu\tau} d\tau \quad (16)$$

and also

$$\phi(\tau) = \int_{-\infty}^{\infty} \text{WS}(\nu) e^{-i2\pi\nu\tau} d\nu \quad (16a)$$

where ν = spatial frequency.

In practice, an approximation of the Wiener spectrum is usually obtained by a direct Fourier transform of the granularity trace itself according to the expression

$$\text{WS}(\nu) = \lim_{x \rightarrow \infty} \frac{1}{2x} \overline{\left| \int_{-x}^x \delta(x) e^{i2\pi\nu x} dx \right|^2} \quad (17)$$

where $\delta(x)$ has the same meaning as before, and the horizontal bar indicates that the average value of several different runs should be taken in order to provide a reasonable value for the approximation.

Returning to the definition of the autocorrelation function [Eq. (14)] it can be seen that the measured values of $\delta(x)$ and $\delta(x + \tau)$ are each actually the true values convolved with the slit response. By the convolution theorem of Fourier transforms, this means that in frequency space

$$\text{WS}(\nu)_{\text{measured}} = \text{WS}(\nu)_{\text{true}} xT(\nu)^2 \quad (18)$$

where $T(\nu)$ is the modulation transfer factor of the measuring system at spatial frequency ν . If $T(\nu)^2$ is divided out of $WS(\nu)_{\text{measured}}$ the underlying net value $WS(\nu)$ is obtained. When this is done for black-and-white film samples exposed to light, the underlying spectrum is found to be flat. For the kinds of samples for which Selwyn's law fails, on the other hand, the net Wiener spectrum is found to contain excess low-frequency power and lack high frequencies. Finally, by the properties of Fourier transforms, the area under the Wiener spectrum curve is the zero-value of the autocorrelation function. But by Eq. (15) this zero-value equals σ^2 for the sample, so that the area under the WS curve gives the rms-granularity of the sample. This can also be seen intuitively, since the value of the WS at special frequency ν is simply the noise power of σ^2 in a spatial frequency band $\nu \pm \Delta\nu/2$.

Doerner³⁹ has shown that the rms-granularity of a negative can be tracked through a printer system in terms of the Wiener spectrum of the negative and the modulation transfer function of the printer system (which includes the MTF of the print stock itself). Doerner's expression is

$$WS(\nu)_p = WS(\nu)_n \gamma_p^2 A(\nu)_{\text{sys}}^2 + WS(\nu)_{\text{ps}} \quad (19)$$

where n and p indicate negative and print respectively, "ps" the print stock itself, and the other symbols have their previous meanings. The granularity/graininess of a print thus depends on the contrast of the print stock and the spatial frequency response of the printer system as well as on the graininess of the negative.

29.17 SHARPNESS AND GRAININESS CONSIDERED TOGETHER

In the foregoing discussion, we have considered the sharpness and graininess aspects of the picture separately. But real photographs frequently suffer from less-than-optimum sharpness and graininess *both*. Bartleson⁴⁰ has studied the subjective quality of such photographs. He concluded that quality was not a linear combination of sharpness and graininess. Instead, "... quality tends to be determined primarily by the poorer of the two attributes. . . . If graininess is high, the print will likely be low in quality regardless of how sharp it may be or, conversely, if sharpness is low, so also will quality be low regardless of how grainy the print appears. . . ."

Bartleson's results are of interest in assessing the quality of electronic images, since, at least at the time of writing, such images exhibit low graininess, but also low sharpness. On the basis of Bartleson's work, such images would be judged to be of low subjective quality. In a comparison of electronic and film imagery published in 1990, Ikenoue and Tabei⁴¹ rated the quality of the former as poor because of the sharpness level.

29.18 SIGNAL-TO-NOISE RATIO AND DETECTIVE QUANTUM EFFICIENCY

Since the information in a photograph is normally carried by the density variation, we may usefully define the output signal-to-noise ratio of the photography by

$$S/N_{\text{out}} = \frac{\Delta\bar{D}}{\sigma} = \frac{\gamma}{\sigma} \cdot \Delta\log H \quad (20)$$

where $\Delta\bar{D}$ is the mean density difference between an element to be detected and its surround, and σ is the rms-granularity of the surround. By Selwyn's law, $\sigma(D)$ varies as the area of the sampling

aperture changes; it is convenient here to take the sampling aperture area A as equal to the area of the image element. Furthermore, for ΔH sufficiently small we may write

$$S/N_{\text{out}} = 0.434 \frac{\gamma}{\sigma} \cdot \frac{\Delta H}{H} \quad (21)$$

Note that γ/σ is a property of the film; it is termed the *detectivity*. $\Delta H/H$, on the other hand, is a property of the object; in fact, it is the object contrast. Practical tests have shown that S/N_{out} should be 4 to 5 if the element is to be detected against its surround, and 8 to 10 if it is to be recognized.⁴²

In 1946, Albert Rose of RCA published a paper⁴³ in which he discussed the performances of the TV pickup tube, the photographic layer, and the human eye on a unified basis. His approach was to compare their performances with that of an "ideal device," that is, a radiation detector whose performance was limited only by the quantum nature of the incoming radiation. Such a perfect detector would report the arrival of every incoming quantum, and add no noise to the signal.

In 1958, R. Clark Jones of Polaroid expanded Rose's work with specific application to photographic layers.⁴⁴ Jones proposed the term *detective quantum efficiency* (DQE) for Rose's performance indicator, and defined it by the expression

$$\text{DQE} = \left[\frac{S/N_{\text{out}}}{S/N_{\text{ideal}}} \right]^2 = \left[\frac{S/N_{\text{out}}}{S/N_{\text{max}}} \right]^2 = \left[\frac{S/N_{\text{out}}}{S/N_{\text{in}}} \right]^2 \quad (22)$$

where S/N_{out} is the signal-to-noise ratio produced by the actual device, and S/N_{ideal} is the ratio that would be produced by the ideal device, given the same input. By the definition of the ideal device, $S/N_{\text{ideal}} = S/N_{\text{max}} = S/N_{\text{in}}$, where S/N_{in} is the signal-to-noise ratio in the input, and is due to the quantum nature of the input. The ratio is squared to make the DQE compatible with various other concepts.

It is easy to derive an expression for the S/N ratio of the input; when this is combined with Eq. (21) the result is

$$\text{DQE} = \left[\frac{0.434\gamma}{\mathcal{G}} \right]^2 \cdot \frac{1}{q} \quad (23)$$

where \mathcal{G} is the Selwyn coefficient and q is the average exposure received by the image in *quanta per unit area*. Since $(1/q)$ is the radiometric speed of the film, it can be seen that the DQE is a performance parameter that combines the gain (gamma), noise, and speed of the layer. As written, the DQE does not involve the sharpness aspect, but this can be included also.⁴⁵

Since Eq. (23) involves standard photographic parameters, it is readily evaluated for a given material. When this is done, it is found that the DQE of typical materials is on the order of 1 to 3 percent, peaking sharply at low densities in the case of black-and-white films. Note that DQE can be calculated for any sensor for which S/N_{out} can be derived. It is interesting to compare the 1 to 3 percent values given above with those of other sensors. Thus, for example, Jones gives a value of 1 percent for the human eye, and 6 percent for an image orthicon tube. On the other hand, a suitable photographic layer recording electrons may approach 100 percent DQE, and a value of 30 percent has been reported for a screened x-ray film.⁴⁶ DQE has also served as a useful approach to considering silver halide mechanisms; see, for example, a paper by Bird, Jones, and Ames that appeared in *Applied Optics* in 1969.⁴⁷

Another figure of merit allied to DQE is the *noise equivalent quanta* (NEQ) which is defined as the number of quanta that a perfect detector would need to produce a record having the same S/N ratio as the system under consideration. It can be shown that

$$q' = \text{DQE} \times q \quad (24)$$

where q is again the number of quanta/unit area used by the real system, and q' the NEQ of a unit area of image.

29.19 RESOLVING POWER

The basic procedure used to measure photographic resolving power follows that used in optics. An American Standard exists.⁴⁸ The standard provides for a suitable test object to be reduced optically onto the material to be tested. (Strictly speaking, the test thus determines the resolution of the lens-film combination, but the resolution capability of the lenses specified is high compared to that of the film.) The developed image is then studied in a microscope to determine the highest spatial frequency in which the observer is “reasonably confident” that the structure of the test pattern can still be detected. Thus, as in optics, the “last resolved” image is a threshold image. However, unlike the optical case, the limit is set not only by the progressive decrease in image modulation as spatial frequency increases, but also by the granular nature of the image.

If an exposure series of the test pattern is made, it is found that the spatial frequency of the limiting pattern goes through a maximum. It is customary to report the spatial frequency limit for the optimum exposure as the resolving power of the film.

Photographic resolving power is now not much measured, but it retains some interest as an example of a signal-to-noise ratio phenomenon. Consider the modulation in the various triplet images in the pattern. By definition,

$$M = \frac{H_{\max} - H_{\min}}{H_{\max} + H_{\min}} = \frac{\Delta H}{2\bar{H}} \quad (25)$$

and from Eq. (21),

$$\Delta D = 0.868 \gamma M \quad (26)$$

so that ΔD within the pattern varies as the modulation, which in turn decreases as the spatial frequency increases. Furthermore, the area of each of the triplets in the pattern (assuming the ISO pattern configuration) = $2.5^2 \lambda^2 = 2.5^2 / \nu^2$. Assuming that Selwyn's law holds, it follows that $\sigma = \mathcal{G}/A^{1/2} = 0.4 \mathcal{G}\nu = C\nu$, where we have lumped the constants. As the spatial frequency increases, the S/N of the triplet decreases because ΔD decreases and the effective rms-granularity increases. The resolving power limit comes at the spatial frequency where the S/N ratio drops to the limit required for resolution. Such a system has been analyzed by Schade.⁴⁹

29.20 INFORMATION CAPACITY

The information capacity, or number of bits per unit area that can be stored on a photographic layer, depends on the size of the point spread function, which determines the smallest element that can be recorded, and on the granularity, which determines the number of gray levels that can be reliably distinguished. The exact capacity level for a given material depends on the acceptable error rate; for one set of fairly stringent conditions, Altman and Zweig⁵⁰ reported levels up to 160×10^6 bits/cm². Jones⁵¹ has published an expression giving the information capacity of films as

$$IC = \frac{1}{2} \iint_0^\infty \log_2 \left(1 + \frac{S(\nu_x, \nu_y)}{N(\nu_x, \nu_y)} \right) d\nu_x d\nu_y \quad (27)$$

where S and N are the spectral distributions of power in the system's signal and noise. For a given spatial frequency, the signal power is given by

$$S = WS_i(\nu) \text{MTF}^2(\nu) \quad (28)$$

where $WS_i(\nu)$ is the value of the Wiener spectrum of the input at ν and $\text{MTF}(\nu)$ is the value of the film's MTF at that frequency. The information capacity thus depends on the frequency response and noise of the system, as would be expected. The matter is discussed by Dainty and Shaw.⁵²

29.21 LIST OF PHOTOGRAPHIC MANUFACTURERS

Agfa Photo Division
 Agfa Corporation
 100 Challenger Road
 Ridgefield, NJ 07660

Eastman Kodak Co.
 343 State Street
 Rochester, NY 14650
 Tel. 1-800-242-2424 for product info.

Ilford Photo Corporation
 70 West Century Boulevard
 Paramus, NJ 07653

Polaroid Corporation
 784 Memorial Drive
 Cambridge, MA 02139
 Tel. 1-800-255-1618 for product info.

E. I. duPont de Nemours and Co.
 Imaging Systems Department
 666 Driving Park Avenue
 Rochester, NY 14613

Fuji Photo Film USA
 555 Taxter Road
 Elmsford, NY 10523

3M Company
 Photo Color Systems Division
 3M Center
 St. Paul, MN 55144-1000

29.22 REFERENCES

1. W. Thomas (ed.), *SPSE Handbook of Photographic Science and Engineering*, Wiley, New York, 1973, sec. 8.
2. B. H. Carroll, G. C. Higgins, and T. H. James, *Introduction to Photographic Theory*, Wiley, New York, 1980, chaps. 6-9.
3. T. H. James (ed.), *The Theory of the Photographic Process*, 4th ed., Macmillan, New York, 1977.
4. J. Sturge, Vivian Walworth, and Alan Shepp (eds.), *Imaging Processes and Materials*, Van Nostrand Reinhold, New York, 1989, chap. 3. (See also Ref. 3, chap. 4.)
5. G. Haist, *Modern Photographic Processing*, vols. I and II, Wiley, New York, 1979.
6. Klaus Hendricks, chap. 20 of Ref. 4.
7. J. H. Altman, F. Grum, and C. N. Nelson, *Phot. Sci. Eng.* **17**:513 (1973).
8. Reference 3, chaps. IV-VII.
9. ANSI/ISO, May 2, 1991. (See also ANSI/ISO, May 1, 1984, which sets forth standardized terminology.)
10. R. W. G. Hunt, *The Reproduction of Colour*, 4th ed., Fountain Press, Tolworth, England, 1987, pp. 247-257.
11. *American National Standard*, ANSI PH2.2, 1984, R1989, PH2.17 1985.
12. C. J. Niederpruem, C. N. Nelson, and J. A. C. Yule, *Phot. Sci. Eng.* **10**:35 (1966).
13. *American National Standard*, ANSI PH3.49, 1971, R1987.
14. H. N. Todd and R. D. Zakia, *Phot. Sci. Eng.* **8**:249 (1964). (See also publications of the American National Standards Institute.)
15. M. Hercher and B. J. Ruff, *J. Opt. Soc. Am.* **57**:103 (1967).
16. W. F. Berg, *Proc. Roy. Soc. of London*, ser. A, **174**:5599 (1940).
17. Reference 2, p. 141.
18. *Scientific Imaging with Kodak Films and Plates*, Publication P315, Eastman Kodak Co., Rochester, N.Y., 1987, p. 63.
19. J. Kapecki and J. Rodgers, *Kirk-Othmer Encyclopedia of Chemical Technology*, 4th ed., vol. 6, Wiley, New York, 1993, p. 965.
20. S. H. Liggero, K. J. McCarthy, and J. A. Stella, *J. Imaging Technology* **10**:1 (1984).

21. Reference 5, chap. 7.
22. E. H. Land, *Phot. Sci. Eng.* **16**:247 (1972).
23. Reference 10, chap. 11 et seq.
24. See, for example, F. R. Clapper in Ref. 3, sec. II, chap. 19.
25. R. M. Evans, *Eye, Film, and Camera in Color Photography*, Wiley, New York, 1959, p. 180.
26. B. J. Thompson, *Progress in Optics VII*, E. Wolf (ed.), North-Holland Publishing Co., Amsterdam, 1969; R. E. Swing, *J. Opt. Soc. Am.* **62**:199 (1972); Dainty and Shaw, Ref. 27, chap. 9.
27. P. Kowalski, *Applied Photographic Theory*, Wiley, New York, p. 77 et seq.
28. C. N. Nelson, Ref. 3, chap. 19.
29. G. C. Higgins and L. A. Jones, *J. Soc. of Mot. Pict. & TV Engrs.* **58**:277 (1952).
30. J. Gasper and J. J. dePalma, Ref. 3, chap. 20.
31. J. C. Dainty and R. Shaw, *Image Science*, Academic Press, London, 1974, chaps. 6 and 7. (See also J. C. Dainty, *Optica Ada* **18**:795, 1971.)
32. R. M. Gendron, *J. Soc. of Mot. Pict. & TV Engrs.* **82**:1009 (1973).
33. E. W. H. Selwyn, *Photography Journal* **75**:571 (1935).
34. C. J. Bartleson, *Photography Journal* **33**:117 (1985).
35. K. F. Stultz and H. J. Zweig, *J. Opt. Soc. Am.* **49**:693 (1959).
36. F. W. Sears, *Optics*, Addison-Wesley, Reading, Mass., 1949, p. 159.
37. *American National Standard*, ANSI PH2.40, 1985, R1991.
38. Reference 31, chap. 8.
39. E. C. Doerner, *J. Opt. Soc. Am.* **52**:669 (1962).
40. C. J. Bartleson, *J. Phot. Sci.* **30**:33 (1982).
41. S. Ikenoue and M. Tabei, *J. Imaging Sci.* **34**:187 (1990).
42. Reference 2, p. 335.
43. A. Rose, *J. Mot. Pict. & TV Engrs.* **47**:273 (1946).
44. R. C. Jones, *Phot. Sci. Eng.* **2**:57 (1958).
45. Reference 31, chap. 8, p. 311 et seq.
46. Reference 4, table 3.2, p. 73.
47. G. R. Bird, R. C. Jones, and A. E. Ames, *Appl. Opt.* **2389** (1969).
48. *American National Standard*, ANSI PH2.33, 1983, R1990.
49. O. H. Schade, *J. Soc. Mot. Pict. & TV Engrs.* **73**:81 (1964).
50. J. H. Altman and H. J. Zweig, *Phot. Sci. Eng.* **7**:173 (1963).
51. R. Clark Jones, *J. Opt. Soc. Am.* **51**:1159 (1961).
52. Reference 31, chap. 10.

John D. Baloga

Imaging Materials and Media
Eastman Kodak Company
Rochester, New York

30.1 INTRODUCTION

Photographic materials are optical devices. Their fabrication and technical understanding encompass the field of optics in the broadest sense. Light propagation, absorption, scattering, and reflection must be controlled and used efficiently within thin multilayer coatings containing tiny light-detecting silver halide crystals and chemistry. Many subdisciplines within classical optics, solid-state theory, and photochemistry describe these processes.

In Chap. 20 of Ref. 1, Altman sweeps broadly through the basic concepts and properties of photographic materials. His brief, high-level descriptions are still generally valid for today's photographic products, and the reader will profit by reviewing that summary. This chapter focuses more sharply on four fundamental photographic technologies intimately related to the broad field of optics, then gives an overview of photographic materials available today.

Section 30.2 discusses the optical properties of multilayer color photographic films and papers. This section outlines the basic structure of the device and describes the principal function of each layer. Light absorption, reflection, scattering, and diffraction properties are discussed in the context of providing minimum optical distortion to the final image.

Silver halide light detection crystals are solid-state devices described by quantum solid-state photophysics and chemistry in addition to crystallography. These devices absorb light to create an electron-hole pair. Delocalization of the electron in the conduction band of the crystal, trapping, reduction of interstitial silver ions, nucleation and growth of clusters of silver atoms, and reversible regression of these phenomena must be controlled and optimized to produce the most sensitive light detectors that also possess long-term stability.

Section 30.3 describes the basic photophysics of silver halides according to our best understanding today. It outlines the solid-state properties most important to image detection and amplification. Surface chemical treatment and internal doping are discussed in the context of providing silver halide emulsions having the highest sensitivity to light.

A color photographic image is formed by high-extinction and high-saturation dyes structurally designed by chemists to optimize their color, permanence, and other useful characteristics. Their properties in both the ground state and photoexcited states are important to color and stability.

Section 30.4 briefly outlines the photochemistry of photographic image dyes. Excited-state properties are described that are important to the photostability of these dyes for image permanence.

Color science guides us to optimize all factors in photographic materials conspiring to render an accurate and pleasing image. These include spectral discrimination during the light detection phase and color correction to compensate for imperfect spectral detection and imperfect image dyes.

Section 30.5 sketches the photophysics and color science of photographic spectral sensitizers with an aim toward describing how modern photographic films sense the world in color nearly as the human eye sees it.

Today there exists a large diversity of photographic films. In addition to different manufacturers' brands, films differ by speed, type, color attributes, format, and a multitude of other factors. This can be confusing. And what are the differences between a consumer film bought in a drugstore and a more expensive professional film?

Section 30.6 gives an overview of the different types of films available today. Differences between high- and low-speed films are described with an understanding of the origins of these differences. Consumer films are compared to professional films and some of the special needs of each type of customer are described. Some general guidelines are offered for film choices among color reversal films, black-and-white films, and color negative films.

In addition to Chap. 20 in Ref. 1, several other texts contain useful information about photographic materials. Kapecki and Rodgers² offer a highly lucid and digestible sketch of basic color photographic film technology. Besides Chap. 20 in Ref. 1, this is a good place for the technically astute but nonpractitioner to gain a high level understanding of basic color photography. The technically detailed treatise by James³ is perhaps the best single comprehensive source containing the science and technology of the photographic system outside photographic manufacturers' internal proprietary libraries. Hunt⁴ provides a good comprehensive treatise on all aspects of color science and additionally gives a good technical review of color photographic materials. Chapter 6 in Ref. 1 contains a useful overview of the theory of light-scattering by particles, Chap. 9 contains a good overview of optical properties of solids, and Chap. 26 on colorimetry provides a good introduction to the basic concepts in that field.

30.2 THE OPTICS OF PHOTOGRAPHIC FILMS AND PAPERS

Introduction

Color photographic materials incorporate design factors that optimize their performance across all features deemed important to customers who use the product. These materials are complex in composition and structure. Color films typically contain over 12 distinct optically and chemically interacting layers. Some of the newest professional color reversal films contain up to 20 distinct layers.* Each layer contributes a unique and important function to the film's final performance. Careful design and arrangement of the film's various layers ultimately determine how well the film satisfies the user's needs.

One very important customer performance feature is film *acutance*, a measure of the film's ability to clearly record small detail and render sharp edges. Because images recorded on color film are frequently enlarged, image structure attributes such as acutance are very important. Magnifications such as those used to place the image on a printed page challenge the film's ability to clearly record fine detail.

The ability of a photographic material to record fine detail and sharp edges is controlled by two factors. The first includes light scatter, diffraction, and reflections during the exposure step. These are collectively called *optical* effects. This factor dominates the film's ability to record fine detail and sharp edges and is therefore critical to the film's design.

*Fujichrome Velvia 50 RVP professional film.

Chemical adjacency effects collectively known as *Eberhard effects* (see Chap. 21 in Ref. 3) are the second factor. In today's color and black-and-white film products these are present during film development and are caused by accumulating development by-products such as iodide ions or inhibitors released by development inhibitor–releasing (DIR) chemicals that restrain further silver development in exposed regions of the image, leading to a chemical unsharp mask effect (pp. 274 and 365 of Ref. 4). These chemical signals also give rise to the film's interlayer interimage effects (HE; see p. 278 of Ref. 4) used for color correction. Film sharpness and HE are strongly related.

This section describes multilayer factors that contribute to light scatter, reflection, diffraction, and absorption in photographic materials. Factors that influence the nonoptical part of film acutance, such as Eberhard effects, were briefly reviewed by Altman.¹ More information about these processes can be found in references cited therein.

Structure of Color Films

Color film structures contain image-recording layers, interlayers, and protective overcoat layers. These layers suspend emulsions and chemistry in a hardened gelatin binder coated over a polyacetate or polyester support. Figure 1 shows the principal layers in a color multilayer film structure.

The overcoat and ultraviolet (UV) protective layers contain lubricants; plastic beads 1 to 5 μm in size called *matte* to impart surface roughness that prevents sticking when the film is stored in roll form; UV-light-absorbing materials; and other ingredients that improve the film's handling characteristics and protect the underlying light-sensitive layers from damage during use and from exposure to invisible UV light. Ultraviolet light is harmful for two reasons: it will expose silver halide emulsions, thereby rendering an image from light invisible to humans, and it promotes photodecomposition of image dyes, leading to dye fade over time (p. 977 of Ref. 2). Visible light exposure must first pass through these overcoats, but they typically do little harm to acutance as they contain nothing that seriously scatters visible light.

The blue-light-sensitive yellow dye imaging layers appear next. Silver halides, with the exception of AgCl—used primarily in color papers—have an intrinsic blue sensitivity even when spectrally sensitized to green or red light. Putting the blue-light-sensitive layers on top avoids incorrect

Overcoat
UV protective layer
Fast yellow layer
Slow yellow layer
Yellow filter layer
Barrier layer
Fast magenta layer
Mid magenta layer
Slow magenta layer
Optional magenta filter layer
Barrier layer
Fast cyan layer
Mid cyan Layer
Slow cyan layer
AHU layer
Plastic support
Optional pelloid AHU layer

FIGURE 1 Simplified diagram showing the key layers in a color photographic film (not drawn to scale).

exposure leading to color contamination because a blue-light-absorbent yellow filter layer, located beneath the blue-sensitive imaging layers, absorbs all blue light that has passed completely through the blue-sensitive layers.

A collection of layers adjacent to one another that contain silver halide sensitized to a common spectral band is called a *color record*. Color records in films are always split into two or three separate layers. This arrangement puts the highest-sensitivity emulsion in the upper (fast) layer, where it gets maximum light exposure for photographic speed, and places low-sensitivity emulsions into the lower layer (or layers), where they provide exposure latitude and contrast control.

One or more interlayers separate the yellow record from the green-light-sensitive magenta record. These upper interlayers filter blue light to avoid blue light exposure color contaminating the red- and green-sensitive emulsion layers below. This behavior is called *punch through* in the trade. These interlayers also contain oxidized developer scavengers to eliminate image dye contamination between color records caused by oxidized color developer formed in one color record diffusing into a different color record to form dye.

The traditional yellow filter material is Carey Lea silver (CLS), a finely dispersed colloidal form of metallic silver, which removes most of the blue light at wavelengths less than 500 nm. The light absorption characteristics of this finely dispersed metallic silver are accounted for by Mie theory.⁵ Unfortunately this material also filters out some green and red light, thus requiring additional sensitivity from the emulsions below. Bleaching and fixing steps in the film's normal processing remove the yellow CLS from the developed image.

Some films incorporate a yellow filter dye in place of CLS. Early examples contained yellow dyes attached to a polymer mordant in the interlayer. A *mordant* is a polymer that contains charged sites, usually cationic, that binds an ionized anionic dye by electrostatic forces. These dyes are removed during subsequent processing steps. Although these materials are free from red and green light absorption, it is a challenge to make them high in blue light extinction to avoid excessive chemical loads in the interlayer with consequent thickening.

Most recently, solid-particle yellow filter dyes⁶ are seeing more use in modern films. These solid dyes are sized to minimize light scatter in the visible region of the spectrum and their absorption of blue light is very strong. They can be made with very sharp spectral cuts and are exceptionally well suited as photographic filter dyes. In some films solid-particle magenta filter dyes⁷ are also used in the interlayers below the magenta imaging layers. Solid-particle filter dyes are solubilized and removed or chemically bleached colorless during the film's normal development process.

Because the human visual system responds most sensitively to magenta dye density, the green-light-sensitive image recording layers are positioned just below the upper interlayers, as close as practical to the source light exposure side of the film. This minimizes green light spread caused when green light passes through the turbid yellow record emulsions. It gives the maximum practical film acutance to the magenta dye record.

A lower set of interlayers separates the magenta record from the red-light-sensitive cyan dye record. These lower interlayers give the same type of protection against light and chemical contamination as do the upper interlayers. The magenta filter dye is absent in some films because red-light-sensitive emulsions are not as sensitive to green light exposure as to blue light exposure.

Located beneath the cyan record, just above the plastic film support, are antihalation undercoat (AHU) layers. The black absorber contained in this bottom layer absorbs all light that has passed completely through all imaging layers, thereby preventing its reflection off the gel-plastic and plastic-air interfaces back into the imaging layers. These harmful reflections cause a serious type of light spread called *halation*, which is most noticeable as a halo around bright objects in the photographic image.

The opacity required in the AHU is usually obtained by using predeveloped black filamentary silver, called *gray gel*, which is bleached and removed during the normal photographic processing steps. Alternatively, in many motion picture films a layer of finely divided carbon suspended in gelatin (*rem jet*) is coated on the reverse side of the film. When placed in this position the layer is called an *AHU pelloid*. It is physically removed by scrubbing just before the film is developed. The newest motion picture films incorporate a black solid-particle AHU filter dye that is solubilized and removed during normal development of the film and does not require a separate scrubbing step.

The overall thickness of a film plays an important role in minimizing harmful effects from light scatter. Because color films are structured with the yellow record closest to the light exposure source, it is especially important to minimize thickness of all film layers in the yellow record and below it, because light scattered in the yellow record progressively spreads as it passes to lower layers. The cyan record shows the strongest dependence on film thickness because red light passes through both yellow and magenta record emulsions en route to the cyan record. Because both emulsions often contribute to red light scatter, the cyan record suffers a stronger loss in acutance with film thickness.

Structure of Color Papers

The optical properties of photographic paper merit special consideration because these materials are coated in very simple structures on a highly reflective white Baryta-coated paper support. *Baryta* is an efficient diffuse reflector consisting of barium sulfate powder suspended in gelatin that produces isotropically distributed reflected light with little absorption.

Photographic papers generally contain about seven layers. On top are the overcoat and UV protective layers, which serve the same functions as previously described for film.

The imaging layers in color papers contain silver chloride emulsions for fast-acting development, which is important to the industry for rapid throughput and productivity. Generally only one layer is coated per color record, in contrast to films, which typically contain two or three layers per color record. The order of the color records in photographic papers differs from that in films due mainly to properties of the white Baryta reflective layer.

Because color paper is photographically slow, exposure times on the order of seconds are common. Most light from these exposures reflects turbidly off the Baryta layer. The imaging layers getting the sharpest and least turbid image are those closest to the reflective Baryta where light spread is least, not those closest to the top of the multilayer as is the case with film.

In addition, the Baryta layer as coated is not smooth. Its roughness translates into the adjacent layer, causing nonuniformities in that layer. Because the human visual system is most forgiving of physical imperfections in yellow dye, the yellow color record must be placed adjacent to the Baryta to take the brunt of these imperfections.

The magenta color record is placed in the middle of the color paper multilayer, as close to the Baryta layer as possible, since sharpness in the final image is most clearly rendered by magenta dye. This leaves the cyan record nearest to the top of the structure, just below the protective overcoats.

The magenta color record in the middle of the multilayer structure is spaced apart from the other two color records by interlayers containing oxidized developer scavengers to prevent cross-record color contamination, just as in films. However, no filter dyes are needed in color paper because silver chloride emulsions have no native sensitivity to light in the visible spectrum.

Light Scatter by Silver Halide Crystals

Because they consist of tiny particles, silver halide emulsions scatter light. Scattering by cubic and cubo-octahedral emulsions at visible wavelengths is most intense when the emulsion dimension ranges from 0.3 to 0.8 μm , roughly comparable to the wavelengths of visible light. This type of scattering is well characterized by Mie⁸ theory and has been applied to silver halides.⁹ We are often compelled to use emulsions having these dimensions in order to achieve photographic speed.

For photographic emulsions of normal grain size and concentration in gelatin layers of normal thickness, multiple scattering predominates. This problem falls within the realm of radiative transfer theory. Pitts¹⁰ has given a rigorous development of radiative transfer theory to the problem of light scattering and absorption in photographic emulsions. A second approach that has received serious attention is the Monte Carlo technique.¹¹ DePalma and Gasper¹² were able to obtain good agreement between their modulation transfer functions (MTFs) determined by a Monte Carlo calculation and experimentally measured MTFs for a silver halide emulsion layer coated at various thicknesses.

Scattering by yellow record emulsions is especially harmful because red and green light must first pass through this record en route to the magenta and cyan records below. All other things being equal, the amount of light scattered by an emulsion layer increases in proportion to the total amount of silver halide coated in the layer.

Color films are constructed with low-scattering yellow record emulsions whenever possible. The amount of silver halide coated in the yellow record is held to a minimum consistent with the need to achieve the film's target sensitometric scale in yellow dye.

High-dye-yield yellow couplers¹³ having high coupling efficiency are very useful in achieving silver halide mass reductions in the yellow record of color reversal films. These provide high yellow dye densities per unit silver halide mass, thereby reducing the amount of silver halide needed to achieve target sensitometry. Some upper-scale (high-density) increase in granularity often results from using these couplers in a film's fast yellow layer, but the benefits of reduced red and green light scatter overcome this penalty, especially because the human visual system is insensitive to yellow dye granularity.

Tabular emulsion grains offer a way to achieve typical photographic speeds using large-dimension emulsions, typically 1.0 to 3.0 μm in diameter, although smaller- and larger-diameter crystals are sometimes used. These do not scatter light as strongly at high angles from normal incidence with respect to the plane of the film as do cubic or octahedral emulsions having comparable photographic speeds. However, diffraction at the edges and reflections off the crystal faces can become a detrimental factor with these emulsions.

Silver halide tabular crystals orient themselves to lie flat in the plane of the gelatin layer. This happens because shear stress during coating of the liquid layer stretches the layer along a direction parallel to the support and also because the water-swollen thick layer compresses flat against the support after it dries by a ratio of roughly 20:1.

Reflections can become especially harmful with tabular emulsion morphologies since light reflecting from the upper and lower faces of the crystal interferes, leading to resonances in reflected light. The most strongly reflected wavelengths depend on the thickness of the tabular crystal. The thickness at which there is a maximum reflectance occurs at fractional multiples of the wavelength¹⁴ given by:

$$t = \frac{\left(m + \frac{1}{2}\right)\lambda}{2n}$$

where t is the thickness of the tabular crystal, λ is the wavelength of light, n is the refractive index of the crystal, and m is an integer.

In an extreme case, tabular emulsions act like partial mirrors reflecting light from grain to grain over a substantial distance from its point of origin. This is called *light piping* by analogy with light traveling through an optical fiber (see Fig. 2). It is especially serious in the cyan record, where red light often enters at angles with respect to perpendicular incidence caused by scattering in upper layers.

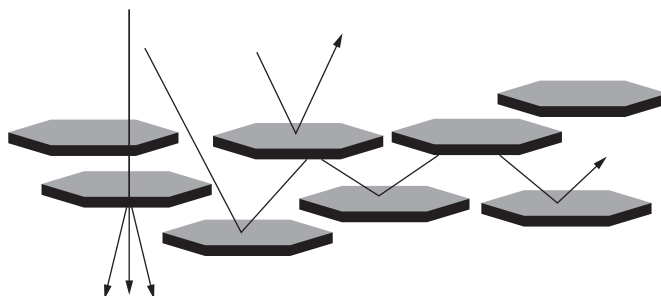


FIGURE 2 Light undergoes reflection, diffraction, and multiple reflections that may lead to light piping in film layers containing tabular silver halide crystals.

Another cause of light piping is the presence of small-grain emulsion contaminants within a tabular grain population. A highly turbid contaminant emulsion, even if present in small amounts, can scatter light at sufficient angles to induce severe light piping through the tabular grain matrix. Modern emulsion science takes great pains to produce tabular emulsions free from highly scattering contaminants.

Gasper¹⁵ treated the problem of scattering from thin tabular silver halide crystals in a uniform gelatin layer. For very thin and large-diameter tabular silver bromide crystals (thickness $<0.06\ \mu\text{m}$, diameter $>1.0\ \mu\text{m}$) light at normal incidence is scattered almost equally in the forward and backward hemispheres. As the grain thickness increases there appears increasing bias for forward scatter in a narrow cone. The efficiencies for scatter and absorption were found to be independent of grain diameter for crystal diameters larger than $\sim 1.0\ \mu\text{m}$. For such crystals, the backscatter and absorption efficiencies are approximately equal to the specular reflectance and absorption of a semiinfinite slab of the same material and thickness.

But the forward scattering efficiency does not approximate the specular transmittance of the semiinfinite slab as a result of interference between the directly scattered forward field and the mutually coherent unscattered field, a process analogous to diffraction that does not occur for the semiinfinite slab with no edge. The specific turbidity depends on grain thickness, but not diameter for diameter $>1.0\ \mu\text{m}$.

A light-absorbing dye is sometimes added to a film to reduce the mean free path of scattered and reflected light leading to an acutance improvement. This improvement happens at the expense of photographic speed, since light absorbed by the dye is not available to expose the silver halide emulsions. However, the trade-off is sometimes favorable if serious light piping is encountered.

Light-absorbing interlayers between color records are sometimes used in color films to help eliminate the harmful effects of reflected light. For example, some films have a magenta filter dye interlayer between the magenta and cyan records. In addition to its usefulness in reducing light punchthrough (green light passing through the magenta record to expose the red-sensitized layers), this filter layer helps eliminate harmful reflections of green light off cyan record tabular emulsions, which bounce back into the magenta record. These reflections, although potentially useful for improving photographic green speed, can be harmful to magenta record acutance.

30.3 THE PHOTOPHYSICS OF SILVER HALIDE LIGHT DETECTORS

Chapter 20 in Ref. 1 gave a very brief sketch of the general characteristics of silver halide crystals used for light detection and amplification in photographic materials. These crystals, commonly termed *emulsions*, are randomly dispersed in gelatin binder layers in photographic films and papers. For photographic applications the most commonly used halide types are AgCl, AgBr, and mixed halide solid solutions of AgCl_xI_y and AgBr_xI_y . In addition, pure phases of AgCl and AgI are sometimes grown epitaxially on AgBr crystals to impart special sensitivity and development properties.

Upon exposure to light, silver halide crystals form a latent image (LI) composed of clusters of three to hundreds of photoreduced silver atoms either within the interior or most usefully on the surface of the crystal where access by aqueous developing agents is easiest.* Higher light exposures result in larger numbers of silver atoms per latent image cluster on average in addition to exposing a larger number of crystals on average.

The detection of light by a silver halide crystal, subsequent conversion to an electron hole pair, and ultimate reduction of silver ions to metallic silver atoms is in essence a solid-state photophysical process. The application of solid-state theory to the photographic process began with Mott and

*Some specialized developing agents can etch into the silver halide crystal to develop the internal latent image (LI), but the most commonly used color negative and color paper developers have little capability to do this. They are primarily surface-developing agents. The color reversal black-and-white developer can develop slightly subsurface LI in color reversal emulsions.

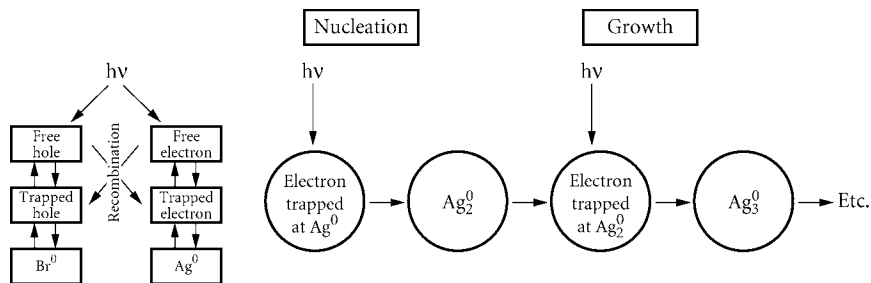


FIGURE 3 Diagram showing the basic features of the nucleation and growth mechanism for formation of a developable latent image site by light exposure of a silver halide crystal.

Gurney¹⁶ and was extended by Seitz.¹⁷ Additional reviews were published^{18–20} as our understanding of these processes grew.

According to the nucleation and growth mechanism,¹⁸ the photographic LI forms through a series of reversible electronic and ionic events (see Fig. 3). An electron is generated in the conduction band either by direct band gap photoexcitation of the crystal in its intrinsic spectral absorption region or by injection from a photoexcited spectral sensitizing dye on the crystal's surface. As this mobile electron travels through the crystal it becomes shallowly and transiently trapped at lattice imperfections carrying a positive charge. The ionized donor may be an interstitial silver ion or a surface defect. The kinked edge of a surface terrace, having a formal charge of $+1/2$ is one such site.

The electron may sample many such traps before becoming more permanently trapped at a preferred site. Trapping of the electron occurs in less than 0.1 ns at room temperature.²¹ The shallowly trapped electron orbits about the trapping site in a large radius. At this stage the formal charge of the trap site becomes $-1/2$ and adjacent surface lattice ions relax, thereby increasing the depth of the trap.²² In octahedral AgBr emulsions this shallow-deep transition occurs within 10 ns at room temperature.²¹

The negative charge formally associated with the trapped electron attracts a mobile interstitial Ag^+ ion, which the trapped electron reduces to an Ag^0 atom. The site's charge reverts to $+1/2$. A second photoelectron is then trapped at the site and reacts with a second interstitial Ag^+ ion to form Ag_2^0 and so on as additional photons strike the crystal. Frenkel defect concentrations—interior plus surface—in the range of $10^{14}/\text{cm}^3$ have been reported for AgBr crystals.¹⁹ Latent image formation depends critically upon the existence of partially charged defects on the surface of the crystal.

Because the single-atom Ag^0 state is metastable and the Ag_2^0 state is not usefully detectable from gross fog by developing solutions (Chap. 4 of Ref. 3), a minimum of three photons must strike the crystal to form a detectable LI site, the second photon within the lifetime of the metastable trapped Ag^0 atom.* The transient existence of one or two silver atom sites has not been directly observed but is strongly inferred by indirect evidence.²⁰

The highest photoefficiencies are achieved after silver halide emulsions are heated with aqueous sulfur and gold salts to produce Ag_2S and AgSAu species on the surface of the crystal.²² Most of the evidence for sulfur sensitization suggests these sites improve electron trapping, perhaps by making trapping lifetimes longer at preexisting kink sites. Gold sensitization accomplishes at least two things: it reduces the size of the latent image center needed to catalyze development and it stabilizes silver atoms formed during the exposure process (Chap. 5 of Ref. 3).

Ideally only a single latent image site forms per exposed silver halide crystal. If more sites nucleate on a crystal, these can compete for subsequent photoelectrons, leading to losses in efficiency related to LI dispersity. One remarkable aspect of this process is that despite the large number of

*A metastable silver atom may dissociate back into an electron-hole pair, and the electron can migrate to another trapping site to form a new silver atom. This may happen many times in principle. The second electron must encounter the silver atom before electron-hole recombination or other irreversible electron loss takes place.

kink sites and chemically sensitized sites on a given crystal, in the most efficient silver halide emulsions only a single LI center forms per grain. The shallow-deep transition that occurs after the initial electron trapping selects one preferred defect out of the plethora of other possibilities for subsequent electron trapping.

The two predominant crystal faces that occur on photographically useful silver halide grains are [100] and [111]. Both surfaces have a negative charge²³ ranging from -0.1 to -0.3 V. A subsurface space charge layer rich in interstitial Ag^+ ions compensates for the charged surface.²⁴ The [100] surfaces are flat with steps containing jogs that are either positively or negatively charged kink sites. The positive kinks on the surface form shallow electron traps.

In contrast, the [111] surface is not well characterized and is believed to be sensitive to surface treatment. Calculations of surface energies suggest that several arrangements may coexist, including microfaceted [100] planes and half layers of all silver or all halide ions in arrangements with hexagonal symmetry.^{24,25}

All silver halides absorb light in the near-UV region of the spectrum. The absorption edge extends to longer wavelengths from AgCl to AgBr to AgI . AgCl and AgBr have indirect band gap energies of 3.29 and 2.70 eV, respectively.¹⁹ Absorption of light produces electrons in a conduction band whose minimum lies at the zone center and holes in the valence band whose maximum is at an L point. Carriers produced at excess energy rapidly thermalize at room temperature since electron mobility is limited by phonon scattering.

Electron mobility in AgBr is about $60 \text{ cm}^2/\text{Vs}$ at room temperature.¹⁹ In AgCl the hole mobility is low, but in AgBr it is substantial—only about a factor of 30 lower than electron mobility.²⁰ Carriers generated by exposure will explore the crystal in a random walk as they scatter off phonons and transiently trap at charged defects. It is expected that electrons can fully sample a $1\text{-}\mu\text{m}$ crystal within less than $1 \mu\text{s}$.

In most commercial photographic materials today, only about 25 percent to 50 percent of the electrons injected into the crystal's conduction band contribute to a developable LI site. The rest are wasted through electron-hole recombination, formation of internal LI that is inaccessible to the developing agent, or formation of multiple nucleation sites termed *dispersity*. The deliberate addition of dopants that act as shallow electron traps reduces the time electrons spend in the free carrier state and thereby limits their propensity to recombine with trapped holes.

Polyvalent transition metal ions are frequently doped into silver halide crystals to limit reciprocity effects, control contrast, and reduce electron hole recombination inefficiencies.²⁶ They act as electron or hole traps and are introduced into the crystal from aqueous hexa-coordinated complexes during precipitation. They generally substitute for $(\text{AgX}_6)^{5-}$ lattice fragments. Complexes of Ru and Ir have been especially useful. Because the dopant's carrier-trapping properties depend on the metal ion's valence state and the stereochemistry of its ligand shell, there are many opportunities to design dopants with specific characteristics. Dopants incorporating Os form deep electron traps in AgCl emulsions with an excess site charge of +3. An effective electron residence lifetime of 550 seconds has been measured for these dopants at room temperature.²⁰ They are used to control contrast at high-exposure regions in photographic papers.

Quantum sensitivity is a measure of the average number of photons absorbed per grain to produce developability in 50 percent of an emulsion population's grains (Chap. 4 in Ref. 27). This microscopic parameter provides a measure of the photoefficiency of a given emulsion; the lower the quantum sensitivity value, the more efficient the emulsion. The quantum sensitivity measurement gives a cumulative measure of latent image formation, detection, and amplification stages of the imaging chain.

Quantum sensitivity has been measured for many emulsions. The most efficient emulsions specially treated by hydrogen hypersensitization yield a quantum sensitivity of about three photons per grain.²⁸ Although hydrogen hypersensitization is not useful for general commercial films because it produces emulsions very unstable toward gross fog, this sensitivity represents an ambitious goal for practical emulsions used in commercial photographic products. The most efficient practical photographic emulsions reported to date have a quantum sensitivity of about five to six photons per grain.²⁹ The better commercial photographic products in today's market contain emulsions having quantum sensitivities in the range of 10 to 20 photons per grain.

30.4 THE STABILITY OF PHOTOGRAPHIC IMAGE DYES TOWARD LIGHT FADE

Introduction

Azomethine dyes are formed in most photographic products by reaction between colorless couplers and oxidized *p*-phenylenediamine developing agents to form high-saturation yellow, magenta, and cyan dyes (Chap. 12 in Ref. 3). Typical examples are shown in Fig. 4. The diamine structural element common to all three dyes comes from the developing agent, where R_1 and R_2 represent alkyl groups. The R group extending from the coupler side of the chromophore represents a lipophilic ballast, which keeps the dye localized in an oil phase droplet. In some dyes it also has hue-shifting properties.

Heat, humidity, and light influence the stability of these dyes in photographic films and papers.³⁰ Heat and humidity promote thermal chemical reactions that lead to dye density loss. Photochemical processes cause dyes to fade if the image is displayed for extended periods of time. Ultraviolet radiation is especially harmful to a dye's stability, which is partly why UV absorbers are coated in the protective overcoat layers of a color film or paper.

Stability toward light is especially important for color papers, where the image may be displayed for viewing over many years. It is less important in color negative film, which is generally stored in the dark and where small changes in dye density can often be compensated for when a print is made. Light stability is somewhat important for color reversal (slide and transparency) film, since a slide is projected through a bright illuminant, but projection and other display times are short compared to values for color papers. A similar situation exists for movie projection films, where each frame gets a short burst of high-intensity light but the cumulative exposure is low.

Evaluation of the stability of dyes toward light is difficult because the time scale of the photochemical reactions, by design, is very slow or inefficient. The quantum yields of photochemical fading of photographic dyes, defined as the fraction of photoexcited dyes that fade, are on the order of 10^{-7} or smaller.^{2,31} Accelerated testing using high-intensity illumination can sometimes give misleading results if the underlying photochemical reactions change with light intensity (p. 266 in Ref. 4).

Given that dyes fade slowly over time, it is best if all three photographic dyes fade at a common rate, thereby preserving the color balance of the image (p. 267 in Ref. 4). This rarely happens. The perceived light stability of color photographic materials is limited by the least stable dye. Historically, this has often been the magenta dye, whose gradual fade casts a highly objectionable green tint to a picture containing this dye. This is most noticeable in images containing memory colors, such as neutral grays and skin tones.

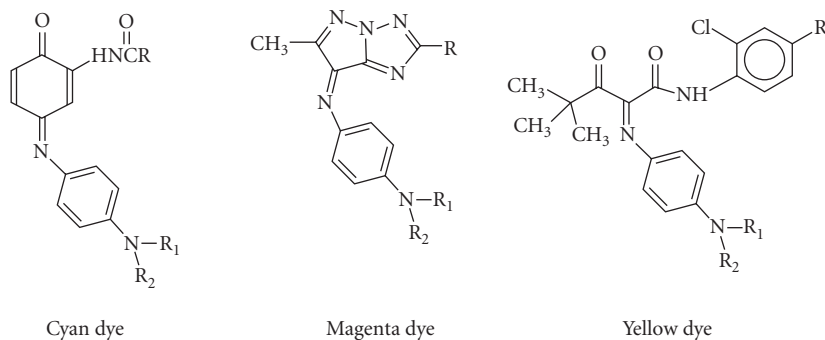


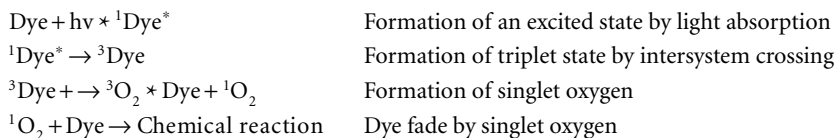
FIGURE 4 Typical examples of azomethine dyes used in photographic materials.

Photochemistry of Azomethine Dyes

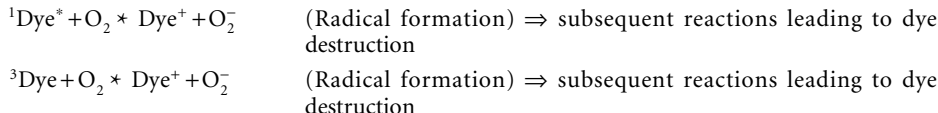
Upon absorption of a photon, a dye becomes reversibly excited to the singlet state. For azomethine dyes, the lifetime of the excited state is estimated to be on the order of picoseconds.^{32–39} Similarly, the lifetime of the triplet state of azomethine dyes, which is expected to form rapidly by intersystem crossing from the singlet excited state, has been estimated to be in the nanosecond range.³⁶ The short lifetime of these species seems to be at the basis of the observed low quantum yields of photochemical fading of azomethine dyes.

The nature of the initial elementary reactions involving excited or triplet states of dyes is not well understood. For some magenta dyes, the fading mechanism was reported to be photooxidative.³⁷ Cyan dye light fade appears to involve both reductive and oxidative steps³⁸ following photoexcitation of the dye.

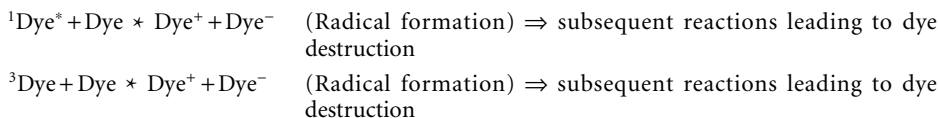
Singlet oxygen has been traditionally postulated to be involved in the oxidative pathway.³⁹ This high-energy species can, in principle, be formed by energy transfer between a triplet dye and molecular oxygen, and can subsequently induce dye destruction, although reaction with singlet oxygen does not always lead to decomposition.⁴⁰



Studies of the photophysical properties of photographic dyes, however, have shown that in general, azomethine dyes are good quenchers of singlet oxygen.⁴¹ This implies that the formation of singlet oxygen by the triplet state of the dyes should be inefficient. An alternative feasible role for molecular oxygen in dye fade involves electron transfer mechanisms.



The chemistry of dye fade may also be started by electron transfer between excited or triplet dye and a dye molecule in its ground state.



Excited State Properties

The few papers that have been published about excited singlet-state properties of azomethine dyes have mainly focused on pyrazolotriazole magenta dyes.^{32–35} These dyes have fluorescence quantum yields on the order of 10^{-4} at room temperature, but increase to ~ 1 in rigid organic glasses at 77 K.⁴⁰ The fluorescence quantum yield is the fraction of photoexcited molecules that emit a quantum of light by fluorescence from an excited state to the ground state having the same multiplicity (usually singlet to singlet). The results at room temperature imply singlet-state lifetimes of only a few picoseconds.⁴⁰

Room-temperature flash photolysis has identified a very short-lived fluorescent state plus a longer-lived nonfluorescent transient believed to be an excited singlet state whose structure is twisted compared to the ground state. This is consistent with the temperature-dependent results that allow rapid conformational change to the nonfluorescent singlet species at room temperature,

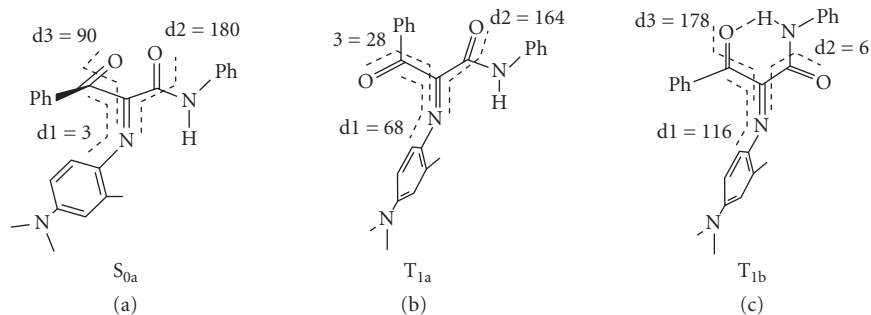


FIGURE 5 (a) Ground state conformation and (b and c) two conformations of the lowest triplet state.

in addition to intersystem crossing to the triplet excited state. But this conformational change is presumably restrained in a cold organic glass matrix, thereby allowing the excited singlet state time to fluoresce back to the ground state.

Investigation of the triplet states of azomethine dyes has proven difficult, and there are no definitive reports of direct observation of the triplet states. There are no reliable reports of phosphorescence from the triplet state of these dyes, although they have been studied by flash photolysis.^{32–35,41–45} Using indirect energy transfer methods, triplet lifetimes have been estimated to be less than 10 ns.³⁶ Similar studies⁴³ gave triplet energies of 166 to 208 kJ mol⁻¹ for yellow dyes, 90 to 120 kJ mol⁻¹ for magenta dyes, and about 90 kJ mol⁻¹ for cyan dyes.

Computational studies have been carried out recently on the ground and lowest triplet states of yellow azomethine dyes.⁴⁵ Consistent with crystallographic studies,^{46,47*} the calculated lowest-energy ground state conformation S_0 is calculated to have the azomethine C=N bond coplanar with the anilide C=O carbonyl, but perpendicular to the C=O carbonyl of the ketone. The energy of the vertical transition triplet state having this conformation is quite high, calculated to be 230 kJ mol⁻¹ above the ground state. Energy minimization on the lowest triplet energy surface starting with the geometry of S_0 results in a relaxation energy of 100 kJ mol⁻¹ and leads to T_{1a} , characterized by a substantial twist around the C=N bond and by increased planarity of the keto carbonyl with the azomethine plane. These conformational changes offer an efficient mechanism for stabilizing the lowest triplet state in these dyes. A second triplet conformer T_{1b} with an energy 25 kJ mol⁻¹ below T_{1a} was also calculated (see Fig. 5).

The low energy values calculated for T_{1a} and T_{1b} relative to S_{0a} (130 and 105 kJ mol⁻¹, respectively) have led to a new interpretation of the reported results⁴³ of energy transfer to S_{0a} . Using the nonvertical energy transfer model of Balzani,⁴⁸ it is shown that observed rates of energy transfer to S_{0a} can be consistent with formation of two distinct triplet states: one corresponding to a higher-energy excited triplet state (T_2) 167 kJ mol⁻¹ above S_{0a} , and the second corresponding to one conformation of the lowest triplet state (T_1) with an energy of 96 kJ mol⁻¹ above S_{0a} , in good agreement with the calculated T_{1b} configuration. The large structural differences between S_{0a} and the conformers of T_1 can explain the lack of phosphorescence in this system.

Light Stabilization Methods

Stabilization methods focus on elimination of key reactants (such as oxygen by various barrier methods), scavenging of reactive intermediates such as free radicals transiently formed during photodecomposition, elimination of ultraviolet light by UV absorbers, or quenching of photoexcited species in the reaction sequence.

*Other unpublished structures exhibiting the same conformation of S_{0a} have been solved at Kodak.

Stabilizer molecules⁴⁹ have been added to photographic couplers to improve the stability of their dyes toward light fade. Although the exact mode by which these stabilizers operate has not been disclosed, it is probable that some quench the photoexcited state of the dyes, thereby preventing further reaction leading to fading; some scavenge the radicals formed during decomposition steps; and some may scavenge singlet oxygen.

Polymeric materials added to coupler dispersions also stabilize dyes. These materials are reported to improve thermal dye decomposition⁵⁰ by physically separating reactive components in a semi-rigid matrix or at least decreasing the mobility of reactive components within that matrix. They may provide benefits for light stability of some dyes by a similar mechanism.

Photographic dyes sometimes form associated complexes or microcrystalline aggregates (p. 322 of Ref. 4). These have been postulated to either improve light stability by electronic-to-thermal (phonon) quenching of excited states, or to degrade light stability by concentrating the dye to provide more available dye molecules for reaction with light-induced radicals. Both postulates may be correct depending upon the particular dyes.

Modern photographic dyes have vastly improved light stabilities over dyes of the past. For example, the magenta dye light stability of color paper dyes has improved by about two orders of magnitude between 1942 and the present time.⁵¹ New classes of dyes^{52,53} have proved especially resistant to light fade and have made memories captured by photographic materials truly archival.

30.5 PHOTOGRAPHIC SPECTRAL SENSITIZERS

Introduction

Spectral sensitizing dyes are essential to the practice of color photography. In 1873, Vogel⁵⁴ discovered that certain organic dyes, when adsorbed to silver halide crystals, extend their light sensitivity to wavelengths beyond their intrinsic ultraviolet-blue sensitivity. Since then, many thousands of dyes have been identified as silver halide spectral sensitizers having various degrees of utility, and dyes exist for selective sensitization in all regions of the visible and near-infrared spectrum. These dyes provide the red, green, and blue color discrimination needed for color photography.

The most widely used spectral sensitizing dyes in photography are known as *cyanine dyes*. One example is shown in Fig. 6. These structures are generally planar molecules as shown.

A good photographic spectral sensitizing dye must adsorb strongly to the surface of the silver halide crystal. The best photographic dyes usually self-assemble into aggregates on the crystal's surface. Aggregated dyes may be considered partially ordered two-dimensional dye crystals.⁵⁵ Blue-shifted aggregates are termed *H-aggregates*, while red-shifted ones, which generally show spectral narrowing and excitonic behavior, are termed *J-aggregates*. Dyes that form J-aggregates are generally most useful as silver halide spectral sensitizers.

A good sensitizing dye absorbs light of the desired spectral range with high efficiency and converts that light into a latent image site on the silver halide surface. The relative quantum efficiency⁵⁶ of sensitization is defined as the number of quanta absorbed at 400 nm in the AgX intrinsic absorption region to produce a specified developed density, divided by the number of like quanta absorbed only by dye within its absorption band. For the best photographic sensitizing dyes this number is only slightly less than unity.

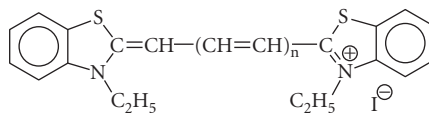


FIGURE 6 Typical cyanine spectral sensitizing dye used in photographic materials.

Adsorption isotherms for J-aggregated cyanine dyes on AgX follow Langmuir behavior with a saturation coverage that indicates closely packed structures.⁵⁷ The site area per molecule derived from the Langmuir isotherms implies dye adsorption along its long axis with heterocyclic rings perpendicular to the surface. Recent scanning tunneling microscope work^{58,59} has convincingly detailed the molecular arrangement within these aggregates, confirming general expectations.

The Photophysics of Spectral Sensitizers on Silver Halide Surfaces

Sensitizing dyes transfer an electron from the dye's excited singlet state to the conduction band of AgX. Evidence for this charge injection comes from extensive correlations of sensitization with the electrochemical redox properties of the dyes.^{22,60,61} In the adsorbed state a good sensitizing dye has an ionization potential smaller than that of the silver halide being sensitized. This allows the energy of the dye excited state to lie at or above the silver halide conduction band even though the energy of the photon absorbed by the dye is less than the AgX band gap. Dye ionization potentials are generally well correlated with their electrochemical oxidation potentials.

The triplet energy level of a typical cyanine dye lies about 0.4 to 0.5 eV (about 35 to 50 kJ mol⁻¹) below the radiationally excited singlet level.⁶² In most cases, this triplet state seems to have little effect on spectral sensitization at room temperature, although it may influence sensitization in special cases.⁶³

Electron transfer takes place by tunneling from the excited-state dye to the silver halide conduction band. Simple calculations verify that penetration of the potential barrier will compete favorably with de-excitation of the dye by fluorescence emission (p. 253 of Ref. 3). The interaction time between an excited spectral sensitizer and silver bromide appears to occur between 10⁻¹³ and 10⁻¹⁰ (p. 237 of Ref. 3). Factors favoring irreversible flow of electrons from dye to silver halide are delocalization within the conduction band, the negative space charge layer on the surface of the crystal, trapping by remote sites on the silver halide surface, and irreversible trapping to form a latent image.

Free electrons from dye sensitization appear in the silver halide conduction band having the same mobility and lifetime as those formed by intrinsic absorption. These electrons participate in latent image formation by the usual mechanism.

After electron transfer, the oxidized dye radical cation becomes the hole left behind. Because there is evidence that a single dye molecule may function repeatedly,⁶⁴ the dye cation "hole" must be reduced again. This can occur by electron tunneling from an occupied site on the crystal's surface, whose energy level is favorable for that process. A bromide ion at a kink is one such site. This leaves a trapped hole on the crystal surface that may be thermally liberated to migrate through the crystal. This hole can lead to conduction band photoelectron loss by recombination. A supersensitizer molecule (discussed later) may also trap the hole.

The formation of the J-aggregate exciton band has significant effects on the excited-state dynamics of the dye. There is the possibility of coherent delocalization of the exciton over several molecules in the aggregate.⁶⁵ The exciton is mobile and can sample over 100 molecules within its lifetime.⁶⁶ This mobility means that the exciton is sensitive to quenching by traps within the aggregate structure. The overall yield of sensitization generally increases to maximum, then decreases somewhat as the aggregate size increases. Optimum sensitizations usually occur below surface monolayer coverage.

Sometimes the spectral sensitivity of a dye is increased by addition of a second substance. If the added sensitivity exceeds the sum of both sensitizers individually, the increase is super-additive and the second substance is called a *supersensitizer*. Maximum efficiency of a supersensitizer often occurs in the ratio of about 1:20 where the supersensitizer is the dilute component.

The supersensitizer may increase the spectral absorption of the sensitizer by inducing or intensifying formation of a J-aggregate. In some cases these changes are caused by a mutual increase in adsorption to the grain surface, as when the sensitizer and supersensitizer are ions of opposite charge.⁶⁷ However, the most important supersensitizers appear to increase the fundamental efficiency of spectral sensitization as measured by the relative quantum yield.

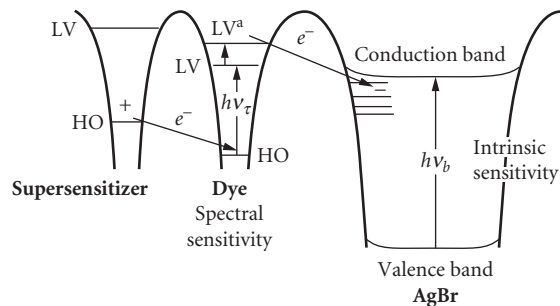


FIGURE 7 Simplified energy level diagram for a spectral sensitizer dye plus a supersensitizer dye adsorbed to a silver bromide surface.

The low concentration of supersensitizer relative to sensitizer suggests trapping of an exciton moving through the quasicrystalline array of J-aggregated sensitizer molecules. One hypothesis is that the supersensitizer molecule traps the exciton, providing a site for facile electron transfer into the silver halide.⁶⁸

Gilman and coworkers^{69–72} proposed that the supersensitization of J-aggregating dyes takes place via hole trapping by the supersensitizer molecules. The exciton moving through the aggregate may be self-trapped at a site adjacent to a supersensitizer molecule. By hypothesis, the highest occupied (HO) electron energy level of the supersensitizer is higher than that of the partially empty HO level of the excited sensitizer molecule. An electron transfers from the supersensitizer to the sensitizer molecule. The sensitizer is reduced to an electron-rich, anion-free radical while the supersensitizer is oxidized to an electron-deficient cation radical of lower energy than the original hole, thereby localizing the hole on the supersensitizer (see Fig. 7).

Following electron transfer from the supersensitizer, the electron-rich free radical anion left on the sensitizer will possess an occupied electron level of higher energy than the excited level of the parent sensitizer. The reduction process thereby raises the level of the excited electron with respect to the conduction band of AgX, with a consequent increase in probability of electron tunneling into the silver halide conduction band.

The various proposed mechanisms of supersensitization are not mutually exclusive. Exciton trapping and hole trapping may operate together. Underlying all mechanisms are the roles of exciton migration throughout the aggregate, its interruption at or near the supersensitizer site, and a more facile transfer of the electron into the AgX conduction band at that localized state.

Spectral sensitizing dyes, especially at high surface coverage, desensitize silver halides in competition with their sensitization ability. Red spectral sensitizers generally desensitize more than green or blue spectral sensitizers do. Desensitization can be thought of as reversible sensitization. For example, a mobile conduction band electron can be trapped by a dye molecule to form a dye radical anion⁷³ or by reduction of a hole trapped on the dye.^{63,74–76} Under normal conditions of film use (not in a vacuum), a dye radical anion may transfer the excess electron irreversibly to an O₂ molecule with formation of an O₂⁻ anion. Either postulate leads to irreversible loss of the photoelectron and consequent inefficiency in latent image formation.

Color Science of Photographic Spectral Sensitizers

Human vision responds to light in the range of 400 to 700 nm. The human eye is known to contain two types of light-sensitive cells termed *rods* and *cones*, so named because of their approximate shapes. Cones are the sensors for color vision. There is strong evidence for three types of cones sometimes termed long (L), middle (M), and short (S) wavelength receptors (Chap. 1 in Ref. 1). The normalized spectral responses of the human eye receptors (Chap. 26, Table 5 in Ref. 1) are

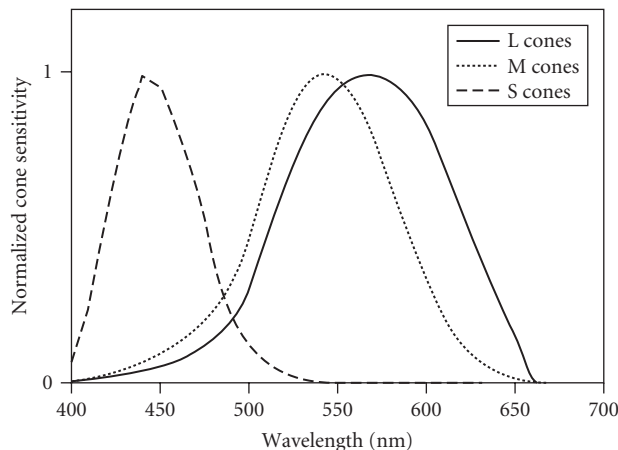


FIGURE 8 Normalized spectral sensitivity of the human eye.

shown in Fig. 8. The ability of humans to distinguish differences in spectral wavelength is believed to rely on difference signals created in the retina and optic nerve and interpreted by the brain from responses to light by the three cone detectors (Chap. 9 in Ref. 4).

The human eye's red sensitivity function peaks at around 570 nm and possesses little sensitivity at wavelengths longer than 650 nm. Also, there exists considerable overlap between red and green sensitivity functions, and to a lesser extent overlap between blue and green and blue and red sensitivity functions. This overlap in sensitivity functions, combined with signal processing in the retina, optic nerve, and brain, allows us to distinguish subtle color differences while simultaneously appreciating rich color saturation.

Photographic films cannot do this. The complex signal processing needed to interpret closely overlapping spectral sensitivity functions is not possible in a photographic material (Chap. 9 in Ref. 4). Some relief from this constraint might be anticipated by considering Maxwell's principle.⁷⁷

Maxwell's principle states that any visual color can be reproduced by an appropriate combination of three independent color stimuli called *primaries*. The amount of each primary per wavelength needed to reproduce all spectral colors defines three color-matching functions (Chap. 6 in Ref. 1) for those primaries. If there exists a set of color-matching functions that minimize overlap between spectral regions, a photographic film with spectral sensitivity like that set of color-matching functions and imaging with the corresponding primary dyes would faithfully reproduce colors.

The cyan, magenta, and yellow dye primaries used in photographic materials lead to theoretical color-matching functions (Chap. 19-II in Ref. 3; Ref. 78)* having more separation than human visual sensitivity, but possessing negative lobes as shown in Fig. 9. This is impossible to achieve in practical photographic films (Chap. 9 in Ref. 4), although the negative feedback from HE effects provides an imperfect approximation. Approximate color-matching functions have been proposed⁸⁰ that contain no negative lobes. But these possess a very short red sensitivity having a high degree of overlap with the green sensitivity, similar to that of the human eye.

It is therefore not possible to build a perfect photographic material possessing the simultaneous ability to distinguish subtle color differences, especially in the blue-green and orange spectral regions, and to render high-saturation colors. Compromises are necessary.

*The theoretical color-matching functions were calculated using the analysis described in Ref. 78, Eq. (9a-c) for block dyes having trichromatic coefficients in Table II, and using the Judd-Vos modified XYZ color matching functions in Chap. 26, Table 2 of Ref. 1.

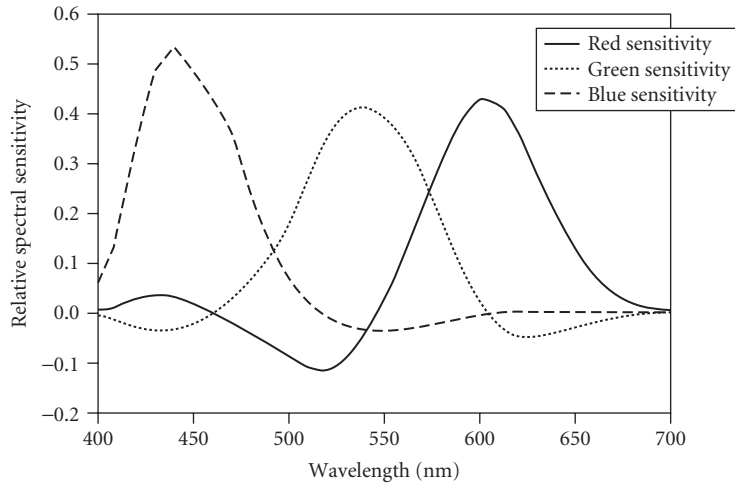


FIGURE 9 Spectral sensitivity for idealized photographic dyes. Note the large negative lobe around 525 nm.

Most photographic materials intended to image colors at high saturation possess spectral sensitivities having small overlap, and in particular possess a red spectral sensitivity centered at considerably longer wavelengths than called for by theoretical color-matching functions. This allows clean spectral detection and rendering among saturated colors but sacrifices color accuracy.

For example, certain *anomalous reflection colors* such as the blue of the lobelia flower possess a sharp increase in spectral reflectance in the long red region not seen by the human eye but detected strongly by a photographic film's spectral sensitivity (see Fig. 10). This produces a serious color

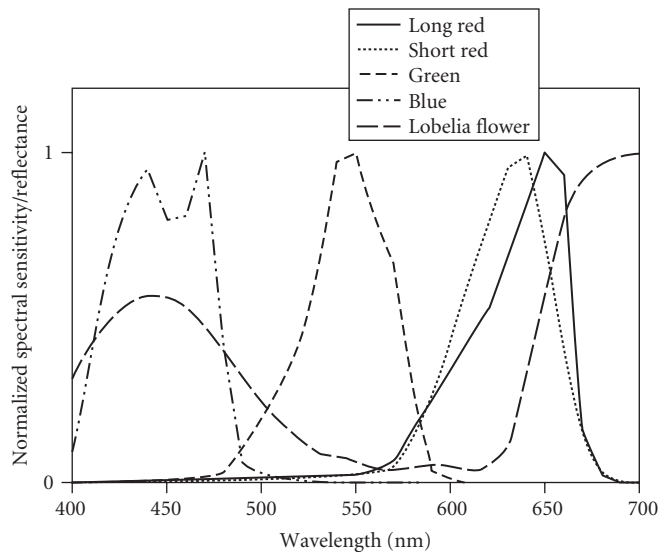


FIGURE 10 Spectral sensitivity for a real photographic film plus reflectance of the lobelia flower.

accuracy error (the flower photographs pink, not blue).⁸⁰ Among other such colors are certain green fabric dyes that image brown or gray. To correct for these hue errors, a special color reversal film was made whose general characteristic involves a short red spectral sensitivity.* This film does a considerably better job of eliminating the most serious hue errors, but unfortunately does not possess high color saturation. It is important in catalogue and other industries where accurate color reproduction is of paramount importance, but most consumers prefer saturated colors.

More recently, color negative films possessing short red spectral sensitivities combined with powerful development inhibitor anchimeric releasing (DIAR) color correction chemistry and masking coupler technology have produced films[†] simultaneously having high color saturation plus accurate color reproduction, at least toward anomalous reflection colors. In addition, these films render improved color accuracy under mixed daylight plus (red-rich) incandescent lighting conditions. These films maintain the high color saturation capability important to consumers but also provide more accurate color hues—although they are not perfect.

A new type of color negative film has been described^{§1} that possesses a fourth color record whose function approximates the large negative lobe in the theoretical red color-matching function of photographic dyes. Several embodiments have been described and are employed in modern films. The most practical utilizes a fourth color record containing silver halide emulsions sensitive to short green light plus a magenta DIR color correction coupler.[‡]

In simple terms, this special color record responds to short green light in the negative lobe region by releasing a development inhibitor that travels into the red-sensitive cyan color record, thereby reducing cyan dye in response to short green light, an effect opposite to that of normal red light exposure. This film's advantages reside in its ability to more accurately record colors in the blue-green spectral region, plus improved color accuracy under certain fluorescent lights that have a pronounced blue-green emission component.

Modern films are moving toward shorter-wavelength red spectral sensitivity. This allows improved color accuracy by shifting the red spectral sensitivity closer to that of the human eye. In addition, films have been designed to crudely approximate the principal negative lobe in the theoretical red color-matching function for photographic dye primaries, thereby producing a more accurate rendition of certain blue-green colors. Although not perfect, these techniques allow modern photographic films to record the world of color in all its richness and subtleties nearly as we ourselves see it.

30.6 GENERAL CHARACTERISTICS OF PHOTOGRAPHIC FILMS

Today there is available a large variety of films including different film speeds, types, and sizes. Many films are offered in both professional and consumer categories. This section describes some general characteristics of films as they relate to film speed and type. Some unique characteristics of professional films are contrasted to consumer films. Finally, color reversal, black-and-white, and color negative films are compared in more detail.

Low-Speed Films Compared to High-Speed Films

Higher granularity, higher sensitivity to ambient background radiation, and higher sensitivity to x rays are fundamental penalties that come with higher photographic speed. High-speed films often suffer additional penalties compared to their lower-speed counterparts, but these additional penalties

*Kodak Ektachrome 100 Film, EPN—a color reversal professional film.

†Kodak Gold 100, Kodak Gold 200, and Kodak Max 400 color negative films.

‡Fuji Realia 100 Film; Fuji Superia 200 and 400 film; Fuji Nexia 200 and 400 film; and Fuji 160 NPL, NPS, and NC professional films. These are all color negative films.

are not intrinsic. They represent conscious choices made by the manufacturer nearly always traceable to the desire to minimize granularity.

As one example, 400-speed color slide film generally has lower color saturation than its 100-speed counterpart. The reduced color saturation comes as a result of lower interlayer interimage effects (IIE) caused by higher iodide content in the mixed silver bromo-iodide crystals, which, in turn, is needed to achieve higher speeds from a smaller sized crystal and is ultimately done to reduce the photographic granularity penalty.

Photographic Speed and Photographic Granularity To understand the intrinsic connection between photographic speed and photographic granularity, recall that light quanta striking a flat surface follow Poisson statistics (Chap. 1 of Ref. 27). Therefore, a uniform light exposure given to a photographic film delivers photon hits Poisson-distributed across the plane of the film. Silver halide crystals are on the order of $1\ \mu\text{m}^2$ in projected area. On average in a film exposure (the product of the intensity of the radiation and the duration) appropriate for an EI100-speed film, the most sensitive silver halide crystals receive:*

E_q (red light; 650 nm)	~48 photons per square micrometer
E_q (green light; 550 nm)	~40 photons per square micrometer
E_q (blue light; 450 nm)	~33 photons per square micrometer

Because the most efficient practical photographic emulsions to date have a quantum sensitivity of about five photons per grain,²⁹ these approximate numbers suggest there is enough light (on average) striking a film during an EI100-speed film exposure to form latent image on efficient $1\text{-}\mu\text{m}^2$ silver halide grains. But an EI800-speed exposure at three stops down begins to push the limits of modern emulsion efficiency at five green photons per square micron.

Silver halide crystals within the film are light collectors. Just as large-area mirrors in telescopes collect more photons and are more sensitive detectors of starlight than smaller mirrors, all other things being equal, larger silver halide crystals likewise collect more photons to become more sensitive light detectors.

Strong correlates exist between silver halide surface area and photographic speed. Photographic speed often linearly follows the average surface area of the silver halide crystal population (p. 969 of Ref. 2): speed $\sim \bar{a}$.

According to the well-known Siedentopf formula (p. 60 in Ref. 27), the root-mean-square photographic granularity at density D is given by $G \sim D\hat{a}$. In this formula, the symbol \hat{a} represents the average cross-sectional area of the image particle projected onto the film plane. To a good approximation, this is proportional (not necessarily linear) to the mass of developed silver regardless of whether the image particle is a dye cloud or silver deposit.

These considerations demonstrate the intrinsic and well-known correlation between photographic speed and photographic granularity. For a given crystal morphology, as the size of the silver halide crystal increases, both surface area and mass increase, leading to a speed plus granularity increase.

At a fundamental level, to lower the granularity of a higher-speed photographic film the sensitivity to light of the silver halide crystals must improve. Then smaller silver halide crystals can be used to attain speed. Factors that contribute to this problem were discussed in Sec. 30.2.

In general, larger silver halide crystals suffer more inefficiencies in the latent-image-forming process because latent-image dispersity, electron-hole recombination, and other irreversible photon waste processes become more problematic as the grain size increases (p. 993 in Ref. 2). Overcoming these inefficiencies at all speeds, but especially at high photographic speeds, is a major challenge of modern photographic science.

Photographic Speed and Sensitivity to High-Energy Radiation High-energy subatomic particles induce developable latent image formation in photographic films. The process by which high-energy

*Analysis based on a typical EI100-speed color reversal film where the exposure falls near the photographically fastest, most sensitive part of the sensitometric curve. This analysis is described on p. 47 of Ref. 27.

charged particles produce latent image usually involves energy transfer to valence band electrons by inelastic scattering. Some of these scattered electrons are promoted into the silver halide conduction band and go on to produce latent image by the usual processes. Higher-speed films are more sensitive to this radiation because their larger silver halide crystals are hit (on average) more frequently and more collisions occur by passage through a larger crystal than through a smaller crystal.

Photographic films are slowly exposed by ambient background radiation (Chap. 23 in Ref. 3) caused by cosmic rays; ubiquitous stone and masonry building materials, which contain trace amounts of radioactive elements; and other low-level radiation sources. These radiation sources emit various types of charged and uncharged high-energy particles and photons including alpha particles, neutrons, γ rays, β particles, muons, and others. This results in a shorter ambient shelf life for very-high-speed films unless special storage precautions are taken.

Photographic emulsions are also sensitive to x rays. The formation of latent image silver by x rays (and γ rays) is due to the action of high-energy electrons emitted during the absorption or inelastic scattering of the electromagnetic radiation by matter (Chap. 23 in Ref. 3). These electrons may enter the conduction band of silver halide and proceed to latent image formation in the usual way. If the primary emitted electron is of sufficiently high energy, it may create secondary electrons of lower energy by inelastic scattering with other electrons. These secondary electrons enter the conduction band.

High-speed films have more sensitivity to x rays than do low-speed films. Some airline travelers ask airport security agents to hand-check their film rather than allowing it to pass through the baggage x-ray scanner. This precaution is prudent for higher-speed film because airport x-ray machines have become stronger emitters recently, especially at airports with special security concerns.

Despite the higher price, inherent penalties, and willful compromises that characterize high-speed films, they offer access to more picture space under low-light conditions. Among the advantages a high-speed film offers are the following:

1. Pictures under low light where flash is not permitted, such as theaters and museums.
2. Pictures under low light using slower lenses such as zoom and telephoto lenses.
3. Pictures under low light when the subject lies beyond the reach of flash.
4. Pictures under low light taken in rapid sequence without waiting for a flash to recharge.
5. Pictures under low light when flash disturbs the subject, as in some nature photography.
6. Pictures of people without flash. Flash is responsible for a major defect known as *red eye*, caused when light from a flash reflects off the retina of the eye.
7. Pictures using a stopped-down camera aperture for improved depth of field, so objects over a wider range of distance from the camera are imaged in focus.
8. Reduced dependence on flash for longer battery life.
9. Low-cost one-time-use cameras (OTUCs) when flash units are not needed.
10. Ability to freeze action using a rapid shutter speed.

For these and other reasons, a high-speed film may be the best choice for many film applications despite its shortcomings compared to lower-speed alternatives.

Professional Film Compared to Amateur Film

Most film manufacturers offer film in a professional category and an amateur (sometimes called *consumer*) category. These differ in many respects.

Consumer film is targeted for the needs of the amateur snap shooter. Drugstores, supermarkets, department stores, and other outlets frequented by nonprofessional photographers sell this film. Because most amateurs are price conscious, discount brands of film are available having some performance compromises that may not be important to some consumers.

To keep the price low, manufacturing tolerances for consumer film performance attributes are wider than those for professional film. In a well-manufactured film these departures from performance

aims are small enough to go unnoticed by most average consumers but large enough to be objectionable to highly discerning professionals.

People want pleasing color in their pictures that renders the main subject attractive to the eye. For most amateurs this means saturated colors enhanced to some extent beyond their true appearance in the original scene. However, this preference is highly selective since over-saturated skin tones are objectionable. Other *memory colors*—those that we recognize and know what they should be, such as green grass and blue sky—must also be rendered carefully to produce the most pleasing color appearance (Chap. 5 in Ref. 4). The best consumer films accommodate these color preferences.

Exposures on a roll of consumer film are sometimes stored for long times. It is not unusual for amateurs to take pictures spanning two successive Christmas seasons on the same roll of film. Thus, latent image stability under a wide range of temperature and humidity conditions is very important to a good consumer film. On the other hand professional photographers develop their film promptly, making this requirement far less demanding in a professional film.

Most amateurs are not skilled photographers. It is advantageous for consumer film to be easy to use and provide good-looking pictures under a wide variety of conditions. Consumer film incorporates more tolerance toward over- and underexposures.

Film for amateurs is offered in three standard formats. The most common format is the 35-mm cassette. Consumer cameras that accept this film have become compact, lightweight, and fully automatic. The camera senses the film speed through the DX coding on the film cassette and adjusts its shutter, aperture, and flash accordingly. Autofocus, autorewind, and automatic film advance are standard on these cameras, making the picture-taking experience essentially point and shoot.

The new Advanced Photo System (APS) takes ease of use several steps further by incorporating simple drop-in film loading and a choice of three popular print aspect ratios from panoramic to standard size prints. A magnetic coating on the film records digital data whose use and advantages are just beginning to emerge.

The third format is 110 film. Although not as popular as 35-mm and APS, it is often found in inexpensive cameras for gift packages. Film for this format comes in a cartridge for easy drop-in loading into the camera.

Professional photographers demand more performance from their film than do typical amateurs. These people shoot pictures for a living. Their competitive advantage lies in their ability to combine technical skill and detailed knowledge of their film medium with an artistic eye for light, composition, and the color of their subject.

Film consistency is very important to professionals, and they pay a premium for very tight manufacturing tolerances on color balance, tone scale, exposure reciprocity, and other properties. They buy their film from professional dealers, catalogues, and in some cases directly from the manufacturer. Professionals may buy the same emulsion number in large quantity, then pretest the film and store it in a freezer to lock in its performance qualities. Their detailed knowledge about how the film responds to subtleties of light, color, filtration, and other characteristics contributes to their competitive advantage.

Professional film use ranges from controlled studio lighting conditions to harsh and variable field conditions. This leads to a variety of professional films designed for particular needs. Some photographers image scenes under incandescent lighting conditions in certain studio applications, movie sets, and news events where the lights are balanced for television cameras. These photographers need a tungsten-balanced film designed for proper color rendition when the light is rich in red component, but deficient in blue compared to standard daylight.

Catalogue photographers demand films that image colors as they appear to the eye in the original scene. Their images are destined for catalogues and advertisement media that must display the true colors of the fabric. This differs from the wants of most amateurs, who prefer enhanced color saturation.

Some professional photographers use high color saturation for artistic effect and visual impact. These photographers choose film with oversaturated color qualities. They may supplement this film by using controlled light, makeup on models, lens gels, or other filtration techniques to control the color in the image.

Professional portrait photographers need film that gives highly pleasing skin tones. Their films typically image at moderate to low color saturation to soften facial skin blotches and marks that most people possess but don't want to see in their portraits. In addition, portrait film features highly

controlled upper and lower tone scale to capture the subtle texture and detail in dark suits and white wedding dresses and the sparkle in jewelry.

Professional photographers image scenes on a variety of film sizes ranging from standard 35-mm size to 120-size, and larger 4×5 to 8×10 sheet-film formats that require special cameras. Unlike cameras made for amateurs, most professional cameras are fully controllable and adjustable by the user. They include motor drives to advance the film rapidly to capture rapid sequence events. Professional photographers often take many pictures to ensure that the best shot has been captured.

Many specialty films are made for professional use. For example, some films are sensitized to record electromagnetic radiation in the infrared, which is invisible to the human eye. Note that these films are not thermal detectors. Infrared films are generally *false colored*, with the infrared-sensitized record producing red color in the final image, the red-light-sensitized record producing green color, and the green/blue-sensitized record producing blue color. Other professional films are designed primarily for lab use, such as duplicating film used to copy images from an original capture film.

Some film applications are important to a professional but are of no consequence to most amateurs. For example, professionals may knowingly underexpose their film and request *push processing*, whereby the film is held in the developing solution longer than normal. This brings up an underexposed image by rendering its midtone densities in the normal range, although some loss in shadow detail is inevitable. This is often done with color transparency film exposed under low-ambient-light conditions in the field. Professional films for this application are designed to provide invariant color balance not only under normal process times but also under a variety of longer (and shorter) processing times.

Multipop is another professional technique whereby film is exposed to successive flashes or *pops*. Sometimes the photographer moves the camera between pops for artistic effect. The best professional films are designed to record a latent image that remains nearly invariant toward this type of exposure.

Consumer film is made for amateur photographers and suits their needs well. Professional film comes in a wide variety of types to fit the varied needs of highly discerning professional photographers who demand the highest quality, specialized features, and tightest manufacturing tolerances—and are willing to pay a premium price for it.

Color Reversal Film

Color reversal films, sometimes called *color transparency* or *color slide* films, are sold in professional and amateur categories. This film is popular among amateurs in Europe and a few other parts of the world but has largely been replaced by color negative film among amateurs in North America. However, it is still very popular among advanced amateurs, who favor it for the artistic control this one-stage photographic system offers them, the ease of proofing and editing a large number of images, and the ability for large-audience presentations of their photographic work. Because the image is viewed by transmitted light, high amounts of image dye in these films can project intense colors on the screen in a dark auditorium for a very striking color impact.

Color reversal film is the medium of choice for most professional commercial photographers. It offers them the ability to quickly proof and select their most favored images from a large array of pictures viewed simultaneously on a light table. The colors are more intense compared to a reflection print, providing maximum color impact to the client. And professional digital scanners have been optimized for color reversal images, making it convenient to digitally manipulate color reversal images and transcribe them to the printed page. Today's professional imaging chain infrastructure has been built around color reversal film, although digital camera image capture is making inroads, especially for low-quality and low-cost applications.

Large-format color reversal film, from 120 size to large sheet formats, is often shot by professionals seeking the highest-quality images destined for the printed page. A larger image needs little or no magnification in final use and avoids the deterioration in sharpness and grain that comes from enlarging smaller images. Additionally, art directors and other clients find it easier to proof larger images for content and artistic appeal.

EI400-speed color reversal films generally possess low color saturation, low acutance, high grain, and high contrast compared to lower-speed color reversal films. The high contrast can lead to some loss in highlight and shadow detail in the final image. Except for very low-light situations, EI100 is

the most popular speed for a color reversal film. It offers the ultrahigh image quality professionals and advanced amateurs demand.

Modern professional color reversal films maintain good color balance under push processing conditions. They also perform very well under a variety of exposure conditions including multipop exposures and very long exposure times.

When professional photographers believe nonstandard push or pull process times are needed, they may request a “clip test” (sometimes also called a “snip test”) whereby some images are cut from the roll of film and processed. Adjustments to the remainder of the roll are based on these preprocessed images. Some photographers anticipate a clip test by deliberately shooting sacrificial images at the beginning of a roll.

Professional color reversal films can be segregated into newer modern films and older *legacy films*. Modern films contain the most advanced technology and are popular choices among all photographers. Legacy films are popular among many professional photographers who have come to know these films’ performance characteristics in detail after many years of experience using them.

Legacy films continue to enjoy significant sales because the huge body of personal technical and artistic knowledge about a film accumulated by a professional photographer over the years contributes to his or her competitive advantage. This inertia is a formidable barrier to change to a new, improved, but different film whose detailed characteristics must be learned.

In some cases, the improved features found in modern films are not important to a particular professional need. For example, a large 4 × 5 sheet film probably needs little or no enlargement, so the poorer grain and sharpness of a legacy film may be quite satisfactory for the intended application. The detailed knowledge the professional has about how the legacy film behaves in regard to light and color may outweigh the image structure benefits in a modern film. Also, all films possess a unique tone scale and color palette. A professional may favor some subtle and unique feature in the legacy film’s attributes.

Kodachrome is a special class of color reversal legacy films. Unlike all other types of color reversal films, this film contains no dye-forming incorporated couplers. The dyes are imbibed into the film during processing so the film’s layers are coated very thin. This results in outstanding image sharpness. Moreover, many photographers prefer Kodachrome’s color palette for some applications.

Kodachrome dyes are noted for their image permanence. Images on Kodachrome film have retained their colors over many decades of storage, making this film attractive to amateurs and professionals alike who want to preserve their pictures. However, modern incorporated coupler color reversal films also have considerably improved dark storage image stability compared to their predecessors.

Table 1 lists some color reversal films available today. This is not a comprehensive list. Most film manufacturers improve films over time, and portfolios change frequently.

TABLE 1 Color Reversal Capture Films

Modern Films	Legacy Films	Kodachrome Films	Tungsten Films
Kodak Ektachrome 100VS (very saturated color)	Fuji Velvia 50 RVP (very saturated color)	Kodachrome 25 PKM	Kodak Ektachrome 64T EPY
Kodak Ektachrome E100S (standard color, good skin tones)	Kodak Ektachrome 64 EPR (standard color)	Kodachrome 64 PKR	Fujichrome 64T RTP
Fuji Astia 100 RAP (standard color, good skin tones)	Kodak Ektachrome 100 Plus EPP (standard color)	Kodachrome 200 PKL	Agfachrome 64T
Fuji Provia 100 RDPII (standard color)	Kodak Ektachrome 100 EPN (accurate color)		Kodak Ektachrome 160T EPT
Agfachrome 100 (standard color)	Kodak Ektachrome 200 EPD (lower color saturation)		Kodak Ektachrome 320T EPJ
Kodak Ektachrome E200 (can push 3 stops)	Agfachrome 200		
	Kodak Ektachrome 400 EPL		
	Fuji Provia 400		
	Agfachrome 400		

Unless a photographer has some particular reason for choosing a legacy color reversal film, the modern films are generally a better overall choice. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Black-and-White Film

Black-and-white (B&W) films and papers are sold through professional product supply chains. Color films and papers have largely displaced B&W for amateur use, with the exception of advanced amateurs who occasionally shoot it for artistic expression.

The image on a B&W film or print is composed of black metallic silver. This image has archival stability under proper storage conditions and is quite stable even under uncontrolled storage. Many remarkable and historically important B&W images exist that date from 50 to over 100 years ago. This is remarkable considering the primitive state of the technology and haphazard storage conditions.

Pioneer nature photographer Ansel Adams reprinted many of his stored negatives long after he stopped capturing images in the field. This artistic genius recreated new expressions in prints from scenes captured many years earlier in his career because his negatives were well preserved.

Ansel Adams worked with B&W film and paper. Many photographic artists work in B&W when color would distract from the artistic expression they wish to convey. For example, many portraits are imaged in B&W.

B&W films pleasingly image an extended range of tones from bright light to deep shadow. They are panchromatically sensitized to render colors into grayscale tones. A B&W film's tone scale is among its most important characteristics.

Photographers manipulate the contrast of a B&W image by push- or pull-processing the negative (varying the length of time in the developer solution) or by printing the negative onto a select contrast paper. Paper contrast grades are indexed from 1 through 5, with the higher index giving higher contrast to the print. Multigrade paper is also available whereby the contrast in the paper is changed by light filtration at the enlarger.

It is generally best to capture the most desired contrast on the negative because highlight or shadow information lost on the negative cannot be recovered at any later printing stage. Paper grades are designed for pleasing artistic effect, not for highlight or shadow recovery.

Legacy films are very important products in B&W photography. Although typically a bit grainier and less sharp than modern B&W films, their tone scale characteristics, process robustness, and overall image rendition have stood the test of time and are especially favored by many photographers. Besides, the B&W photographic market is small so improvements to this line of films occur far less frequently than do improvements to color films.

Unlike the case with color films, which are designed for a single process, a photographer may choose among several developers for B&W films. Kodak Professional T-Max developer is a good choice for the relatively modern Kodak T-Max B&W film line. Kodak developer D-76 is also very popular and will render an image with slightly less grain and slightly less speed. It is a popular choice for legacy films. Kodak Microdol-X developer is formulated to give very low grain at the expense of noticeable speed loss. Kodak developer HC110 is popular for home darkrooms because of its low cost and ease of use, but it renders a grainier image compared to other developers. Other manufacturers, notably Ilford Ltd., carry similar types of B&W developers.

A new type of B&W film has emerged that is developed in a standard color negative process. This incorporated-coupler 400-speed film forms a black-and-white image from chemical dyes instead of metallic silver. Among its advantages are very fine grain at normal exposures in the commonly available color negative process. Its shortcomings are the inability to control contrast on the negative by push or pull processing and its objectionable high grain when underexposed compared to a comparably underexposed 400-speed legacy film. And many photographers prefer the tone scale characteristics found in standard B&W films.

Table 2 lists some B&W films available today. This is not a comprehensive list. Modern films generally feature improved grain and sharpness compared to legacy films. However, the legacy films are favorite choices among many photographers because of their forgiving process insensitivity, the

TABLE 2 Black and White Films

Modern Films	Legacy Films	Specialty Films
Kodak T-Max 100 Professional (fine grain, high sharpness)	Ilford Pan F Plus 50	Kodak Technical Pan 25 (fine grain and very high sharpness)
Ilford Delta 100 Pro (fine grain, high sharpness)	Kodak Plus-X 125	Kodak IR (infrared sensitive)
Kodak T-Max 400 Professional (finer grain compared to legacy films)	Ilford FP4 Plus 125	Ilford SFX 200 (extended long red sensitivity but not IR)
Ilford Delta 400 Pro (finer grain compared to legacy films)	Kodak Tri-X 400	Ilford Chromogenic 400 (incorporated couplers, color negative process)
Fuji Neopan 400 (finer grain compared to legacy films)	Ilford HP5 Plus 400	Kodak Professional T400 CN (incorporated couplers, color negative process)
Kodak T-Max P3200 Professional (push process to high speed, high grain)		
Ilford Delta P3200 Pro (push process to high speed, high grain)		

experience factor of knowing their detailed behavior under many conditions, and the characteristic look of their images. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Color Negative Film

Today's color negative film market is highly segmented. It can be most broadly divided into consumer films and professional films. The consumer line is further segmented into 35-mm film, 110 format film (a minor segment whose image size is 17×13 mm), single-use cameras, and the new Advanced Photo System (APS). Each segment serves the needs of amateur photographers in different ways. The basic films are common across segments, with the main difference being camera type.

Consumer color negative film is a highly competitive market with many film and camera manufacturers offering products. In general, films for this marketplace offer bright color saturation that is most pleasing to amateurs. The higher-speed films show progressively more grain than do the lower-speed films. The highest-speed films are less often used when enlargement becomes significant, such as in APS and 110 format, but are prevalent in single-use cameras.

Consumer 35-mm Films Films in 35-mm format form the core of all consumer color negative products. These same basic films are also found in single-use cameras, 110 format cameras, and APS cameras. Film speeds include 100, 200, 400, and 800. All manufacturers offer films at 100, 200, and 400 speed, but only a few offer films at speeds of 800 and above.

The 400-speed film is most popular for indoor shots under low light and flash conditions. Lower-speed films are most often used outdoors when light is plentiful.

The best consumer films feature technology for optimized color rendition including the high color saturation pleasing to most consumers, plus accurate spectral sensitization for realistic color rendition under mixed lighting conditions of daylight plus fluorescent and incandescent light, which is often present inside homes. Technology used for accurate spectral sensitization was briefly described for Fujicolor Superia and Kodak Gold films in Sec. 30.5.

Kodak consumer films segment into Kodak Gold and Kodak Royal Gold films. Kodak Gold consumer films emphasize colorfulness, while Kodak Royal Gold films emphasize fine grain and high sharpness.

Films made in 110 format generally fall into the 100- to 200-speed range because the enlargement factors needed to make standard-size prints place a premium on fine grain and high sharpness.

Single-Use Cameras The single-use camera marketplace is extremely competitive because of the high popularity of this film and camera system among amateurs. These low-cost units can be bought at all consumer outlets and are especially popular at amusement parks, theme parks, zoos, and similar family attractions. After the film is exposed, the photographer returns the entire film plus camera unit for processing. Prints and negatives are returned to the customer while the camera body is returned to the manufacturer, repaired as needed, reloaded with film, and repackaged for sale again. These camera units are recycled numerous times from all over the world.

Single-use cameras come in a variety of styles including a waterproof underwater system, low-cost models with no flash, regular flash models, an APS style offering different picture sizes, and a panoramic style. Single-use cameras typically contain 800-speed film, except for APS and panoramic models, which usually contain 400-speed film due to enlargement demands for finer grain.

Advanced Photo System (APS) APS is the newest entry into consumer picture taking. The size of the image on the negative is 29×17 mm, about 60 percent of the image size of standard 35-mm film (image size 36×24 mm). This puts a premium on fine grain and high sharpness for any film used in this system. Not all film manufacturers offer it. Speeds offered are 100, 200, and 400, with the higher speeds being the most popular.

Some manufacturers offer a B&W film in APS format. The film used is the incorporated-coupler 400-speed film that forms a black-and-white image from chemical dyes and is developed in a standard color negative process.

The APS camera system features simple drop-in cassette loading. Unlike the 35-mm cassette, the APS cassette contains no film leader to thread into the camera. When loaded, the camera advances the film out of the cassette and automatically spools it into the camera, making this operation mistake-proof. Three popular print sizes are available; panoramic, classic, and high-definition television (HDTV) formats. These print sizes differ in their width dimension.

A thin, nearly transparent magnetic layer is coated on the APS film's plastic support. Digital information recorded on this layer, including exposure format for proper print size plus exposure date and time, can be printed on the back of each print. Higher-end cameras offer prerecorded titles—for example, "Happy Birthday"—that can be selected from a menu and placed on the back of each print. Additional capabilities are beginning to emerge that take more advantage of this magnetic layer and the digital information it may contain.

The APS system offers easy mid-roll change on higher-end cameras. With this feature, the last exposed frame is magnetically indexed so the camera "remembers" its position on the film roll. A user may rewind an unfinished cassette to replace it with a different cassette (different film speed, for example). The original cassette may be reloaded into the camera at a later time. The camera will automatically advance the film to the next available unexposed frame.

Markings on the cassette indicate whether the roll is unexposed, partially exposed, fully exposed, or developed. There is no uncertainty about whether a cassette has been exposed or not. Unlike the 35-mm system, where the negatives are returned as cut film, the negatives in the APS system are rewound into the cassette and returned to the customer for compact storage. An index print is given to the customer with each processed roll so that each numbered frame in the cassette can be seen at a glance when reprints are needed.

Compact film scanners are available to digitize pictures directly from an APS cassette or 35-mm cut negatives. These digitized images can be uploaded into a computer for the full range of digital manipulations offered by modern photo software. Pictures can be sent over the Internet, or prints can be made on inkjet printers. For the best photo-quality prints, special photographic ink cartridge assemblies are available with most inkjet printers to print onto special high-gloss photographic-quality paper.

Table 3 summarizes in alphabetical order many brands of 35-mm consumer film available today worldwide. Some of these same basic films appear in single-use cameras, APS format, and 110 format, although not all manufacturers listed offer these formats. This is not a comprehensive list. Because this market is highly competitive, new films emerge quickly to replace existing films. It is not unusual for a manufacturer's entire line of consumer films to change within three years. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

TABLE 3 Consumer 35-mm Color Negative Films

Manufacturer and Brand Name					
Agfacolor HDC Plus	100	200	400		
Fujicolor Superia	100	200	400	800	
Ilford Colors	100	200	400		
Imation HP	100	200	400		
Kodak Gold	100	200	Max400	Max800	
Kodak Royal Gold	100	200	400		1000
Konica Color Centuria	100	200	400	800	
Polaroid One Film	100	200	400		

Professional Color Negative Film Professional color negative film divides roughly into portrait and wedding film and commercial and photojournalism film. Portrait and wedding films feature excellent skin tones plus good neutral whites, blacks, and grays. These films also incorporate accurate spectral sensitization technology for color accuracy and excellent results under mixed lighting conditions. The contrast in these films is about 10 percent lower than in most consumer color negative films for softer, more pleasing skin tones, and its tone scale captures highlight and shadow detail very well.

The most popular professional format is 120 size, although 35-mm and sheet sizes are also available. Kodak offers natural color (NC) and vivid color (VC) versions of Professional Portra film, with the vivid color having a 10 percent contrast increase for colorfulness and a sharper look. The natural color version is most pleasing for pictures having large areas of skin tones, as in head and shoulder portraits.

Commercial films emphasize low grain and high sharpness. These films offer color and contrast similar to those of consumer films. Film speeds of 400 and above are especially popular for photojournalist applications, and the most popular format is 35 mm. These films are often push processed.

Table 4 summarizes in alphabetical order many brands of professional color negative film available today worldwide. This is not a comprehensive list. Because this market is highly competitive, new films emerge quickly to replace existing films. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Silver halide photographic products have enjoyed steady progress during the past century. This chapter has described most of the technology that led to modern films. Most noteworthy among these advances are

1. Efficient light management in multilayer photographic materials has reduced harmful optical effects that caused sharpness loss and has optimized beneficial optical effects that lead to efficient light absorption.
2. Proprietary design of silver halide crystal morphology, internally structured halide types, transition metal dopants to manage the electron-hole pair, and improved chemical surface treatments has optimized the efficiency of latent image formation.

TABLE 4 Professional Color Negative Films

Manufacturer and Brand Name (Commercial Film)					
Agfacolor Optima II Prestige	100	200	400		
Fujicolor Press			400	800	
Kodak Professional Ektapress	PJ100		PJ400	PJ800	
Konica Impresa	100	200			3200 SRG
Manufacturer and Brand Name (Portrait Film)					
Agfacolor Portrait		160 XPS			
Fujicolor		160 NPS	400 NPH	800 NHGII	
Kodak Professional Portra		160 NC	400 NC		Pro 1000
		160 VC	400 VC		
Konica Color Professional		160			

3. Transition metal dopants and development-modifying chemistry have improved process robustness, push processing, and exposure time latitude.
4. New spectral sensitizers combined with powerful chemical color correction methods have led to accurate and pleasing color reproduction.
5. New image dyes have provided rich color saturation and vastly improved image permanence.
6. New film and camera systems for consumers have made the picture-taking experience easier and more reliable than ever before.

Photographic manufacturers continue to invest research and product development resources in their silver halide photographic products. Although digital electronic imaging options continue to emerge, consumers and professionals can expect a constant stream of improved silver halide photographic products for many years to come.

30.7 REFERENCES

1. M. Bass, E. Van Stryland, D. Williams, and W. Wolfe (eds.), *Handbook of Optics*, vol. 1, 2d ed., McGraw-Hill, New York, 1995.
2. J. Kapecki and J. Rodgers, *Kirk-Othmer Encyclopedia of Chemical Technology*, vol. 6, 4th ed., John Wiley & Sons, New York, 1993.
3. T. H. James (ed.), *The Theory of the Photographic Process*, 4th ed., Macmillan, New York, 1977.
4. R. W. G. Hunt, *The Reproduction of Colour*, Fountain Press, Tolworth, England, 1987.
5. E. Klein and H. J. Metz, *Photogr. Sci. Eng.* **5**:5 (1961).
6. R. E. Factor and D. R. Diehl, U. S. Patent 4,940,654, 1990.
7. R. E. Factor and D. R. Diehl, U. S. Patent 4,855,221, 1989.
8. G. Mie, *Ann. Physik.* **25**:337 (1908); M. Kerker, *The Scattering of Light and Other Electromagnetic Radiation*, Academic Press, New York, 1969.
9. D. H. Napper and R. H. Ottewill, *J. Photogr. Sci.* **11**:84 (1963); *J. Colloid Sci.* **18**:262 (1963); *Trans. Faraday Soc.* **60**:1466 (1964).
10. E. Pitts, *Proc. Phys. Soc. Lond.* **67B**:105 (1954).
11. J. M. Hammersley and D. C. Handscomb, *Monte Carlo Methods*, John Wiley & Sons, New York, 1964.
12. J. J. DePalma and J. Gasper, *Photogr. Sci. Eng.* **16**:181 (1972).
13. M. Motoki, S. Ichijima, N. Saito, T. Kamio, and K. Mihayashi, U. S. Patent 5,213,958, 1993.
14. Miles V. Klein, *Optics*, John Wiley & Sons, New York, 1970, pp. 582–585.
15. J. Gasper, Eastman Kodak Company, unpublished data.
16. N. F. Mott and R. W. Gurney, *Electronic Processes in Ionic Crystals*, Oxford University Press, London, 1940.
17. F. Seitz, *Rev. Mod. Phys.* **23**:328 (1951).
18. J. F. Hamilton, *Adv. Phys.* **37**:359 (1988).
19. A. P. Marchetti and R. S. Eachus, *Advances in Photochemistry*, vol. 17, John Wiley & Sons, New York, 1992, p. 145.
20. R. S. Eachus, A. P. Marchetti, and A. A. Muentner, *Ann. Rev. Phys. Chem.* **50**:117 (1999).
21. R. J. Deri, J. P. Spoonhower, and J. F. Hamilton, *J. Appl. Phys.* **57**:1968 (1985).
22. T. Tani, *Photographic Sensitivity*, Oxford University Press, Oxford, UK, 1995, p. 235.
23. L. M. Slifkin and S. K. Wonnell, *Solid State Ionics* **75**:101 (1995).
24. J. F. Hamilton and L. E. Brady, *Surf. Sci.* **23**:389 (1970).
25. R. C. Baetzold, Y. T. Tan, and P. W. Tasker, *Surf. Sci.* **195**:579 (1988).
26. R. S. Eachus and M. T. Olm, *Annu. Rep. Prog. Chem., Sect. C Phys. Chem.* **83**:3 (1986).
27. J. C. Dainty and R. Shaw, *Image Science*, Academic Press, London, 1974.

28. R. K. Hailstone, N. B. Liebert, M. Levy, R. T. McCleary, S. R. Girolmo, D. L. Jeanmaire, and C. R. Boda, *J. Imaging Sci.* **3**(3) (1988).
29. J. D. Baloga, "Factors in Modern Color Reversal Films," *IS&T 1998 PICS Conference*, p. 299 (1998).
30. R. J. Tuite, *J. Appl. Photogr. Eng.* **5**:200 (1979).
31. W. F. Smith, W. G. Herkstroeter, and K. I. Eddy, *Photogr. Sci. Eng.* **20**:140 (1976).
32. P. Douglas, *J. Photogr. Sci.* **36**:83 (1988).
33. P. Douglas, S. M. Townsend, P. J. Booth, B. Crystall, J. R. Durrant, and D. R. Klug, *J. Chem. Soc. Faraday Trans.* **87**:3479 (1991).
34. F. Wilkinson, D. R. Worrall, and R. S. Chittock, *Chem. Phys. Lett.* **174**:416 (1990).
35. F. Wilkinson, D. R. Worrall, D. McGarvy, A. Goodwin, and A. Langley, *J. Chem. Soc. Faraday Trans.* **89**:2385 (1993).
36. P. Douglas, S. M. Townsend, and R. Ratcliffe, *J. Imaging Sci.* **35**:211 (1991).
37. P. Egerton, J. Goddard, G. Hawkins, and T. Wear, *Royal Photographic Society Color Imaging Symposium*, Cambridge, UK, September 1986, p. 128.
38. K. Onodera, T. Nishijima, and M. Sasaki, *Proceedings of the International Symposium on the Stability and Conservation of Photographic Images*, Bangkok, Thailand, 1986.
39. Y. Kaneko, H. Kita, and H. Sato, *Proceedings of IS & T's 46th Annual Conference*, 1993, p. 299.
40. R. J. Berry, P. Douglas, M. S. Garley, D. Clarke, and C. J. Winscom, "Photophysics and Photochemistry of Azomethine Dyes," *IS & T 1998 PICS Conference*, p. 282 (1998).
41. W. F. Smith, W. G. Herkstroeter, and K. I. Eddy, *J. Am. Chem. Soc.* **97**:2164 (1975).
42. W. G. Herkstroeter, *J. Am. Chem. Soc.* **95**:8686 (1973).
43. W. G. Herkstroeter, *J. Am. Chem. Soc.* **97**:3090 (1975).
44. P. Douglas and D. Clarke, *J. Chem. Soc. Perkin Trans.* **2**:1363 (1991).
45. F. Abu-Hasanayn, *Book of Abstracts, 218th ACS National Meeting*, New Orleans, LA, August 22–26, 1999.
46. K. Hidetoshi et al. *International Congress of Photographic Science meeting*, Belgium, 1998.
47. N. Saito and S. Ichijima, *International Symposium on Silver Halide Imaging*, 1997.
48. V. Balzani, F. Bolletta, and F. Scandola, *J. Am. Chem. Soc.* **102**:2152 (1980).
49. R. Jain and W. R. Schleigh, U. S. Patent 5,561,037, 1996.
50. O. Takahashi, H. Yoneyama, K. Aoki, and K. Furuya, "The Effect of Polymeric Addenda on Dark Fading Stability of Cyan Indoaniline Dye," *IS & T 1998 PICS Conference*, p. 329 (1998).
51. R. L. Heidke, L. H. Feldman, and C. C. Bard, *J. Imag. Tech.* **11**(3):93 (1985).
52. S. Cowan and S. Krishnamurthy, U. S. Patent 5,378,587 (1995).
53. T. Kawagishi, M. Motoki, and T. Nakamine, U. S. Patent 5,605,788 (1997).
54. H. W. Vogel, *Berichte* **6**:1302 (1873).
55. J. E. Maskasky, *Langmuir* **7**:407 (1991).
56. J. Spence and B. H. Carroll, *J. Phys. Colloid Chem.* **52**:1090 (1948).
57. A. H. Hertz, R. Danner, and G. Janusonis, *Adv. Colloid Interface Sci.* **8**:237 (1977).
58. M. Kawasaki and H. Ishii, *J. Imaging Sci. Technol.* **39**:210 (1995).
59. G. Janssens, J. Gerritsen, H. van Kempen, P. Callant, G. Deroover, and D. Vandenbroucke, *The Structure of H-, J-, and Herringbone Aggregates of Cyanine Dyes on AgBr(111) Surfaces*. Presented at ICPS 98 International Conference on Imaging Science, Antwerp, Belgium (1998).
60. P. B. Gilman, *Photogr. Sci. Eng.* **18**:475 (1974).
61. J. Lenhard, *J. Imaging Sci.* **30**:27 (1986).
62. W. West, "Scientific Photography," in *Proceedings of the International Conference at Liege, 1959*, H. Sauvenier (ed), Pergamon Press, New York, 1962, p. 557.
63. T. Tani, *Photogr. Sci. Eng.* **14**:237 (1970).
64. J. Eggert, W. Meidinger, and H. Arens, *Helv. Chim. Acta.* **31**:1163 (1948).
65. J. M. Lanzafame, A. A. Muentner, and D. V. Brumbaugh, *Chem. Phys.* **210**:79 (1996).

66. A. A. Muentner and W. Cooper, *Photogr. Sci. Eng.* **20**:121 (1976).
67. W. West, B. H. Carroll, and D. H. Whitcomb, *J. Phys. Chem.* **56**:1054 (1952).
68. R. Brunner, A. E. Oberth, G. Pick, and G. Scheibe, *Z. Elektrochem.* **62**:146 (1958).
69. P. B. Gilman, *Photogr. Sci. Eng.* **11**:222 (1967).
70. P. B. Gilman, *Photogr. Sci. Eng.* **12**:230 (1968).
71. J. E. Jones and P. B. Gilman, *Photogr. Sci. Eng.* **17**:367 (1973).
72. P. B. Gilman and T. D. Koszelak, *J. Photogr. Sci.* **21**:53 (1973).
73. H. Sakai and S. Baba, *Bull. Soc. Sci. Photogr. Jpn.* **17**:12 (1967).
74. B. H. Carroll, *Photogr. Sci. Eng.* **5**:65 (1961).
75. T. Tani, *Photogr. Sci. Eng.* **15**:384 (1971).
76. T. A. Babcock, P. M. Ferguson, W. C. Lewis, and T. H. James, *Photogr. Sci. Eng.* **19**:49 (1975).
77. E. J. Wall, *History of Three Color Photography*, American Photographic Publishing Company, Boston, MA, 1925.
78. A. C. Hardy and F. L. Wurzburg Jr., "The Theory of Three Color Reproduction," *J. Opt. Soc. Am.* **27**:227 (1937).
79. M. L. Pearson and J. A. C. Yule, *J. Color Appearance* **2**:30 (1973).
80. S. G. Link, "Short Red Spectral Sensitizations for Color Negative Films," *IS & T 1998 PICS Conference*, p. 308 (1998).
81. Y. Nozawa and N. Sasaki, U. S. Patent 4,663,271 (1987).

IMAGE TUBE INTENSIFIED ELECTRONIC IMAGING

C. Bruce Johnson

*Johnson Scientific Group Inc.
Phoenix, Arizona*

Larry D. Owen

*NuOptics International
Phoenix, Arizona*

31.1 GLOSSARY

B_s	phosphor screen brightness, photometric units
CCDs	charge-coupled devices
CIDs	charge-injection devices
E_i	image plane illuminance, lux
E_s	scene illuminance, lux
e	electronic charge, coulombs
FO	fiberoptic
FOV	field-of-view, degrees
fc	illuminance, photometric, foot candles = lm/ft^2
f_N	spatial Nyquist frequency, cycle/mm
f_{fto}	limiting resolution at fiberoptic taper output
F_{si}	input window signal flux
ftL	luminance, photometric (brightness), foot Lamberts = lm/ft^2
G_m	VMCP electron gain, e/e
HVPS	high-voltage power supply
II	image intensifier
LLL	low-light-level
lx	illuminance, photometric, lux = lm/m^2
M_{fot}	magnification of fiberoptic taper
MCP	microchannel plate
MTF	modulation transfer function, 0 to 1.0
N_{essa}	number of stored SSA electrons per input photoelectron, e/photon
N_f	total number of frames, #
N_p	number of photoelectrons, #

$N_{ps}(\lambda)$	number of photons per second, photon/s
PDA	photodiode arrays
P	phosphor screen efficiency, photon/eV
$P_p(\lambda)$	radiometric power spectral distribution, W
QLI	quantum limited imaging
Q_{ssa}	stored SSA charge per input photoelectron from the photocathode, C
R_s	scene reflectance, ratio
R_{sn}	signal-to-noise ratio, ratio
$S(\lambda)$	absolute spectral sensitivity, mA/W
$S(f)$	squarewave response versus frequency, cycles/mm
SIT	silicon-intensifier-target vidicon
SNR	signal-to-noise ratio
sb	luminance, photometric (brightness), stilbs = cd/cm ²
SSA	silicon self-scanned array
T_f	filter transmission, 0 to 1.0
T_{fot}	transmission of fiberoptic taper, 0 to 1.0
T_n	lens T-number = $FN/\sqrt{\tau_0}$
T_{ssa}	transmission of fiberoptic window on the SSA, 0 to 1.0
V_a	phosphor screen, actual applied voltage, V
V_d	phosphor screen, "dead-voltage," V
V_m	VMCP applied potential, V
V_s	MCP-to-screen applied potential, V
$Y(\lambda)$	quantum yield (electrons/photon), percent
Y_k	quantum yield, photoelectrons/photon
Y_{ssa}	SSA quantum yield, e/photon
τ_e	the exposure period, s
τ_i	CCD charge integration period, s
τ_o	lens transmission, 0 to 1.0
Φ_p	photon flux density, photon/m ² /s

31.2 INTRODUCTION

It is appropriate to begin our discussion of image tube intensified (II) electronic imaging with a brief review of natural illumination levels. Figure 1 illustrates several features of natural illumination in the range from full sunlight to overcast night sky conditions. Various radiometric and photometric illuminance scales are shown in this figure. Present silicon self-scanned array (SSA) TV cameras, having frame rates of 1/30 to 1/25 s, operate down to about 0.5 lx minimum illumination.

The generic term *self-scanned array* is used here to denote any one of several types of silicon solid-state sensors available today which are designed for optical input. Among these are charge-coupled devices (CCDs), charge-injection devices (CIDs), and photodiode arrays (PDAs). Vol. II, Chaps. 32, "Visible Array Detectors," and 33, "Infrared Detector Arrays," contain detailed information on these types of optical imaging detectors. Specially designed low-light-level (LLL) TV cameras making use of some type of image intensifier must be used for lower exposures, i.e., lower illumination and/or shorter exposures.

The fundamental reason for using an II SSA camera instead of a conventional SSA camera is that low-exposure applications require the low-noise optical image amplification provided by an II

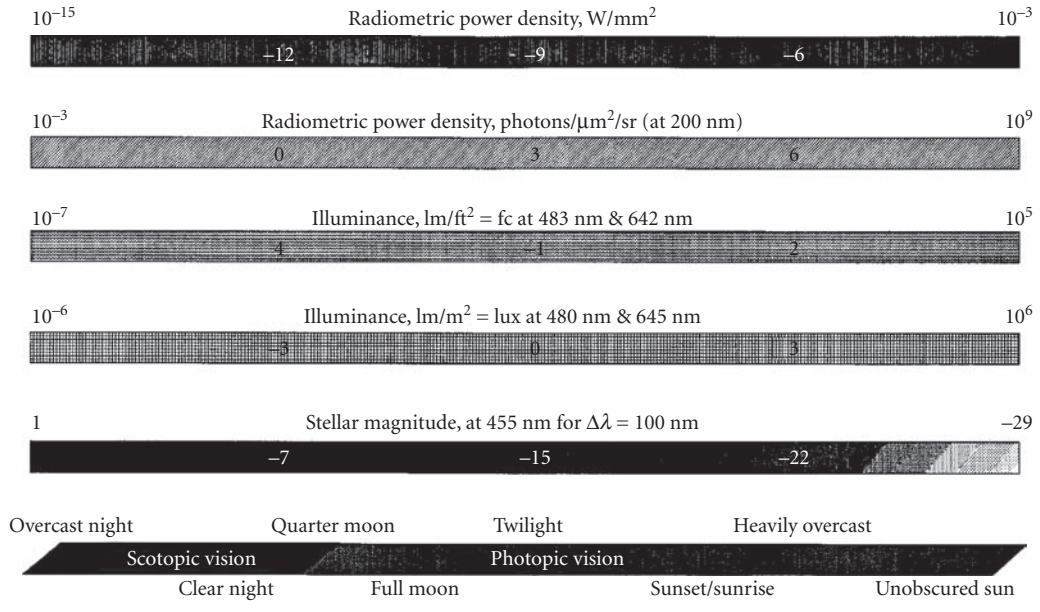


FIGURE 1 Various optical illumination ranges.

to produce a good signal-to-noise ratio from the SSA camera. Other important applications arise because of the ability to electronically shutter IIs as fast as 1 ns or less and the higher sensitivity of IIs in certain spectral regions. The following sections deal with the optical interface between the object and the II SSA, microchannel plate proximity-focused IIs, and II SSA detector assemblies. By using auto-iris lenses and controlling both the electronic gain and gating conditions of the II, II SSA cameras can provide an interscene dynamic range covering the full range of twelve orders of magnitude shown in Fig. 1. Several applications for II SSAs are discussed later in the chapter under “Applications.”

31.3 THE OPTICAL INTERFACE

It is necessary to begin our analysis of II SSA cameras with a brief discussion of the various ways to quantify optical input and exposure. Two fundamental systems are used to specify input illumination: radiometric and photometric. These systems are briefly described, and the fundamentals of optical image transfer are discussed. Detailed aspects of radiometry, photometry, and optical image transfer are discussed in Vol. II, Chap. 34, “Radiometry and Photometry” and Vol. I, Chap. 4, “Transfer Function Techniques.” However, enough information is presented in this chapter to allow the reader to properly design, analyze, and apply II SSA imaging technology for a wide variety of practical applications.

Quantum Limited Imaging Conditions

Quantum limited imaging (QLI) conditions exist in a wide variety of applications. An obvious one is that of LLL TV imaging at standard frame rates, i.e., 33-ms exposure periods, under nighttime illumination conditions. For example, under full moonlight input faceplate illumination conditions,

only ~ 1000 photons enter a $10 \times 10 \mu\text{m}^2$ image pixel in a 33-ms frame period. Assuming a quantum yield of 10 percent, an average of only 100 electrons is generated, and the maximum SNR achievable in each pixel and each frame is only $\sqrt{100} = 10$. Alternatively, under full unobscured sunlight input faceplate illumination conditions, an electronically gated camera with gatewidth limited exposure period of 10 ns produces a total of $(1\text{E}9 \text{ photons}/\mu\text{m}^2/\text{s})(10 \lambda \times 10 \mu\text{m}^2)(10 \text{ ns}) = 1000$ photons, or the same SNR as for the LLL operating conditions noted above. These are both clearly QLI operating conditions. II SSA camera technology is used to obtain useful performance in both of these types of applications. Without the use of an II, a bare SSA does not meet the requirements for useful SNR under these conditions.

Radiometry

The unit of light flux in the radiometric system is the watt. The watt can be used anywhere in the optical spectrum to give the number of photons per second (N_{ps}) as a function of wavelength (λ). Since the photon energy $E_p(\lambda)$ is

$$E_p = \frac{hc}{\lambda} \quad (1)$$

where h is Planck's constant and c is the velocity of light in vacuum, the radiometric power $P_p(\lambda)$, in watts, is given by

$$P_p(\lambda) = \left(\frac{hc}{\lambda}\right) \cdot N_{\text{ps}}(\lambda) \quad (2)$$

or

$$P_p(\lambda) = (2 \cdot 10^{-25}) \cdot \frac{N_{\text{ps}}(\lambda)}{\lambda} \quad (3)$$

where N_{ps} is the number of photons per second. Alternatively, the photon rate is given by

$$N_{\text{ps}}(\lambda) = (5 \times 10^{24}) \lambda P_p(\lambda) \quad \text{photons/s} \quad (4)$$

For example, one milliwatt of 633-nm radiation from an He-Ne laser is equivalent to $(5\text{E}24)(633\text{E}-9)(1\text{E}-3) = 3.2\text{E}15$ photons/s.

Radiometric flux density, in W/m^2 , represents a photon rate per unit area, and radiometric exposure per unit area is the product of the flux density times the exposure period. The active surface of a photoelectronic detector produces a current density in response to an optical flux density input, while a total signal charge is produced per unit area in the same detector during a given exposure period.

Rose¹ has shown that all types of optical detectors, e.g., photographic, electronic, or the eye, are subject to the same fundamental limits in terms of signal-to-noise ratio (R_{sn}), optical input, and exposure period. In summary, the noise in a measured signal of N_p photoelectrons during a fixed exposure period is $\sqrt{N_p}$, so that

$$R_{\text{sn}} = \sqrt{N_p} \quad (5)$$

The brightness (B_s) of a scene that produces this signal in a square pixel of dimensions ($y \cdot y$), as a result of the optical transfer and conversion from the source to the detector, possibly through a medium that absorbs, scatters, and focuses photons, is

$$B_s = \frac{C \cdot N_p}{y^2} \quad (6)$$

where C is a constant. In terms of signal-to-noise ratio,

$$B_s = \frac{C \cdot R_{sn}^2}{y^2} \quad (7)$$

Thus, for twice the signal-to-noise ratio, the scene brightness must be increased four times, or the throughput of the optical system must be quadrupled, etc. Also, if the pixel size is reduced by a factor of two, the same changes in scene brightness or optical throughput must be made in order to maintain the same signal-to-noise ratio. Under QLI conditions, higher resolution necessarily requires more input flux density for equal signal-to-noise ratio, and higher resolution inherently implies less sensitivity. The Rose limit should be used often as a proof check on design and performance estimates of LLL and other QLI imaging systems.

As an example, assume a simple imaging situation such as a single pixel, e.g., a star in the night-time sky, and an II SSA camera having an objective lens of diameter D_o . Also assume that the starlight is filtered, to observe only a narrow wavelength band, and that the photon flux density from the star is Φ_p (photon/m²/s). The number of photoelectrons produced at the photocathode of the II SSA detector (N_p) is given by

$$N_p = \Phi_p \cdot T_f \cdot \left(\frac{\pi D_o^2}{4} \right) \cdot \tau_o \cdot Y_k \cdot \tau_c \quad (8)$$

where T_f is the filter transmission, τ_o is the lens transmission, Y_k is the quantum yield of the window/photocathode assembly in the II SSA camera, and τ_c is the exposure period. Note that the II SSA camera parameters which determine the rate of production of signal photoelectrons are filter transmission, lens diameter, quantum yield, and exposure period. The key one is of course the lens diameter, and not lens f -number, for this kind of imaging; it is important, however, for extended sources such as terrestrial scenes.

Photometry and the Camera Lens

A lens on the II SSA camera is used to image a scene onto the input window/photocathode assembly of the II SSA. The relationship between the scene (E_s) and II SSA image plane (E_i) illuminances in lux (lx) is

$$E_i = \frac{\pi \cdot E_s \cdot R_s \cdot \tau_o}{(4 \cdot FN^2 \cdot (m+1)^2)} \quad (9)$$

where R_s is the scene reflectance, τ_o is the optical transmission of the lens, FN is the lens f -number, and m is the scene-to-image magnification. If E_s is in foot-lamberts, then the π is dropped and E_i is in footcandles.

Alternatively, Eq. (9) becomes

$$E_i = \frac{E_s \cdot R_s}{(4 \cdot T_n^2 \cdot (m+1)^2)} \quad (10)$$

using the T-number of the lens, where

$$T_n = \frac{FN}{\sqrt{\tau_o}} \quad (11)$$

The sensitivity of an II is usually given in two forms, i.e., “white-light” luminous sensitivity, in units of $\mu\text{A}/\text{lm}$, and absolute spectral sensitivity, in units of A/W as a function of wavelength, as discussed later in the section “Input Window/Photocathode Assemblies” in Sec. 31.4.

Example: A scene having an average reflectance of 50 percent receives LLL “full-moon” illumination of $1.0\text{E} - 2$ fc. If a lens having a T-number of 3.0 is used, and the scene is at a distance of 100 m

from a lens with a focal length of 30 mm, what is the input illumination at the II SSA? Since the distance to the scene is much longer than the focal length of the lens, the magnification is much smaller than unity and m can be neglected. Thus,

$$E_i = \frac{E_s R_s}{(4 \cdot T_n^2)} \quad (12)$$

For the given values, the input illumination at the II SSA is bound to be $E_i = (1.0E - 2 \text{ fc}) / (0.50)/(4(3.0)^2) = 1.4E - 4 \text{ fc}$.

General Considerations

It is of prime importance in any optoelectronic system to couple the maximum amount of signal input light into the primary detector surface, e.g., the window/photocathode assembly of an II SSA. In order to achieve the maximum signal-to-noise ratio, the modulation transfer function of the input optic and the spectral sensitivity of the II SSA must be carefully chosen. As shown in Fig. 2, the spectral sensitivity of a silicon SSA is much different than that of a Gen-3 image intensifier tube.

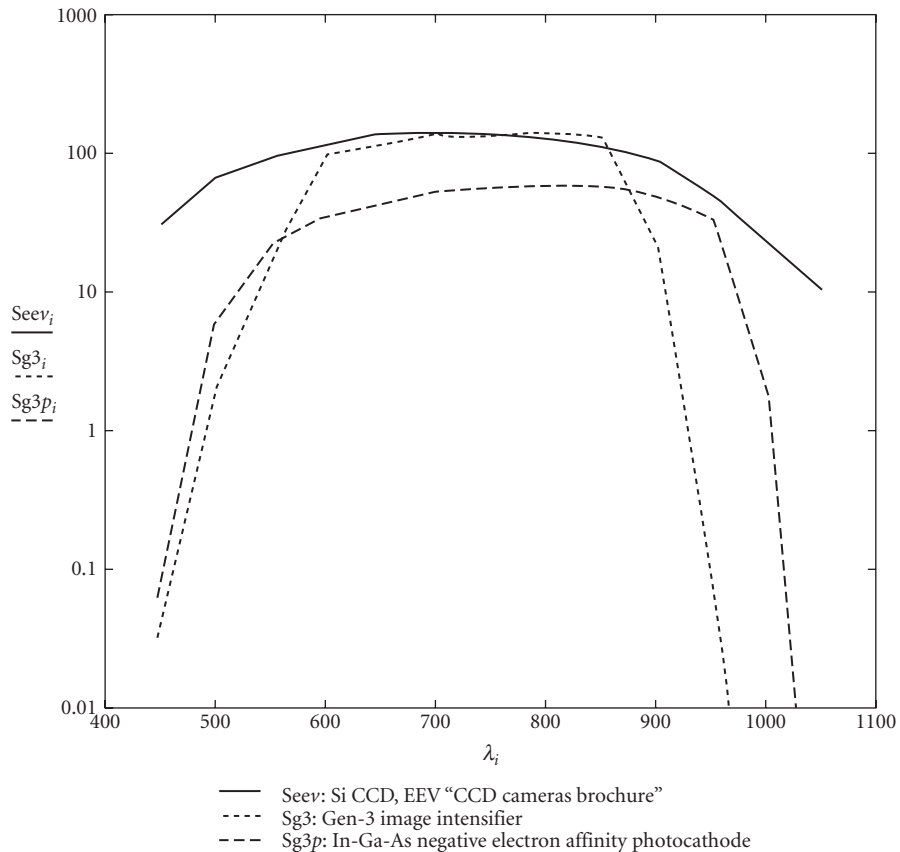


FIGURE 2 Absolute spectral sensitivity S (mA/W) versus wavelength λ (nm) of a frame-transfer type of CCD, a gen-III image intensifier, and an II having an In-Ga-As negative electron affinity photocathode.

Thus, an optimized objective lens design for a CCD will be much different than that for an II SSA. The dynamic range characteristics are also very different, since IIs will handle seven orders-of-magnitude interscene dynamic range, using a combination of II gain control and electronic duty-cycle gating, while SSAs will only provide about two orders of magnitude.²

Several factors must be considered if the overall system resolution and sensitivity are to be optimized. For example, the spectral responses of many optical input SSAs and/or lenses used in commercial cameras have been modified by using filters to reduce the red and near-ir responses to give more natural flesh tones. In an II SSA the filter may have little effect if the filter is on the SSA. The filter should not be used in the objective lens for the II SSA since a major portion of the signal will be filtered out. If a color SSA is to be used in an intensified system using relay lens coupling, sacrifice of both sensitivity and resolution will result. This is due to the matrix color filter used in these SSA chip designs. Most of the signal will go into green bandpass filter elements, and very little will go into the blue and red elements. The color matrix filter is usually bonded to the surface of the SSA chip; thus these SSA types are not used for fiberoptically coupled II SSAs.

The ideal objective lens design for an II SSA needs to be optically corrected over the spectral range of sensitivity of the II and the spectral range of interest. For special-purpose photosensitivity covering portions of the uv, blue, or near-ir spectral regions, appropriate adjustments must be made in the lens design. Although they may be adequate for many applications, it is very seldom that a commercial CCTV lens is optimized for nighttime illumination, or other LLL or QLI conditions.

Another very important part of an optimized II SSA camera design is to make the proper choice of II and SSA formats. This subject is discussed in detail later under "Fiberoptic-Coupled II/SSAs," under Sec. 31.5. The input of the II SSA system is the II, and the most likely choice will be one with an 18-mm active diameter, since the widest choice of II features is available in this size. Image intensifiers are also available having 25- and 12-mm active diameters, but these are generally more expensive. Regarding the SSA standard format sizes, the standard commercial TV formats are named by a longtime carryover from the days when vidicons were used extensively. Thus 2/3-, 1/2-, and 1/3-in format sizes originally referred to the diameters of the vidicon envelope and not the actual image format.

31.4 IMAGE INTENSIFIERS

An image intensifier (II) module, when properly coupled to an SSA camera, produces a low-light-level electronic imaging capability that is extremely useful across a broad range of application areas, including spectral analysis, medical imaging, military cameras, nighttime surveillance, high-speed optical framing cameras, and astronomy. An immediate advantage of using an II is that its absolute spectral sensitivity can be chosen from a wide variety of window/photocathode combinations to yield higher sensitivity than that of a silicon SSA. Since recently developed IIs are very small, owing to the use of microchannel plate (MCP) electron multipliers, the small size of a solid-state SSA camera is not severely compromised. In summary, advantages of using MCP IIs are

- Long life
- Low power consumption
- Small size and mass
- Rugged
- Very low image distortion
- Linear operation
- Wide dynamic range
- High-speed electronic gating, e.g., a few nanoseconds or less

An image intensifier can be thought of as an active optical element which transforms an optical image from one intensity level to another, amplifying the entire image at one time, i.e., all pixels are

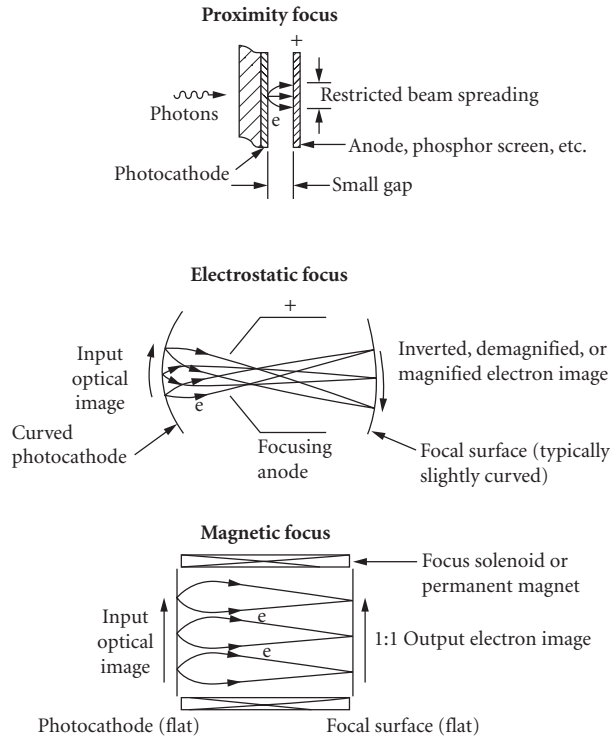


FIGURE 3 Electron lenses.

amplified in parallel and relatively independent of each other. In most cases, the resultant output image is more intense than that of the input image. The level of image amplification depends on the composite efficiency of all the conversion steps of the process involved in the image intensification operation and the basic definition of amplification. The term *image intensifier* is generally used to refer to a device that transforms visible and near-visible light into brighter visible images. Devices which convert nonvisible radiation, e.g., uv or ir, into visible images are generally referred to as *image converters*. For simplicity we refer to both types of image amplifiers/converters as “IIs” in this chapter.

Three general families of IIs exist, shown schematically in Fig. 3, that are based upon the three kinds of electron lenses used to extract the signal electrons from the photocathode, namely,

- Proximity focus IIs
- Electrostatic focus IIs
- Magnetic focus IIs

The first image tubes used a “proximity-focus” electron lens.³ Having inherently low gain and resolution, the proximity-focus lens was dropped in favor of electrostatic focus and magnetic focus IIs. The so-called Generation-0 and Generation-1 image tubes made for the U.S. Army used electrostatically focused IIs. The input end of the silicon-intensifier-target (SIT) vidicon also made use of electrostatic focusing. Magnetic focusing was used extensively in the old TV camera tubes, e.g., image orthicons, image isocons, and vidicons, and also for large-active-area and high-resolution IIs for specialized military and scientific markets.

With the development of the MCP, which was achieved for the U.S. Army’s Generation-2 types of night-vision devices, it became practical to use a proximity-focused electron lens again to meet the

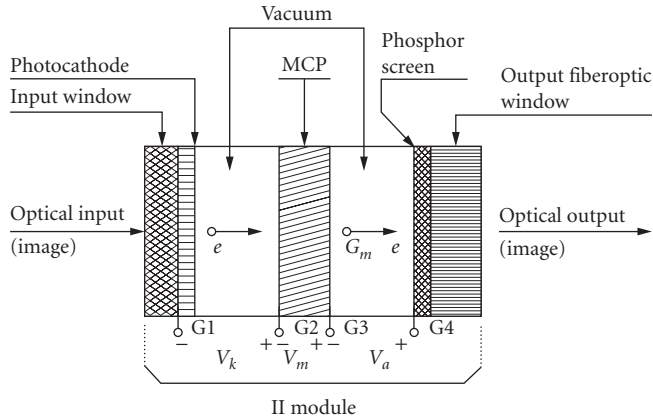


FIGURE 4 Schematic design of a proximity-focused MCP image intensifier tube module.

needs for extremely small and low-mass IIs. These “Gen-2” tubes are being used extensively for military night-vision applications, e.g., night-vision goggles for helicopter pilots, individual soldier helmet mounted night-vision goggles, etc. The most recently developed “Gen-3” IIs have higher sensitivity and limiting resolution characteristics than Gen-2 IIs, and they are used in similar night-vision systems.

Both the Gen-2 and Gen-3 types of IIs are available for use as low-noise, low-light-level amplifiers in II SSA cameras. In addition, by choosing special input window/photocathode combinations outside the military needs for Gen-2 and Gen-3 devices, a very wide range of II SSA spectral sensitivities can be achieved, well beyond silicon’s range. For II SSA camera applications, we will focus our attention exclusively on the use of proximity-focused MCP IIs because of their relative advantages over other types of IIs.

The basic components of a proximity-focused MCP II are shown schematically in Fig. 4. This type of II contains an input window, a photocathode, a microchannel plate, a phosphor screen, and an output window. The *photocathode* on the vacuum side of the *input window* converts the input optical image into an electronic image at the vacuum surface of the photocathode in the II. The *microchannel plate* (MCP) is used to amplify the electron image pixel-by-pixel. The amplified electron image at the output surface of the MCP is reconverted to a visible image using the *phosphor screen* on the vacuum side of the *output window*. This complete process results in an output image which can be as much as 20,000 to 50,000 times brighter than what the unaided eye can perceive. The input window can be either plain transparent glass, e.g., Corning type 7056, fiberoptic, sapphire, fused-silica, or virtually any optical window material that is compatible with the high-vacuum requirements of the II. The output window can be glass, but it is usually fiberoptic, with the fibers going straight through or twisted 180° for image inversion in a short distance.

A block diagram of a generalized high-voltage power supply (HVPS) used to operate the II is given in Fig. 5. For dc operation, the basic HVPS provides the following typical voltages:

$$\begin{aligned} V_k &= 200 \text{ V} \\ V_m &= 800 \text{ V for an MCP} \\ &\quad (V_m = 1600 \text{ V for a VMCP}) \\ &\quad (V_m = 2400 \text{ V for a ZMCP}) \\ V_a &= 6000 \text{ V} \end{aligned}$$

For high-speed electronic gating of the II, the photocathode is normally gated off by holding the G1 electrode a few volts positive with respect to the G2 electrode. Then, to gate the tube on and off

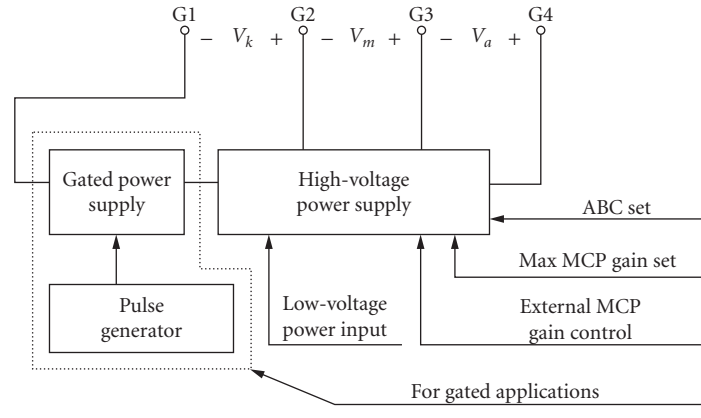


FIGURE 5 MCP image intensifier high-voltage supply.

for a short period, a pulse generator is used to control the output of the gated power supply to the normal gated on condition, i.e., $V_k = 200$ V with the polarity as shown in Fig. 5.

The dc HVPSs for IIs draw very little power, and they can be operated continuously using two AA cells, e.g., 3-V input voltage, for about 2 days. These dc HVPSs are available in small flat-packs or wraparound versions. Gated HVPSs, excluding the pulse generator, are generally at least two times larger than their dc counterparts.

In operation, an input image is focused onto the input window/photocathode assembly, producing a free-electron image pattern which is accelerated across the cathode-to-MCP gap by an applied bias voltage V_k . Electrons arriving at the MCP are swept into the channels, causing secondary electron emission gain due to the potential V_m applied across the MCP input and output electrodes. Finally, the amplified electron image emerging from the output end of the MCP is accelerated by the voltage V_a applied across the MCP-to-phosphor screen gap so that they strike an aluminized phosphor screen on a glass or FO output window with an energy of about 6 keV. This energy is sufficient to produce an output image which is many times brighter than the input image. The brightness gain of the MCP II is proportional to the product of the window/photocathode sensitivity to the input light, the gain of the MCP, and the conversion efficiency of the phosphor-screen/output-window assembly. Each of these key components and/or assemblies is discussed in more detail in the following sections of this section.

Input Window/Photocathode Assemblies

The optical spectral range of sensitivity of an II, or the II SSA that it is used in, is determined by the combination of the optical transmission properties of the window and the spectral sensitivity of the photocathode. In practice, a photocathode is formed on the input window in a high-vacuum system to produce the window/photocathode assembly as shown in Fig. 6. This assembly is then vacuum-sealed onto the II body assembly, and the finished II is then removed from the vacuum system. This type of photocathode processing is called *remote processing* (RP), because the alkali metal generators, antimony sources, and/or other materials used to form the photocathode are located outside of the vacuum II tube. Since there is no room for these photocathode material generators, remote processing must be used for MCP IIs. Also, IIs made using remote processing are found to have significantly less spurious dark current emission than the older Gen-0 and Gen-1 types of IIs having internally processed photocathodes.

The short wavelength cutoff of a window/photocathode assembly is determined by the optical transmission characteristic of the window, i.e., its thickness and material composition. The absolute

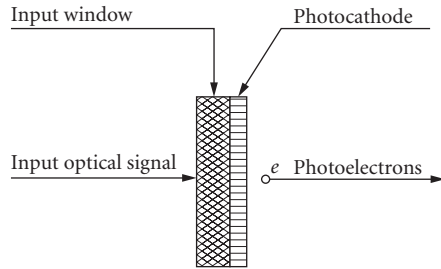


FIGURE 6 Input window/photocathode assembly.

spectral sensitivity of the photocathode determines the midrange and long wavelength cutoff characteristics of the assembly. Photocathode materials having longer wavelength cutoffs also have lower bandgap energies and generally higher thermionic emission than photocathodes with shorter wavelength cutoffs.

The spectral quantum efficiencies of various window/photocathode combinations are shown in Fig. 7 for comparison. Useful spectral bands range from the uv to the near-ir, depending upon the particular combination chosen. This figure shows the spectral sensitivity advantages that can be achieved with II SSAs. Other advantages are discussed throughout this chapter.

Note that the window/photocathode spectral quantum efficiency $[Y(\lambda)]$ curves given in Fig. 7 represent the ratio of the average number of photoelectrons produced per input photon as a function of wavelength λ . Alternatively, window/photocathode response can be specified in terms of

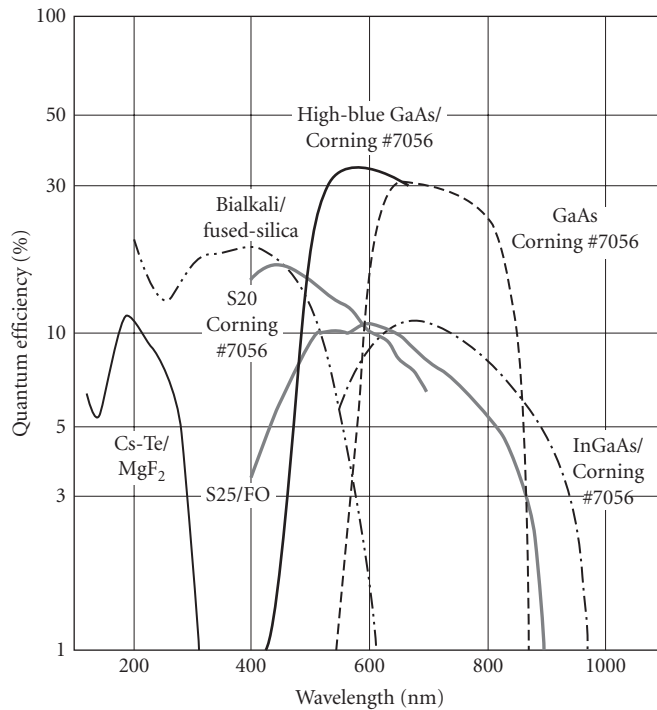


FIGURE 7 Window/cathode spectral quantum efficiencies.

absolute spectral sensitivity [$S(\lambda)$], or defined as the ratio of photocathode current per watt incident as a function of wavelength. These two parameters are related by the convenient equation

$$Y(\lambda) = \frac{124 \cdot S(\lambda)}{\lambda} \tag{13}$$

where Y is the quantum yield in percent, S is the absolute sensitivity in mA/W, and λ is the wavelength in nm.

Microchannel Plates

The development of the microchannel plate (MCP) was a revolutionary step in the art of making IIs. Although developed for and used in modern military passive night-vision systems, MCP IIs are being used today in nearly all II SSA cameras.

An MCP is shown schematically in Fig. 8. Microchannel plates are close-packed-hexagonal arrays of channel electron multipliers. With a voltage V_m applied across its input and output electrodes, the MCP produces a low-noise gain G_m , e.g., a small electron current (I_{in}) from a photocathode produces an output current $G_m I_{in}$. In addition to its function as a low-noise current amplifier, the MCP retains the current density pattern or “electron image” from its input to output electrodes. It is also possible to operate two MCPs (VCMP) or three MCPs (ZMCP) in face-to-face contact to achieve electron gains as high as about $1E7$ e/e in an II tube, as shown in Fig. 9. Other general characteristics of these types of MCP assemblies are also given in Fig. 9.

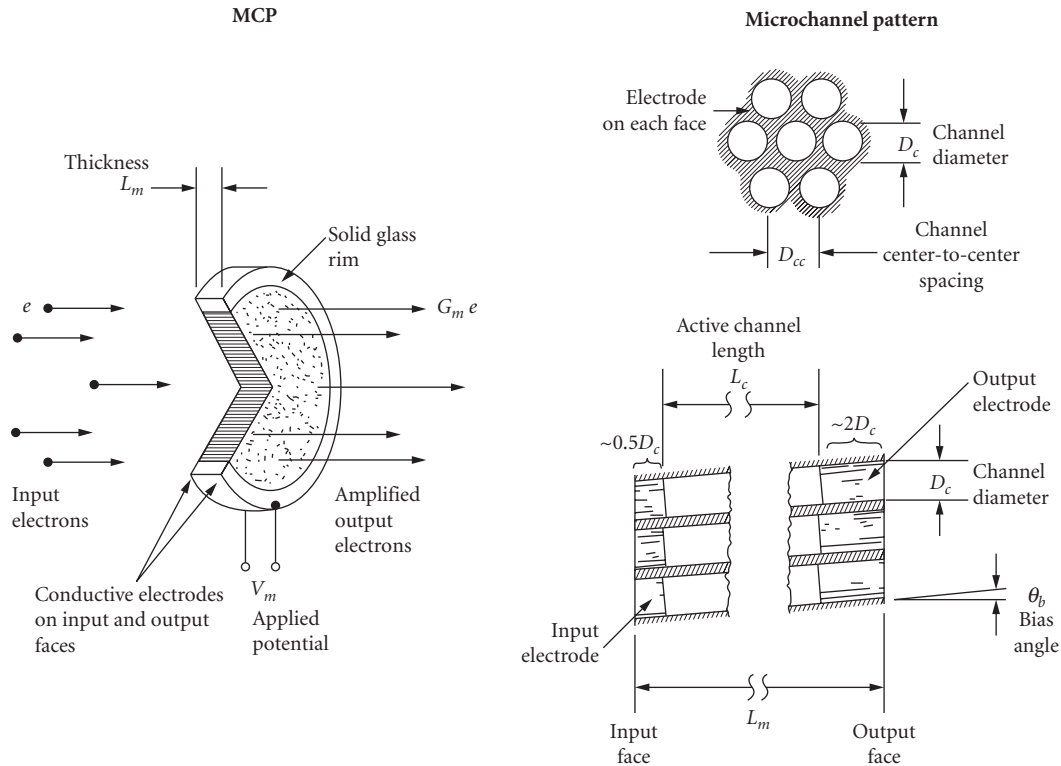


FIGURE 8 MCP parameters.




	V_m max (kV)	G_m max (e/e)	Pulse height distribution (% FWHM)	Relative limiting resolution (units)
MCP: 	1.0	1E3	Negative exponential	1.00
VMCP: 	2.0	1E5	120	0.71
ZMCP: 	3.0	1E7	80	0.50

FIGURE 9 General characteristics of MCPs, VMCPs, and ZMCPs.

The approximate limiting spatial resolutions of MCPs depend upon the channel center-to-center spacings, as follows:

Channel Diameter (μm)	Channel Center-to-Center Spacing (μm)	Approximate Limiting Resolution (lp/mm)
4	6	83
6	8	63
8	10	50
10	12	42
12	15	33

As shown in Fig. 8, MCPs are made to have channel axes that make a “bias angle” (θ_b) with respect to the normal to its input and output faces. This bias angle improves electron gain and reduces noise factor by reducing “boresighting” of electrons into the channels. The MCP bias current or “strip current” (I_s) that results from the voltage applied to the MCP sets an upper limit to the maximum linear dynamic range of the MCP. Generally, when the output current density of the MCP is in excess of about 10 percent of the strip current density, the MCP ceases to remain a linear amplifier. Conventional MCPs have strip current densities of about $1 \mu\text{A}/\text{cm}^2$, and recent high-output-technology MCPs (HOT MCPs)⁴ have become available that have strip current densities as high as about $40 \mu\text{A}/\text{cm}^2$. Electron-gain characteristics of MCP assemblies are given approximately by the equation and associated parameters shown in Table 1.

TABLE 1 MCP Gain Equation and Gain Parameters

$G_m(V_m) = \left(\frac{V_m}{V_c}\right)^g$			
Type	V_c (V)	g (units)	(L_m/D_c) (units)
MCP	350	8.5	40
MCP	530	13	60
VMCP	700	17	80
ZMCP	1050	25	120

Power noise factors for conventional MCPs, used in Gen-2 IIs, and “filmed-MCPs,” used in Gen-3 IIs, are approximately 2.0 and 3.5, respectively. Detailed information on MCP gain, noise factors, and other parameters are given by Eberhardt.⁵ Note that MCP gain is a strong function of the channel length-to-diameter ratio. The parameter V_c in the gain equation is the “crossover” voltage for the channel, i.e., it is the MCP applied voltage at which the gain is exactly unity.

Phosphor Screens

Output spectral and temporal characteristics of a wide variety of screens are given in an Electronic Industries Association publication.⁶ The phosphor materials covered in this publication are listed in Table 2. Both the old “P-type” and the new two-letter phosphor designations are given in this table. Any of these phosphor screen materials can be used in proximity-focused MCP IIs. However, one very commonly used phosphor is the type KA (P20) because it has a high conversion efficiency, its output spectral distribution matches the sensitivity of a silicon SSA reasonably well, it is fast enough for conventional 1/30-s frame times, it has high resolution, and it is typically used in direct-view night-vision IIs.

The three main components of an aluminized phosphor-screen/output-window assembly, of the type used in a proximity focused MCP II, are shown schematically in Fig. 10. An aluminum film electrode is deposited on the electron input side of the phosphor to accelerate the MCP output to high energy, e.g., about 6 keV, and to increase the conversion efficiency of the assembly by reflecting light toward the output window. The phosphor itself is deposited on the glass or fiberoptic output window.

Decay times, or persistence, and relative output spectral distributions for a variety of phosphor types are given in Fig. 11. Key phosphor assembly parameters that should be accounted for in the design of MCP II SSAs are MCP-to-phosphor applied potential (V_a), effective “dead-voltage” resulting from electron transmission losses in the aluminum film, phosphor screen energy input-to-output conversion efficiency, optical transmission of the glass or fiberoptic window, sine-wave MTF of the assembly, phosphor persistence, and output spectral distribution.

Before specifying the use of a particular phosphor, the operational requirements of the II SSA camera should be reviewed. The phosphor persistence should be short compared to the SSA frame

TABLE 2 Worldwide Phosphor-Type Designation System*

P1	GJ	P20	KA	P38	LK
P2	GL	P21	RD	P39	GR
P3	YB	P22	X(XX)	P40	GA
P4	WW	P23	WG	P41	YD
P5	BJ	P24	GE	P42	GW
P6	WW	P25	LJ	P43	GY
P7	GM	P26	LC	P44	GX
P10	ZA	P27	RE	P45	WB
P11	BE	P28	KE	P46	KG
P12	LB	P29	SA	P47	BH
P13	RC	P31	GH	P48	KH
P14	YC	P32	GB	P49	VA
P15	GG	P33	LD	P51	VC
P16	AA	P34	ZB	P52	BL
P17	WF	P35	BG	P53	KJ
P18	WW	P36	KF	P55	BM
P19	LF	P37	BK	P56	RF
				P57	LL

*Cross reference: old-to-new designations.

Source: Adapted from Electronic Industries Association Publication, no. 116-A, 1985.

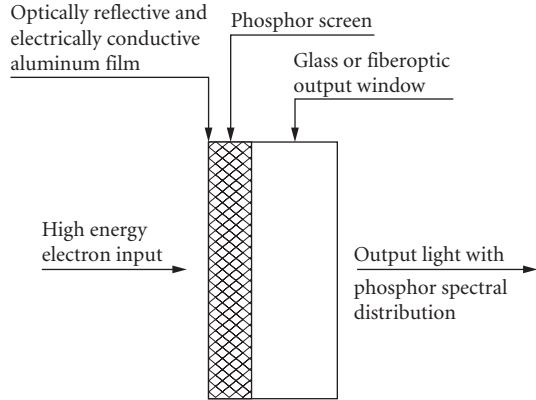


FIGURE 10 Aluminized phosphor screen and window assembly.

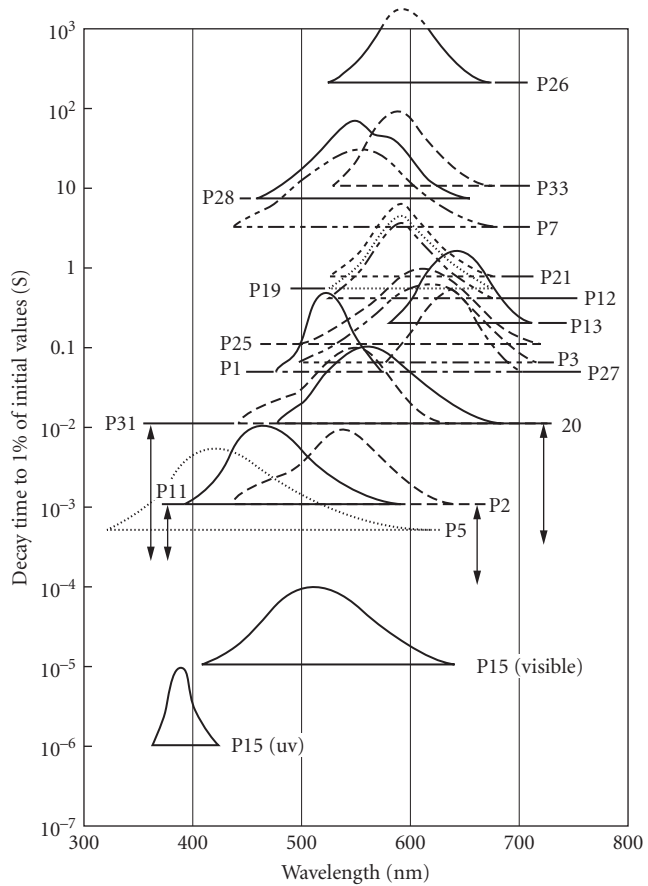


FIGURE 11 Phosphor screen decay times and spectral outputs. (Reprinted with permission from United Mineral and Chemical Co.)

time to minimize image smear due to rapidly moving objects. Also, the absolute conversion efficiency of the phosphor assembly and its relative output spectral distribution should be spectrally matched⁷ to the sensitivity of the SSA for maximum coupling efficiency.

Typical absolute spectral response characteristics, i.e., the phosphor spectral efficiency (radiated watts per nanometer per watt excitation) as a function of wavelength, of aluminized phosphor screens are given in Ref. 7. The associated phosphor screen efficiencies are also given in this reference in three different ways:

- Typical quantum yield factor: photons out per eV input
- Typical absolute efficiency: radiated watts per watt excitation
- Typical luminous equivalent: radiated lumens per radiated watt

For example, a type KA(P20) aluminized phosphor-screen/glass window assembly is found to have its peak output at 560 nm and a typical quantum yield factor of 0.063 photons/eV. Thus, an electron which leaves the MCP and strikes the assembly with 6 keV of energy, and for a “dead-voltage” of 3 kV, approximately $(6 - 3) \text{ keV} \times 0.063 \text{ photons/eV} = 190$ photons will be produced at the output.

Proximity-Focused MCP IIs

By combining the image transfer and conversion properties of the three major proximity-focused MCP II assemblies discussed earlier, i.e.,

- Input window/photocathode
- Microchannel plate
- Phosphor screen/output window

the operational characteristics of the II itself, as shown in Fig. 4, can be determined.

For example, consider an II CCD application for a space-based astronomical telescope that requires more than 10 percent quantum yield at 200 nm, but minimum sensitivity beyond 300 nm. It is desired that the top end of the dynamic range be at an input window signal flux (F_{si}) of 1000 photon/pixel/s at 250 nm. Let the CCD have a 1-in vidicon format, i.e., an active area of $11.9 \times 8.9 \text{ mm}^2$, with 325 vertical columns and 244 horizontal rows of pixels. The limiting resolution of even a dual-MCP (VMCP) image tube has a limiting resolution that is significantly higher than the horizontal pixel spatial Nyquist frequency (f_N) in the CCD, so that the pixel size at the input to the II will be essentially the same as that of the CCD. Let us rough-in an II design by making the following additional assumptions:

MCP-to-phosphor applied potential (V_a)	6000 V
Phosphor screen type	KA (P20)
Phosphor screen/window-quantum yield (P_q)	0.06 photon/eV
Phosphor screen dead voltage (V_d)	3000 V
CCD charge integration period (τ_i)	33 ms
CCD pixel full-well charge	1 pC = $6.3E6 \text{ e}$

An II with an 18-mm active diameter can be used, since the diagonal of the CCD active area is 14.9 mm. From Fig. 7, the $\text{MgF}_2/\text{Cs-Te}$ window/photocathode assembly will be chosen, having a quantum yield (Y_k) of 0.12 at 200 nm, to meet the spectral sensitivity requirements.

Let's now proceed to estimate the required gain of the MCP structure, decide what kind of an MCP structure to use, and determine its operating point. A first-order estimate of the stored pixel charge (Q_{ccd}) for the given input signal flux density is

$$Q_{ccd} = F_{si} \cdot Y_k \cdot G_m \cdot (V_a - V_d) \cdot P_q \cdot Y_{ccd} \cdot \tau_i \quad (14)$$

Since

$$F_{si} = 1000 \text{ photon/pixel/s}$$

$$Y_k = 0.10 \text{ e/photon}$$

$$Y_{ccd} = 0.3 \text{ e/photon}$$

it is found that $Q_{ccd} = G_m (178 \text{ e/pixel})$. Setting this charge equal to the full-well pixel charge gives $G_m = 6.3E6 \text{ e/pixel}/(178 \text{ e/pixel}) = 3.5E4 \text{ e/e}$. This MCP assembly gain is easily satisfied by using a VMCP. From Table 1, it is found that the gain of a VMCP is given approximately by $G_m = (V_m/700)^{17} = 3.5E4 \text{ e/e}$. Solving for V_m gives $V_m = 1300 \text{ V}$.

Thus, a first-order estimate for the general requirements to be placed in the II to do the job is as follows:

Active diameter	18 mm
Quality area	(11.9 × 8.9 mm)
Input window/photocathode	Fused-Silica/Cs-Te
MCP assembly	VMCP
Aluminized phosphor screen assembly	KA/FO window

Coupling this II to the specified FO input window CCD, e.g., by using a suitable optical cement, will meet the specified objective. Other parameters like the dark count rate per pixel as a function of temperature, the DQE of the II CCD, cosmetic, uniformity of sensitivity, and other specifications will have to be considered as well before completing the design.

Recent “Generations” of MCP IIs The most impressive improvement in direct-view night-vision devices has come with the advent of Gen-3 technology. The improvement, which is most apparent at very low light levels, is mainly due to the use of GaAs as the photocathode material. At higher light levels, e.g., half-moon to full-moon conditions, the Gen-2+ gives somewhat better performance. Key to the detection of objects under LLL conditions is the efficiency of the photocathode; the Gen-3 sensitivity is typically a factor of 3 higher. Also, the spectral response of Gen-3 matches better to the night sky spectral illumination. This equates to being able to see at almost one decade lower scene illumination with Gen-3. A summary of proximity-focused MCP image intensifier general characteristics is given in Table 3.

TABLE 3 Summary of Proximity-Focused MCP Image Intensifier General Characteristics

Minimum Active Diameter (mm)	Input Window Material*	Spectral Sensitivity Range (nm)	MCP Assembly Type	Temperature Rating		Output Window Material	Minimum Limiting Resolution (lp/mm)	Technology Type
				Storage (°C)	Operating (°C)			
11.3	FS	160–850	MCP	–55, +65	–20, +40	FO	25	Gen-2
12.0	FS, G, FO	160–900	MCP, VMCP, ZMCP	–57, +65	–51, +45	FO, G	45, 29, 20	Gen-2
17.5	FS, G, FO	600–900	MCP, VMCP, ZMCP	–57, +65	–51, +45	FO, G	45, 25, 20	Gen-2
17.5	G, FO	500–1100	MCP, VMCP, ZMCP	–57, +95	–51, +52	FO, G	45, 25, 20	Gen-3
25.0	FS, G, FO	160–900	MCP	–57, +65	–51, +45	FO, G	40	Gen-2
25.0	G, FO	500–1100	MCP	–57, +95	–51, +52	FO, G	40	Gen-3

*FS = fused silica; G = Corning #7056 glass; FO = fiberoptic.

Technology Type	Options	
	Photocathode	Phosphor
Gen-2	All but GaAs, InGaAs	Wide selection
Gen-3	GaAs, InGaAs	Wide selection

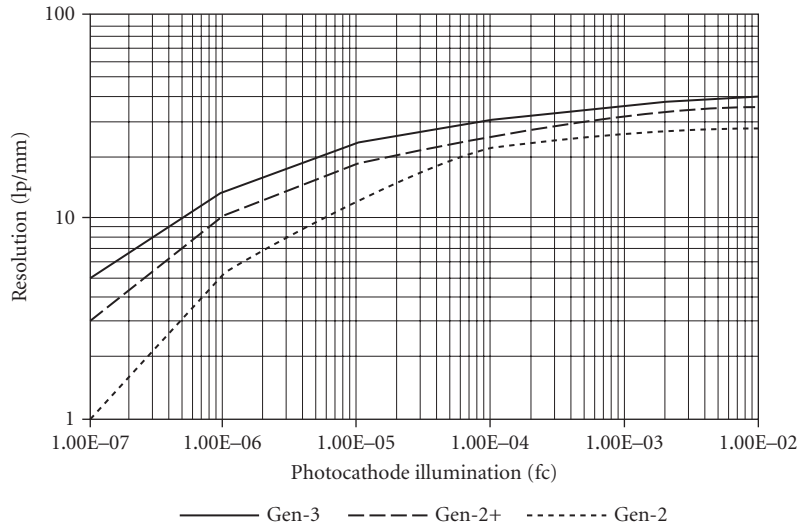


FIGURE 12 Image intensifier tube resolution curves.

For systems design work, it is useful to know the approximate characteristics of the three most recent generations in terms of II resolution versus photocathode illumination. The resolution transfer curves shown in Fig. 12 give the II resolution, observable by the eye, as a function of input illumination for Gen-2, Gen-2+, and Gen-3 IIs. These curves do not include system optics degradations, except in the sense that a human observer made the resolution measurements using a 10-power eyepiece in viewing the output image of the II.

Improved Performance Gen-2 IIs Recent enhancements in the dynamic range performance of Gen-2 IIs for direct-view applications have been made which also benefit II SSA camera performance. Improvement goals were to increase both the usable output brightness and the LLL gain of Gen-2 IIs. Night-vision devices are normally used at light levels ranging from full moon to just below quarter-moon, or in dark city environments with ample scattered light. It is important to have good contrast over as wide a light-level range as possible. To get this extended dynamic range, the gain should be held nearly constant to as high a level as possible, for improved contrast at the high-light levels. Any gain improvement should be attained with little or no increase in noise, to ensure good performance at the minimum light levels. Reducing the objective lens f -number as low as possible also improves system performance and gain. However, f -number reduction by itself may create problems in the system dynamic range if the II and its power supply assembly is not appropriately adjusted to match the optical throughput.

Figure 13 shows the extended dynamic range of a Gen-2+ II and power supply assembly, as compared to the typical MIL-SPEC Gen-2 assembly. Increasing the gain in a standard Gen-2 assembly by increasing the gain control voltage, i.e., the MCP voltage, will not give the same benefits as the Gen-2+. Ideally, a change of one unit in input brightness should result in a proportional output brightness change. The increased near-linear gain range up to higher-output light levels in the Gen-2+ improves the contrast at the higher levels. Brightness limiting begins reducing the gain to hold the output brightness constant after the automatic brightness control (ABC) limit of the power supply is reached. The increased gain of the Gen-2+ improves the performance at the lowest-light levels as well.

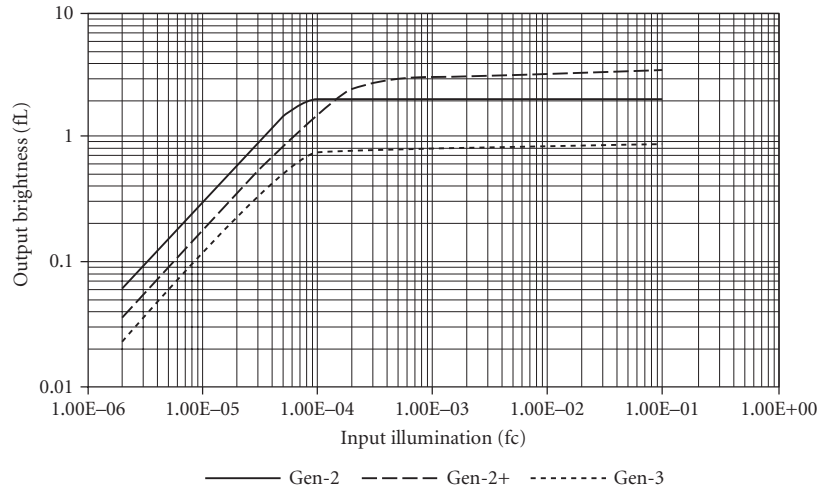


FIGURE 13 Output versus input transfer characteristics of Gen-2, Gen-2+, and Gen-3 II/power supply assemblies.

31.5 IMAGE INTENSIFIED SELF-SCANNED ARRAYS

There are several reasons to consider using an II SSA instead of an SSA alone. One obvious reason is to achieve LLL sensitivity. Figure 14 shows the limiting resolution versus faceplate illumination characteristic of CID camera operating in the unintensified and intensified modes.⁸ It is seen that

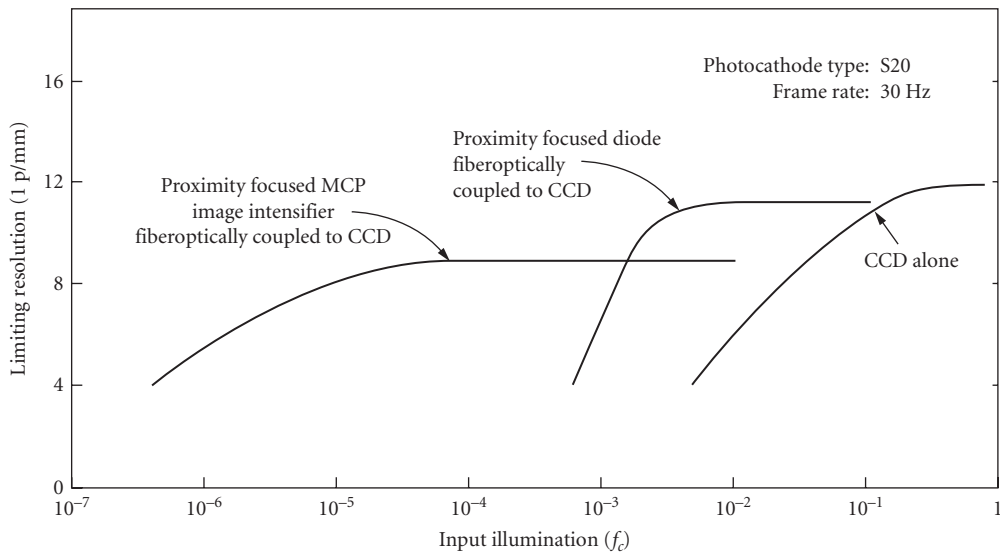


FIGURE 14 Resolution versus input illumination characteristics of a conventional optical input CCD camera and the same camera fiberoptically coupled to an MCP image intensifier tube. (From Ref. 8.)

LLL sensitivity is achieved by coupling the CID to an image intensifier tube, albeit at the expense of reduced high-light resolution. Other reasons for using an II SSA are

- High-speed electronic gating, down to a few nanoseconds, for framing cameras, LADAR, smoke and fog penetration
- Improved spectral sensitivity
- Use in a TV camera system that operates automatically under lighting conditions ranging from nighttime to full daylight conditions.
- High-sensitivity and high-speed-gated optical multichannel analyzers (OMAs)

Fiberoptic-Coupled II/SSAs

Figure 15 shows a schematic design of a fiberoptically (FO) coupled II SSA assembly. These designs are modular, since an II module is optically coupled to an SSA module. Virtually any type of image tube can be optically coupled to an SSA. The fiberoptically coupled design shown in Fig. 15 requires the use of an II having a fiberoptic output window and an SSA having an SSA input window. A fiberoptic taper, instead of a simple unity magnification FO window, is also generally required to efficiently couple the output of the II into the SSA, and this is shown in Fig. 15 as a separate module. The various fiberoptic modules are joined at interfaces 1, 2, and 3, using optical cement, optical grease, immersion oil, or "air." For the highest-resolution image transfer across these interfaces, it is necessary that the gap length at each interface be kept short, and the numerical aperture of the fiberoptic windows should be kept as low as possible, consistent with the SNR and gain requirements. It has been shown⁹ that the first interface can be eliminated by making the fiberoptic taper part of the II and depositing the phosphor screen directly onto it, and interface 3 can also be eliminated by coupling the fiberoptic taper directly to the SSA. The properties of the image transfer and conversion components shown in Fig. 15 can be used to estimate the overall performance characteristics of the fiberoptically coupled II SSA camera.

The terminology used to define SSA image format sizes derives from the earlier vidicon camera tube technology. The mass, volume, and power requirements of vidicon cameras are much larger than SSA cameras. Vidicons also have image distortion and gamma characteristics which must be accounted for, whereas SSAs and II SSSAs using proximity-focused IIs are nearly distortion-free with

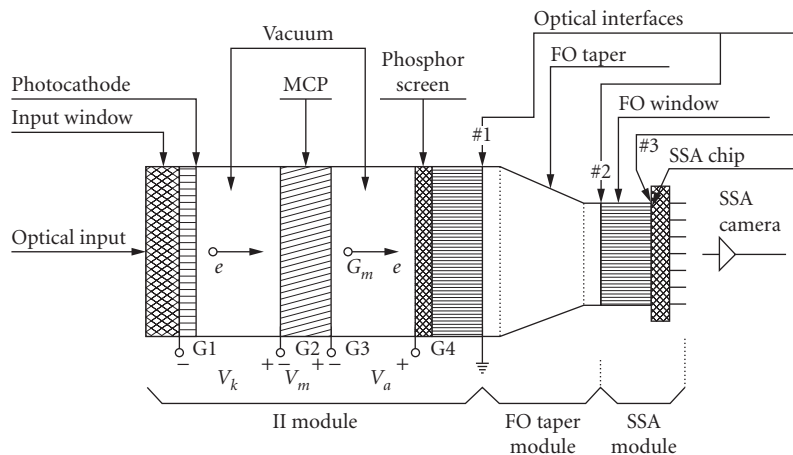


FIGURE 15 Schematic design of fiberoptically coupled II SSA assembly.

TABLE 4 Comparison of Basic Image Intensifier Diameters, SSA Format Sizes, Matching Fiberoptic Taper Magnifications, and Limiting Resolutions at the Fiberoptic Taper Output Surface (for 45 lp/mm Intensifier)

Image Intensifier Active Dia. (mm)	SSA Format		Diagonal (mm)	(M_{fot}) FOT Magnification	(f_{ito}) Limiting Resolution at FOT Output (lp/mm)
	Vidicon (in)	(mm)			
25	1	11.9 × 8.9	14.9	0.596	76
25	2/3	8.8 × 6.6	11.0	0.440	102
18	1	11.9 × 8.9	14.9	0.828	54
18	2/3	8.8 × 6.6	11.0	0.611	74
18	1/2	6.5 × 4.85	8.1	0.451	100
12	2/3	8.8 × 6.6	11.0	0.917	49
12	1/2	6.5 × 4.85	8.1	0.676	67
12	1/3	4.8 × 3.6	6.0	0.500	90

linear, i.e., unity gamma, input/output transfer characteristics over wide intrascene dynamic ranges. Table 4 gives the basic II active diameters, SSA format sizes, SSA active-area diagonal lengths, fiberoptic taper magnifications (M_{fot}) required to couple II outputs to the SSAs, and limiting resolutions (f_{ito}) at the fiberoptic taper output. Figure 16 shows schematically the relative sizes of the standard active diameters of IIs and the standard SSA formats.

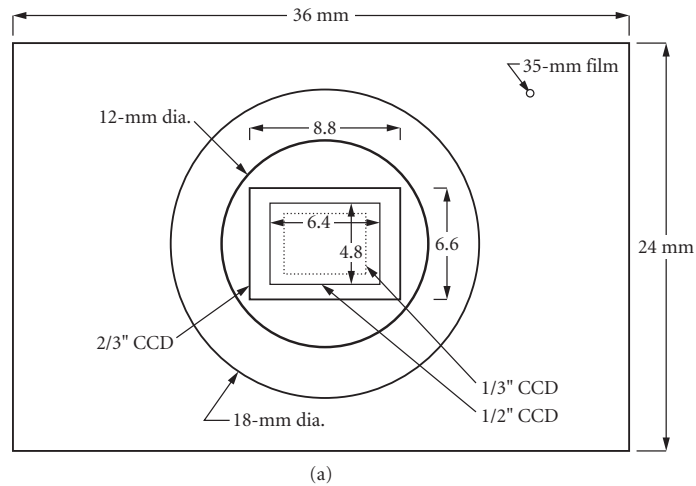


FIGURE 16a Typical 35-mm film, image intensifier and SSA formats.

	18-mm Φ image tube	12-mm Φ image tube	1" CCD	2/3" CCD	1/2" CCD	1/3" CCD	35-mm film	Units
Diagonal	18.0	12.0	14.9	11.0	8.1	6.0	43.3	mm
Vertical	10.8	7.2	8.9	6.6	4.9	3.6	24.0	mm
Horizontal	14.4	9.6	11.9	8.8	6.5	4.8	36.0	mm
Area	155.5	69.1	105.9	58.1	31.5	17.3	864.0	mm ²

(b)

FIGURE 16b Typical dimensions for image intensifiers and SSAs using 3:4 format.

The present limiting resolution range of MCP IIs is 36 to 51 lp/mm. In an II SSA, the resolution of the II should be matched, in some sense, to that of the SSA. For example, it is unwise to use a low-resolution II and fiberoptic lens combination with a much higher resolution CCD.

Lens-Coupled II SSAs

Figure 17 is a schematic design for a lens-coupled II SSA assembly. The differences between this design and the fiberoptic-coupled II SSA design described earlier are that the output window of the II can be either fiberoptic or glass, and a lens is used instead of an FO taper to couple the output optical image from the II directly into a conventional optical input SSA, i.e., no FO window is required at the SSA. Although the lens-coupling efficiency is lower, its image distortion and resolution performance is superior to the FO-coupled design. Also, the chance for possible adverse rf interference at the sensitive input to the SSA camera from the II high-voltage power supply is less than for the lens-coupled design.

Parameters to Specify Typical parameters to specify for an MCP II SSA detector assembly, using either fiberoptic or lens-coupling, are as follows:

- Sensitivity
 - White-light (2856 K) ($\mu\text{A}/\text{lm}$)
 - Spectral sensitivity (mA/W versus nm)
 - Sensitivity (mA/W at specified wavelength)
- EBI (lm/cm^2 at 23°C)
- MCP applied potential for 10 K fl/fc luminous gain (V)
- Horizontal resolution at specified input illumination (TVL)
- Shades-of-gray (units)
- Cosmetic properties
 - Uniformity (percent)
 - Bright spots (number allowable in format zone) Dark spots (number allowable in format zone)
- Burn-in (procedure)

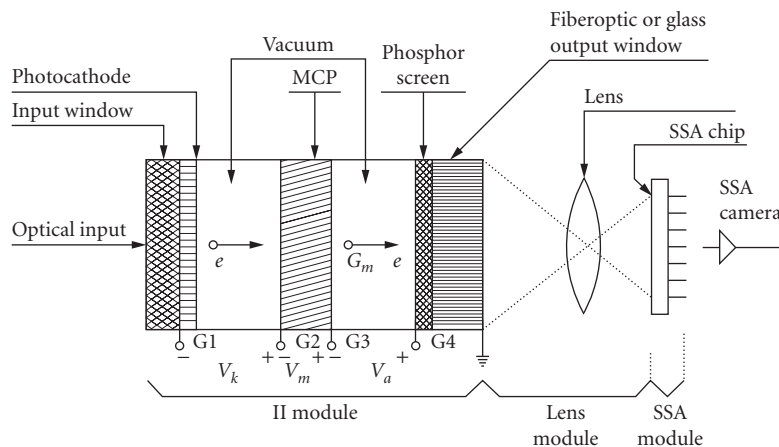


FIGURE 17 Schematic design of lens-coupled II SSA assembly.

- Mechanical specifications
- Dimensions (interface drawing)
- Mass (g)
- Environmental (specified)

Electron-Bombarded SSA

Since the early work by Abraham et al.¹⁰ which showed the feasibility of achieving useful electron gain by electron bombardment (EB) of a silicon diode in a photomultiplier tube, several attempts have been made to achieve similar operation using an SSA specially designed for EB input, instead of optical input. The charge gain (G_{eb}) resulting from the electron bombardment is given by

$$G_{eb} = \frac{(V_a - V_d)}{3.6} \quad (15)$$

where V_a is the acceleration voltage and V_d is the “dead-voltage” of the EBSSA. It was quickly found that successful CCD operation could not be obtained by simply bombarding the normal optical input side of the chip with electrons, because interface states soon form which prevent readout of the chip and other problems. By thinning a CCD chip to 10 to 15 μm from the “backside” and operating in a backside EB-mode, useful performance is achieved. In this way, 100 percent of the silicon chip is sensitive to incident photoelectrons, and it becomes technically feasible to make EBSSA cameras.

Proximity Focused EBSSAs A proximity-focused EBSSA is shown schematically in Fig. 18. In this design, the input light enters the window/photocathode assembly to generate the signal photoelectrons which are accelerated to about 10-keV energy and bombard the thinned backside of the EBSSA. Note that no MCP, no MCP-to-screen gap, no phosphor screen/output window assembly, and no fiberoptic or lens coupling is used to transfer the electronic image to the SSA for readout. Thus, higher limiting resolution is attainable. Also, the power noise factor associated with the EBSSA gain process is lower than that of MCP devices, and image lag is eliminated because no phosphor is used. Early work on proximity-focused EBDDs was done by Barton et al.,¹¹ Williams,¹² and Cuny et al.¹³ By 1979, a 100×160 pixel TI CCD was used in this type of detector and put into a miniature TV camera. With an acceleration voltage of $V_a = 15$ kV, an electron gain of 2000 was achieved, along with a Nyquist limited resolution of 20 lp/mm. Recent advances have brought this technology closer

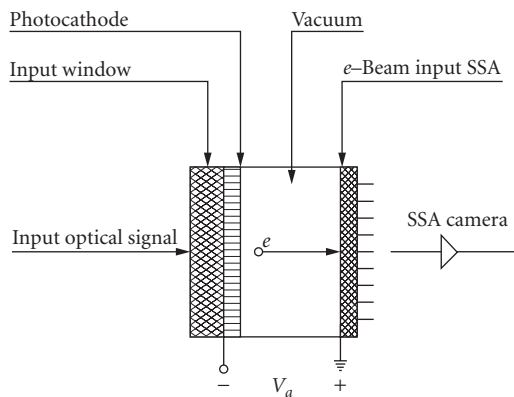


FIGURE 18 Electron bombarded SSA (EBSSA).

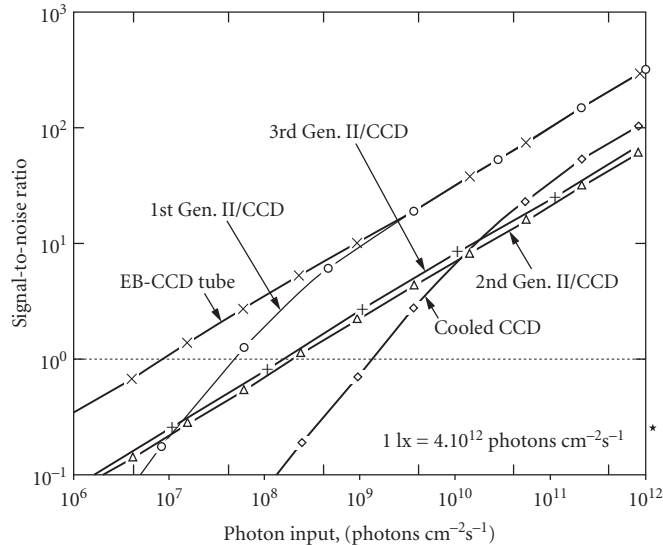


FIGURE 19 Comparison of the signal-to-noise ratio of various optoelectronic imagers versus the photon input. (From Ref. 14.)

to extensive usage possibilities. Richard et al.¹⁴ have compared the SNR characteristics of an EB CCD tube, various other types of II CCDs, and bare CCDs. Their results are shown in Fig. 19.

In order to achieve its full performance capabilities, the energy of the bombarding electrons must be absorbed by the active silicon SSA material, photoelectrons must not be lost, the exposure of the EBSSA to high-energy electrons should not cause a life problem, and it must be possible to read out the stored charge pattern in the SSA. It is found that recombination phenomena at the EB-input face can be reduced with a p^+ passivation layer, e.g., by using $3E17 \text{ cm}^{-3}$ boron doping, which reduces back-diffusion of signal electrons, front-diffusion of “dark” charges from the rear face, reduced diffusion length, separation of holes and electrons by the built-in electric field, and higher surface conductivity, thus better voltage stability, at the rear face.

Internally processed (IP) and remotely processed (RP) or “transfer” photocathodes have been used in EBSSAs. It is generally found that the internal processing produces consistently higher-background and spurious noise problems due to field emission from tube body parts and the photocathode. Both types of photocathode processes have yielded long-life EBCCD detectors.

Proven applications to date for EBSSA detectors:

- Photon-counting wavefront sensor (adaptive optics), European Space Organization 3.6-m telescope at La Silla, Chile
- NASA, Goddard Space Flight Center, Oblique Imaging EB CCD uv sensitive camera

Advantages of EBSSA cameras over MCP II-based II-SSAs:

- No image lag
- Higher resolution
- Single photoelectron detection per frame per pixel
- Higher DQE

Digital II SSA Cameras Consider a photon-counting imaging detector consisting of an MCP image intensifier tube (II) that is fiberoptically coupled to a silicon solid-state self-scanned array

(SSA) chip in a TV camera. Incoming photons at wavelength λ pass through the input window of the II and produce an average quantum yield of Y_k photoelectrons per photon at the photocathode. The resulting photoelectrons (e) are accelerated into the MCP electron multiplier assembly. Amplified output electrons from this low-noise electron multiplier are accelerated into an aluminized phosphor screen on the output window of the II. The number of output photons from the II per photoelectron is proportional to the electron gain in the MCP (G_m), the effective electron bombardment energy at the phosphor screen (ϵV_s), and finally the electron-input to photon-output conversion efficiency (P) at the phosphor screen. As discussed earlier the optical transmission of the input window and the actual quantum yield of the photocathode are usually factored together in the average quantum yield parameter Y_k , and the optical transmission of the output window is also normally factored together with the actual conversion efficiency of the phosphor screen in the screen efficiency parameter P .

The output photon pulse from the II, resulting from the single detected input photon, is coupled into the SSA via the fiberoptic taper, which matches the output size of the II to the size of the SSA, and a fiberoptic window on the SSA. This photon pulse is then converted to an electron signal charge packet (Q_{ssa}) at the SSA. The number of electrons stored per pixel in the SSA depends upon the area of the photon pulse at the SSA, the spatial distribution of photons in this pulse, and the area per pixel in the SSA. Thus, in addition to the above II factors, the stored charge in the SSA per photoelectron is also proportional to the optical transmissions of the FO taper (T_{fot}) and SSA window (T_{ssa}), and the quantum yield of the SSA (Y_{ssa}).

By using two or three conventional MCPs in cascade, i.e., VMCPs or ZMCPs, the gain can be made so large that it completely overrides any normal room-temperature thermal dark current in an SSA at a conventional RS-170 rate. In this photon-counting mode of operation, a charge signal above a preset threshold value is looked for. When it is found in a given pixel, a "1" is stored in memory for that pixel's address, "0s" are stored in pixel addresses where this condition is not met, and the entire frame is read out. By reading out a total of N_f frames, the dynamic range can be made as high as N_f if the dark count rate is negligible. Thus, photon-counting imaging can achieve a very large dynamic range.

Another advantage of photon-counting imaging is that the image resolution can also be made very high by centroiding the detected charge packets in the SSA. Since the performance of a centroiding camera depends upon the signal-processing algorithm, this will not be analyzed here. Instead, the reader is referred to several references in which centroiding is discussed.¹⁵

Let us next calculate the stored charge and number of stored electrons in a photon-counting II SSA per photoelectron. Assume that a proximity-focused VMCP II is coupled to the SSA with a fiberoptic taper. For our analysis, some typical values will be used for the operating voltage and gain of a VMCP: the acceleration voltage between the VMCP and the phosphor screen, the efficiency of an aluminized type KA (P20) phosphor screen, the optical transmissions of a fiberoptic taper and an SSA fiberoptic window, and the quantum yield of an SSA.

Definitions for the parameters that will be used are summarized as follows:

Q_{ssa}	stored SSA charge per input photoelectron from the photocathode
Y_k	photocathode quantum yield, e/photon
G_m	VMCP electron gain, e/e
V_m	VMCP applied potential, V
V_s	MCP-to-screen applied potential, V
V_d	phosphor screen "dead-voltage"
P	phosphor screen efficiency, photon/eV
T_{fot}	transmission of fiberoptic taper
T_{ssa}	transmission of fiberoptic window on the SSA
Y_{ssa}	quantum yield of SSA, e/photon
e	electron charge, 1.6E-19 C
N_{essa}	number of stored SSA electrons per input photoelectron

Using these definitions, the general equation for the charge stored in the SSA per input photoelectron is given by Eq. (16).

$$Q_{ssa} = e \cdot G_m \cdot (V_s - V_d) \cdot P \cdot T_{\text{tot}} \cdot T_{ssa} \cdot Y_{ssa} \quad (16)$$

Thus, Q_{ssa} is given by the product of the VMCP gain, the effective electron bombardment energy at the aluminized phosphor screen, the conversion efficiency of the phosphor screen assembly, the transmissions of the FO taper and the SSA's FO window, and finally the quantum yield of the SSA.

For

$$V_m = 1380 \text{ V}$$

the VMCP electron gain is

$$G_m(V_m) = \left(\frac{V_m}{700 \text{ V}} \right)^{17}$$

$$G_m(V_m) = 1 \times 10^5$$

By using the following values for the additional parameters

$$V_s = 5500 \text{ V}$$

$$V_d = 2500 \text{ V}$$

$$P = 0.06 \text{ photon/eV}$$

$$T_{\text{tot}} = 0.6$$

$$T_{ssa} = 0.8$$

$$Y_{ssa} = 0.5 \text{ e/photon}$$

it is found that the stored charge per photoelectron is

$$Q_{ssa} = 7 \times 10^{-13} \text{ C}$$

and that the number of electrons stored in the SSA per input photoelectron is

$$N_{\text{essa}} = Q_{ssa} / e$$

$$N_{\text{essa}} = 4 \times 10^6 \text{ electrons}$$

Since the full-well or saturation charge for an SSA pixel is on the order of 1 pC, this VMCP II/SSA assembly is seen to qualify as a photon-counting imaging detector.

Modulation Transfer Function and Limiting Resolution The modulation transfer function (MTF) of an II SSA camera is determined by a convolution of the individual MTFs of the camera lens, II, II-to-SSA-coupling fiberoptic or lens, fiberoptic-to-fiberoptic interfaces, fiberoptic-to-SSA interface, SSA, etc. There are several ways to determine the MTF of an existing II SSA camera.

For example, the II SSA camera can be focused on a spatial frequency burst pattern, i.e., a periodic pattern of black and white bars in which the spatial frequency of the bars increases in one direction. Alignment of the pattern's bars with pixel columns and readout of the modulation of the spatial pattern in the pixel row direction gives the squarewave MTF of the camera $S(f)$, where f is the spatial frequency in cycles/mm. Conversion of this square-wave MTF to a sine-wave MTF is accomplished by using the Fourier transform at a given frequency:

$$T(f) = \left(\frac{\pi}{4} \right) \cdot \left(S(f) + \frac{S(3f)}{3} - \frac{S(5f)}{5} + \frac{S(7f)}{7} + \dots \right) \quad (17)$$

By calculating several values of $T(f)$ from the known square-wave function, the sine-wave MTF can be determined and used for camera system optical image transfer analysis. The limiting resolution is often taken to be the spatial frequency value for this sine-wave MTF at which the modulation drops to a few percent.

Also, to use the above example as an illustration, the MTF of an II SSA is a function of the direction in which the spatial frequency burst pattern is aligned. Self-scanned array pixels are not generally square, and the distances between centers of pixels in the horizontal and vertical directions are not generally the same. Thus, the corresponding MTFs in the horizontal and vertical directions are different, and the MTF is a function of the angle that the burst pattern makes with the rows and columns of pixels.

A convenient specification of the spatial frequency response of an SSA or an II SSA camera is the Nyquist frequency (f_N), defined to be the reciprocal of twice the distance between the pixels. For example, if an SSA has rows of pixels spaced on 20- μm center-to-center, then the horizontal Nyquist spatial frequency is

$$f_{N,h} = \frac{1}{(2 \cdot (0.02 \text{ mm}))} = 25 \text{ cycles/mm}$$

Assuming an II limiting resolution (f_{II}) of 32 cycle/mm, and assuming that the MTFs are gaussian, then an estimated value for the limiting resolution of the II SSA camera is

$$\left(\frac{1}{f_{\text{cam}}^2} \right) = \left(\frac{1}{f_{II}^2} \right) + \left(\frac{1}{f_{\text{ssa}}^2} \right) \quad (18)$$

or

$$f_{\text{cam}} = \frac{f_{II} \cdot f_{\text{ssa}}}{\sqrt{(f_{II}^2) + (f_{\text{ssa}}^2)}} \quad (19)$$

For the values used in this example, $f_{\text{cam}} = 20$ cycle/mm. Although this gaussian estimate is convenient to use, a more exact estimate can be made by multiplying the various component sine-wave MTFs to find the II SSA camera's MTF, and from this its limiting resolution can also be found.

For example, actual MTF measurements⁹ on an II SSA camera, having a proximity-focused 18-mm active diameter MCP II fiberoptically coupled to a CCD with an $m = 8 \text{ mm}/18 \text{ mm} = 0.44$ magnification taper, showed that the MTF of the CCD, referred to the II input, is given by $T_{\text{CCD}}(f) = \exp - (f/9.0)^{1.4}$, where f is the spatial frequency in cycles/mm. The Nyquist frequency of the CCD was $f_N = 20$ cycle/mm, and the MTF of the complete camera was found to be $T_{\text{II SSA}}(f) = \exp - (f/6.3)^{1.1}$, both referred to the II input.

31.6 APPLICATIONS

In time, it is expected that most of the quantum-limited and LLL TV applications will use some form of II SSA camera, instead of an intensified vidicon-based camera. A few of the major application areas for II SSA cameras are highlighted in this section.

Optical Multichannel Analyzers

Optical multichannel analyzers (OMAs) are instruments used to measure optical radiation in linear patterns, e.g., spectra or two-dimensional images. Photographic film, single-channel photomultiplier tubes, and TV camera tubes, e.g., vidicons, have been replaced by SSAs, e.g., CCDs, and II SSAs, e.g., image tube intensified CCDs or photodiode arrays (PDAs). Four distinct application areas exist for OMAs using II SSA detectors:¹⁶

Application	Detector*	Time
	Type	Resolution
Spectroscopy	IILPDA	50 ms
Time-resolved pulsed laser spectroscopy	GIILPDA	<5 ns
Time-resolved pulsed laser imaging spectroscopy	GISSCCD	<5 ns/spectrum
Time-resolved imaging	GISSCCD	<5 ns/spectrum

*IILPDA = image intensified photodiode array; GIILPDA = gated image intensified photodiode array; GISSCCD = gated image intensified charge-coupled device.

In each of these types of systems, the incoming radiation is converted to charge packets that are stored in each pixel of the SSA. Each line and/or field of pixels is read out by a suitable camera electronics, and the pixel charge values are stored in a computer for subsequent processing and analysis. The sensitivity of the II is chosen for best performance over the range of wavelengths being investigated. The dynamic range of OMAs can be as high as 18 bits. In comparison with the older single-element scanning system, modern II SSA-based OMAs acquire spectra up to 1000 times faster and/or with higher SNR during a given measuring period. Three common OMA applications are Raman spectroscopy, multiple input spectroscopy, and small-angle light scattering.

Range Gating and LADAR

Range gating is becoming increasingly important because the required technology now exists at an affordable cost and the signal-to-noise ratio improvement is much higher than nongated conventional TV imaging. A gated laser sends out a laser pulse of only a few nanoseconds duration while the II SSA camera is gated off. No reflection from scattering or reflection in the medium between the camera and the object is allowed to be registered in the II SSA camera.¹⁷ The II SSA camera is gated on only at the moment when the light packet from the object returns to the camera and, after exposing for the duration of the outgoing pulse, it is returned to a gated-off condition. By repeating this process and controlling the image-storing conditions, high-contrast images having high signal-to-noise ratios can be achieved. Another obvious advantage of the range-gated system is that the time-of-flight between pulse output and receipt gives the range to the object being viewed, thus leading to the realization of a *laser detection and ranging* (LADAR) system.

Microchannel plate image tubes offer high-speed gating and spectral response advantages to LADAR systems. They can be electronically gated to a few nanoseconds, i.e., providing distance resolutions of a few feet for LADAR systems. Present MCP IIs offer the user a broad range of spectral sensitivity, including near-ir imaging at 1060 nm, so that powerful and efficient lasers may be used for optimum LADAR system performance. Also, the ruggedness, extremely small size, and low power drain characteristics of II SSA cameras make them very attractive for LADAR applications for spacecraft and unmanned autonomous vehicles.

Day/Night Cameras

Full day-night interscene dynamic range capability, while maintaining a high signal-to-noise ratio, is achievable using an II SSA camera in conjunction with an auto-iris camera lens.¹⁸ The principle of operation of this type of camera can be described as follows. Assume that operation begins at the lowest light level to be encountered. The MCP voltage in the II is set to operate with high SNR for the camera lens, an auto-iris lens, set to its lowest f -number. As the light level is increased, the system operates over two orders-of-magnitude dynamic range. For four orders-of-magnitude higher light level inputs, the f -setting of the auto-iris lens is increased by a feedback circuit, driven by the peak-to-peak video output signal from the SSA camera. The effective exposure of the SSA is next automatically reduced as the light level increases by four-and-one-half orders-of-magnitude, again to

maintain a high SNR, by duty-cycle gating the MCP II. Finally, another two-and-one-half orders-of-magnitude in input light level are accommodated by reducing the gain of the MCP, i.e., by reducing its applied operating potential. The total interscene dynamic range achievable with this type of automatically controlled day/night camera is 13 orders-of-magnitude. In addition to its wide dynamic range capability, this type of camera is also able to make rapid narrow-band spectral samples across a wide spectral range, e.g., from the uv to the near-ir.

Mosaic II SSA Cameras

For very high amounts of image information throughput, multiple SSAs are used to read out large area IIs. For example, a 75-mm active diameter MCP II can be coupled fiberoptically to four individual SSA cameras. The fiberoptic couplers are made to butt against each other at the output of the II, and their output ends are optically coupled to the SSAs. Parallel readout of the SSAs is then accomplished, giving the advantage of high-resolution readout without the disadvantage of having to use a wide bandwidth video electronic system or lower frame rates.

One such II SSA camera is designed for x-ray radiology image input.¹⁹ The x-ray input image is converted to a visible light image at a scintillator screen that is in optical contact with the 6-in-diameter fiberoptic window. This scintillator/window assembly is, in turn, coupled to the input of a 6-in-diameter proximity-focused diode II, i.e., without an MCP, for modest light gain and good image quality. Six fiberoptic tapers in a 2×3 matrix couple the images from the six adjacent output sections of the II to six CCD cameras which operate in parallel to continuously read out the converted x-ray image.

Other Applications

Other applications for II SSA cameras are the following:

- Semiconductor circuit inspection
- Astronomical observations
- X-ray imaging
- Coronary angiography
- Mammography
- Nondestructive testing
- Multispectral video systems

Active Imaging Active imaging is becoming increasingly important because the required technology now exists at an affordable cost and the signal-to-noise ratio improvement is much higher than nonactive conventional TV imaging. Two types of active imaging presently exist, i.e., “line-scanned” and “range-gated.”

In a line-scanned imaging system, a narrow beam from a cw laser is raster-scanned across the object to be viewed, and the resulting reflected light is collected by a lens and detector assembly which receives and measures light from the illuminated field-of-view (FOV). Large FOV scenes can be scanned in a short period of time, which is a major advantage of this system. The signal from the detector is finally processed by a video electronics system and displayed, for direct viewing, or image processed as required. Limiting-resolution is set by the beam diameter achievable at the object. Thus systems operating in space; atmospheric and underwater environments have significantly different limiting-resolution characteristics.

In a range-gated type of system, a gated laser sends out a laser pulse of only a few nanoseconds duration while the TV camera is gated off. No reflection from scattering or reflection in the medium between the camera and the object is allowed to be registered in the TV camera. The TV camera is

gated on only at the moment when the light packet from the object returns to the camera and, after exposing for the duration of the outgoing pulse, the TV camera is returned to a gated-off condition. By repeating this process and controlling the image-storing conditions, high-contrast images having high signal-to-noise ratios can be achieved.

Another obvious advantage of the range-gated system is that the time of flight between pulse output and receipt gives the range to the object viewed, thus leading to the realization of a laser detection and ranging (LADAR) system. MicroChannel plate image tubes offer high-speed gating and spectral response advantages to LADAR systems. They can be electronically gated to a few nanoseconds, i.e., distance resolutions of a few feet, and in some parts of the optical spectrum they offer high sensitivity.

31.7 REFERENCES

1. A. Rose, "A Unified Approach to the Performance of Photographic Film, Television Pickup Tubes, and the Human Eye," *J. Soc./Motion Picture Engrs.* **47**:273–294 (1946).
2. D. E. Caudle, "Dynamic Range Enhancement Techniques for Gated, Solid-State Intensified Cameras," *SPIE* **1155**:104–109 (1990).
3. G. Hoist, J. H. de Boer, and C. F. Veenemans, *Physica* **1**:297 (1934).
4. HOT MCP™ is a trademark of Galileo Electro-Optics Corp.
5. E. H. Eberhardt, "An Operational Model for MicroChannel Plate Devices," *IEEE Trans. Nucl. Sci.* **NS-28**:712–717 (1981).
6. "Optical Characteristics of Cathode Ray Tubes," EIA Publication, no. 116-A, 1985.
7. H. P. Westman, ed., *Reference Data for Radio Engineers*, 5th ed., Howard W. Sams & Co, Inc., Indianapolis, 1969, pp. 16–34–16–37.
8. J. J. Cuny, T. F. Lynch, and C. B. Johnson, "Proximity Focused Image Tube Intensified Charge Injection Device (CID) Camera for Low Light Level Television," *SPIE* **203**:75–79 (1979).
9. G. M. Williams, Jr., "A High Performance LLLTV CCD Camera for Nighttime Pilotage," *SPIE* **1655**:14–32 (1992).
10. J. M. Abraham, L. G. Wolfgang, and C. N. Inskeep, "Application of Solid-State Elements to Photoemission Devices," *Adv. EEP* **22B**:671 (1966).
11. J. B. Barton, J. J. Curry, and D. R. Collins, "Performance Analysis of EBS-CCD Imaging Tubes/Status of ICCD Development," *Proc. Int. Conf. Apps. of CCDs*, San Diego, Calif., 1975.
12. J. T. Williams, "Test Results on Intensified Charge Coupled Devices," *SPIE* **78**:78–82 (1976).
13. J. J. Cuny, T. F. Lynch, and C. B. Johnson, "Small Intensified Charge Injection Device Cameras for Low Light Television," *MEDE '79 Conf. Proc.*, Interavia SA, Publishers, 1979, pp. 836–845.
14. J. C. Richard, M. Vittot, and J. C. Rebuffie, "Recent Developments and Applications of Electron-Bombarded CCD in Imaging," *Proc. Conf. on Photoelectronic Image Devices*, SEP91, London, IOP Pub. Ltd., 1992.
15. A. Boksenberg and D. E. Burgess, "An Image Photon Counting System for Optical Astronomy," *Adv. EEP* **33B**:835–849 (1972).
16. H. W. Messinger, "Modern Optical Multichannel Analyzers Capture Images Without Film," *Laser Focus World*, March 1992, pp. 91–94.
17. R. A. Sturz and D. E. Caudle, "Capabilities of New Cost-Effective Near Infrared Imaging," *Adv. Imaging*, April 1993, pp. 60–62.
18. D. E. Caudle, "Low Light Level Imaging Systems Application Considerations and Calculations," *SPIE* **1346**:54–63 (1990).
19. H. Roehrig et al., "High Resolution X-Ray Imaging Device," *SPIE* **1072**:88–99 (1989).

Timothy J. Tredwell

*Sensor Systems Division
Imager Systems Development Laboratory
Eastman Kodak Company
Rochester, New York*

32.1 GLOSSARY

A	area of the pixel
C_{FD}	total capacitance of the floating diffusion in a CCD output
C_g	gate capacitance per unit area
C_r	readout line capacitance
J_D	total dark current per unit area
J_s	surface generation current
L_e	diffusion length of electrons in silicon
L_p	diffusion length of holes in silicon
N_A	p -type dopant concentration in silicon
N_D	n -type dopant concentration in silicon
N_e	total number of electrons collected in a pixel
$N(0)$	number of photons entering the silicon
$N(x)$	number of photons remaining a distance x below the surface
n_i	intrinsic carrier concentration in silicon
$P(x)$	probability that an electron-hole pair generated a distance x from the surface will be collected before recombination
q	electron charge
S_o	surface recombination velocity
T_{int}	integration time of light in an image sensor
$T(\lambda)$	transmission of light
V_{bi}	built-in voltage for a silicon pn junction
$W(V)$	width of the depletion layer at a given bias voltage in the MOS capacitor or junction-photodiode
$\alpha(\lambda)$	absorption coefficient of light

ϵ_s	silicon dielectric constant
μ_e	electron mobility in silicon
μ_p	hole mobility in silicon
τ_o	depletion-layer lifetime
τ_p	minority carrier hole lifetime in silicon
Φ_s	electrostatic potential at the silicon-silicon dioxide, also called surface potential

32.2 INTRODUCTION

Since the invention of the image sensor in 1964, solid state image sensors have advanced in resolution, sensitivity, and image quality to the point where they have replaced other methods of converting visible light to electronic signals in nearly all imaging applications. There are two types of image sensors: *area image sensors*, which are used in cameras, and *linear sensors*, which are used in scanning applications. Cameras using area image sensors dominate the camcorder, video and broadcast, machine vision, scientific, and medical fields. Area image sensors for camcorder applications are typically 400,000 picture elements in resolution, 60 dB in dynamic range, and have noise levels of a few tens of electrons. Area image sensors for scientific applications may have resolutions of over six million elements, dynamic ranges exceeding 80 dB, and noise levels approaching a single electron. Scanners employing linear solid-state image sensors dominate facsimile, document scanner, digital copier, and film scanner applications. Linear sensors range from 2000-element monochrome arrays with 40 dB of dynamic range used in facsimile applications to 8000 or more element trilinear arrays with 80 dB of dynamic range for high-performance color scanning applications.

The steps involved in image sensing consist of (1) converting the incoming photons to charge at picture element (pixel), and (2) transferring that charge to an output amplifier and converting the charge to a voltage or current signal which can be sensed by circuits external to the sensor. The image-sensing elements are described first, followed by readout elements. Sensor architectures for area and linear sensors are described next.

32.3 IMAGE SENSING ELEMENTS

There are four basic types of structures which are used for image sensing: the junction photodiode, the photocapacitor, the pinned (p^+np) photodiode, and the photoconductor.* The first three are generally fabricated in single-crystal silicon as part of the image sensor; the photoconductor is usually fabricated from amorphous silicon deposited over the image sensor. The photoconversion process begins with the absorption of a photon in silicon resulting in the generation of a single electron-hole pair. The absorption of light at a particular wavelength is given by

$$N(x, \lambda) = N(0)e^{-\alpha(\lambda)x} \quad (1)$$

where $N(x)$ is the number of photons remaining at a distance x below the surface, $N(0)$ is the number of photons entering, and $\alpha(\lambda)$ is the absorption coefficient.^{1,2} The absorption coefficient is shown in Fig. 1 as a function of wavelength λ for single-crystal silicon, doped polycrystalline silicon, and hydrogenated amorphous silicon. The absorption depth is defined as the inverse of the absorption coefficient [$d(\lambda) = 1/\alpha(\lambda)$]. In single-crystal silicon, the absorption depth is 0.4 μm in the blue (450 nm), 1.5 μm

*For some scientific applications in which high quantum efficiency and fill factor are essential, the silicon wafer is thinned to 10 μm or less in thickness and illuminated from the backside. The frontside contains an area charge coupled device, which is used to collect the photogenerated carriers.

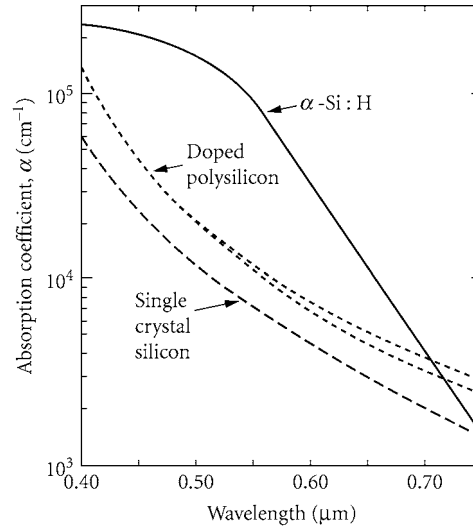


FIGURE 1 Absorption coefficient for light in single-crystal silicon, heavily doped polycrystalline silicon, and hydrogenated amorphous silicon as a function of wavelength.

in the green (550 nm), and 3.0 μm in the red (640 nm). In the infrared, the absorption depth increases to 10.5 μm at 800 nm. Beyond 1100 nm, the absorption is virtually zero because the photon energy is less than the 1.1-eV silicon bandgap.

Junction Photodiode

The junction photodiode is one of the most common image-sensing elements. The physical structure and band diagram of the junction photodiode are shown in Fig. 2a for a *p*-type substrate. The *n*-type region is formed by ion implantation or diffusion of phosphorous or arsenic to a depth of 2000 to 10,000 \AA into the *p*-type silicon. The *n*-type dopant region is usually graded, with the highest concentration at the surface. The gradient in *n*-type dopant concentration results in a gradient in electrostatic potential which accelerates photogenerated carriers (holes in the *n*-type region) away from the surface. This reduces loss of photogenerated carriers to surface recombination. The photodiode is typically operated with a reverse bias V of 1 to 5 V. A depletion layer is formed between the *n*- and *p*-type regions.* The width of the depletion layer $W(V)$ for an abrupt *n + p* junction is given by

$$W(V) = \sqrt{\frac{2\epsilon_s(V + V_{bi})}{qN_A}} \quad (2)$$

where q is the electronic charge, ϵ_s is the silicon dielectric constant, N_A is the *p*-type dopant concentration, and V_{bi} is the built-in voltage given by

$$V_{bi} = \frac{kT}{q} \ln \frac{N_A}{n_i} \quad (3)$$

*For a detailed review of the device physics of junction diodes and MOS capacitors, see Sze, *Physics of Semiconductor Devices*, Wiley, New York, 1969.

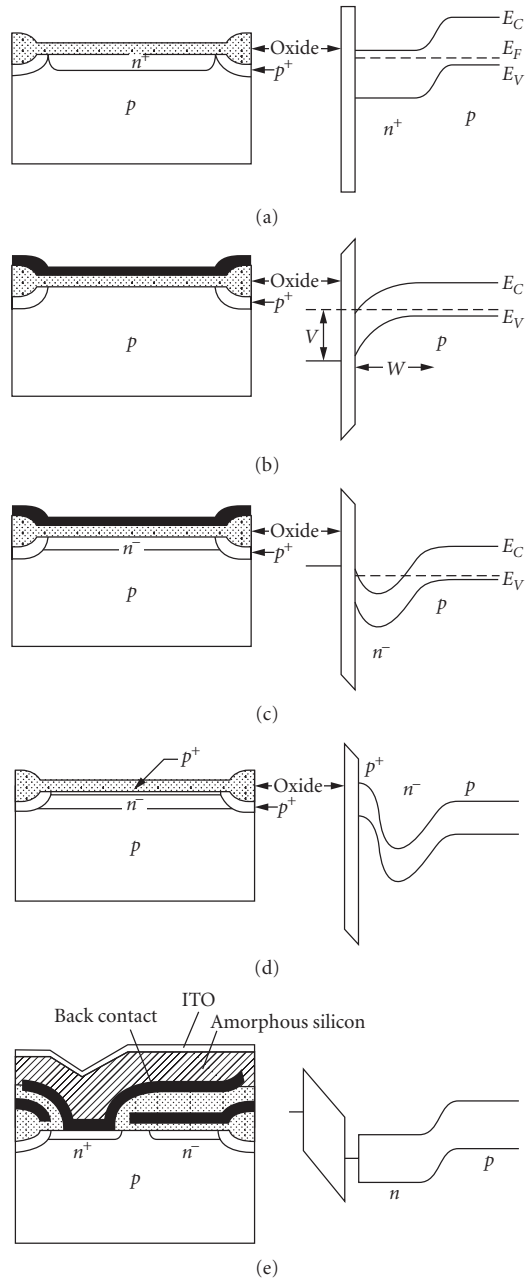


FIGURE 2 Cross-sectional diagrams and band diagrams for (a) junction photodiode; (b) surface-channel MOS capacitor; (c) buried-channel MOS capacitor; (d) pinned or hole-accumulated diode; and (e) amorphous silicon photoconductor.

where n_i is the intrinsic carrier concentration. For silicon doped at $1 \times 10^{16} \text{ cm}^{-3}$, the depletion width at 5-V reverse bias is 0.8 μm .

If the diode is illuminated, some of the photons will be absorbed in the n -region, some in the depletion layer, and the remainder in the p -type substrate. The quantum efficiency $\eta(\lambda)$ is the ratio of the charge collected to the number of photons incident on the diode (i.e., 100-percent quantum efficiency refers to one electron-hole pair collected for every incident photon). The quantum efficiency depends on three factors: transmission $T(\lambda)$ of light through the overlying layers into silicon, absorption of light in silicon, and the probability $P(x)$ that an electron-hole pair generated a distance x from the surface will be collected before recombination:

$$\eta(\lambda) = T(\lambda) \int_{x=0}^{x=\infty} (1 - e^{-\alpha(\lambda)x}) P(x) dx \quad (4)$$

The transmission of light through the overlying layers into the silicon can be calculated using standard multilayer interference models.

The collection of the photogenerated charge takes place by two processes: drift and diffusion. Drift is the movement of electrons and holes due to an electric field. Even for small electric fields, the transport of carriers by drift will dominate diffusion. This is the case in the depletion region. Outside the depletion region, such as in the p -type substrate, there is no electric field, and carrier transport occurs by diffusion. For example, a photon absorbed in the p -type region will excite an electron from the valence band to the conduction band. The electron will move in a three-dimensional random walk until it recombines or encounters the edge of the depletion region, where it is swept across the junction by the electric field. The probability that an electron at a distance x' from the junction can diffuse to the junction before recombining is given by

$$P(x') = e^{-x'/L_e} \quad \text{where} \quad L_e = \sqrt{\frac{KT}{q}} \mu_e \tau_e \quad (5)$$

in which L_e is the diffusion length, μ_e is the electron mobility, and τ_e is the electron lifetime (typically $\sim 1 \mu\text{s}$). For a 1- μs lifetime, the electron diffusion length in p -type silicon is 50 μm . Electrons generated from photons absorbed at less than the diffusion length from the junction have a high probability of collection.*† Similarly, photons absorbed in the n -type layer create electron-hole pairs. The holes must travel by diffusion to the junction in order to be collected. The probability that a hole a distance x' from the junction can diffuse to the junction is given by

$$P(x') = e^{-x'/L_p} \quad \text{where} \quad L_p = \sqrt{\frac{KT}{q}} \mu_p \tau_p \quad (6)$$

in which L_p is the diffusion length, μ_p is the hole mobility, and τ_p is the hole lifetime. For a 1- μs lifetime, the hole diffusion length in n -type silicon is 30 μm . Since the n -type region in an $n + p$ junction is less than 1 μm thick, the holes have no difficulty diffusing through the n -type region to the junction unless the n -type region is so heavily damaged or so heavily doped that the hole lifetime is

*In image sensors, the collection of photogenerated carriers can be complicated by a variety of factors. The doping concentration may not be uniform on either the n - or p -sides. On the n -side, the dopant concentration is designed to decrease from the surface to the junction, building in a potential gradient for holes away from the surface and toward the junction, preventing surface recombination. On the p -side, the dopant concentration may not be uniform owing to the use of epitaxial layers or wells diffused into silicon. Additionally, the carrier lifetime may not be uniform. Impurity gettering, a process used in many image sensors to remove metallic contaminants will leave a region of crystalline defects in the silicon starting 20 to 50 μm beneath the silicon surface. The defects result in a very short electron lifetime in this region. Finally, diffusion takes place in three dimensions; a carrier absorbed beneath a given pixel in an array will diffuse laterally by the same amount it diffuses vertically. This can cause it to be collected in adjacent pixels.

†See, for example, Lavine et al., "Steady State Photocarrier Collection in Silicon Imaging Devices," *IEEE Transactions on Electron Devices* ED-30:1123–1133 (Sept. 1983).

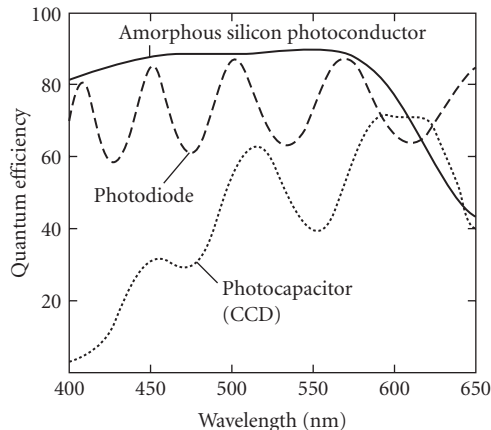


FIGURE 3 Quantum efficiency as a function of wavelength for photocapacitor, photodiode, and amorphous silicon photoconductor.

very short. More typically, loss of quantum efficiency on the n -side results from recombination at the surface.

The quantum efficiency of a junction photodiode is illustrated as a function of wavelength in Fig. 3. In the ultraviolet, the light is absorbed very near the surface and some of the photo-generated charge can be lost to surface recombination. In the 420- to 700-nm region, most of the light is absorbed close to the junction and is easily collected. The structure in the quantum efficiency illustrated in Fig. 3 results from the multilayer reflections of light in the oxide layer which overlies the photodiodes used in image sensors. The positions of the minima and maxima depend on the thicknesses and indices of refraction of the layers overlying the silicon. Beyond 800 nm, some of the photons are absorbed sufficiently deep in silicon that the carriers recombine before they can diffuse to the junction; this results in a decrease in quantum efficiency at longer wavelengths.

In most image sensing applications, the junction diode at each picture element is used not only to collect the photogenerated carriers but also to store the carriers until they can be read out. In an imaging array, each photodiode would be reset to a reverse bias V_r by a MOS gate. The capacitance of the diode of area A for an abrupt $n + p$ junction is given by $C(V) = \epsilon_s A / [W(V)]$ where $W(V)$ is the depletion width.

When a photogenerated carrier is collected by the junction, it is stored on the junction until it is read out. Storage of a carrier will cause the voltage on the junction to decrease by $q/C(V)$. When the photogenerated charge is removed from the junction during readout, the junction voltage is restored to its original value.

If so much photogenerated charge is stored on the photodiode that the voltage drops to zero before the charge can be read out, additional charge cannot be collected and will diffuse into the p -type substrate. This condition is referred to as *saturation*. The diffusion of excess charge into neighboring photosites is called *blooming*.

One of the difficulties encountered in using the junction photodiode in image sensor applications is *image lag*. The combination of the capacitance of the photodiode and the channel resistance of the MOS transfer gate used to read out the diode give rise to a time constant for transferring the photogenerated charge from the diode to the readout structure. As a result, not all the charge can be completely drained from the diode during the short reset times typically used in imaging applications. The remaining charge is drained in successive readouts, causing an afterimage. This effect is called *image lag*.³

MOS Capacitor

The MOS capacitor consists of the silicon substrate (taken to be p -type in this section), a thin layer of silicon dioxide (typically 200 to 1000 Å thick), and an electrode (typically polycrystalline silicon doped heavily n -type with phosphorous). The physical structure and the band diagrams of the surface channel MOS capacitor are shown in Fig. 2*b* for a p -type substrate. If the gate of the MOS capacitor is biased positive, a depletion layer is created in the p -type silicon substrate. The depth of the depletion layer depends on the substrate doping, gate voltage, and oxide thickness. The calculation of the depth of the depletion layer is somewhat more complex than the photodiode* and depends on the electrostatic potential at the surface, called the surface potential Φ_s . For values typical of image sensors ($N_A = 1 \times 10^{15}$, 5-V gate bias, 500-Å oxide thickness) the depletion layer is 2.4 μm deep. On the edges of the photocapacitor (see Fig. 2*b*) is a heavily doped p -type region overlaid by a thick (2000- to 5000-Å) oxide layer. This is called the field or *channel stop* region. Because of the heavier p -type doping, the field is not depleted by the voltage on the gate. The channel stops confine the electrons to the channel region.

If the MOS capacitor is illuminated, a fraction of the light will be reflected, a fraction will be absorbed in the polysilicon, and the rest will be transmitted into the silicon substrate. The absorption coefficient of heavily doped polysilicon is shown as a function of wavelength in Fig. 1. The absorption coefficient is $4 \times 10^4 \text{ cm}^{-1}$ at 450 nm and $1.2 \times 10^4 \text{ cm}^{-1}$ at 550 nm. For a 3000-Å-thick polysilicon electrode, less than 30 percent of the blue light and less than 70 percent of the green light is transmitted through the polysilicon. The photons entering the silicon are absorbed, either in the depletion layer of the MOS capacitor or in the undepleted p -type silicon beneath the depletion layer. Those photogenerated electrons created in the undepleted p -type region move by diffusion until they are captured by the depletion layer or they recombine. The electrons are held at the silicon-SiO₂ interface, where they remain until they are read out. The quantum efficiency as a function of wavelength is illustrated in Fig. 3. The efficiency is low in the blue owing to absorption in the polysilicon. The structure in the quantum efficiency as a function of wavelength is due to multilayer interference in the polysilicon-oxide-silicon stack.

Charge is stored in the MOS capacitor at the silicon-silicon dioxide interface as a layer of sheet-charge only a few hundred angstroms thick. As more photogenerated charge is added, the surface potential decreases. If sufficient photogenerated charge is added, the surface potential becomes zero and no additional charge can be stored. This condition is saturation.

The capacitance per unit area on which photogenerated charge is stored is the parallel capacitance of the oxide and the depletion layer:

$$C^{-1} = \left(\frac{t_{\text{ox}}}{\epsilon_{\text{ox}}} + \frac{W(\Phi_s)}{\epsilon_{\text{si}}} \right) \quad (7)$$

where the depletion width $W(\Phi_s)$ is give by

$$W(\Phi_s) = \sqrt{\frac{2\epsilon_{\text{si}}\Phi_s}{qN_A}} \quad (8)$$

in nearly all cases, the oxide capacitance is the dominant term. As a result, the storage capacity of the MOS capacitor is significantly larger than the junction diode in which the charge is stored only on the depletion capacitance.

A variant of the surface-channel photocapacitor is the buried-channel photocapacitor. The structure and band diagram are shown in Fig. 2*c*. In this device, a lightly-doped n -type region is diffused or implanted into the silicon surface early in the fabrication process. This n -type region is sufficiently lightly doped that it is fully depleted. The n -type dopant in the buried channel results in

*To determine the depletion depth, the surface potential Φ_s must be calculated which depends on the voltage on the gate of the MOS capacitor, the oxide thickness, substrate doping, and weakly on temperature. See reference on p. 32.3.

a band diagram with a potential minimum, or well, for electrons just below the surface.* This well is separated from the surface by a few thousand angstroms in distance and about 1 V in potential. When the buried-channel photocapacitor is illuminated, the electrons collect in the buried channel and do not contact the surface. The primary purpose of the buried channel is to prevent electrons from being trapped by interface states at the silicon-silicon dioxide interface.

The photocapacitor has several advantages over the junction photodiode. These include higher storage capacity per unit area, zero lag readout, and generally lower dark current. The principal disadvantage is the low quantum efficiency in the blue. In some applications, a transparent electrode, such as indium-tin-oxide (ITO) may be substituted to improve the blue response.⁴ ITO has very low absorption over the visible (420 to 750 nm) and can be deposited sufficiently conductive for use as a gate electrode in an image sensor.

Another approach to achieving high quantum efficiency is the thinned backside-illuminated charge-coupled device (CCD). In this approach, the CCD (which is an array of MOS capacitors) is fabricated on the frontside of a silicon wafer. The wafer is then thinned from the backside to 10 μm or less in thickness. The backside is passivated to prevent surface recombination. Photons entering the backside are absorbed in silicon beneath the MOS capacitors (usually buried channel). The photogenerated carriers diffuse to the capacitors, where they are held until they are read out. This device has quantum efficiency similar to the photodiode in the visible. However, because the silicon is thin, some of the photons at wavelengths beyond 700 nm will not be absorbed and so the quantum efficiency falls off beyond 700 nm. Owing to their extreme complexity in process and their extremely fragile design, backside-illuminated image sensors are limited to special scientific applications, especially astronomy.

Pinned Photodiode

The third type of photosensitive element is the pinned (p^+np) photodiode.⁵ This is sometimes called the *hole accumulation diode*, or HAD. This element combines the best features of the photodiode and photocapacitor, offering the high blue response of the photodiode with the high charge capacity, zero lag, and low dark current of the buried-channel photocapacitor. The pinned photodiode consists of a very shallow (<2000 Å) P^+ layer overlying an n -type buried-channel region. The structure and band diagrams are shown in Fig. 2d. The p^+ surface layer, which contacts the p^+ channel stop region on the sides, holds the electrostatic potential at the surface at 0 V. When the photodiode is illuminated, the photogenerated electrons are held in the n -type buried-channel region just below the surface.

The quantum efficiency of the pinned photodiode is nearly identical to that of the photodiode shown in Fig. 2. Because there is no overlying polysilicon electrode, the blue response is very high, similar to that of the photodiode. Because the buried-channel region can be completely emptied, the pinned photodiode does not have the lag of a normal junction (n^+p or p^+n) photodiode. The pinned photodiode has been the most widely used image-sensing element in interline area CCDs used for camcorders and for industrial and medical cameras. The pinned photodiode is also used in some linear image sensors, particularly where low image lag is critical.

Photoconductor

The last type of photosensitive device is the photoconductor. The most common material for the photoconductor is hydrogenated amorphous silicon, although other material systems have been explored. Amorphous silicon photoconductors have been used for two types of image sensors: area image

*For a review of the calculation of electrostatic potential and charge capacity of buried-channel MOS capacitors and charge-coupled devices, see B. C. Burke, G. Lubberts, E. A. Trabka, and T. J. Tredwell, *IEEE Transactions on Electron Devices* ED-31(4):423 (April 1984).

sensors, where it is deposited on top of an area array to improve fill factor (i.e., the proportion of the picture element which is photosensitive), and contact linear sensors, where it is deposited on large ceramic or glass substrates to fabricate very long line or linear arrays.

The structure of an amorphous silicon photoconductor on a CCD image sensor and the corresponding band diagrams are shown in Fig. 2e. The hydrogenated amorphous silicon photoconductor consists of a back electrode, an undoped amorphous silicon layer approximately 1 μm thick, and a transparent top electrode.⁶ Additional doped amorphous silicon or silicon nitride layers may be added to the amorphous silicon front or back surfaces to improve performance. Photons absorbed in the amorphous silicon generate electron-hole pairs. Photogenerated electrons and holes are quickly swept to the back and front electrodes, respectively, because of the high electric field in the photoconductor. When the amorphous silicon is used as part of an area image sensor, the electrons on the back electrode can be transferred into the readout element; when used as part of a contact linear array, the voltage change can be amplified and read out through a multiplexer.

The quantum efficiency of an amorphous silicon photoconductor is shown in Fig. 3. Owing to the wider bandgap of the amorphous silicon, photons of wavelength greater than about 650 nm are not absorbed. The advantage of the amorphous silicon is the high quantum efficiency across the visible wavelengths and the ability to fabricate devices either on top of area CCDs for higher fill factor or to process on glass or ceramic for very large linear sensors. The disadvantages include charge trapping at defects in the amorphous silicon and low carrier mobility, both of which lead to field-dependent nonlinear response and image lag when used in an image sensor. Improvements in material and device technology have mitigated many of these disadvantages at the expense of process and device complexity.⁷

Antiblooming in Charge-Sensing Elements

Blooming in image sensors occurs when the charge generated in an image-sensing element exceeds its capacity. If no method is provided to remove this excess charge, it will be injected either onto the readout element (CCD or MOS readout line) or into the substrate. If the excess charge is injected onto the readout element, it will usually appear as a bright line in the image. If it is injected into the substrate, the charge can diffuse in a circular pattern and be collected by neighboring elements.

There are two basic types of antiblooming circuits: lateral and vertical.^{8,*} In lateral antiblooming, illustrated in Fig. 4a and b, a MOS antiblooming gate and an antiblooming drain are provided adjacent to each image-sensing element. Excess charge on the sensing element overflows the antiblooming gate onto the antiblooming drain. The antiblooming drains of all elements on the array are connected and the current sunk in a bias supply.

In a vertical antiblooming structure, illustrated in Fig. 4c and d for a pinned photodiode sensing element, the image-sensing element is fabricated in a shallow, lightly doped *p*-well. A large (10- to 30-V) bias is applied to the *n*-type silicon substrate, causing the *p*-well underneath the photodiode to become completely depleted. Once the charge on the diode exceeds its capacity, the excess charge flows over the saddle-point in the *p*-well, under the diode, and into the substrate. The substrate is connected to a bias supply which sinks the blooming current. The vertical antiblooming structure has the advantage of requiring no additional silicon area, so that antiblooming is achieved with no reduction in fill factor. Lateral antiblooming, on the other hand, requires additional silicon area in each pixel, thereby reducing fill factor. The vertical antiblooming has the disadvantage of lower quantum efficiency at wavelengths longer than about 500 nm. Because any light absorbed below the photodiode or photocapacitor is drained into the substrate, the quantum efficiency of photodiodes or photocapacitors with vertical antiblooming is reduced.

*Other types of antiblooming are occasionally used in image-sensing arrays. One of these is charge pumping, in which a MOS gate is repeatedly clocked in order to cause excess charge held underneath it to recombine at interface states at the silicon-silicon dioxide interface. In charge pumping, the MOS gate is pulsed sufficiently negative to cause accumulation, resulting in the interface states filling with holes. When the gate is pulsed out of accumulation, excess electrons can recombine with the holes.

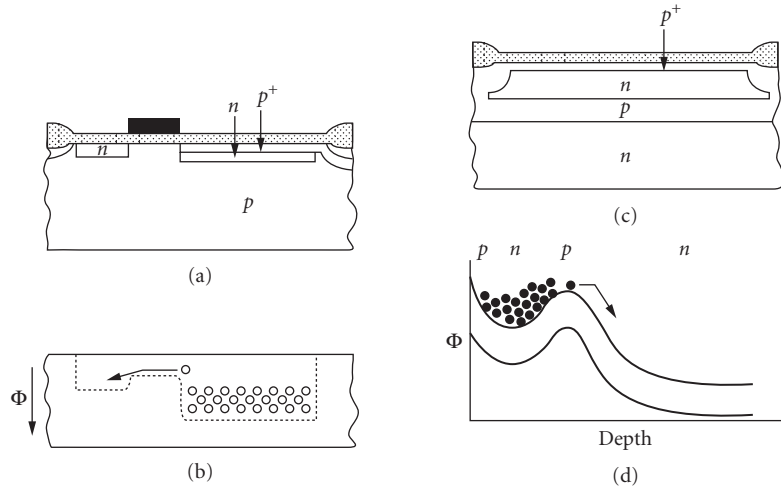


FIGURE 4 Antiblooming methods for image sensors: (a) cross section of lateral antiblooming structure; (b) illustration of electrostatic potential and charge overflow in lateral antiblooming; (c) cross section of vertical antiblooming; and (d) band diagram and illustration of charge overflow in vertical antiblooming.

Dark Current in Photosensing Elements

Signal in photosensing elements is the result of collection of electron-hole pairs generated by the absorption of light. However, charge is generated at each photosensing element even in the absence of light. This generation is called dark current and is a result of thermal generation of electron-hole pairs. The thermal generation occurs at defects, such as impurities or crystalline defects, in the bulk of silicon and at surface states at the silicon-silicon dioxide interface.

There are four sources of dark current: (1) diffusion current, which is the thermal generation of carriers in the undepleted n - and p -type regions; (2) depletion layer generation current, which occurs in the depletion layer of a diode or MOS capacitor; (3) surface generation current, which is the generation of electron-hole pairs at interface states at the Si-SiO₂ interface; and (4) leakage, which refers to generation at extended defects such as impurity clusters or stacking faults, particularly in the presence of a large electric field. These sources of dark current are illustrated for a photodiode and a photocapacitor in Fig. 5. The generation of the electron-hole pairs in both diffusion current and depletion-layer generation-recombination current occurs almost exclusively at impurity sites. Impurities with energy levels near midgap, such as gold, copper, and iron, are particularly effective in the thermal generation of charge. The depletion-layer generation current is given by*

$$J_g(V) = \frac{qn_i W(V)}{\tau_o} \quad (9)$$

where $W(V)$ is the width of the depletion layer at a given bias voltage in the MOS capacitor or junction-photodiode, n_i is the intrinsic carrier concentration, and τ_o is the depletion-layer lifetime, for carriers.

*The expressions here are for the generation current. The total current is the sum of the generation and recombination current given by

$$J_{gr}(V) = (qn_i W(V) / \tau_o) (e^{qV/2kT} - 1)$$

However, for diodes under 0.2 V or more of reverse bias ($qV/kT > 8$), such as would be the case in an image sensor, the recombination current is negligible.

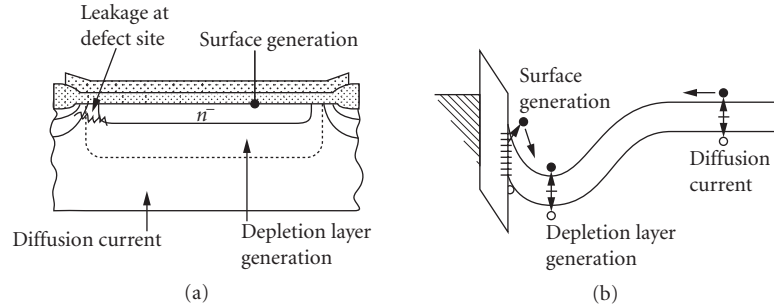


FIGURE 5 (a) Cross section of buried-channel MOS capacitor and (b) band diagram of buried-channel MOS capacitor illustrating mechanisms for dark current generation.

Typical values of τ_0 in clean MOS processes would be 1 to 10 ms and values of J_{gr} would be 30 to 300 pA/cm².

The diffusion generation current for a *p*-type region is given by*

$$J_{\text{diff}} = \frac{qn_i^2 L_e}{N_a \tau_e} \quad (10)$$

Since the intrinsic carrier concentration n_i depends on temperature as $e^{-E_g/2kT}$, depletion-layer generation current and diffusion current will have different temperature dependences. Depletion-layer generation current will increase as $e^{-E_g/2kT}$, which corresponds to a doubling in dark current for every 9 to 11°C increase in temperature near room temperature. Diffusion current will increase as $e^{-E_g/2kT}$, which corresponds to a doubling every 4.5 to 5.5°C increase in temperature.

The surface generation current J_s occurs almost exclusively at regions where the depletion layer intersects the Si-SiO₂ interface, such as the surface region between the *n*⁺- and *p*-regions around a photodiode or in the depleted surface under a MOS capacitor. It is given by

$$J_s = \frac{qn_i s_o}{2} \quad (11)$$

where s_o is the surface recombination velocity. A typical value of J_s for a MOS surface would be 100 pA/cm².^{9†} The surface generation depends on temperature in the same manner as the depletion-layer current. In very clean MOS processes (i.e., low concentration of metallic impurities in silicon), the surface current is often the largest component of the overall dark current.

Leakage current occurs at extended defects in silicon, such as impurity clusters and stacking faults, particularly when these defects are in a depletion layer and so are subject to a high electric field. While there is no single analytical expression for the current generated by an extended defect,

*The full expression for the diffusion current from the undepleted *p*-region in a diode is

$$J_{\text{diff}} = \frac{qn_i^2 L_e}{N_a \tau_e} (e^{qV/kt} - 1)$$

However, for reverse biases more than 0.2 V ($qV/kt > 8$), as would be the case in a photodiode in an image sensor, only the generation term is important.

†In calculating the total surface current, the current density J_s is multiplied by the area of the depleted surface. For a *pn* junction diode, this would be the area of depleted surface region around the diode between the *n*- and *p*-regions (i.e., approximately the junction perimeter times the depletion-layer width); for a MOS capacitor it would be the entire area under the MOS capacitor unless sufficient sheet charge of electrons had formed at the surface to invert the surface. The surface recombination velocity is given by $S_o = N_{st} v_{th} \sqrt{\sigma_n \sigma_p}$ where N_{st} is the density of surface states near midgap, v_{th} is the thermal velocity (2.0×10^7 cm/s) and σ_n and σ_p are the electron and hole cross sections for midgap surface states. See reference on p. 32.3 for a complete explanation.

they are characterized by very high values of dark current ($\gg 1$ nA/cm²), a very strong dependence on the electric field, and, in regions of high electric field, only a very slight dependence on temperature. In an image sensor, they are visible as “bright spots” in a few isolated pixels against the otherwise low level of background thermal generation.¹⁰

Values for the dark current vary widely because of variation in the amount of impurities in the silicon; the dark current can range from 0.01 nA/cm² in very high quality image sensors to >10 nA/cm² in sensors with significant metallic contamination. The total number of electrons N_e collected in a pixel is

$$N_e = JAT_{\text{int}}/q \quad (12)$$

where J is the total dark current per unit area, A is the area of the pixel, and T_{int} is the integration time. For a 1/2-in format CCD such as would be used in a camcorder, typical values would be a dark current of 0.5 nA/cm², a pixel area of 100 μm^2 , and an integration time of 1/30 second; the number of thermally generated electrons would be 105 electrons per pixel. Image sensors developed for scientific purposes might have a dark current 5 to 10 times lower at room temperature. These same devices might also be operated below room temperature to reduce the thermal generation to levels below one electron per pixel.

There are two types of noise associated with the charge generated by the dark current: shot noise and pattern noise. The shot noise due to the dark current is the square root of the number of dark electrons in a pixel. Pattern noise is due to pixel-to-pixel variations in the dark current and is often highly correlated between neighboring pixels. A numerical value for the dark pattern noise is typically obtained by using the standard deviation of the pixel values in the dark from a large region of an imager, where the values are obtained by averaging over many frames to eliminate shot noise and other temporal noise sources.*

32.4 READOUT ELEMENTS

The readout element transfers the charge from the image-sensing element (photodiode, photocapacitor, or photoconductor) to the output of the image sensor. In a linear sensor, the readout is only in one direction. In an area sensor, both x and y readout is required. There are two basic types of readout elements: charge-coupled devices, or CCDs, and x - y addressed photodiode arrays (typically called MOS arrays owing to the MOS transistors used in the addressing of the pixels). CCDs are by far the most widely used readout elements owing to their very low noise. Nearly all consumer camcorders, facsimile machines, scanners, copiers, and professional and scientific cameras utilize CCDs. Applications of MOS arrays are largely restricted to those where addressing of an individual pixel or subarray is required.

Charge-Coupled Device (CCD)

CCD Operation The CCD works by moving packets of charge physically at or near surface of silicon from the image-sensing element to an output, where the charge packet is converted into a voltage. The CCDs is formed by an array of overlapping MOS capacitors.^{11–13} There are a number of different types of CCDs, including four-phase, three-phase, two-phase, virtual (single-) phase, and ripple-clocked CCDs. The number of phases refers to the number of separately clocked elements

*A histogram of the dark current values of the pixels is rarely gaussian. In most cases, it exhibits an extended tail at high dark current values due to pixels with crystalline defects. The histogram may also exhibit quantization due to pixels with integral number of impurities (i.e., 1, 2, or 3 gold atoms); the quantization may even be used as a signature of the impurity present. See, for example, McColgin et al., “Effects of Deliberate Metal Contamination on CCD Image Sensors,” *Materials Research Society Symposium Proceedings*, 262: 769 (1992) and “Dark Current Quantization in CCD Image Sensors,” *Proc. 1992 International Electron Device Meeting*, Washington, D.C., 113–116 (1992).

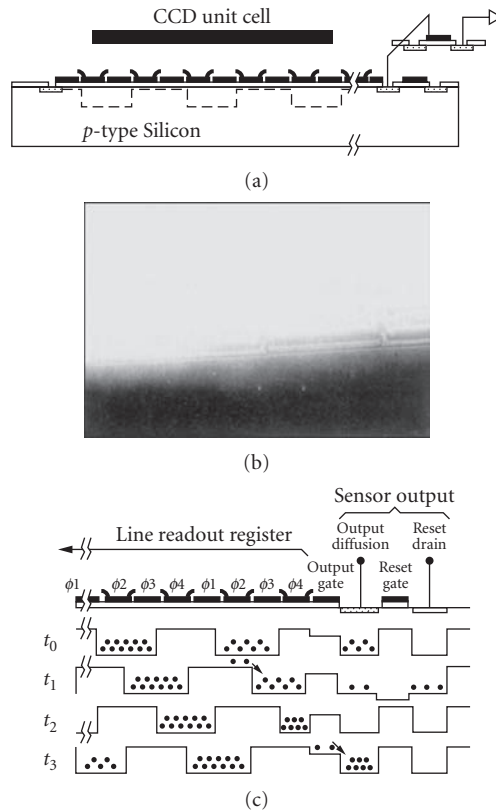


FIGURE 6 (a) Cross-sectional diagram of unit cell of a four-phase CCD; (b) scanning electron micrograph of a unit cell of a four-phase CCD; and (c) illustration of charge transfer along a four-phase CCD, showing both transfer along the register and at the sensor output.

within a single stage of the CCD. These are described later. For understanding the principle of charge transfer in a CCD, consider the four-phase CCD illustrated in Fig. 6a.

The four-phase CCD physically comprises a silicon substrate (assumed to be *p*-type for purposes of illustration), a gate oxide 300 to 1000 Å thick, and overlapping polysilicon electrodes 1000 to 5000 Å thick which have been heavily doped with phosphorus to lower their resistance. For a four-phase CCD, two levels of electrode are required. The first level is deposited, then defined photolithographically to form two phases (ϕ_1 and ϕ_3). A thin (500-Å) oxide is grown over the first level of polysilicon to insulate it from the second level. The second level of polysilicon is then deposited, doped with phosphorus, and defined to form the other two phases (ϕ_2 and ϕ_4). An electron micrograph of a four-phase CCD with six-micron long gates is shown in Fig. 6b.

The process of charge transfer in a four-phase CCD is illustrated in Fig. 6c. In order to hold a packet of electrons, two adjacent gates (ϕ_2 and ϕ_3 , for example) would be held at a high positive potential ($\sim +5$ V) while the other two phases would be held at a low potential (~ 0 V). A depletion layer, or well, is formed under ϕ_2 and ϕ_3 , allowing electrons to be held at or below the surface. The other two phases, ϕ_1 and ϕ_4 , serve as potential barriers, keeping the charge packet under ϕ_2 and ϕ_3 . To transfer the electrons through silicon, the electrode ahead of the charge packet (ϕ_4) is clocked

positive and the electrode behind (ϕ_2) is clocked negative. The electrons move along the silicon surface following the positive potential. This is repeated through all four phases, during which the charge packet is moved forward one pixel.

The CCD may be either surface-channel or buried-channel. In a surface-channel CCD (Figs. 2b and 6a) the electrons are held at the silicon surface. Surface-channel CCDs are rarely used owing to the trapping of electrons at interface states at the silicon surface. At the silicon-SiO₂ interface there is a density of states of about $1 \times 10^{10} \text{cm}^{-2} \text{eV}^{-1}$. These states can trap electrons from one charge packet and reemit the electrons into a later charge packet. This results in transfer inefficiency. In the buried-channel CCD (Figs. 2c and 6d), a lightly doped *n*-type layer is formed in silicon. This *n*-type layer results in a potential well for electrons just below the surface rather than at the surface (Fig. 2b). As a result, the electrons remain separated from the interface and are not trapped by interface states. This results in the ability to transfer charge from one stage to another with very high efficiency. Virtually all CCDs for image-sensing applications utilize buried-channel CCDs. The clock voltages used to drive CCD gates typically swing by 5 to 8 V between high and low levels.

For buried-channel CCDs, transfer rates of up to 20 MPix/s can usually be achieved without special design considerations. Above this rate, special attention must be paid to optimizing the electric field along the direction of transfer to speed up charge transfer. Transfer rates exceeding 50 MPix/s have been achieved in optimized CCD designs.¹⁴

CCD Output At the output of the CCD is a circuit to convert the charge packets into a voltage signal. By far the most common type of output circuit is the floating diffusion with source follower amplifier. The floating diffusion output is shown schematically in Fig. 7a and a photograph of the end of a CCD shift register with the first stage of the amplifier is shown in Fig. 7b. The floating diffusion output consists of an output gate (OG), a floating diffusion, a reset gate (RG), and a reset drain. The floating diffusion is an *n*⁺-type region between the output gate and the reset gate. The floating diffusion is connected to the gate of a source-follower amplifier. The output gate is held at a fixed dc voltage, creating a barrier potential over which the packet of electrons can be transferred onto the floating diffusion when the last phase of the CCD register is clocked to its low (~ 0 V) voltage.

The sequence of events in the CCD output is shown in Fig. 6c. As the charge packet is transferred along the CCD, it arrives at the last phase (ϕ_4) before the output gate (time t_2 in Fig. 6c). When ϕ_4 is clocked low (time t_3), the packet of electrons is transferred over the output gate onto the floating diffusion (time t_0). The voltage of the floating diffusion changes by an amount

$$V = Nq/C \quad (13)$$

where N is the number of electrons and C is the total capacitance of the floating diffusion itself, the interconnect to the source-follower amplifier and the input capacitance of the amplifier. In typical CCDs, this capacitance would be in the 10- to 50-fF range. Because this capacitance is so small, there is a large change in voltage for a small change in charge. The charge-to-voltage conversion is an important parameter in CCD design; for a 10-fF capacitance the charge-to-voltage conversion is 16 $\mu\text{V}/\text{electron}$. After the change in voltage has been sensed by the amplifier, the charge packet must be removed from the floating diffusion before the next packet arrives. This is achieved by clocking the reset gate positive (time t_1 in Fig. 6c), allowing the charge packet to flow from the floating diffusion to the reset drain. The reset drain is held at a constant positive voltage, typically ~ 10 V. The reset drain is then turned off and the floating diffusion is prepared to accept the next packet of electrons.

The change in voltage from the floating diffusion is typically buffered on-chip by a source-follower (Fig. 7c). A two-stage source follower is most often used. The first stage utilizes very small-dimensioned FET transistors in order to minimize the input capacitance. The second uses much larger FET transistors in order to achieve sufficient drive current to overcome off-chip capacitances such as package and lead capacitances on the circuit board. The source-follower amplifier is typically designed to have a bandwidth on the order of ten times the CCD pixel rate. For high-pixel-rate applications (> 10 MPix/s), three-stage source-follower amplifiers are employed.

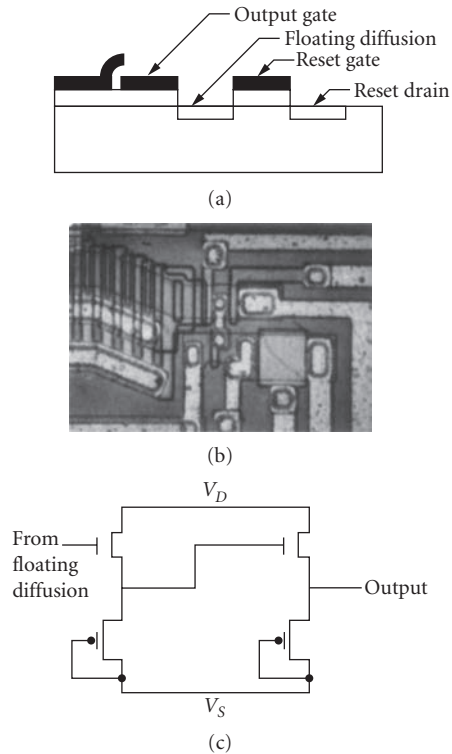


FIGURE 7 (a) Cross section of the floating diffusion output for a CCD; (b) photomicrograph of a floating diffusion output including the first stage of the source-follower amplifier; and (c) two-stage source-follower amplifier with on-chip loads.

For special purpose applications, such as CCD signal processing, nondestructive readout is required. Floating gate output amplifiers are used in these applications. These amplifiers are similar to the floating diffusion except that the floating diffusion is replaced by a MOS gate which is connected to the MOS amplifier. Other outputs, such as buried-channel JFET structures, have seen limited use in very low noise applications.¹⁵

Types of CCDs In addition to the four-phase CCD described here, there are a number of other types of CCD shift registers, including three-phase, two-phase, and virtual phase. The different types are illustrated in Fig. 8. The four-phase CCD (Fig. 8a) was described previously. The three-phase (Fig. 8b) consists of three different layers of polysilicon electrodes. The charge is normally held under one of the three; during transfer, the gate ahead of the charge packet is clocked positive and the gate holding the charge is clocked negative in order to transfer the charge one gate ahead. The three-phase CCD has the advantage of a shorter unit cell than the four-phase but at the expense of additional processing complexity (i.e., a third polysilicon layer).

In the two-phase CCD, each phase has a barrier and a well region. The barrier is formed by doping the barrier region slightly less n -type than the well region, making the electrostatic potential in the barrier region a few volts lower than the well region for the same gate voltage. Electrons will flow over the barrier region and be held in the well region. There are two methods of fabricating a two-phase CCD. In the first method (Fig. 8c), two separate gates are used for each phase; one gate receives

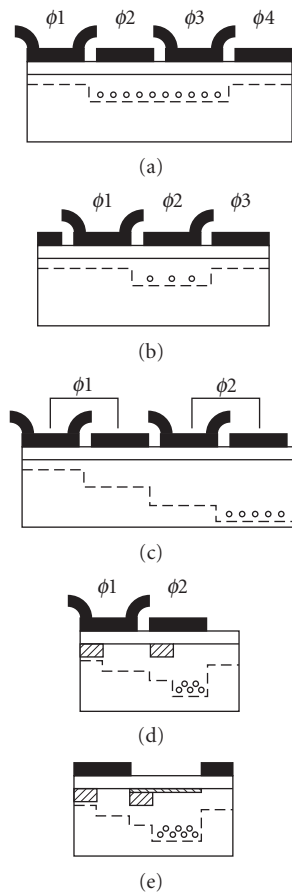


FIGURE 8 Types of CCD registers: (a) four-phase; (b) three-phase; (c) pseudo-two-phase; (d) true two-phase; and (e) virtual phase. In the pseudo-two-phase CCD, each phase consists of two polysilicon gates, one of which is offset in potential from the other by an implanted threshold adjustment. In the true two-phase CCD, each phase consists of a single polysilicon gate in which the implanted threshold-adjust region is formed underneath a portion of the gate.

a threshold-adjust ion implantation in order to create the barrier region. The two gates, however, are connected and thus driven at the same voltage. In the other method (Fig. 8d), the barrier and well regions are created under a single polysilicon gate. The two-phase CCD is very commonly used for three reasons. First, it requires only two clocks (ϕ_1 and ϕ_2), simplifying drive circuitry. Second, the clocks are complementary. This reduces clock feedthrough. Third, the two-phase has a higher horizontal density, especially the two-phase with barrier and well regions under the same gate.

The virtual phase CCD^{16,17} (Fig. 8e) consists of one gate with both barrier and well regions within it and a second region, the virtual phase, in which a shallow high-dose ion implantation is used to create a heavily doped surface region which pins the surface potential at 0 V. The virtual phase has both barrier and well regions within it and acts the same as a polysilicon gate held at a constant voltage. Charge is first transferred from the clocked phase into the virtual phase, then the charge is transferred from the virtual phase into the next clocked phase. The virtual phase CCD has a number of advantages, including higher quantum efficiency than two-polysilicon-level area CCD sensors and simpler clocking. The disadvantages of the virtual phase include the need for larger voltage

swings on the single-clock and high-clock feedthrough into the output due to the lack of complementary clocks.

CCD Characteristics The four major performance parameters of a CCD shift register and output amplifier are charge-handling capacity, charge-transfer efficiency, charge-to-voltage conversion ratio, and noise. The charge capacity is the number of electrons which can be held and transferred in the CCD shift register. As charge is added to a shift register, a point is reached at which excess charge cannot be held; the excess charge either overflows into adjacent pixels or overflows into the bulk beneath the CCD or, in the case of a buried-channel CCD, overflows the barrier to the Si-SiO₂ surface. The charge capacity is a function of the device design, device layout, and CCD process. Figure 9a shows the electrostatic potential and charge distribution in a buried-channel CCD at three levels of charge: approximately one-quarter of saturation, at saturation, and beyond saturation. Below saturation, the electrons fill the center of the buried channel in the region of largest potential, separated from the surface by approximately 0.2 μm in distance and about 500 mV in potential. As more charge is added, the electron distribution spreads toward the surface and the potential barrier to the surface drops. Beyond saturation, the electrons contact the surface directly, resulting in charge-transfer inefficiency and blooming to neighboring pixels. The charge capacity of most CCD shift registers is of the order 1×10^{12} electrons/cm² of area in which the charge is held. In most CCD cells, the charge is held in only a fraction of the total cell area both along and across the CCD register. Typically, area CCD image sensors are designed with CCD charge capacities in the range of 50,000 to 200,000 electrons. Linear CCD image sensors, because of the larger amount of silicon area available for the CCD shift register, often are designed for 100,000 to 1,000,000 electrons.

The second major performance parameter for CCD shift registers is charge transfer efficiency. The transfer of charge from one stage to the next is neither instantaneous nor complete, limiting both transfer rate and the total number of stages in the CCD. There are two intrinsic mechanisms governing charge transfer: drift and diffusion. Drift is the movement of charge in the presence of an electric field. There are two origins of the electric field seen by an electron during charge transfer:

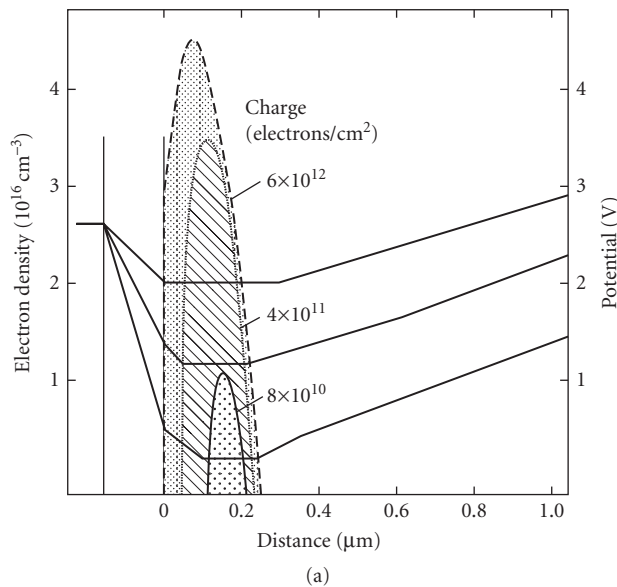
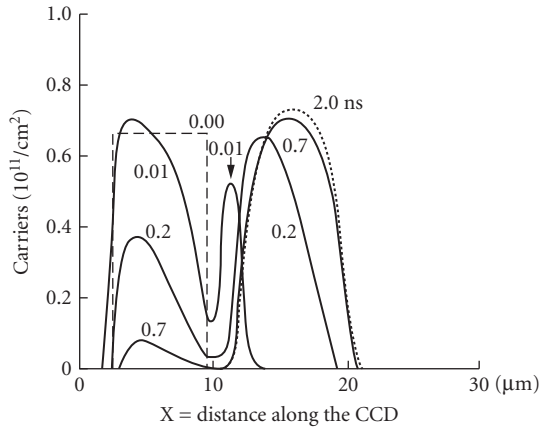
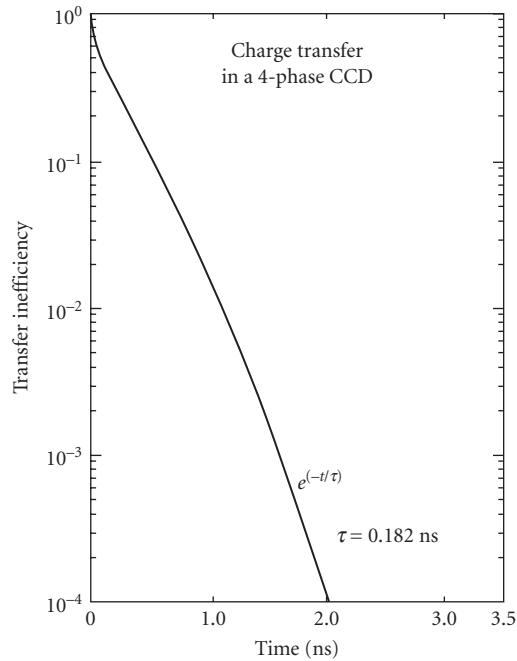


FIGURE 9 (a) Charge distribution and electrostatic potential as a function of distance from the surface for a buried-channel CCD for three levels of charge density.



(b)



(c)

FIGURE 9 (b) Charge density as a function of distance along a CCD for various times following the beginning of charge transfer from one 8- μm stage to another and (c) charge transfer inefficiency as a function of time from the start of charge transfer for the same example. (*Continued*)

the self-induced field resulting from the other electrons under the gate and the externally induced or fringing field. During the early stages of charge transfer, the self-induced field is large and is the dominant factor. After the charge concentration under the gate has decreased to a low value, the remainder of the charge transfer will be governed either by diffusion or by drift due to externally induced or fringing fields.

For a surface-channel CCD, the self-induced fields can be estimated from the formula:

$$E = -\frac{q}{C} \frac{dN}{dx} \quad (14)$$

where C is the gate capacitance per unit area and $N(x)$ is the density of electrons as a function of distance x from the edge of the electrode. In the early stages of charge transfer, both N and dN/dx are large. As the transfer proceeds, both N and dN/dx become small and the charge transfer is governed by fringing fields and diffusion.

Fringing fields are due to the two-dimensional nature of the electrostatic potential. If the charge is being transferred from gate 1 to 2, the potential will change smoothly between the two gates. The effect of the potential from one gate will typically extend 1 to 3 μm into a neighboring gate. The charge within the range of the fringing field will move by drift to the neighboring gate. Charge out of the range of the fringing fields will move by diffusion.

Figure 9b shows an example of charge transfer calculated for charge transfer from one 8- μm -long CCD gate into an adjacent gate. The charge density is shown as a function of distance along the CCD at several times after the start of charge transfer. At the start of the transfer, all the charge is under the first gate. At 0.01 ns into the transfer, the charge has moved from the edge of the first gate into the second gate, as a result of self-induced drift. By 0.2 ns, approximately half the charge has been transferred. By 0.7 ns, over 90 percent of the charge has moved into the second gate, leaving a residual in the first. At this point, the self-induced drift is sufficiently small that drift due to fringing fields and diffusion are the dominant mechanisms. Figure 9c shows the charge transfer inefficiency as a function of time for this example. Two slopes are evident in the transfer inefficiency. In the first 0.5 ns, the charge is transferred rapidly owing to the self-induced drift. For times longer than 0.7 ns, the transfer is due to fringing fields in this example.

For CCDs with longer gates or lower fringing fields than the example above, the final ~10 percent of the charge must transfer by diffusion. Charge transfer by diffusion follows an exponential time dependence.

$$N(t) = N(0)e^{-t/\tau_{\text{diff}}} \quad (15)$$

where, the time constant τ_{diff} for diffusion is

$$\tau_{\text{diff}} = \frac{4L_g^2}{\pi^2 D} \quad (16)$$

Here, L_g is one gate length and $D = KT\mu_e/q$.

For electrons at room temperature, $D = 25.8 \text{ cm}^2/\text{s}$. For an 8- μm long gate, the diffusion time constant is ~10 ns. To achieve a transfer inefficiency below 2×10^{-5} per transfer, 11 time-constants, or 110 ns in this example, are required. For this reason CCDs are typically designed with gate lengths shorter than ~8 μm and to build-in fringing fields.

In the simplest case for very low levels of transfer inefficiency, the total transfer inefficiency in a CCD register is the product of the number of stages N in the register and the transfer inefficiency per stage. Each stage will require two or more transfers. Virtual phase and two-phase require two transfers per stage, three-phase requires three transfers, and four-phase requires four transfers. The inefficiency per stage, however, will likely depend in a complicated manner on the amount of charge in the charge packet, the amount of charge in preceding charge packets, the voltages, and frequency of operation. The charge in the preceding packets will affect the filling of both bulk and interface traps.

In addition to the intrinsic sources of transfer inefficiency, there are a variety of extrinsic sources. These include surface and bulk traps and potential wells and barriers. The traps and the potential obstacles hold back an amount of charge from a charge packet. The charge is reemitted to later charge packets. The inefficiency due to traps depends on whether the traps have been filled by preceding charge packets and the emission time constants of the traps, and so is not modeled in a simple manner.

The intrinsic sources of noise in CCD shift registers include dark current and output amplifier noise. In addition, other sources of noise not intrinsic to the CCD itself include noise due to clock feedthrough from the CCD clocks to the output signal and noise in external electronics. The generation of dark current in CCD shift registers is the same as described at the end of Sec. 32.3 for photosensing elements. Associated with the dark current are both shot noise and pattern noise. The magnitude of the pattern noise in a CCD shift register is reduced over that of a single element since the charge packet averages the dark current over many pixels as it is transferred to the output.

The noise associated with the CCD output consists of the Johnson or thermal noise, the $1/f$ noise of the output amplifier, and the kTC noise associated with resetting the floating diffusion. The kTC noise is due to uncertainty in the amount of charge remaining on the floating diffusion following reset owing to thermal fluctuations in the reset gate. The rms noise σ_n in the number of electrons caused by kTC noise is given by

$$\sigma_n = \sqrt{\frac{kT}{q}} C_{\text{FD}} \quad (17)$$

where C_{FD} is one total floating diffusion capacitance. For a 10-fF floating diffusion capacitance, the rms noise is ~ 40 electrons.

The kTC noise may be eliminated entirely by use of a signal-processing technique called correlated double sampling (CDS). In correlated double sampling, the output level is sampled before the charge packet is transferred onto the floating diffusion and again after transfer of the charge packet (times t_2 and t_0 respectively in Fig. 6c). The two values are subtracted either by an analog circuit or by digital subtraction. Any uncertainty in the voltage level of the floating diffusion following reset is subtracted and thus the kTC noise eliminated. The output amplifier low-frequency noise, called $1/f$ noise because of the inverse frequency ($1/f$) shape of the noise power spectrum, is also reduced but not eliminated by correlated double sampling. The total noise of a CCD output amplifier is in the range of 7 to 40 rms electrons per pixel depending on the amplifier design and the operating speed. Values of a single rms electron or less have been obtained for very slow pixel rates under cooled conditions.¹⁸

MOS Readout

The other major category of readout structures is the MOS readout. The individual light-sensing elements (photodiodes, photocapacitors, or photoconductors) at each pixel are connected to a readout line by means of a transfer gate. Each pixel along the readout line is addressed separately by addressing circuitry. When a particular pixel is addressed, the transfer gate is turned on and the charge transferred from the pixel to the readout line. An amplifier at the end of the readout line senses the change in voltage or current resulting from the charge transfer.

Typically, the pixels would be addressed serially along the line. The first pixel would be addressed, causing the charge from the image-sensing element to be transferred onto the readout line. The voltage change or current would be sensed, the readout line reset to its original voltage if necessary, and the next pixel addressed. This is different from a CCD. In the CCD, charge from all pixels is transferred into the CCD register simultaneously. Individual pixels or groups of pixels cannot easily be addressed in a random fashion by the CCD, but this random addressing can be accomplished readily by the MOS readout.

There are several types of MOS readout devices. These include the CID,¹⁹ the AMI,²⁰ and the CMD²¹ in addition to the normal MOS array. The CID has no readout line. Each pixel consists of two overlapping gates, one controlled by a row address and the other by a column address.¹¹ When neither a row nor a column of a particular pixel is being addressed, the photogenerated charge is held under both gates and can be transferred between them. When a row of a pixel is addressed, the charge transfers onto the column gate. Then both row and column are addressed, the charge is injected into the substrate, and the current sensed. The CID is not widely used in visible imaging applications because the charge conversion sensitivity is very poor and noise very high compared to the CCD or the other MOS architectures.

In the amplified MOS imager (AMI), the image-sensing element at each pixel consists of a phototransistor rather than a simple photodiode. The photogenerated charge is stored on the gate of the MOS transistor.¹⁹ When a particular pixel is addressed, the photogenerated charge modulates the transistor current. This current amplification at each pixel helps to overcome many of the noise and speed limitations of conventional MOS arrays.

MOS readout differs in an important way from CCD readout. In MOS readout, the charge is transferred from a single pixel onto a readout line and the change in voltage or current in the readout line is sensed. In a CCD, the charge packets are kept intact while being transferred physically to a low-capacitance output. The lower sensitivity of a simple MOS array can be illustrated as follows. The change in voltage on the readout line is given by $V = Nq/C$ where N is the number of electrons, and C is the readout line capacitance. Because the readout line covers the full length of the array, its capacitance is in the picofarad range (typically 2 to 10 pF depending on design and process). This compares to the 10-fF capacitance for the CCD output. As a result, the voltage swings on the readout line are very small (16 nV/electron for a 10-pF readout line capacitance). This leads to a high sensitivity to clock noise due to capacitive feedthrough of the row and column address clocks onto the readout line. The feedthrough may be many times larger than the signal in most MOS sensors. Once the charge has been transferred onto the readout line, it is sensed either by a current-sensitive amplifier or by a voltage-sensitive amplifier, followed by a reset of the readout line to its original voltage.

CCD readout has the advantages of very high sensitivity and low noise. However, CCD readouts are limited in charge-handling capacity, while MOS readouts are capable of carrying very large amounts of charge and so are not as limited on the high end of the dynamic range. However, because the MOS readout line has much higher capacitance than the CCD, the sensitivity is lower and the noise is higher. Another difference is in the readout architecture. The CCD readout is essentially serial and not suited to random readout or partial-array readout. The MOS array, however, can be addressed in a manner similar to a memory, making it well-suited to pixel or partial-array addressing.

32.5 SENSOR ARCHITECTURES

Solid-state image sensors are classified into two basic groupings: linear and area. Linear sensors include single-line arrays, multilinear arrays for color scanning, and time delay and integrate (TDI) arrays for low-light-level scanning. Area sensor architectures include the frame transfer CCD, the interline transfer CCD, and various forms of MOS x - y addressed arrays.

Linear Image Sensor Arrays

Linear sensors are used almost exclusively in scanning systems for scanning documents, film, and three-dimensional still objects. There are two basic classes of scanning systems: contact scanners and reduction scanners. These are illustrated in Figs. 10*a* and *b*. In reduction scanners (Fig. 10*a*), the sensor is smaller than the document to be scanned; lenses are used to image the document onto the sensor. In contact scanners (Fig. 10*b*), the sensor is the same width as the item to be scanned, usually a document. Relay optics is used between the sensor and the document. Selfoc lenses (Fig. 10*c* and *d*) and roof-mirror-lens arrays are the two types of relay optics used most frequently.

There are three basic architectures for linear sensing arrays: MOS line arrays, CCD linear and multilinear sensors, and time-delay and integrate (TDI) sensors. These architectures are illustrated in Fig. 11. The MOS array is used most often in contact scanning applications where material or processing problems make CCD arrays impractical. These applications include arrays fabricated from polysilicon or amorphous silicon on nonsilicon substrates, arrays covering large distances, or arrays requiring special processing (such as logarithmic amplification) at each pixel. The CCD linear and multilinear arrays are used most often in reduction scanning where wide dynamic range

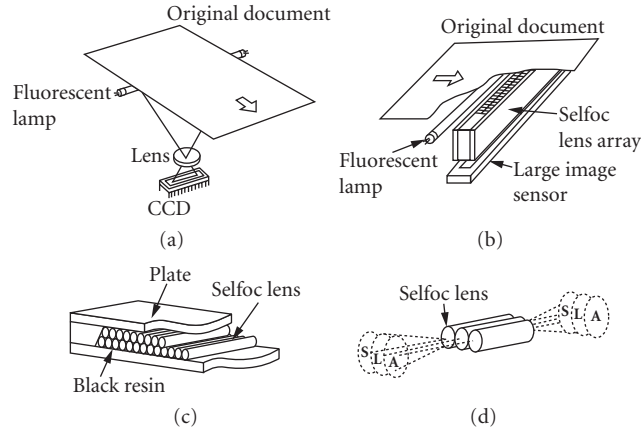


FIGURE 10 (a) Reduction document scanner using linear image sensor, in which the image on the page is reduced by a lens onto the line array; (b) contact document scanner in which the length of the image sensor is the same as the document width and one-to-one relay optics is used to transfer the image from the document to the array; (c) diagram of a self-foc lens array used as transfer optics between document and array in a contact scanner; and (d) operation of a self-foc lens array.

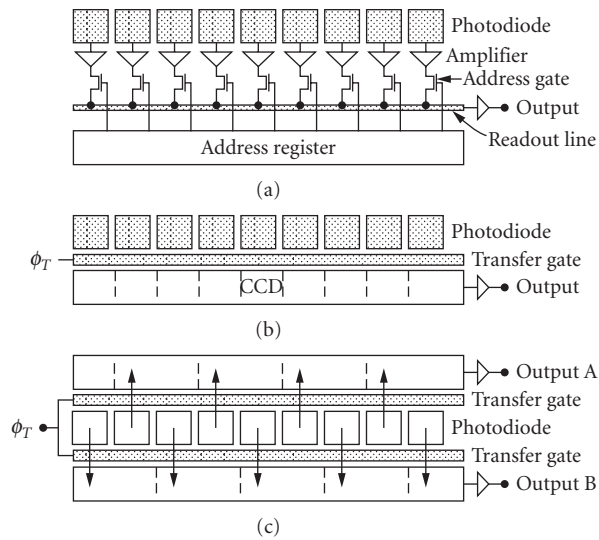


FIGURE 11 Architectures for linear image sensors: (a) MOS line array consisting of photodiodes, preamplifier, MOS switches addressed by an address register, and a readout line with amplifier; (b) linear CCD image sensor consisting of photo-diodes, transfer gate, and CCD readout; and (c) linear CCD image sensor with two CCD output registers, one for the odd diodes and the other for the even diodes, for higher horizontal pitch.

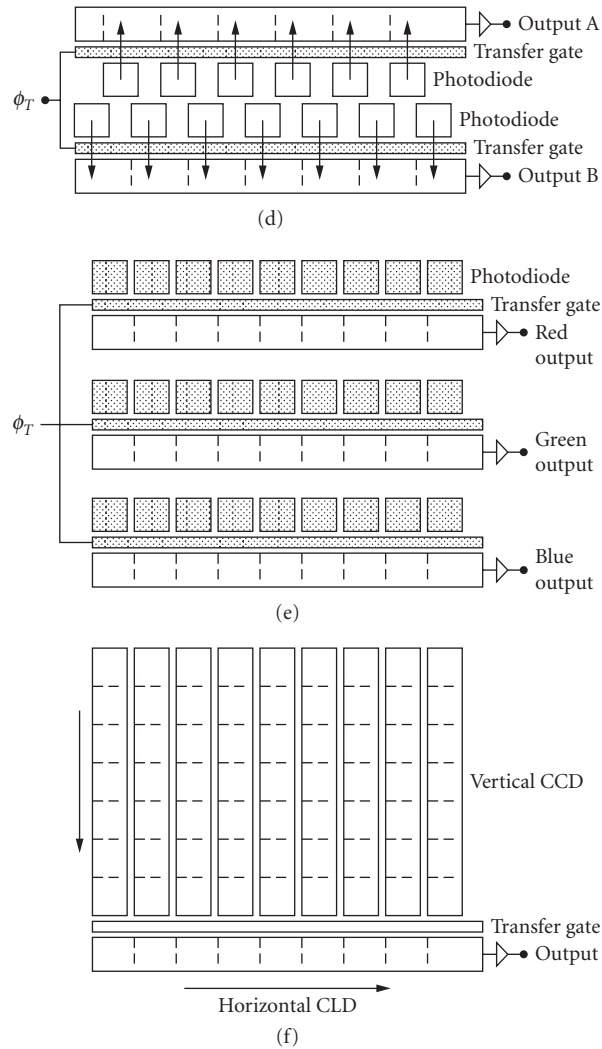


FIGURE 11 (d) Staggered linear CCD image sensor with two rows of photodiodes offset by one-half pixel to increase horizontal sampling; (e) trilinear CCD in which three CCDs are fabricated on the same silicon die, each with its own color filter; and (f) time-delay and integrate (TDI) array, in which the charge in the vertical registers is clocked in phase with the motion of the scene or document being imaged in order to increase signal-to-noise. (*Continued*)

and small pixel size are required. However, contact arrays are also often realized by butting multiple CCD arrays end-to-end. TDI arrays are used in very low-light-level scanning applications where integration over many lines is required to achieve adequate signal-to-noise ratio.

The MOS linear array (Fig. 11a) consists of individual photosensing elements, an amplifier, an address switch, and a readout line and amplifier. The photosensing element is usually a photodiode,

although photoconductors are used in contact scanning arrays fabricated from amorphous silicon. Since the charge generated in the diode is generally very small (in the hundreds to thousands of electrons), a simple amplifier is usually placed at each pixel to drive the high-capacitance readout line. The use of amplification at each pixel can allow some signal-processing functions, such as logarithmic amplification, clipping, triggering and latching, etc., to be performed at the pixel. A MOS switch placed after the amplifier allows each pixel to be addressed in sequence; the switch is driven by an address register. At the end of the readout line is an amplifier which may buffer and/or amplify the signal. The MOS array has the advantage of process simplicity and the ability to perform signal processing at each pixel; it has the disadvantage of low signal level (because of the large readout line capacitance) and pattern noise introduced by feedthrough from the switching transistor.

The linear CCD image sensor is the most often used architecture for scanning applications owing to its low noise, high sensitivity, wide dynamic range, and small pixel pitch. Figure 11*b* shows the simplest type of linear CCD, in which a single row of photodiodes is connected to a single CCD register via a transfer gate. In operation, the signal is integrated on the image-sensing element (generally a photodiode) for a line time. The horizontal CCD is then stopped, the transfer gate opened, and the charge transferred from all the photodiodes simultaneously to the CCD. The transfer time is typically a few microseconds. The transfer gate is then closed, integration resumed for the next line, and the CCD clocked to read out the charge packets. Many arrays also feature antiblooming for situations where the light level may not be controlled, as well as electronic shuttering, which allows an integration time on the photodiodes to be less than the readout time of the CCD.

For linear-sensing applications requiring a higher pixel density, a double-sided readout is often used (Fig. 11*c*). In this architecture, a CCD array is placed on either side of the line of photodiodes and charge transfer from the diodes alternates between the top and bottom CCDs. This architecture uses lower horizontal clock rates and a higher pixel pitch, since the diode pitch can usually be made smaller than the CCD pitch. The charge packets from the two arrays may be multiplexed into one output if desired. The disadvantage of this architecture is the slight differences between even and odd pixels, due to slight differences in the two outputs (or slight differences due to multiplexing the two registers).

Another architecture which is used to further decrease the sampling pitch is the staggered linear array (Fig. 11*d*). In the staggered array, two rows of photodiodes are offset by a half pixel. The two rows are read out by CCD arrays. The first array is delayed (usually in a digital line store) and then combined with the second to form a double-density scan.

Multilinear arrays (Fig. 11*d*) have been developed for color scanning applications. In this architecture, several (usually three) linear arrays are combined on the same silicon die separated by a distance equivalent to an integral number of scan lines. Color filters (either integral or in close proximity) are aligned over the arrays. External digital line delays are used to realign the three arrays. Separate electronic shuttering may be provided for each array in order to adjust for differences in intensity in each of the bands.

The third major class of line arrays is time-delay and integrate, or TDI, arrays.^{22,23} The TDI architecture is shown in Fig. 11*f*. TDI arrays are used when inadequate signal-to-noise ratio from a single-line array requires averaging over multiple lines. Applications for TDI arrays include high-speed document scanners and space-based imaging systems. Instead of a single row of photodiodes, the TDI array utilizes CCD stages in the vertical dimension which are clocked synchronously with the movement of the document to be scanned. The signal level in a TDI array increases linearly with the number of stages. The noise level, however, increases most as the square root of the number of stages.

Area Image Sensor Arrays

There are three major classes of area image sensor architectures: MOS diode arrays, frame-transfer CCDs, and interline-transfer CCDs. Within each there are a number of variations. CCDs have come to dominate the majority of applications owing to their higher sensitivity. However, MOS arrays are still used for specialized applications where addressability or high readout rate is important.

TABLE 1 Image Area Dimensions and Pixel Dimensions for Various Format Image Sensors

The format name is given in inches based on historic image tube formats. The pixel dimensions are based on a 484×768 pixel array.

Optical format	1 in	2/3 in	1/2 in	1/3 in	1/4 in
Active area (4:3 aspect ratio)					
Height (mm)	9.6	6.6	4.8	3.3	2.4
Width (mm)	12.8	8.8	6.4	4.4	3.2
Diagonal (mm)	16	11	8	5.5	4
Pixel dimensions ($484 \text{ lines} \times 768 \text{ pixels}$)					
Height (μm)	19.8	13.6	9.9	6.8	4.9
Width (μm)	16.7	11.5	8.3	5.7	4.2

Historically, the physical dimensions of the active imaging areas of CCD arrays for consumer and commercial applications are specified by the size of the vidicon tube which it replaces. The common format sizes include 1/4, 1/3, 1/2, and 1 in. In most cases, the aspect ratio is 4:3, reflecting the television standard. Table 1 lists the formats and the corresponding dimensions of the imaging area of the array.

The standards for the number of vertical lines in arrays for consumer and professional video are usually based on the corresponding television standards, including NTSC, PAL, and the various HDTV standards. NTSC has 484 active lines, PAL has 575, the Japanese HDTV standard has 1035, and the European HDTV standard has 1150. In all cases the lines are interlaced; i.e., odd lines are read out in one field and even lines in the next. Historically, for sensors for NTSC television, the number of pixels horizontally in a line has been associated with multiples of the color subcarrier frequency; common horizontal pixel counts include 384, 576, and 768. Sensors for the Japanese HDTV standard typically have 1920 horizontal pixels. The pixels are rectangular rather than square. Table 1 lists approximate pixel dimensions in micrometers for a few common formats for a $484(V) \times 768(H)$ pixel image sensor.

For industrial, scientific, graphics electronic photography, digital television, and multimedia applications, however, the nonsquare pixels and interlaced readout of sensors based on television standards are significant disadvantages. Image sensors designed for these industrial and commercial applications typically have square pixels and progressive scan readout. In interlaced NTSC readout, for example, the first field consists of lines 1, 3, 5 \dots 483 and is read out in the first 1/60 second of a frame. The second field consists of lines 2, 4, 6 \dots 484 and is read out in the second 1/60 second of a frame. The resulting temporal and spatial displacement between the two fields is undesirable for these applications. In addition, digital compression of interlaced scan moving images, especially with motion estimation, is difficult and also introduces artifacts.

In progressive scan readout each line is read out sequentially. There is no even or odd field, only a single frame. As a result there are no temporal and spatial sampling displacement differences. However, for a given resolution and frame rate, the readout rate of a progressive scan image sensor is double that of an interlaced scan. In addition, for these applications, the number of pixels is often based on powers of 2: such as 512×512 or 1024×1024 —facilitating memory mapping and image processing.

MOS Area Array Image Sensors The architecture of MOS area arrays is illustrated in Fig. 12.²⁴ It consists of the imaging array, vertical and horizontal address registers, and output amplifiers. The pixel of a MOS array consists of an image-sensing element (photodiode, photocopacitor, phototransistor, or photoconductor), a row-address gate, and a vertical readout line. The row-address gate is bussed horizontally across the array and is driven from a row-address register on the side(s) of the array. At the start of a line, a single row is addressed, causing the charge from all the photodiodes in

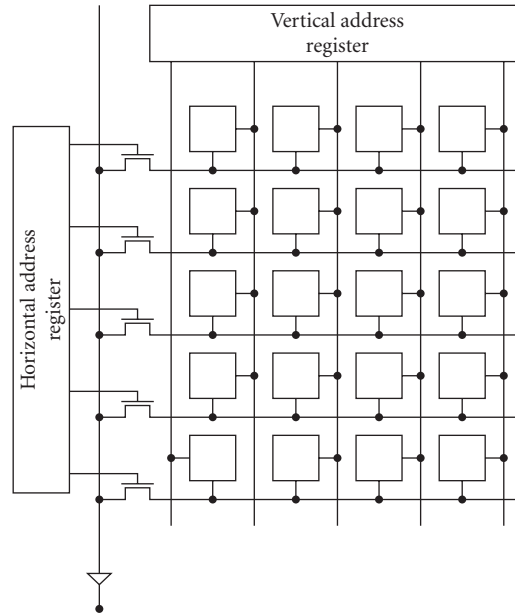


FIGURE 12 MOS photodiode array, consisting of vertical and horizontal address registers and readout line.

a row to be transferred onto the vertical readout line. Horizontal address gates are placed at the bottom of the vertical readout line. Buffer amplifiers to drive the horizontal readout line may also be placed at the end of the vertical readout line. The horizontal address register then serially addresses each vertical readout line, sequentially turning on the horizontal address gates. After the horizontal addressing is completed, the readout lines may be reset and precharged and the next row addressed.

There is a wide range of variations on this basic architecture. An example is the charge modulation device, or CMD,^{25,26} array in which a phototransistor is placed at each pixel site in order to achieve amplification and to achieve high currents to drive the capacitance of the vertical and horizontal readout lines. Other examples include arrays with sophisticated charge collection circuits at the end of each vertical readout line.

Frame Transfer CCD Image Sensors CCD area arrays fall into two categories: frame transfer and interline transfer. The simplest form of frame transfer CCD is the full-frame type, shown in Fig. 13a. A photograph of a few pixels of a frame transfer CCD is shown in Fig. 14a. The array consists of a single image area composed of vertical CCDs and a single horizontal register with an output amplifier at its end. In this architecture, the pixel consists of a single stage of a vertical CCD. This type of device requires an external shutter. When the shutter is opened, the entire surface of the sensor is exposed and the charge is collected in the CCD potential wells at each pixel. After the shutter is closed, the sensor is read out a row at a time by clocking a row of the vertical register into the horizontal register, then clocking the horizontal register to read out the row through the output amplifier. For higher readout rates dual horizontal CCDs are used in parallel. The full-frame CCD has the advantage of progressive scan readout high fill factor, very low noise, and wide dynamic range. However, it requires an external shutter. It is most often used in still electronic photography, scientific, industrial, and graphics applications.

For motion imaging applications, a shutter is not practical. In order to overcome the need for a shutter, frame transfer CCDs incorporate a storage area in addition to imaging area. For interlaced video applications, this storage area is sufficiently large to hold a field (242 lines in NTSC television).

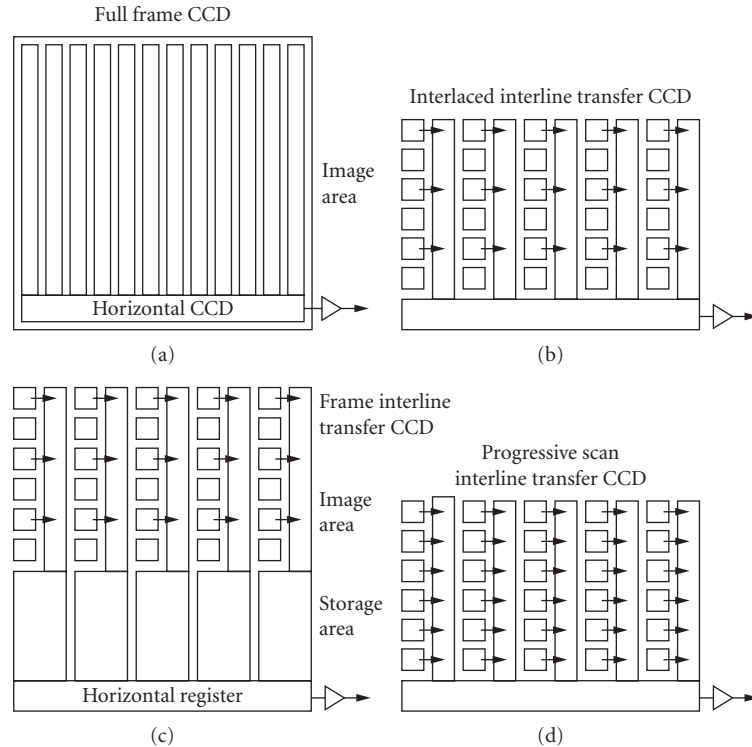
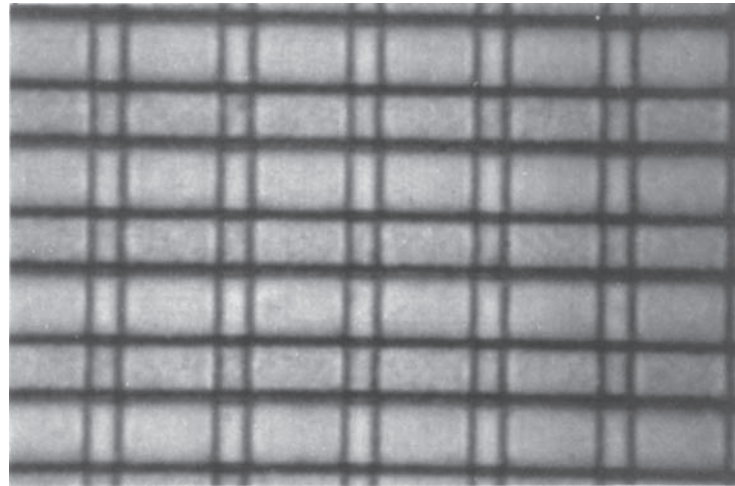


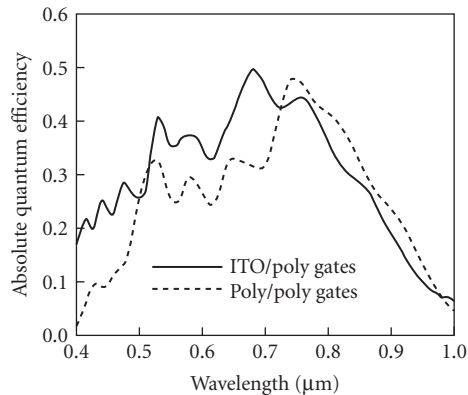
FIGURE 13 (a) Architecture of full-frame CCD; (b) architecture of interlaced interline transfer CCD; (c) architecture of frame interline transfer CCD, in which a storage area is provided to reduce smear during readout; and (d) progressive scan interline transfer CCD, in which every photodiode is read out into the vertical CCD simultaneously.

The image area consists of vertical CCDs. For interlaced NTSC video, there are 242 pixels vertically in the image area. Interlace is achieved by changing the gates under which integration is performed in the even and odd fields. For example, for a four-phase CCD, integration would be performed under phases 1 and 2 in one field and 3 and 4 in another, effectively shifting the sampling area by half a pixel in each field. The storage area consists also of 242 pixels vertically. The device operation is as follows. The image area integrates for a field time and the photogenerated charge held in the vertical CCDs. The vertical CCDs are then clocked in order to rapidly transfer the charge from the image area into the storage area. Because the sensor is still under illumination, this transfer time must be much shorter than the integration time. This transfer typically requires 0.2 to 0.5 ms. The storage area is then read out a row at a time by transferring a row into the horizontal register and clocking this register. While this readout is occurring, the image area is integrating the next field. The most significant disadvantage of frame transfer CCDs is the image smear caused by illumination during the transfer from the image to the storage area. This smear can be on the order of 3 percent.

In both frame transfer and full-frame devices, the light must pass through the polysilicon electrodes before being absorbed in silicon. Owing to the high absorption of short wavelengths in the polysilicon, the quantum efficiency in the blue is only about 20 percent and in the green is about 50 percent. Figure 14b shows the quantum efficiency for a full-frame image sensor with polysilicon gates. Three approaches have been used to improve the efficiency: the virtual phase CCD, transparent electrodes, and backside illumination. In the virtual phase CCD (see “Types of



(a)



(b)

FIGURE 14 (a) Photograph of a few pixels in the image area of a full-frame CCD and (b) quantum efficiency as a function of wavelength for a full-frame CCD with polysilicon and with transparent indium-tin-oxide electrodes.

CCDs” in Sec. 32.4), the second polysilicon electrode is replaced with a very shallow highly doped p^+ layer, very similar to the pinned photodiode (Fig. 2d). Since there is no electrode over this phase to absorb the light, higher quantum efficiency is achieved, particularly at wavelengths less than 500 nm. Backside illumination provides nearly unity quantum efficiency but requires that the sensor be thinned to less than 10 μm . Owing to its cost, backside thinning is employed only in sensors for very specialized scientific or aerospace applications. Indium-tin-oxide, or ITO, is the most commonly used transparent electrode.²⁷ Usually it is substituted for the second level of polysilicon. Figure 14b also shows the quantum efficiency of a full-frame device in which ITO has been substituted for one of the polysilicon gates.

Interline Transfer CCD Image Sensors The interline transfer CCD is fundamentally different from the frame-transfer CCD in that, in addition to the vertical CCD, the pixel also contains a separate image-sensing element (photodiode, pinned photodiode, photocapacitor, or photoconductor) and a

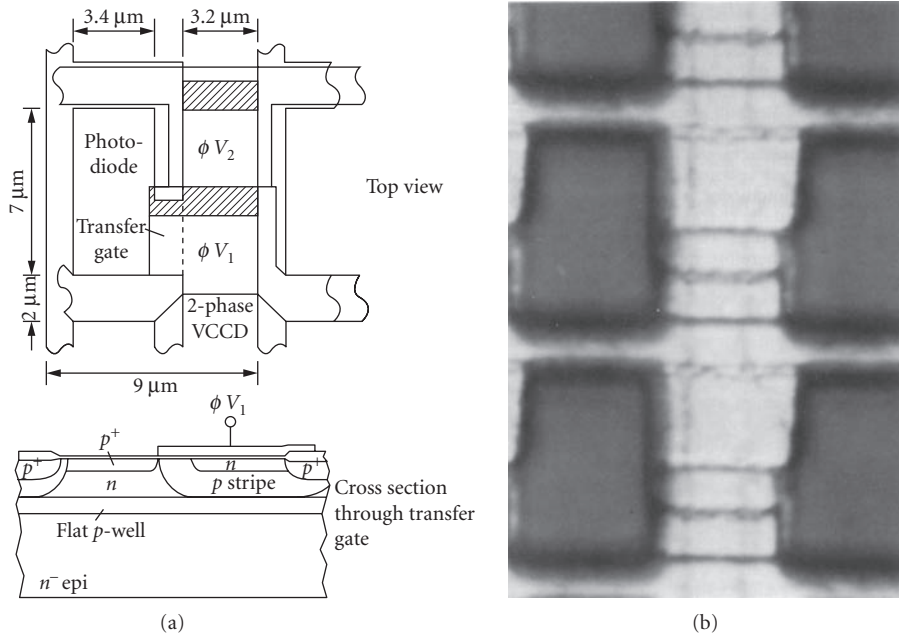


FIGURE 15 (a) Diagram of a pixel of an interline transfer CCD and (b) photograph of a pixel from an interline transfer CCD.

transfer region between the photodiode and the vertical CCD. Interline CCDs with photodiodes for sensing elements are considered first. Figure 13*b* illustrates the architecture of an interline transfer CCD and Fig. 15 illustrates a pixel of the CCD. The CCD and the transfer region between the diode and CCD are covered with a light shield (typically aluminum, although metal silicides are also used). The light shield prevents any light from entering the vertical CCD registers, allowing them to be read out while the sensor is illuminated. When the sensor is illuminated, the photogenerated charge is held on the photodiode. During the vertical retrace interval at the end of a field, the photogenerated charge is transferred into the vertical CCD by clocking the CCD gate over the transfer region. Once the charge has been transferred from the diodes into the vertical CCDs, the diodes resume integrating and the vertical CCDs are clocked in order to transfer a row at a time from the image area into the horizontal CCD.

For consumer and most commercial applications, interlaced interline CCDs are used to be consistent with television standards. In interlaced interline CCDs there is one vertical CCD stage for every two photodiodes. During the retrace time before the first field the charge from the odd rows of photodiodes is transferred into the vertical CCDs; the charge is then transferred out a row at a time into the horizontal CCD. During the retrace time before the second field the charge from the even rows of photodiodes is transferred into the vertical CCDs; once again, the charge is transferred out a row at a time into the horizontal CCD. In NTSC television the fields are approximately 1/60 second long; in PAL the field time is 1/50 second.

Because the vertical CCDs in an interline CCD are covered by a light shield, very little stray light is absorbed in the CCD. However, due to light scattering under the lightshield and lateral diffusion of photogenerated electrons, some charge can reach the vertical registers as they are read out during a field. This results in smear. For consumer applications the level of smear is not noticeable. However, for especially demanding applications such as television broadcast cameras, a field storage area is added to the bottom of the image area. This architecture is called the frame interline transfer, or FIT CCD²⁸ (Fig. 13*c*). Following transfer of the photogenerated charge from the diodes

into the vertical CCDs, the vertical CCDs are clocked to rapidly shift the charge from the image area into the storage area. This transfer typically takes less than 0.5 ms, reducing smear 30-fold. One line at a time is transferred from the storage area to the horizontal register and the readout.

For still electronic photography, scientific, computer-related, graphics and professional applications, interlaced video is not desirable. For these applications a progressive scan architecture is utilized. In a progressive scan interline CCD there is a full CCD stage for every photodiode.²⁸ Following integration, the photogenerated charge from all the photodiodes is transferred into the vertical CCDs. The photodiodes resume integration and the charge packets from the vertical registers are transferred into the horizontal a row at a time. The progressive scan interline CCD requires twice the vertical CCD density and is therefore more complicated to fabricate. However, progressive scan readout provides many advantages in image quality for both motion and still imaging.

Nearly all interline CCDs used in camcorder, broadcast camera, or commercial applications utilize vertical antiblooming in order to prevent blooming when the sensor is illuminated beyond saturation. A cross-section diagram of an interline CCD with vertical antiblooming is shown in Fig. 15*b*. The device illustrated uses a pinned photodiode photosensing element. The CCD is built on an *n*-type silicon substrate. A *p*-well is formed about 2 μm deep in the *n*-type silicon. An *n*-type buried-channel is then formed followed by a *p*+ surface layer. The *n*-type substrate is reverse biased with respect to the *p*-well. When the photodiode is illuminated above saturation, the excess electrons spill out of the *n*-type buried channel and into the substrate. Owing to the vertical overflow, however, photogenerated carriers from photons absorbed below the *p*-well are drawn into the substrate and are not collected by the diode. Thus the quantum efficiency of these devices falls rapidly at wavelengths beyond 550 nm. Figure 16 shows the quantum efficiency of an interline CCD with vertical antiblooming as a function of wavelength. The internal quantum efficiency of the photodiode is nearly 100 percent to about 550 nm, after which it decreases. However, since the photodiode only occupies about 20 percent of the pixel, and this aperture is typically reduced further by the light shield in order to eliminate optical scattering into the vertical CCD, the actual efficiency of the device is only about 15 percent. The photoresponse is linear at low signal levels but becomes nonlinear at charge levels near saturation.²⁹

Two major approaches are used to improve the fill factor (and therefore the quantum efficiency) of interline CCDs. One uses microlens arrays to focus the light incident on a pixel onto the photodiode and the other a vertical integration of image sensors with amorphous silicon photoconductors to achieve a high fill factor. Figure 17 shows a microlens array on top of an interline CCD.³⁰ The

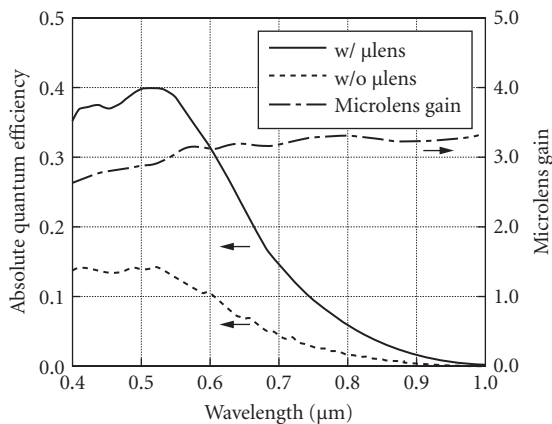


FIGURE 16 Quantum efficiency as a function of wavelength for an interline transfer CCD with and without microlens array.

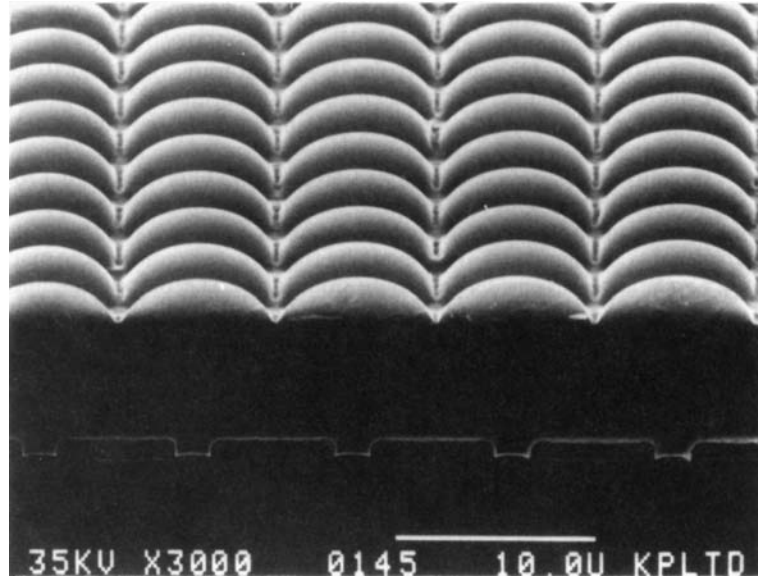


FIGURE 17 Scanning electron micrograph of microlens array on top of an interline transfer CCD.

lenses are formed by coating a spacer layer on the wafer of CCD devices followed by a lens-forming layer. The lens-forming layer is patterned and then reflowed to form the lens arrays. The quantum efficiency of an interline CCD with and without a microlens array is shown in Fig. 18. A 3-fold improvement in quantum efficiency is achieved because light from nearly the entire pixel area is focused onto the photodiode.

The structure of the photoconductor and its band diagram in the second approach⁷ are illustrated in Fig. 2e. The amorphous silicon is about 1 μm thick; owing to its high absorption coefficient in the visible, it can achieve nearly 100 percent internal quantum efficiency. The back contact of the

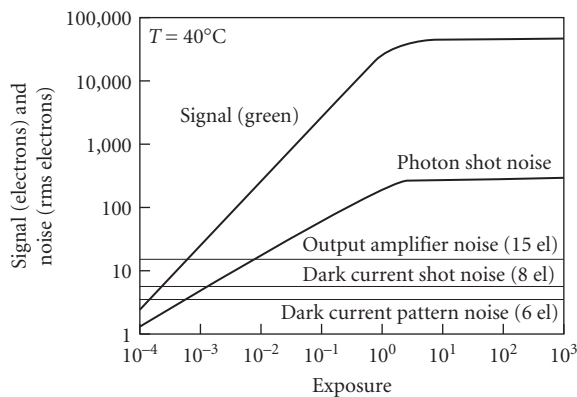


FIGURE 18 Quantum efficiency of interline CCD with and without microlens array.

photoconductor contacts the diode in the interline CCD. Light absorbed in the amorphous silicon generates electron-hole pairs. The amorphous silicon is biased such as to create a high field which sweeps out the electrons to the back contact, where they can be stored on the diode until they are transferred into the CCD. Because the photoconductor is fabricated on top of the CCD pixel, it can have nearly 100 percent fill factor. In addition, because of the wide bandgap of the amorphous silicon, its dark current (due to thermal generation of carriers) is often lower than the single-crystal silicon. However, the amorphous silicon photoconductors suffer difficulties related to charge trapping. Owing to impurities and dangling or strained bonds, there is a high density of traps in amorphous silicon. These can trap charge carriers and reemit them at a later time, causing image lag. Because the trapping and detrapping is field-dependent, it can also result in nonlinear response.

CCD Performance In all CCDs, both interline and frame transfer, the signal output is linear with the illumination except at levels approaching saturation. This linear response is in marked contrast to image tubes, which exhibit highly nonlinear response. For interline CCDs, the total charge capacity ranges from 100,000 electrons for larger cells (such as the $13.6(\text{V}) \times 11.6(\text{H})\text{-}\mu\text{m}$ cell typical for a 2/3-in format) to 20,000 electrons or less for smaller cells (such as the $6.8 \times 5.8\text{-}\mu\text{m}$ cell in a 1/3-in format $484 \times 768\text{-pixel}$ CCD).

The principal noise sources in both interline and full-frame image sensors are dark current pattern and shot noises and output amplifier noise. To illustrate the contributions, CCDs have dark current levels at 40°C of less than 1 nA/cm^2 . For a 2/3-in format sensor, this would correspond to about 320 electrons per pixel dark level. The corresponding shot noise would be 18-rms electrons and the pattern noise would typically be at a similar level. However, because of the nonrandom nature of pattern noise, its appearance is considerably more noticeable. The output amplifier noise would also be at about the 15-rms electron level. Thus, the overall noise for this example would be in the 30-electron range and the dynamic range would be over 2000 for a charge capacity of 90,000 electrons. For scientific applications where the sensor can be cooled and the readout performed at a lower frequency, noise levels less than 5 electrons can be achieved and even subelectron noise has been reported.¹⁸

Color Imaging

Silicon based CCDs are monochrome in nature. That is, they have no natural ability to determine the varying amounts of red, green, and blue (RGB) illumination presented to the photodetectors. There are three techniques to extract color information.

1. *Color Sequential* (Fig. 19)—A color image can be created using a CCD by taking three successive exposures while switching in optical filters having the desired RGB transmittances. This approach is normally used only to provide still images of stationary scenes. The resulting image is then reconstructed off-chip. The advantage to this technique is that resolution of each color can remain that of the CCD itself. The disadvantage is that three exposures are required, reducing frame times by more than a factor of three. Color misregistration can also occur due to subject or camera motion. The filter switching assembly also adds to the mechanical complexity of the system.
2. *Three-Chip Color* (Fig. 19)—Three-chip color systems use an optical system to split the scene into three separate color images. A dichroic prism beam splitter is normally used to provide RGB images. Color images can then be detected by synchronizing the outputs of the three CCDs. The disadvantage to such a system is that the optical complexity is very high and registration between sensors is difficult.
3. *Integral Color Filter Arrays (CFA)* (Fig. 19)—Instead of performing the color filtering off-chip, filters of the appropriate characteristics can be fabricated above individual photosites.^{31,32} This approach can be performed during device fabrication using dyed (e.g., cyan, magenta, yellow) photoresists in various patterns. The major problem with this approach is that each pixel is sensitive to only one color. Off-chip processing is required to “fill in” the missing color information between pixels.^{33,34}

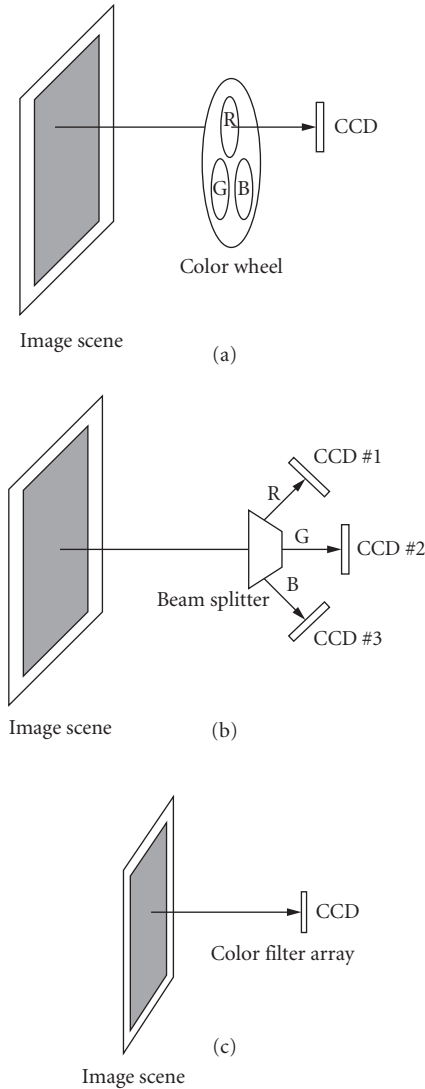


FIGURE 19 Methods of color separation in cameras with area image sensors: (a) color sequential using a color wheel; (b) a prism with dichroic beam splitters and three image sensors; and (c) single-chip image sensor with integral color filter array.

In order to minimize size, weight, and cost, most consumer color camcorders use a CCD sensor with an integral CFA. The photosites are covered with individual color filters—for example, a red, green, and blue striped filter, or a green, magenta, cyan, and yellow mosaic filter. Some popular CFA patterns are shown in Fig. 20. Because each photosite can sense only one color, the color sampling is not coincident. For example, a blue pixel might be seeing a white line, while nearby red and green pixels are seeing a dark line in the scene. As a result, high-frequency luminance edges can be aliased into bright color bands. These color bands depend not only on the color filter pattern used, but also

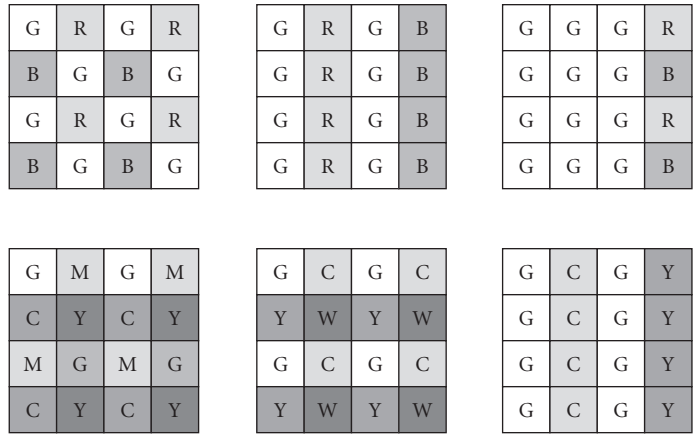


FIGURE 20 Common color filter array patterns where, R = Red, G = Green, B = Blue, Y = Yellow, M = Magenta, C = Cyan, and W = White.

on the optical prefilter and CFA interpolation algorithm. The color bands are caused by aliasing, which is a property of any sampled system. Aliasing occurs when the frequency of the input signal is greater than the Nyquist limit of one-half the sampling frequency. If the input frequency is well below the Nyquist limit, there are many samples per cycle. This allows the input to be reconstructed perfectly, if a proper reconstruction filter is used. When the input frequency is greater than the Nyquist limit, there are less than two samples per cycle. The sample values now define a new curve, which has a frequency lower than the input frequency. In effect, the high frequency takes on the alias of a lower frequency. Aliasing is a particular problem with color sensors, since the sampling phase is different for the different color photosites. Therefore, the aliased signal has different phases for different colors. This creates the color bands.

The color aliasing is reduced by using an optical anti-aliasing or “blur” filter, positioned in front of the color CCD sensor.³⁵ Blur filters are typically made of birefringent quartz, with the crystal axis oriented at a 45° angle, as shown in Fig. 21. In this orientation, the birefringent quartz exhibits the double refraction effect. An unpolarized input ray emerges as two polarized output rays, labelled *o*- and *e*-rays. The output ray separation is proportional to the filter’s thickness, *T*. A 1.5-mm-thick plate will give a separation of about 9 μm. Figure 21 shows a simple “two-spot” filter. More complex filters use three or more pieces of quartz cemented in a stack.

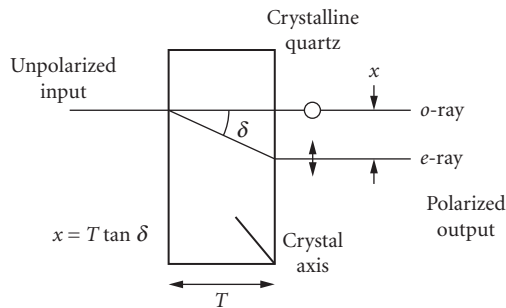


FIGURE 21 Birefringent blur filter used to reduce aliasing in single-chip color image sensors.

32.6 REFERENCES

1. G. Lubberts and B. C. Burkey, "Optical and Electrical Properties of Heavily Phosphorous-Doped Epitaxial Silicon Layers," *J. Appl. Phys.*, **55**(3):760–763 (Feb. 1, 1984).
2. O. S. Heavens, *Optical Properties of Thin Solid Films*, Dover, New York, 1965.
3. N. Teranishi et al., "An Interline CCD Image Sensor with Reduced Image Lag," *IEEE Trans. Electron Devices* **ED-31**(12):1829–1833 (Dec. 1984).
4. S. Kosman et al., "A Large Area 1.3 Megapixel Full-Frame CCD Image Sensor with a Lateral Overflow Drain and a Transparent Gate Electrode," *Proc. IEEE International Electron Device Meeting*, Washington, D.C. (Dec. 1990).
5. B. Burkey et al., "The Pinned Photodiode for an Interline Transfer CCD Image Sensor," *Proc. 1984 International Electron Device Meeting*, Washington, D.C. (Dec. 1986), p. 28.
6. M. Sasaki, H. Ihara, and Y. Matsunaga, "A 2/3-in 400k-Pixel Sticking-Free Stact CCD Image Sensor," *IEEE Trans. Electron Devices*, **ED-28**(11) (Nov. 1993).
7. H. Yamashita et al., "A 2/3-Inch 2-Megapixel Stack CCD Imager," *Proceedings of the 1994 International Solid State Circuits Conference*, p. 224.
8. E. G. Stevens et al., "A 1-Megapixel, Progressive-Scan Image Sensor with Antiblooming Control and Lag-Free Operation," *IEEE Trans. Electron Devices*, **ED-38**(5) (May 1991).
9. D. K. Schroder, "The Concept of Generation and Recombination Lifetimes in Semiconductors," *IEEE Trans. Electron Devices* **ED-29**(8) (Aug. 1982), p. 17.
10. J. Van der Spiegel and G. J. Declerck, *Solid State Electronics* **27**:147 (1984).
11. Carlo, H. Sequin and Michael F. Tompsett, *Charge Transfer Devices Advances in Electronics and Electron Physics*, suppl. 8, Academic Press, Inc., New York, 1975.
12. M. J. Howes and D. V. Morgan, (eds.), *Charge Coupled Devices and Systems*, Wiley, Chicago, 1979.
13. D. F. Barbe et al., *Charge-Coupled Devices*, Springer-Verlag, New York, 1980.
14. M. Azuma et al., "A Fixed Pattern Noise Free 2/3", 1.3 M-Pixel CCD Image Sensor for HDTV Camera System," *Proc. 1991 IEEE International Solid-State Circuits Conference*, San Francisco, 212–213 (Feb. 1991).
15. Y. Matsunaga and S. Ohsawa, "1/3 Inch Interline Transfer CCD Image Sensor with Negative Feedback 94 dB Dynamic Range Charge Detector," *Proc. 1991 IEEE International Solid-State Circuits Conference*, San Francisco, 210–211 (Feb. 1991).
16. J. Janesick, "Open Pinned-Phase CCD Technology," *Proc. SPIE*, vol. 1159, San Diego (Nov. 1989).
17. J. Hyncek, "Virtual Phase Technology: A New Approach to Fabrication of Large-Area CCDs," *IEEE Trans. Electron Devices* **ED-28**(5):483–489 (May 1981).
18. J. Janesick et al., "Fano-Noise-Limited CCDs," *Optical and Optoelectric Applied Science and Engineering Symposium: X-Ray Instrumentation in Astronomy*, San Diego (Aug. 14–19, 1988).
19. J. Carbone et al., "New Low Noise Random Access, Radiation Resistant and Large Format Charge Injection Device (CID) Imagers," *Proc. of SPIE Conference 1900* (Jan. 31–Feb. 4, 1993).
20. M. Sugawara et al., "An Amplified MOS Imager Suited for Image Processing," *Proc. 1994 IEEE International Solid-State Circuits Conference*, San Francisco, 228–229 (Feb. 1994).
21. M. Ogata et al., "A Small Pixel CMD Image Sensor," *IEEE Trans. Electron Devices*, **ED-38**(5):1005–1010 (May 1991).
22. D. H. McCann, M. H. White, A. P. Turley, and R. A. Frosch, "Time Delay and Integration Detectors Using Charge Transfer Devices," U.S. Patent Number 4,280,141, July 1981.
23. Thompson et al., "Time-Delay-and-Integration Charge-Coupled Devices Using Tin Oxide Gate Technology," *IEEE Trans. Electron Devices* **ED-25**(2):132–134 (Feb. 1978).
24. S. Ohba et al., "MOS Area Sensor: Part 2: Low-Noise MOS Area Sensor with Antiblooming Photodiodes," *IEEE Trans. Electron Devices* **ED-27**(8):1682–1687 (Aug. 1980).
25. T. Nakamura et al., "A New MOS Image Sensor Operating in a Non-Destructive Readout Mode," *Proc. 1986 IEEE International Electron Device Meeting*, Washington, D.C., 353–356 (Dec. 1986).
26. Nomoto et al., "A 2/3 Inch 2M-Pixel CMD Image Sensor with Multi-Scanning Functions," *Proc. 1993 IEEE International Solid-State Circuits Conference*, San Francisco, 196–197 (Feb. 1993).

27. D. H. McCann et al., "Buried Channel CCD Imaging Arrays with Tin Oxide Transparent Gates," *IEEE International Solid State Circuits Conference*, pp. 30, 31, 261, 262 (Feb. 1978).
28. K. Harada, "A 2/3-Inch, 2M-Pixel Frame Interline Transfer CCD HDTV Image Sensor," *Proc. 1992 IEEE International Solid-State Circuits Conference*, San Francisco, 170–171 (Feb. 1992).
29. E. Stevens, "Photoresponse Nonlinearity of Solid State Image Sensors with Antiblooming Protection," *IEEE Trans. Electron Devices* **ED-38**(2):229–302 (Feb. 1991).
30. A. Weiss et al., *J. Electrochem. Soc.* 133:110C (1986).
31. P. Dillion et al., "Fabrication and Performance of Color Filter Arrays for Solid-State Imagers," *IEEE Trans. Electron Devices* **ED-25**:97–101 (1978).
32. P. Dillion, D. Lewis, and F. Kaspar, "Color Imaging System Using a Single CCD Array," *IEEE Trans. Electron Devices*, **ED-25**:102–107 (Feb. 1978).
33. K. A. Parulski, et al., "A Digital Color CCD Imaging System Using Custom VLSI Circuits," *IEEE Trans. Consumer Electronics* **35**(3):382–389 (Aug. 1989).
34. K. A. Parulski, "Color Filters and Processing Alternatives for One-Chip Cameras," *IEEE Trans. Electron Devices* **ED-32**(8):1381–1389 (Aug. 1985).
35. J. E. Greivenkamp, "Color Dependent Optical Prefilter for the Suppression of Aliasing Artifacts," *Appl. Opt.* **29**(5):676–684, (Feb. 1990).

INFRARED DETECTOR ARRAYS

Lester J. Kozłowski

*Altasens, Inc.
Westlake Village, California*

Walter F. Kosonocky*

*New Jersey Institute of Technology
University Heights
Newark, New Jersey*

33.1 GLOSSARY

A_{det}	detector area
A_I	gate modulation current gain (ratio of integration capacitor current to load current)
A_V	amplifier voltage gain
C_{amp}	amplifier capacitance
C_{det}	detector capacitance
C_{FIS}	fill-and-spill gate capacitance for a Tompsett type CCD input
C_{fb}	CTIA feedback or Miller capacitance
C_{gd}	FET gate-drain overlap capacitance
C_{gs}	FET gate-source capacitance
C_L	CTIA band-limiting load capacitance
C_{out}	sense node capacitance at the CCD output
C_T	effective feedback (transcapacitance) or integration capacitance for a capacitive transimpedance amplifier
$C_{T\lambda}$	spectral photon contrast
cte	charge transfer efficiency
$D_{\lambda\text{pk}}^*$	peak detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
D_{bb}^*	blackbody detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
D_{th}^*	thermal detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
e^-	electron
E_g	detector energy gap
$f/\#$	conventional shorthand for the ratio of the focal length of a lens to its diameter
f_{chop}	chopper frequency

*Deceased.

f_{frame}	display frame rate
f_{knee}	frequency at which the $1/f$ noise intersects the broadband noise
f_s	spatial frequency (cycles/radian)
$g_{m, \text{LOAD}}$	gate transconductance of the load FET in the gate-modulated input circuit
g_m	gate transconductance of a Field Effect Transistor
h	Planck's constant
I_D	FET drain current
I_{det}	detector current
I_{photo}	detector photocurrent
k	Boltzmann constant
K_{amp}	amplifier FET noise spectral density at 1 Hz
K_{det}	detector noise spectral density at 1 Hz
K_{FET}	FET noise spectral density at 1 Hz
L	length-to-width ratio of a bar chart (always set to 7)
MRT	minimum resolvable temperature (K)
MTF	modulation transfer function for the optics, detector, readout, the integration process, or the composite sensor
n	detector junction ideality or diffusion constant
$N_{\text{amp}, 1/f}$	number of noise carriers for one integration time due to amplifier FET $1/f$ noise
$N_{\text{amp}, \text{white}}$	number of noise carriers for one integration time due to amplifier FET white noise
N_c	number of photo-generated carriers integrated for one integration time
n_{det}	detector junction ideality or diffusion constant
NE ΔT	noise equivalent temperature difference (K)
n_{FET}	subthreshold FET ideality
N_{FPA}	composite (total) FPA noise in carriers
$N_{\text{KTC, channel}}$	CTIA broadband channel noise in carriers
$N_{\text{load, white}}$	number of noise carriers for one integration time due to CTIA load FET white noise
N_{os}	display overscan ratio
N_{PHOTON}	shot noise of photon background in carriers
NSD	noise spectral density of a detector or field effect transistor; the $1/f$ noise is often specified by the NSD at a frequency of 1 Hz
N_{sf}	source follower noise
N_{ss}	serial scan ratio
q	electron charge in coulombs
Q_B	photon flux density (photons/cm ² -s) incident on a focal plane array
Q_D	charge detected in a focal plane array for one integration time
Q_{max}	maximum charge signal at saturation
R_{det}	detector resistance
R_{LOAD}	gate modulation load resistance
R_0	detector resistance at zero-bias resistance
R_0A	detector resistance-area product at zero-bias voltage
R_r	detector resistance in reverse-bias resistance
S/N	signal-to-noise ratio
SNR _T	target signal-to-noise ratio
S_V	readout conversion factor describing the ratio of output voltage to detected signal carriers
T	operating temperature

tce	thermal coefficient of expansion
TCR	thermal coefficient of resistance for bolometer detectors
T_D	time constant for correlated double sampling process normally set by Nyquist rate
t_{int}	integration time
U	residual nonuniformity
V_{br}	detector reverse-bias breakdown voltage, sometimes defined as the voltage where $R_r = R_0$
V_D	FET drain voltage
V_{det}	detector bias voltage
V_{DS}	FET drain-to-source voltage
V_G	FET gate voltage
v_n	measured rms noise voltage
ΔA_I	gate modulation current gain nonuniformity
Δf	noise bandwidth (Hz)
ΔI_{photo}	differential photocurrent
ΔT	scene temperature difference creating differential photocurrent ΔI_{photo}
ΔV_S	signal voltage for differential photocurrent ΔI_{photo}
Δx	horizontal detector subtense (mradian)
Δy	vertical detector subtense (mradian)
η	detector quantum efficiency
η_{BLIP}	percentage of BLIP
$\eta_{\text{inj, DI}}$	injection efficiency of detector current into the source-modulated FET of the direct injection input circuit
η_{inj}	injection efficiency of detector current
η_{noise}	injection efficiency of DI circuit noise into integration capacitor
η_{pc}	quantum efficiency of photoconductive detector
η_{pv}	quantum efficiency of photovoltaic detector
λ_c	detector cutoff wavelength (50 percent of peak response, μm)
σ_{det}	noise spectral density of total detector noise including photon noise
$\sigma_{\text{input, ir}}$	noise spectral density of input-referred input circuit noise
σ_{LOAD}	noise spectral density of input-referred load noise
$\sigma_{\text{mux, ir}}$	noise spectral density of input-referred multiplexer noise
σ_{VT}	rms threshold voltage nonuniformity across an FPA
τ_{amp}	amplifier time constant (s)
τ_{eye}	eye integration time (s)
τ_o	optical transmission
ω	angular frequency (radians)
$\langle e_{\text{amp}} \rangle$	buffer amplifier noise for buffered direct injection circuit

33.2 INTRODUCTION

Infrared sensors have been available since the 1940s to detect, measure, and image the thermal radiation emitted by all objects. Due to advanced detector materials and microelectronics, large scanning and staring focal plane arrays (FPA) with few defects are now readily available in the short wavelength infrared (SWIR; 1 to 3 μm), medium wavelength infrared (MWIR; ≈ 3 to 5 μm), and long wavelength infrared (LWIR; ≈ 8 to 14 μm) spectral bands. We discuss in this chapter the disparate FPA technologies, including photon and thermal detectors, with emphasis on the emerging types.

IR sensor development has been driven largely by the military. Detector requirements for missile seekers and forward looking infrared (FLIR) sensors led to high-volume production of photoconductive (PC) HgCdTe arrays starting in the 1970s. Though each detector requires direct connection to external electronics for purposes of biasing, signal-to-noise ratio (SNR) enhancement via time delay integration (TDI), and signal output, the first-generation FPAs displaced the incumbent Pb-salt (PbS, PbSe) and Hg-doped germanium devices, and are currently being refined using custom analog signal processing,¹ laser-trimmed solid-state preamplifiers, etc.

Size and performance limitations of first-generation FLIRs necessitated development of self-multiplexed FPAs with on-chip signal processing. Second-generation thermal imaging systems use high-density FPAs with relatively few external connections. Having many detectors that integrate longer, low-noise multiplexing and on-chip TDI (in some scanning arrays), second-generation FPAs offer higher performance and design flexibility. Video artifacts are suppressed due to the departure from ac-coupling and interlaced raster scan, and external connections are minimized. Fabricated in monolithic and hybrid methodologies, many detector and readout types are used in two basic architectures (staring and scanning). In a monolithic FPA, the detector array and the multiplexing signal processor are integrated in a single substrate. The constituents are fabricated on separate substrates and interconnected in a hybrid FPA.

FPAs use either photon or thermal detectors. Photon detection is accomplished using intrinsic or extrinsic semiconductors and either photovoltaic (PV), photoconductive (PC), or metal insulator semiconductor (MIS) technologies. Thermal detection relies on capacitive (ferro- and pyroelectric) or resistive bolometers. In all cases, the detector signal is coupled into a multiplexer and read out in a video format.

Infrared Applications

Infrared FPAs are now being applied to a rapidly growing number of civilian, military, and scientific applications such as industrial robotics and thermography (e.g., electrical and mechanical fault detection), medical diagnosis, environmental and chemical process monitoring, Fourier transform IR spectroscopy and spectroradiometry, forensic drug analysis, microscopy, and astronomy. The combination of high sensitivity and passive operation is also leading to many commercial uses. The passive monitoring provided by the addition of infrared detection to gas chromatography-mass spectroscopy (GC-MS), for example, yields positive chemical compound and isomer detection without sample alteration. Fusing IR data with standard GC-MS aids in the rapid discrimination of the closely related compounds stemming from drug synthesis. Near-IR (0.7 to 0.1 μm) and SWIR spectroscopy and fluorescence are very interesting near-term commercial applications since they pave the way for high-performance FPAs in the photochemical, pharmaceutical, pulp and paper, biomedical, reference quantum counter, and materials research fields. Sensitive atomic and molecular spectroscopies (luminescence, absorption, emission, and Raman) require FPAs having high quantum efficiency, low dark current, linear transimpedance, and low read noise.

Spectral Bands

The primary spectral bands for infrared imaging are 3 to 5 and 8 to 12 μm because atmospheric transmission is highest in these bands. These two bands, however, differ dramatically with respect to contrast, background signal, scene characteristics, atmospheric transmission under diverse weather conditions, and optical aperture constraints. System performance is a complex combination of these and the ideal system requires dual band operation. Factors favoring the MWIR include its higher contrast, superior clear-weather performance, higher transmissivity in high humidity, and higher resolution due to $\sim 3 \times$ smaller optical diffraction. Factors favoring the LWIR include much-reduced background clutter (solar glint and high-temperature countermeasures including fires and flares have much-reduced emission), better performance in fog, dust, and winter haze, and higher immunity to atmospheric turbulence. A final factor favoring the LWIR, higher

S/N ratio due to the greater radiance levels, is currently moot because of technology limitations. Due to space constraints and the breadth of sensor applicability, we focus on target/background metrics in this section.

The signal collected by a visible detector has higher daytime contrast than either IR band because it is mainly radiation from high-temperature sources that is subsequently reflected off earth-based (ambient temperature; ≈ 290 K) objects. The high-temperature sources are both solar (including the sun, moon, and stars) and synthetic. Since the photon flux from high- and low-temperature sources differs greatly at visible wavelengths from day to night, scene contrasts of up to 100 percent ensue.

Reflected solar radiation has less influence as the wavelength increases to a few microns since the background radiation increases rapidly and the contrast decreases. In the SWIR band, for example, the photon flux density from the earth is comparable to visible room light (10^{13} photons/cm²-s). The MWIR band ($\sim 10^{15}$ photon flux density) has lower, yet still dynamic, daytime contrast, and can still be photon-starved in cold weather or at night.

The net contribution from reflected solar radiation is even lower at longer wavelengths. In the LWIR band, the background flux is equivalent to bright sunlight ($\approx 10^{17}$ photons/cm²-s). This band thus has even lower contrast and much less background clutter, but the “scene” and target/background metrics are similar day and night. Clear-weather performance is relatively constant.

Depending on environmental conditions, however, IR sensors operating in either band must discern direct emission from objects having temperatures very near the average background temperature (290 K) in the presence of the large background and degraded atmospheric transmissivity. Under conditions of uniform thermal soak, such as at diurnal equilibrium, the target signal stems from minute emissivity differences.

The spectral photon incidence for a full hemispheric surround is

$$Q = \tau_{cf} \int_{\lambda_1}^{\lambda_2} Q_{\lambda}(\lambda) d\lambda \quad (1)$$

if a zero-emissivity bandpass filter having in-band transmission τ_{cf} , cut-on wavelength λ_1 , and cutoff wavelength λ_2 is used (zero emissivity obtained practically by cooling the spectral filter to a temperature where its self-radiation is negligible). The photon flux density, Q_B (photons/cm²-s), incident on a focal plane array is

$$Q_B = \frac{1}{4(f/\#)^2 + 1} Q \quad (2)$$

where $f/\#$ is the conventional shorthand for the ratio of the focal length to the diameter (assumed circular) of the limiting aperture or lens. The cold shield $f/\#$ limits the background radiation to a field-of-view consistent with the warm optics to eliminate extraneous background flux and concomitant noise. The background flux in the LWIR band is approximately two orders of magnitude higher than in the MWIR.

The spectral photon contrast, $C_{T\lambda}$, is the ratio of the derivative of spectral photon incidence to the spectral photon incidence, has units K^{-1} , and is defined

$$C_{T\lambda} = \left(\frac{\partial Q}{\partial T} \right) / \left(\frac{Q}{T} \right) \quad (3)$$

Figure 1 is a plot of $C_{T\lambda}$ for several MWIR subbands (including 3.5 to 5, 3.5 to 4.1, and 4.5 to 5 μm) and the 8.0 to 12 μm LWIR spectral band. The contrast in the MWIR bands at 300 K is 3.5 to 4 percent compared to 1.6 percent for the LWIR band. While daytime MWIR contrast is even higher due to reflected sunlight, an LWIR FPA offers higher sensitivity if it has the larger capacity needed for storing the larger amounts of photogenerated (due to the higher background flux) and detector-generated carriers (due to the narrow bandgap). The photon contrast and the background flux are key parameters that determine thermal resolution as will be described later under “Performance Figures of Merit.”

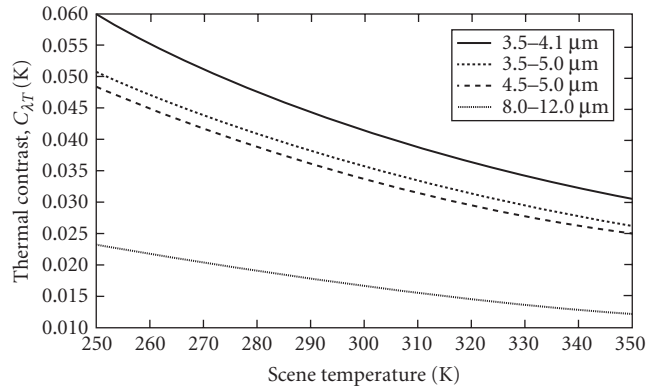


FIGURE 1 Spectral photon contrast in the MWIR and LWIR.

Scanning and Staring Arrays

The two basic types of FPA are scanning and staring. The simplest scanning device consists of a linear array as shown in Fig. 2a. An image is generated by scanning the scene across the strip. Since each detector scans the complete horizontal field-of-view (one video raster line) at standard video frame rates, each resolution element or pixel has a short integration time and the total detected charge can usually be accommodated.

A staring array (Fig. 2b) is the two-dimensional extension of a scanning array. It is self-scanned electronically, can provide enhanced sensitivity, and is suitable for lightweight cameras. Each pixel

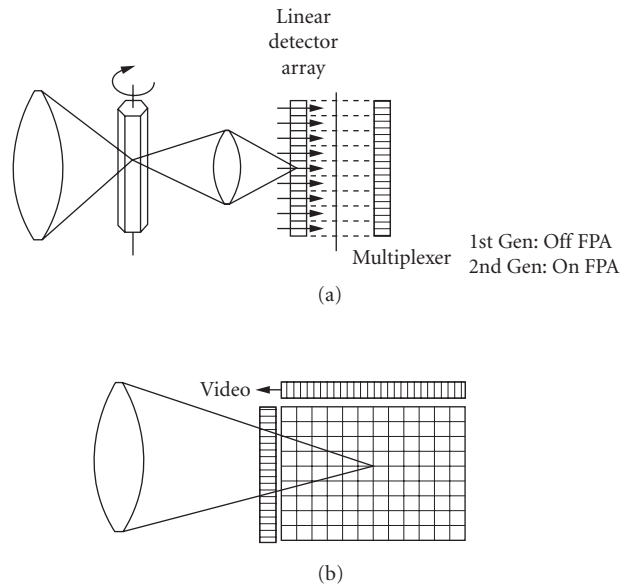


FIGURE 2 Scanning (a) and staring (b) focal plane arrays.

is a dedicated resolution element, but synchronized dithering of sparsely populated arrays is sometimes used to enhance the effective resolution, minimize spatial aliasing, and increase the effective number of pixels. Although theoretically charge can be integrated for the full frame time, the charge-handling capacity is inadequate at terrestrial LWIR backgrounds.

Detectors

Infrared detectors convert IR photons and energy to electrical signals. Many types are used in FPAs (as shown in Fig. 3²) including photon and thermal detectors that address diverse requirements spanning operating temperatures from 4 K to room temperature. Figure 4 compares the quantum efficiencies of several detector materials.

Intrinsic detectors³ usually operate at higher temperatures than extrinsic devices, have higher quantum efficiencies, and dissipate less power. Backside-illuminated devices, consisting of an absorbing epitaxial layer on a transparent substrate, are used in hybrid FPAs and offer the advantages of nearly 100 percent active detector area, good mechanical support, and high quantum efficiency. The most popular intrinsic photovoltaics are HgCdTe and InSb. These detectors are characterized by their quantum efficiency (η), zero-bias resistance (R_0), reverse-bias resistance (R_r), junction ideality or diffusion constant n , excess noise (if any) versus bias, and reverse-bias breakdown voltage (V_{br}), which is sometimes defined as the voltage where $R_r = R_0$.

PtSi is a photovoltaic Schottky barrier detector (SBD) which is the most mature for large monolithic FPAs. IR detection is via internal photoemission over a Schottky barrier (0.21 to 0.23 eV). Characteristics include low (≈ 0.5 percent for broadband 3.5 to 5.0 μm) but very uniform quantum efficiency, high producibility that is limited only by the Si readout circuits, full compatibility with VLSI technology, and soft spectral response with peak below 2 μm and zero response just beyond 6 μm . Internal photoemission dark current requires cooling below 77 K.

HgCdTe is the most popular intrinsic photoconductor, and various linear arrays in several scanning formats are used worldwide in first-generation FLIRs. For reasons of producibility and cost, HgCdTe photoconductors have historically enjoyed a greater utilization than PV detectors despite the latter's higher quantum efficiency, higher D^* by a factor of $(2\eta_{pv}/\eta_{pc})^{1/2}$, and superior modulation transfer function (MTF). Nevertheless, not all photoconductors are good candidates for FPAs due to their low detector impedance. This includes the intrinsic materials InSb and HgCdTe.

The most popular photoconductive material system for area arrays is doped extrinsic silicon (Si: x ; where x is In, As, Ga, Sb, etc.), which is made in either conventional or impurity band conduction [IBC or blocked impurity band (BIB)] technologies. Early monolithic arrays were doped-Si devices, due primarily to compatibility with the silicon readout. Extrinsic photoconductors must be made relatively thick (up to 30 mils; doping density of IBCs, however, minimizes this thickness requirement but does not eliminate it) because they have much lower photon capture cross section than intrinsic detectors. This factor adversely affects their MTF in systems having fast optics.

Historically, Si:Ga and Si:In were the first mosaic focal plane array PC detector materials because early monolithic approaches were compatible with these dopants. Nevertheless, problems in fabricating the detector contacts, early breakdown between the epitaxial layer and the detector material (double injection), and the need for elevated operating temperatures helped force the general move to monolithic PtSi and intrinsic hybrids.

The most advanced extrinsic photoconductors are IBC detectors using Si:As and Si:Ga.⁴ These have reduced recombination noise (negating the $\sqrt{2}$ superiority in S/N that PV devices normally have) and longer spectral response than standard extrinsic devices due to the higher dopant levels. IBC detectors have a unique combination of PC and PV characteristics, including extremely high impedance, PV-like noise (reduced recombination noise since IBC detectors collect carriers both from the continuum and the "hopping" impurity band), linear photoconductive gain, high uniformity, and superb stability. The photo-sensitive layer in IBCs is heavily doped to achieve hopping-type conduction. A thin, lightly doped ($10^{10}/\text{cm}^2$) silicon layer blocks the hopping current before

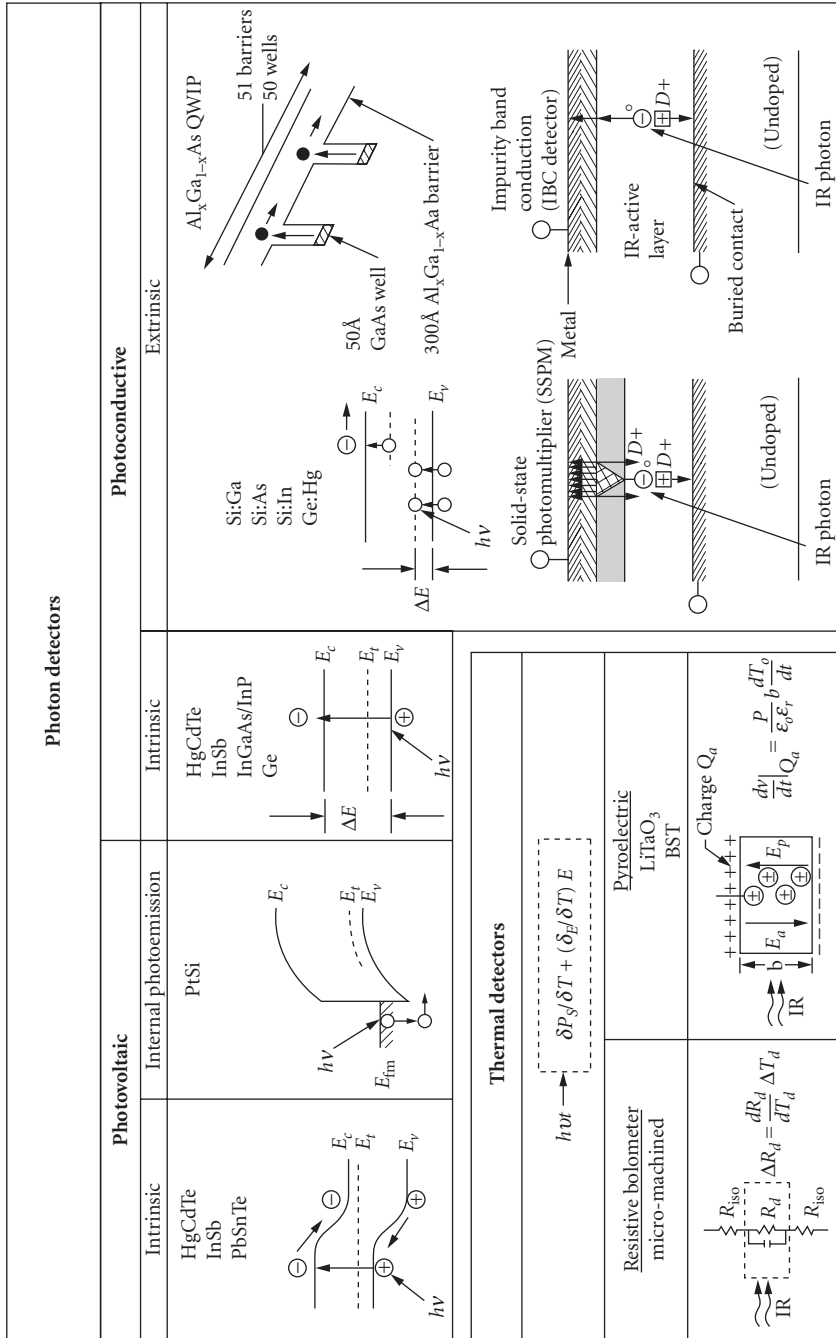


FIGURE 3 Photon and thermal detectors.

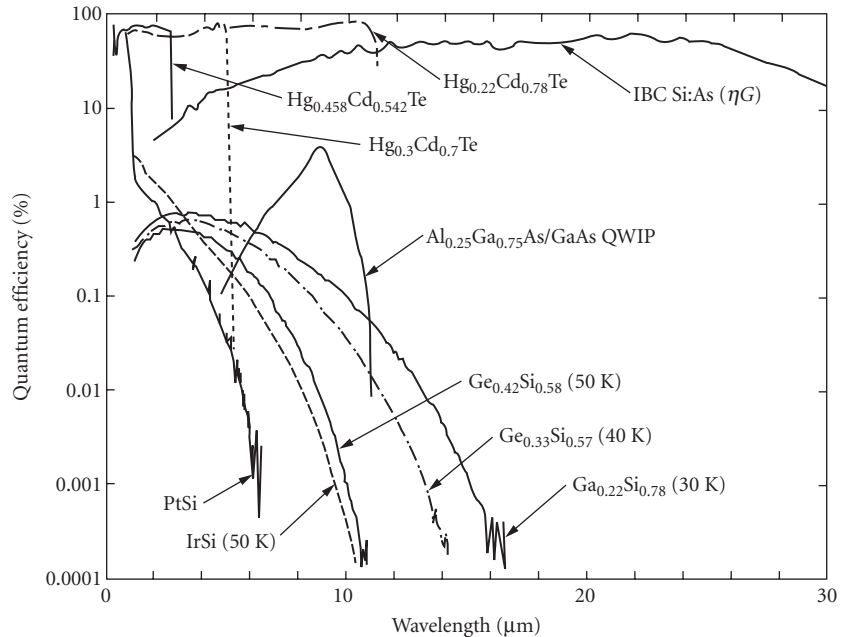


FIGURE 4 Quantum efficiency versus wavelength for several detector materials.

it reaches the device electrode to reduce noise. Specially doped IBCs (see cross-sectional views in Fig. 3) operate as solid-state photomultipliers (SSPM) and visible light photon counters (VLPC) in which photoexcited carriers are amplified by impact ionization of impurity-bound carriers.⁵ The amplification allows counting of individual photons at low flux levels. Standard SSPMs respond from 0.4 to 28 μm .

An alternative custom tunable detector is the GaAs/AlGaAs quantum well infrared photodetector (QWIP). Various QWIP photoconductive⁶ and photovoltaic⁷ structures are being investigated as low-cost alternatives to II-VI LWIR detectors like HgCdTe. Infrared detection in the typical PC QWIP is via intersubband or bound-to-extended-state transitions within the multiple quantum well superlattice structure. Due to the polarization selection rules for transitions between the first and second quantum wells, the photon electric field must have a component parallel to the superlattice direction. Light absorption in *n*-type material is thus anisotropic with zero absorption at normal incidence. The QWIP detector's spectral response is narrowband, peaked about the absorption energy. The wavelength of peak response can be adjusted via quantum well parameters and can be made bias-dependent.

Various bolometers, both resistive and capacitive (pyroelectric), are also available. Bolometers sense incident radiation via energy absorption and concomitant change in device temperature in both cooled moderate-performance and uncooled lower-performance schemes. Much recent research, which was previously highly classified, has focused on both hybrid and monolithic uncooled arrays and has yielded significant improvements in the detectivity of both resistive and capacitive bolometer arrays. The resistive bolometers currently in development consist of a thin film of a temperature-sensitive resistive material film which is suspended above a silicon readout. The pixel support struts provide electrical interconnect and high thermal resistance to maximize pixel sensitivity. Recent work has focused on the micromachining necessary to fabricate mosaics with low thermal conductance using monolithic methodologies compatible with silicon.

Capacitive bolometers sense a change in elemental capacitance and require mechanical chopping to detect incident radiation. The most common are pyroelectric detectors. J. Cooper⁸ suggested the use of pyroelectric detectors in 1962 as a possible solution for applications needing a low-cost IR FPA with acceptable performance. These devices have temperature-dependent spontaneous polarization. Ferroelectric detectors are pyroelectric detectors having reversible polarization. There are over a thousand pyroelectric crystals, including several popularly used in hybrid FPAs; e.g. lithium tantalate (LiTaO_3), triglycine sulfate (TGS), and barium strontium titanate⁹ (BaSrTiO_3).

33.3 MONOLITHIC FPAs

A monolithic FPA consists of a detector array and the readout multiplexer integrated on the same substrate. The progress in the development of the monolithic FPAs in the last two decades has been strongly influenced by the rapid advances in the silicon VLSI technology. Therefore, the present monolithic FPAs can be divided into three categories reflecting their relationship to the silicon VLSI technology. The first category includes the “complete” monolithic FPAs in which the detector array and the readout multiplexer are integrated on the same silicon substrate using processing steps compatible with the silicon VLSI technology. They include the extrinsic Si FPAs reported initially in the 1970s,¹⁰ FPAs with Schottky barrier,¹¹ heterojunction detector FPAs, and microbolometer FPAs.¹²

The second category will be referred to here as the “partial monolithic” FPAs. This group includes narrowband detector arrays of HgCdTe ¹³ and InSb ¹⁴ integrated on the same substrate only with the first level of multiplexing, such as the row and column readout from a two-dimensional detector array. In this case the multiplexing of the detected signal is completed by additional silicon IC chips usually packaged on the imager focal plane.

The third category represents “vertically integrated” photodiode (VIP) FPAs. These FPAs are functionally similar to hybrids in the sense that a silicon readout multiplexer is used with the narrow-bandgap HgCdTe detectors. However, while in the hybrid FPA the completed HgCdTe detector array is typically connected by pressure contacts via indium bumps to the silicon multiplex pads; in the case of the vertically integrated FPAs, HgCdTe chips are attached to a silicon multiplexer wafer and then the fabrication of the HgCdTe photodiodes is completed including the deposition and the definition of the metal connections to the silicon readout multiplexer.

In the following sections we will review the detector readout structures, and the main monolithic FPA technologies. It should also be noted that most of the detector readout techniques and the architectures for the monolithic FPAs were originally introduced for visible silicon imagers. This heritage is reflected in the terminology used in the section.

Architectures

The most common structures for the photon detector readout and architectures of monolithic FPAs are illustrated schematically in Figs. 5 and 6.

MIS Photogate FPAs: CCD, CID, and CIM Most of the present monolithic FPAs use either MIS photogates or photodiodes as the photon detectors. Figure 5 illustrates a direct integration of the detected charge in the potential well of a MIS (photogate) detector for a charge coupled device (CCD) readout in (a), a charge injection device (CID) readout in (b), and a charge-integration matrix (CIM) readout in (c). The unique characteristic of the CCD readout is the complete transfer of charge from the integration well without readout noise. Also in a CCD FPA the detected charge, Q_D , can be transferred via potential wells along the surface of the semiconductor that are induced by clock voltages but isolated from the electrical pickup until it is detected by a low-capacitance (low kTC noise) on-chip amplifier. However, because of the relatively large charge transfer losses ($\sim 10^{-3}$ per transfer) and a limited charge-handling capacity, the use of nonsilicon CCD readout has

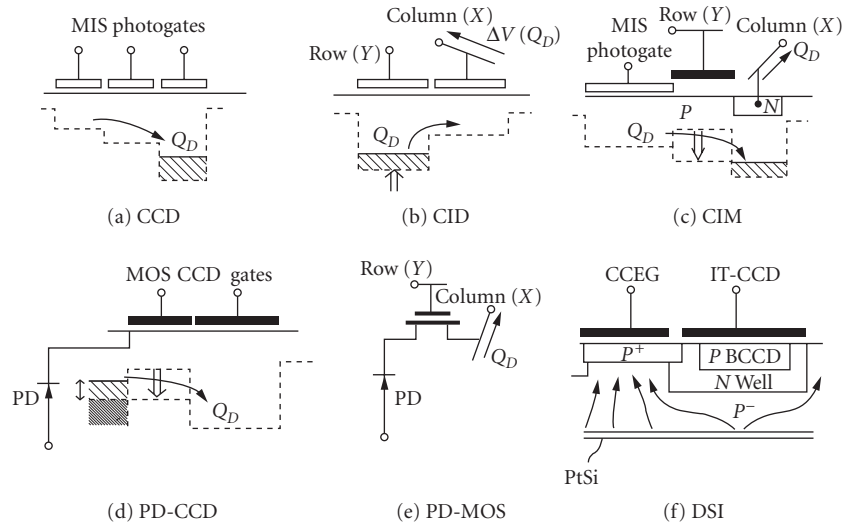


FIGURE 5 Photodetector readout structures.

been limited mainly to HgCdTe TDI FPAs. Such TDI imager is shown in Fig. 6a as a frame transfer (FT) type CCD area imager performing a function of a line sensor with the effective optical integration time increased by the number of TDI elements (CCD stages) in the column CCD registers. In this imager the transfer of the detected charge signal between CCD wells of the vertical register is adjusted to coincide with the mechanical motion of the image.

In the FT-CCD TDI FPA, the vertical registers perform the functions of charge detection and integration as well as transfer. The detected image is transferred one line at a time from the parallel vertical registers to the serial output registers. From there it is transferred at high clock rate to produce the output video. Similar TDI operation can also be produced by the interline-transfer (IT) CCD architecture shown in Fig. 6e. However, in the case of IT-CCD readout, the conversion of infrared radiation into charge signal photodetection is performed by photodiodes.

In the CID readout, see Fig. 5b, the detected charge signal is transferred back and forth between the potential wells of the MIS photogates for nondestructive X - Y addressable readout, $\Delta V(Q_D)$, that is available at a column (or a row) electrode due to the displacement current induced by the transfer of the detected charge signal, Q_D . At the end of the optical integration time, the detected charge is injected into the substrate by driving both MIS capacitors into accumulation.

CID FPAs with column readout for single-output-port and parallel-row readout are illustrated schematically in Fig. 6b and c, respectively. Another example of a parallel readout is the CIM FPA shown in Fig. 6d. The parallel readout of CID and CIM FPAs is used to overcome the inherent limitation on charge-handling capacity of these monolithic FPAs by allowing a short optical integration time with fast frame readout and off-chip charge integration by supporting silicon ICs.

Silicon FPAs: IT-CCD, CSD, and MOS FPAs The monolithic FPAs fabricated on silicon substrate take advantage of well-developed silicon VLSI process technology. Therefore, silicon ($E_g = 1.1$ eV), which is transparent to infrared radiation having wavelength longer than $1.0 \mu\text{m}$, is often used to produce monolithic CCD and MOS FPAs with infrared detectors that can be formed on silicon substrate. In the 1970s there was great interest in the development of monolithic silicon FPAs with extrinsic Si:In and Si:Ga photoconductors. However, since the early 1980s most progress was reported on monolithic FPAs with Schottky-barrier photodiodes, GeSi/Si heterojunction photodiodes, vertically integrated photodiodes, and resistive microbolometers. With the exception of the

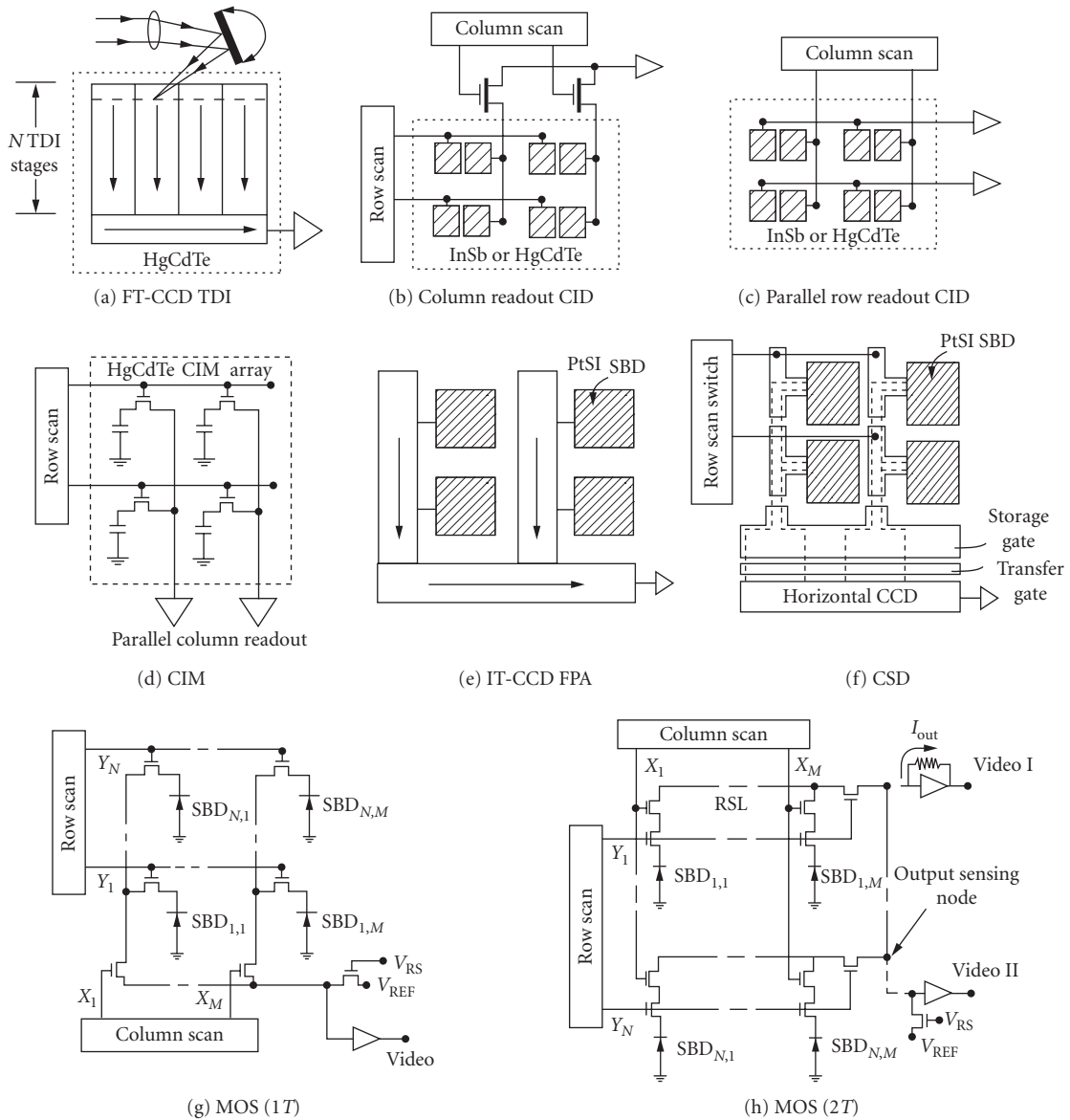


FIGURE 6 Monolithic FPA architectures.

resistive microbolometers, in the form of thin-film semiconductor photoresistors formed on micro-machined silicon structures, all of the above infrared detectors can be considered to be photodiodes and can be read out either by IT-CCD or MOS monolithic multiplexer.

The photodiode (PD) CCD readout, see Fig. 5d, is normally organized as the interline-transfer (IT) CCD staring FPA, shown in Fig. 6e. The IT-CCD readout has been used mostly for PtSi Schottky-barrier detectors (SBDs). The operation of this FPA consists of direct integration of the detected charge signal on the capacitance of the photodiode. At the end of the optical integration time, a frame readout is initiated by a parallel transfer of the detected charge from the photodiodes to the

parallel vertical CCD registers. From there the detected image moves by parallel transfers one line at a time into the horizontal output register for a high-clock-rate serial readout.

The design of the Schottky-barrier IT-CCD FPA involves a trade-off between fill factor (representing the ratio of the active detector area to the pixel area) and maximum saturation charge signal (Q_{mix}). This trade-off can be improved by the charge sweep device (CSD) architecture, shown in Fig. 6f, that has also been used as a monolithic readout multiplexer for PtSi and IrSi SBDs.

In CSD FPA the maximum charge signal is limited only by the SBD capacitance since its operation is based on transferring the detected charge signal from one horizontal line corresponding to one or two rows of SBDs (depending on the type of the interlacing used) into minimum geometry vertical CCD registers. During the serial readout of the previous horizontal line, the charge signal is swept into a potential well under the storage gate by low-voltage parallel clocking of the vertical registers. Then, during the horizontal blanking time, the line charge signal is transferred in parallel to the horizontal CCD register for serial readout during the next horizontal line time.

The main advantage of the silicon CCD multiplexer is relatively low readout noise, from a few electrons to the order of several tens of electrons (depending on video rate, sense capacitance, and CCD technology), so that a shot-noise-limited operation can be achieved at relatively low signal levels. But as the operating temperature is lowered below 60 K, the charge transfer losses of buried-channel CCDs (BCCDs) become excessive due to the freeze-out of the BCCD implant. Therefore, detectors requiring operation at 40 K or lower are more compatible with MOS readout device technology.

Photodiode (PD) MOS readout (see Fig. 5e) represents another approach to construction of an X-Y addressable silicon multiplexer. These types of monolithic MOS multiplexers used for readout of PtSi SBDs are illustrated in Fig. 6g and h.

A single-output-port FPA with one MOSFET switch per detector, MOS (1T), is shown in Fig. 6g. During FPA readout, the vertical scan switch transfers the detected charge signal from one row of detectors to the column lines. Then the column lines are sequentially connected by MOSFET switches to the output sense line under the control of the horizontal scan switch. The main limitation of the MOS (1T) FPA is a relatively high readout noise (on the order of 100 electrons/pixel) due to sensing of small charge signals on large-capacitance column lines. This readout noise can be decreased with a row readout MOS (2T) FPA having two MOSFET switches per detector. In this case, low readout noise can be achieved using current sensing, it is limited by the noise of the amplifier, and for voltage sensing it can be reduced by correlated double sampling (CDS).¹⁵ A readout noise of 300 rms electrons/pixel was achieved at Sarnoff for a 640×480 low-noise PtSi MOS (2T) FPA designed with row buffers and 8:1 multiplexing of the output lines;¹⁶ 2T MOS FPA read noise of $60e^-$ was later achieved by reading via capacitive transimpedance amplifier column buffers in $0.5 \mu\text{m}$ CMOS technology.¹⁷

An alternative form of the MOS (1T) FPA architecture is an MOS FPA with parallel column readout for fast frame operation. This silicon VIP FPA, resembling CIM architecture in Fig. 6d, can be used with HgCdTe vertically integrated PV detectors.

Direct-Charge-Injection Silicon FPAs All of the silicon monolithic FPAs thus far described use separately defined detectors. A direct-charge-injection type monolithic silicon FPA with a single detector surface is a PtSi direct Schottky injection (DSI) imager that is made on thinned silicon substrate having a CCD or MOS readout on one side and PtSi SBD charge-detecting surface on the other side.¹⁸ A cross-sectional area of one pixel of this FPA for IT-CCD readout is shown in Fig. 5f. In the operation of this imager, the *p*-type buried-channel CCD formed in an *n*-well removes charge from a P^+ charge-collecting electrode that in turn depletes a high-resistivity *p*-type substrate. Holes injected from the PtSi SBD surface into the *p*-type substrate drift through the depleted *p*-type substrate to the P^+ charge-collecting electrode. The advantages of the DSI FPA include 100 percent fill factor, a large maximum charge due to the large capacitance between the charge collecting electrode and the overlapping gate, and that the detecting surface does not have to be defined. A 128×128 IT-CCD PtSi DSI FPA was demonstrated;¹⁹ however, the same basic structure could also be used with other internal photoemission surfaces such as IrSb or Ge:Si.

Microbolometer FPAs A microbolometer FPA for uncooled applications consists of thin-film semiconductor photoresistors micromachined on a silicon substrate. The uncooled IR FPA is fabricated as an array of microbridges with a thermoresistive element in each microbridge. The

resistive microbolometers have high thermal coefficient of resistance (TCR) and low thermal conductance between the absorbing area and the readout circuit which multiplexes the IR signal. As each pixel absorbs IR radiation, the microbridge elemental resistance changes accordingly with its temperature.

Metal films have traditionally been used to make the best bolometer detectors because of their low $1/f$ noise. These latest devices use semiconductor films of 500 Å thickness having TCR of 2 percent per °C. The spacing between the microbridge and the substrate is selected to maximize the pixel absorption in the 8- to 14- μm wavelength range. Standard photolithographic techniques pattern the thin film to form detectors for individual pixels. The thin film TCR varies over an array by ± 1 percent, and produces responsivity of 70,000 V/W in response to 300 K radiation. This has been sufficient to yield 0.1°C NE ΔT with an $f/1$ lens. Potential low-cost arrays at prices similar to that of present large IC memories are possible with this technology.

Scanning and Staring Monolithic FPAs

In an earlier section we have reviewed the available architectures for the construction of two-dimensional scanning TDI, scanning, and staring FPAs. The same basic readout techniques, however, are also used for line-sensing imagers with photodiodes and MIS photogate detectors. For example, a line-scanning FPA corresponds to a vertical column CCD with the associated photodiodes of an IT-CCD, a column readout CID, a MOS (IT) FPA with only one row of detectors. However, since the design of a staring FPA is constrained by the size of the pixels, there is more space available for the readout multiplexer of a line of detectors. Therefore, design of the monolithic silicon multiplexer for a line detector array may also resemble the complexity of silicon multiplexer for hybrid FPAs.

33.4 HYBRID FPAs

Hybrid FPAs are made by interconnecting, via either direct or indirect means, a detector array to a multiplexing readout. Several approaches are pursued in both two- and three-dimensional configurations. Hybrids are typically made by either epoxying detector material to a processed silicon wafer (or readout) and subsequently forming the detectors and electrical interconnects by, for example, ion-milling; by mating a fully processed detector array to a readout to form a “two-dimensional” hybrid;²⁰ or by mating a fully processed detector array to a stack of signal processors to form a three-dimensional stack (*Z*-hybrid or “3D-IC”). The detector is usually mounted on top of the multiplexer and infrared radiation impinges on the backside of the detector array. Indium columns typically provide electrical and, often in conjunction with various epoxies, mechanical interconnect.

Hybrid methodology allows independent optimization of the detector array and the readout. Silicon is the preferred readout material due to performance and the leveraging of the continuous improvements funded by commercial markets. Diverse state-of-the-art processes and lithography are hence available at a fraction of their original development cost.

Thermal expansion match In a hybrid FPA, the detector array is attached to a multiplexer which can be of a different material. In cooling the device from room temperature to operating temperature, mechanical strain builds up in the hybrid due to the differing coefficients of thermal expansion. Hybrid integrity requires detector material that has minimum thermal expansion mismatch with silicon. Based on this criterion, the III-V and II-VI detectors are favored over Pb-salts. Silicon-based detectors are matched perfectly to the readout; these include doped-Si and PtSi. The issue of hybrid reliability has prompted the fabrication of II-VI detectors on alternative substrates to mitigate the mismatch. HgCdTe, for example, is being grown on sapphire (PACE-I),²¹ GaAs (liftoff techniques are available for substrate thinning or removal), and silicon in addition to the lattice-matched Cd(Zn)Te substrates. Detector growth techniques²² include liquid phase epitaxy and vapor phase epitaxy (VPE). The latter includes metal organic chemical vapor deposition (MOCVD) and molecular beam epitaxy (MBE).

Hybrid Readout

Hybrid readouts perform the functions of detector interface, signal processing, and video multiplexing.²³ The hybrid FPA readout technologies include

- Surface channel charge coupled device (SCCD)
- Buried channel charge coupled device (BCCD)
- x - y addressed switch-FET (SWIFET) or direct readout (DRO) FET arrays
- Combination of MOSFET and CCD (MOS/CCD)
- Charge-injection device (CID)

Early hybrid readouts were either CIDs²⁴ or CCDs, and the latter are still popular for silicon monolithics. However, x - y arrays of addressed MOSFET switches are superior for most hybrids for reasons of yield, design flexibility, simplified interface, and direct leveraging of Moore's Law for ongoing improvements and cost reduction. The move to the FET-based, direct readouts is key to the dramatic improvements in staring array producibility and is a consequence of the spin-off benefits from the silicon memory markets. DROs are fabricated with high yield and are fully compatible with advanced processes that are available at captive and commercial foundries. We will thus focus our discussion on these families. Though not extensively, CCDs are still sometimes used in hybrid FPAs.²⁵

Nonsilicon readouts Readouts have been developed in Ge, GaAs, InSb, and HgCdTe. The readout technologies include monolithic CCD, charge-injection device (CID), charge-injection matrix (CIM), enhancement/depletion (E/D) MESFET (GaAs), complementary heterostructure FET (C-HFET; GaAs),²⁶ and JFET (GaAs and Ge). The CCD, CID, and CIM readout technologies generally use MIS detectors for monolithic photon detection and signal processing.

The CID and CIM devices rely on accumulation of photogenerated charge within the depletion layer of a MIS capacitor that is formed using a variety of passivants (including CVD and photo-SiO₂, anodic SiO₂, and ZnS). A single charge transfer operation then senses the accumulated charge. Device clocking and signal readout in the CIDs and CIMs relies on support chips adjacent to the monolithic IR FPA. Thus, while the FPA is monolithic, the FPA assembly is actually a multichip hybrid.

CCDs have been demonstrated in HgCdTe and GaAs.²⁷ The n -channel technology is preferred in both materials for reasons of carrier mobility and device topology. In HgCdTe, for example, n -MOSFETs with CVD SiO₂ gate dielectric have parameters that are in good agreement with basic silicon MOSFET models. Fairly elaborate circuits have been demonstrated on CCD readouts, e.g., an on-chip output amplifier containing a correlated double sampler (CDS).

GaAs has emerged as a material that is very competitive for niche applications including IR FPAs. Since GaAs has very small thermal expansion mismatch with many IR detector materials including HgCdTe and InSb, large hybrids are possible, and VPE detector growth capability suggests future development of composite monolithic FPAs. The heterostructure (H-)MESFET and C-HFET technologies are particularly interesting for IR FPAs because low $1/f$ noise has been demonstrated; noise spectral densities at 1 Hz of as low as 0.5 $\mu\text{V}/\sqrt{\text{Hz}}$ for p -HIGFET and 2 $\mu\text{V}/\sqrt{\text{Hz}}$ for the enhancement H-MESFET²⁸ at 77 K have been achieved. The H-MESFET has the advantage of greater fabrication maturity (16 K SRAM and 64 \times 2 readout demonstrated), but the C-HFET offers lower power dissipation.

Direct Readout Architectures The DRO multiplexer consists of an array of FET switches. The basic multiplexer has several source follower stages that are separated at the cell, row, and column levels by MOSFET switches which are enabled and disabled to perform pixel access, reset, and multiplexing. The signal voltage from each pixel is thus direct-coupled through the cascaded source follower architecture as shown, for example, in Fig. 7. Shift registers generate the various clock signals; a minimum of externally supplied clocks is required. Since CMOS logic circuitry is used, the clock levels do not require precise adjustment for optimum performance. The simple architecture also gives high functional yield even in readout materials less mature than silicon.

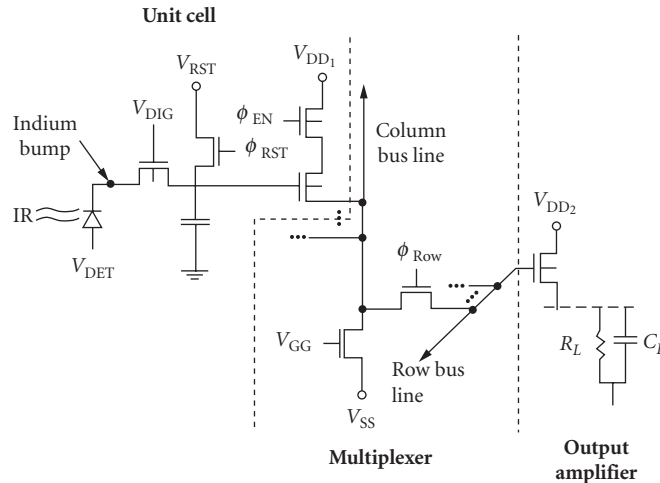


FIGURE 7 Direct readout schematic (shown with direct injection input).

Owing to the relatively low internal impedances beyond the input circuit, multiplexer noise is usually negligible. The inherent dynamic range is often >100 dB and the FPA dynamic range is limited only by the output-referred noise of the input circuit and the maximum signal excursion. The minimum read noise for DROs in imaging IR FPAs is typically capacitor reset noise. Correlated double samplers are thus used to suppress the reset noise for highest possible SNR at low integrated signal level.

In addition to excellent electrical characteristics, the DRO has excellent electro-optical properties including negligible MTF degradation and no blooming. Crosstalk in DRO-based FPAs is usually detector-limited since the readouts typically have low (<0.005 percent) electrical crosstalk. DROs also have higher immunity to clock feedthrough noise due to their smaller clock capacitances. Substrate charge pumping, which causes significant FET backgating²⁹ and transconductance degradation in SCCDs, is low in DROs.

X-Y addressing and clock generation Both static and dynamic shift registers are used to generate the clock signals needed for cell access, reset, and pixel multiplexing. Static registers offer robust operation and increased hardness to ionizing radiation in trade for increased FET count and preference for CMOS processes. Dynamic registers use fewer transistors in NMOS or PMOS processes, but require higher voltages, have lower maximum clocking rate, and must be carefully designed to avoid being affected by incident radiation.

Dynamic shift registers use internal bootstrapping to regenerate the voltage at each tap. The circuit techniques limit both the lowest and highest clock rates and require fine-tuning of the MOSFET design parameters for the specific operating frequency. More importantly, the high internal voltages stress conventional CMOS processes.

Electronically scanned staring FPAs The inability to integrate photogenerated charge for full staring frame times is often handled by integration time management. Since the photon background for the full 8- to 12- μm spectral band is over two orders of magnitude larger than the typical MWIR passband, and since the LWIR detector dark current is several orders of magnitude larger than similarly sized MWIR devices, LWIR FPA integration duty cycle can be quite poor. It is sometimes prudent to concede the limited duty cycle by electronically scanning the staring readout. Electronic scanning refers to a modified staring FPA architecture wherein the FPA is operated like a scanning FPA but without optomechanical means. Sensitivity is enhanced beyond that of a true scanning FPA by, for example, using multiple readout bus lines to allow integration times longer than one row time. The sharing reduces circuit multiplicity and frees unit cell real estate to share circuitry and

larger integration capacitance, and to use the otherwise parasitic bus capacitance to further increase capacity. More charge can thus be integrated even though the duty cycle is preset to $1/N$, where N is the number of elements on each common bus.

Time delay integration scanning FPAs While no longer extensively used for staring readouts, SCCDs are used in scanning readouts to incorporate on-focal plane TDI since they have higher dynamic range than FET bucket brigades. Dynamic range >72 dB and as high as 90 dB are typically achieved with high TDI efficacy. Two architectures dominate. In one, the CCD is integrated adjacent to the input circuit in a contiguous unit cell. In the second, the input circuit is segregated from the CCD in a sidecar configuration. The latter offers superior cell-packing density and on-chip signal processing in trade for circuit complexity. Figure 8 shows the schematic circuit for a channel of a scanning readout having capacitive transimpedance amplifier input circuit (discussed earlier), common TDI channel bus, and fill-and-spill³⁰ input to a sidecar SCCD TDI. This scheme integrates CMOS and CCD processes for much on-chip signal processing in very fine orthoscan pitch.³¹

The readout conversion factor, i.e., volts out per electrons in, for the sidecar CTIA scheme is

$$S_v = \frac{\Delta V}{e^-} = \frac{C_{F/S}}{C_T} A_{V_1} A_{V_2} \frac{q}{C_{out}} \quad (4)$$

where $C_{F/S}$ is the fill-and-spill gate capacitance, C_T is the integration/feedback capacitance, A_{V_x} characterizes the various source follower gains, and C_{out} is the sense node capacitance at the CCD output. The ratio of $C_{F/S}$ to C_T sets a charge gain that allows design-tailoring for managing dynamic range or lowering input-referred noise. High charge-gain yields read noise that is limited by the input circuit and not by the transfer noise³² of the high-carrier-capacity SCCD.

MOSFET bucket brigades are also used as TDI registers since simpler, all-MOS designs and processes can be used. Advantages include compatibility with standard MOS and CMOS, and capability for external clocking using specific CMOS-compatible clock levels. The latter potential advantage is mitigated in the sidecar TDI scheme by appropriately sizing the SCCD registers and the charge gain to yield the desired CCD clock levels, for example. Disadvantages include higher TDI register noise due to kTC noise being added at each transfer and limited signal excursion.

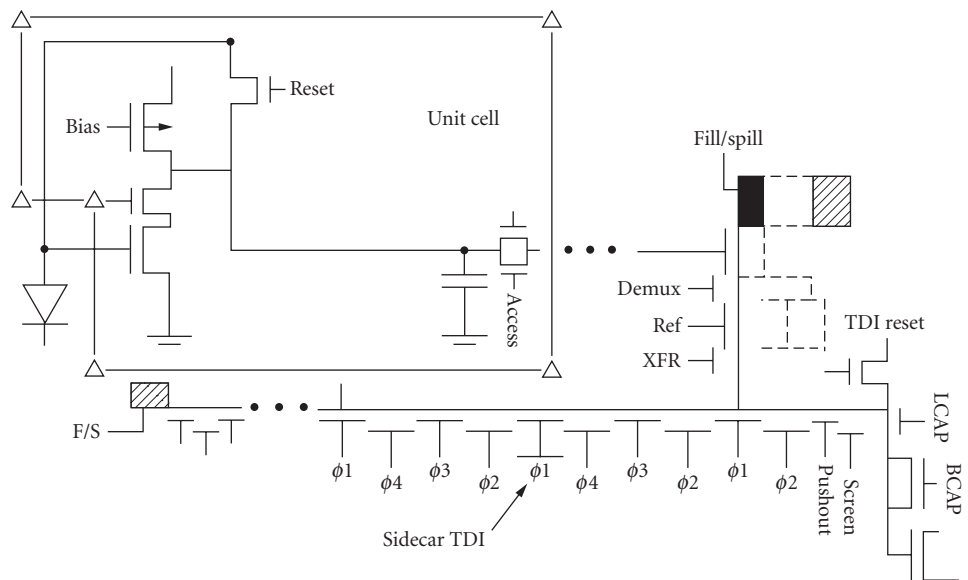


FIGURE 8 Sidecar TDI with capacitive transimpedance amplifier input circuit.

Output circuits Output circuitry is usually kept to a minimum to minimize power dissipation. Circuit design thus tends to focus on the trades between voltage-mode and current-mode output amplifiers, although on-chip signal processing is now at system-on-chip (SoC) level of sophistication, including low-speed A/D conversion, switched-capacitor filtering,³³ and on-chip nonuniformity correction. Voltage-mode outputs offer better S/N performance across a wider range in backgrounds for a given readout transimpedance. Current-mode outputs offer wider bandwidth and better drive capability at higher clock frequencies.

Detector Interface: Input Circuit After the incoming photon flux is converted into a signal by the detector, it is coupled into the readout via a detector interface circuit.³⁴ Signal input is optical in a monolithic FPA, so signal conditioning is limited. In hybrid FPAs and some composite material monolithics, the signal is injected electrically into the readout. The simplest input schemes offering the highest mosaic densities include direct detector integration (DDI) and direct injection (DI). More complex schemes trade simplicity for input impedance reduction [buffered direct injection (BDI) and capacitive transimpedance amplification (CTIA)], background suppression (e.g., gate modulation), or ultralow read noise with high speed (CTIA). We briefly describe the more popular schemes and their performance. Listed in Table 1 are approximate performance-describing equations for comparing the circuits schematically shown in Fig. 9.

Direct detector integration Direct detector integration (Fig. 9a), also referred to as source follower per detector (SFD), is used at low backgrounds and long frame times (frame rates typically ≤ 15 Hz in large staring arrays). Photocurrent is stored directly on the detector capacitance, thus requiring the detector to be heavily reverse-biased to maximize dynamic range. The changing detector voltage modulates the gate of a source follower whose drive FET is in the cell and whose current source is common to all the detectors in a column or row. The limited cell area constrains the source followers' drive capability and thus the bandwidth.

The DDI unit cell typically consists of the drive FET, cell enable transistor(s), and reset transistor(s). A detector site is read out by strobing the appropriate row clock, thus enabling the output source follower. The DDI circuit is capable of read noise as low as a few electrons per pixel.

Direct injection Direct injection (Fig. 9b) is perhaps the most widely used input circuit due to its simplicity and high performance. The detector directly modulates the source of a MOSFET. The direct coupling requires that detectors with p -on- n polarity, as is the case with InSb and most photovoltaic LWIR detectors, interface p -type FETs (and vice versa) for carrier collection in the integration capacitor. In surface channel CCDs, the input transistor's drain is virtual, as formed by a fully enhanced well, and often doubles as the integration capacitor.

Practical considerations, including limited charge-handling capacity, constrain the DI input to operation with high-impedance MWIR or limited cutoff ($\lambda_c \leq 9.5 \mu\text{m}$) LWIR detectors. The associated background photocurrent for the applications where direct injection can be used mandates that the DI FET operate subthreshold.³⁵ The subthreshold gate transconductance, g_m , is independent of FET geometry:³⁶

$$g_m = \left(\frac{\partial I_D}{\partial V_G} \right) \Big|_{V_{DS}=\text{constant}} = \frac{q \left(\eta_{\text{inj}} \left\{ I_{\text{photo}} + \frac{V_{\text{det}}}{R_{\text{det}}} - I_{\text{det}_0} \left(e^{(qV_{\text{det}}/nkT)} - 1 \right) \right\} \right)}{nkT} \approx \frac{qI_D}{nkT} \quad (5)$$

The injection efficiency, η_{inj} , of detector current into the DI FET is

$$\eta_{\text{inj, DI}} = \frac{g_m R_{\text{det}}}{1 + g_m R_{\text{det}}} \left[\frac{1}{1 + \frac{j\omega C_{\text{det}} R_{\text{det}}}{1 + g_m R_{\text{det}}}} \right] \quad (6)$$

where R_{det} and C_{det} are the detectors' dynamic resistance and capacitance, respectively. Poor DI circuit bandwidth occurs at low-photon backgrounds due to low g_m .

TABLE 1 Focal Plane Array Performance

Input Circuit	Percentage of BLIP	Detector Noise	Input-Referred Circuit Noise	Input-Referred MUX Noise	Transimpedance
Direct Detector Integration	$\left[\frac{(\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}})^{1/2}}{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}} + qR_{\text{det}}^2} + N_{\text{f}}^2 + \frac{2kT\tau_{\text{int}}}{qR_{\text{int}}} \right]^{1/2}$	$I_{\text{det}}^2 = \left[\frac{4kT}{R_{\text{det}}} + 2qI_{\text{det}} \right] \Delta f + \int \left(\frac{K_{\text{det}}}{f} \right) df$	$N_{\text{sf}} \approx \sqrt{2} \left[\int_{\nu}^{\Delta f} V^2(f) \frac{(1 - \cos 2\pi f T_D)}{[1 + (2\pi f T_D)^2]} df \right]^{1/2}$ $S_{\nu} = \left(\frac{C_{\text{det}} A_{\nu}}{q} \right)^{-1}$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{A_{\nu}^2} kTC_{\text{mux}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{det}}} \right) A_{\nu}$
Direct Injection	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\frac{\Delta f [1 + \omega^2 C_{\text{det}}^2 R_{\text{det}}^2 \left(\frac{8}{3} kTg_m + \frac{K_{\text{FET}}}{f\alpha} \right)]}{g_m^2 R_{\text{det}}^2} \right] df$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_{\nu}$
Buffered Direct Injection	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\eta_{\text{noise}}^2 \left(\frac{8}{3} kTg_m + \frac{K_{\text{FET}}}{f\alpha} \right) + A_{\text{amp}}^2 \left(\epsilon_{\text{amp}}^2 \right) \right] df$ $\eta_{\text{noise}} = \frac{1 + j\omega R_{\text{det}} C_{\text{int}}}{1 + (1 + A_{\nu}) g_m R_{\text{det}} + j\omega R_{\text{det}} (1 + A_{\nu}) C_{\text{int}}}$ $A_{\text{amp}} = \left(\frac{g_m}{R_{\text{det}}} \right) \frac{(1 + j\omega R_{\text{det}} C_{\text{det}})}{1 + (1 + A_{\nu}) g_m R_{\text{det}} + j\omega [C_{\text{det}} + (1 + A_{\nu}) C_{\text{int}}] R_{\text{det}}}$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_{\nu}$
Chopper-Stabilized BDI	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{load}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\eta_{\text{noise}}^2 \left(\frac{8}{3} kTg_m \right) + A_{\text{amp}}^2 \left(\epsilon_{\text{amp}} \right) \right] df$ $\eta_{\text{noise}} = \left(\frac{g_m}{1 + (1 + A_{\nu}) g_m R_{\text{det}} + j\omega R_{\text{det}} C_{\text{int}}} \right)$ $A_{\text{amp}} = \left(\frac{g_m}{R_{\text{det}}} \right) \frac{(1 + j\omega R_{\text{det}} C_{\text{det}})}{1 + (1 + A_{\nu}) g_m R_{\text{det}} + j\omega [C_{\text{det}} + (1 + A_{\nu}) C_{\text{int}}] R_{\text{det}}}$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_{\nu}$
Gate Modulation (FET Load)	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{load}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\frac{1}{A_{\nu}^2} \left(2qI_{\text{input}} + \frac{K_{\text{FET,input}}}{f\alpha} \right) \right] df$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{A_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{A_{\nu} t_{\text{int}}}{C_{\text{int}}} \right) A_{\nu}$
Capacitive Transimpedance Amplifier	$\left[\frac{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}}}{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}} + \frac{2kT\tau_{\text{int}}}{qR_{\text{det}}}} + N_{\text{amp,lif}}^2 + N_{\text{amp,white}}^2 + N_{\text{load,white}}^2 \right]^{1/2}$	$N_{\text{amp,lif}} \approx \frac{C_{\text{det}} K_{\text{amp}} \sqrt{2}}{q} \left[\frac{5\tau_{\text{int}}}{\tau_{\text{amp}}} \right]$ $N_{\text{amp,white}} \approx \frac{C_{\text{det}}}{q} \sqrt{\frac{8}{3} \frac{kT}{g_m \tau_{\text{amp}}}}$ $N_{\text{load,white}} \approx \frac{C_{\text{det}}}{qA_{\nu} \text{amp}} A_{\nu}^2 \sqrt{\frac{2kT}{C_L}}$		$I_{\text{mux,ir}}^2 = kTC_{\text{input}} \Delta f$	$Z_T = \frac{\tau_{\text{int}}}{C_T} A_{\nu, \text{sf}}$ $C_T = \frac{(C_{\text{gd}} + C_{\text{gs}} + C_{\text{det}}) + A_{\nu} (C_{\text{fb}} + C_{\text{gd}})}{A_{\nu}}$

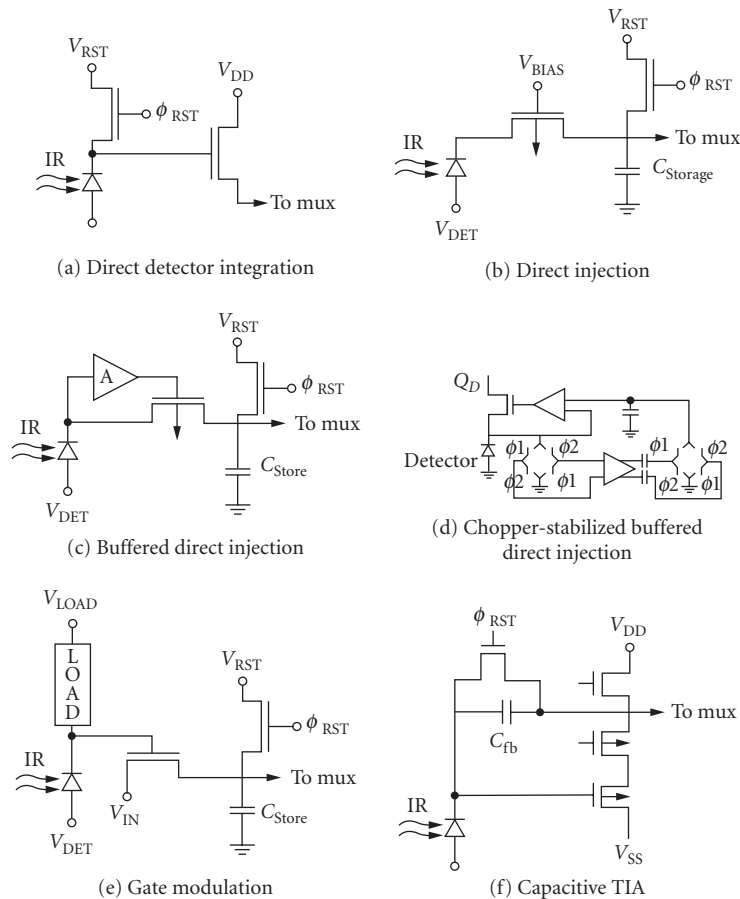


FIGURE 9 Hybrid FPA detector interface circuits.

The injection efficiency varies across an FPA due to FET threshold, detector bias, and detector resistance nonuniformity. Changes in detector current create detector bias shifts since the input impedance is relatively high. In extreme cases a large offset in threshold gives rise to excess detector leakage current and $1/f$ noise, in addition to fixed pattern noise. The peak-to-peak threshold voltage nonuniformity spans the range from ≈ 1 mV for silicon p -MOSFETs to over 125 mV for some GaAs-based readouts.

Depending upon the interface to the multiplexing bus, the noise-limiting capacitance, C_{input} , is approximately the integration capacitance or the combined integration and bus capacitance. Some direct-injection cells are thus buffered with a source follower (see Fig. 7). Omitting the source follower reduces the readout transimpedance (due to charge splitting between the integration capacitor and thus bus line capacitance) in trade for larger integration capacitance since more unit-cell real estate is available.

Increasing the pixel density has required a continuing reduction in cell pitch. Figure 10 plots the charge-handling capacity as functions of cell pitch and minimum gate length for representative DI designs using the various minimum feature lengths. Also plotted is the maximum capacity assuming the cell is composed entirely of integration capacitor ($225 \text{ \AA} \text{ iO}_2$). A $27\text{-}\mu\text{m}$ DI cell, fabricated in $1.25\text{-}\mu\text{m}$ CMOS, thus had similar cell capacity as an earlier $60\text{-}\mu\text{m}$ DI cell in $3\text{-}\mu\text{m}$ CMOS. Limitations on cell real estate, operating voltage, and the available capacitor dielectrics nevertheless

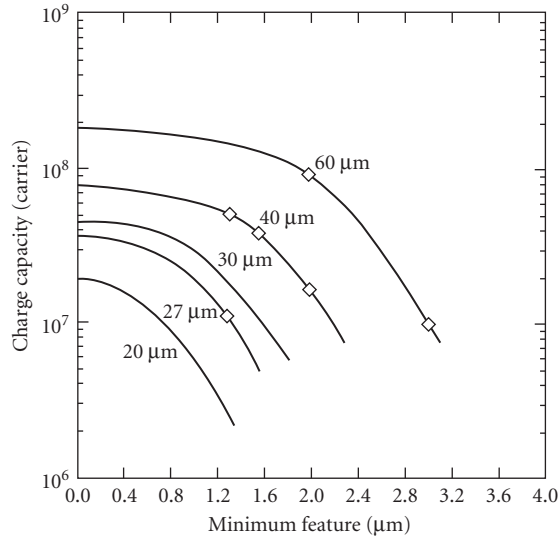


FIGURE 10 Direct injection charge-handling capacity versus cell pitch for standard CMOS processes.

dictate maximum integration times that are often shorter than the total frame time. This duty cycling equates to degradation in detective quantum efficiency.

Buffered direct injection A significant advantage of the source-coupled input is MOSFET noise suppression. This suppression is implied in the η_{BLIP} expression shown in Table 1. When injection efficiency is poor, however, MOSFET noise becomes a serious problem along with bandwidth. These deficiencies are ameliorated via buffered direct injection (BDI),³⁷ wherein a feedback amplifier with open-loop gain $-A_v$ (Fig. 9c) is added to the DI circuit. The buffering increases injection efficiency to near-unity, increases bandwidth by orders of magnitude, and suppresses the DI FET noise.

BDI has injection efficiency

$$\eta_{\text{inj}} \equiv \frac{g_m R_{\text{det}} (1 + A_v)}{1 + g_m R_{\text{det}} (1 + A_v) + j\omega R_{\text{det}} [(1 + A_v) C_{\text{amp}} + C_{\text{det}}]} \quad (7)$$

where C_{amp} is the Miller capacitance of the amplifier. Circuit bandwidth is maximized by lowering the amplifiers' Miller capacitance to provide detector-limited frequency response that is lower than that possible with DI by the factor $(1 + A_v)$.

The noise margin of the BDI circuit is superior to DI, even though two additional noise sources associated with the feedback amplifier are added. The dominant circuit noise stems from the drive FET in the amplifier. The noise power for frequencies less than $1/(2\pi R_{\text{det}} C_{\text{det}})$ is directly proportional to the detector impedance. Amplifier noise is hence a critical issue with low impedance ($<1 \text{ M}\Omega$) detectors including long wavelength photovoltaics operating at temperatures above 80 K.

Chopper-stabilized buffered direct injection The buffered direct injection circuit has high $1/f$ noise with low-impedance detectors. While the MOSFET $1/f$ noise can be suppressed somewhat by reverse-biasing the detectors to the point of highest resistance, detector $1/f$ noise may then dominate. Other approaches include enlarging MOSFET gate area, using MOS input transistors in the lateral bipolar mode, or using elaborate circuit techniques such as autozeroing and chopper stabilization. Chopper stabilization is useful if circuit real estate is available, as in a scanning readout. Figure 9d shows a block diagram schematic circuit of chopper-stabilized BDI.

Chopper stabilization refers to the process of commutating the integrating detector node between the inverting and noninverting inputs of an operational amplifier having open-loop gain, A_V . This chopping process shifts the amplifier's operating frequency to higher frequencies where the amplifier's noise is governed by white noise, not $1/f$ noise. At chopping frequencies $f_{\text{chop}} \gg f_k$, the equivalent low-frequency input noise of the chopper amplifier is equal to the original amplifier white-noise component.³⁸ The amplifier's output signal is subsequently demodulated and filtered to remove the chopping frequency and harmonics. This scheme also reduces the input offset nonuniformity by the reciprocal of the open-loop gain, thereby generating uniform detector bias. Disadvantages include high circuit complexity and the possibility of generating excess detector noise via clock feedthrough-induced excitation of traps, particularly with narrow bandgap photovoltaic detectors.

Gate modulation Signal processing can be incorporated in a small unit cell by using a gate-modulated input structure (c.f. MOSFET load gate modulation in Fig. 9e).³⁹ The use of an MOSFET as an active load device, for example, provides dynamic range management via automatic gain control and user-adjustable background pedestal offset since the detector current passes through a load device with resistance R_{LOAD} . The differential gate voltage applied to the input FET varies for a change in photocurrent, ΔI_{photo} , as

$$\Delta V_G = R_{\text{LOAD}} \eta_{\text{inj, DI}} \Delta I_{\text{photo}} \quad (8)$$

The current injected into the integration capacitor is

$$I_{\text{input}} = g_m R_{\text{LOAD}} \eta_{\text{inj, DI}} \Delta I_{\text{photo}} \quad (9)$$

The ratio of I_{input} to I_{photo} is the current gain, A_I , which is

$$A_I = \frac{g_m}{g_{m, \text{LOAD}}} \eta_{\text{inj, DI}} \quad (10)$$

The current gain can self-adjust by orders of magnitude depending on the total detector current. Input-referred read noise of tens of electrons has thus been achieved with high-impedance SWIR detectors at low-photon backgrounds. The same circuit has also been used at LWIR backgrounds with LWIR detectors having adequate impedance for good injection efficiency.

The current gain expression suggests a potential shortcoming for imaging applications since the transfer characteristic is nonlinear, particularly when the currents in the load and input FETs differ drastically. In conjunction with tight specifications for threshold uniformity, pixel functionality can be decreased, dynamic range degraded, and imagery dominated by spatial noise. The rms fractional gain nonuniformity (when operating subthreshold) of the circuit is approximately

$$\frac{\Delta A_I}{A_I} = \frac{q \sigma_{VT}}{n_{\text{FET}} kT} \quad (11)$$

where σ_{VT} is the rms threshold nonuniformity. At 80 K, state-of-the-art σ_{VT} of 0.5 mV for a 128 × 128 FPA, and $n = 1$, the minimum rms nonuniformity is ≈ 7 percent.

Capacitive transimpedance amplifier (CTIA) Many CTIAs have been successfully demonstrated. The most popular approach uses a simple CMOS inverter⁴⁰ for feedback amplification (Fig. 9f). Others use a more elaborate differential amplifier. The two schemes differ considerably with respect to open-loop voltage gain, bandwidth, power dissipation, and cell real estate. The CMOS inverter-based CTIA is more attractive for high-density arrays. The latter is sometimes preferred for scanning or Z-hybrid applications where real estate is available primarily to minimize power dissipation.⁴¹

In either case, photocurrent is integrated directly onto the feedback capacitor of the transimpedance amplifier. The minimum feedback capacitance is set by the amplifier's Miller capacitance and defines the maximum circuit transimpedance. Since the Miller capacitance can be made very small (<5 fF), the resulting high transimpedance yields excellent margin with respect to

downstream system noise. The transimpedance degrades when the circuit is coupled to large detector capacitances, so reducing pixel size serves to minimize read noise and the circuits' attractiveness will continue to increase in the future.

The CTIA allows extremely small currents to be integrated with high efficiency and tightly regulated detector bias. The amplifier open-loop gain, $A_{V,amp}$, ranges from as low as on the order of ten to higher than several thousand for noncascoded inverters, and many thousands for cascoded inverter and differential amplifier designs. The basic operating principle is to apply the detector output to the inverting input of a high-gain CMOS differential amplifier operated with capacitive feedback. The feedback capacitor is reset at the detector sampling rate. The noise equivalent input voltage of the amplifier is referred to the detector impedance, just as in other circuits.

The CTIA's broadband channel noise sets a lower limit on the minimum achievable read noise, is the total amplifier white noise, and can be approximated for condition of large open-loop gain by

$$N_{kTC,channel}^2 = \frac{kTC_{fb}}{q^2} \left[\frac{C_{det} + C_{fb}}{C_L + \frac{C_{fb}C_{det}}{C_{fb} + C_{det}}} \right] \quad (12)$$

where C_{fb} is the feedback capacitance including the Miller capacitance and integration capacitance and C_L is the output load capacitance. This expression provides an intuitive formula for minimizing noise: detector capacitance must be low (i.e., minimize detector shunting capacitance which reduces closed-loop gain) and output capacitance high (i.e., limit bandwidth). Of the three amplifier noise sources listed in Table 1, amplifier $1/f$ noise is often largest. For this reason, the CMOS inverter-based CTIA has best performance with p -on- n detectors and p -MOSFET amplifier FET due to the lower $1/f$ noise.

33.5 PERFORMANCE: FIGURES OF MERIT

In the early days of infrared technology, detectors were characterized by the noise equivalent power (NEP) in a 1-Hz bandwidth. This was a good specification for single detectors, since their performance is usually amplifier limited. The need to compare detector technologies for application to different geometries and the introduction of FPAs having high-performance on-board amplifiers and small parasitics necessitated normalization to the square root of the detector area for comparing S/N. R. C. Jones⁴² thus introduced detectivity (D^*), which is simply the reciprocal of the normalized NEP and has units $\text{cm} \cdot \sqrt{\text{Hz}}/\text{W}$ (or Jones).

While D^* is well-suited for specifying infrared detector performance, it can be misleading to the uninitiated since the D^* is highest at low background. An LWIR FPA operating at high background with background-limited performance (BLIP) S/N can have a D^* that is numerically lower than for a SWIR detector having poor S/N relative to the theoretical limit. Several figures of merit that have hence proliferated include other ways of specifying detectivity: e.g., thermal D^* (D_{th}^*), blackbody D^* (D_{bb}^*), peak D^* (D_{pk}^*), percentage of BLIP (%BLIP or η_{BLIP}), and noise equivalent temperature difference ($NE\Delta T$). Since the final output is an image, however, the ultimate figure of merit is how well objects of varying size are detected and resolved in the displayed image. The minimum resolvable temperature (MRT) is thus a key benchmark. These are briefly discussed in this section.

Detectivity (D^*)

D^* is the S/N ratio normalized to the electrical bandwidth and detector area. In conjunction with the optics area and the electrical bandwidth, it facilitates system sensitivity estimation. However, D^* can be meaningless unless the test conditions, including magnitude and spectral distribution of the flux source (e.g., blackbody temperature), detector field-of-view, chopping frequency (lock-in amplifier),

background temperature, and wavelength at which the measurement applies. D^* is thus often quoted as “blackbody,” since the spectral responsivity is the integral of the signal and background characteristics convolved with the spectral response of the detector. D_{bb}^* specifications are often quantified for a given sensor having predefined scene temperature, filter bandpass, and cold shield $f/\#$ using a generalized expression.

Peak detectivity is sometimes preferred by detector engineers specializing in photon detectors. The background-limited peak detectivity for a photovoltaic detector is

$$D_{\lambda_{\text{pk}}}^* = \sqrt{\frac{\eta}{2Q_B}} \frac{\lambda_{\text{pk}}}{hc} \quad (13)$$

and refers to measurement at the wavelength of maximum spectral responsivity. For detector-limited scenarios, such as at higher operating temperatures or longer wavelengths (e.g., $\lambda_c > 12 \mu\text{m}$ at operating temperature less than 78 K or $\lambda_c > 4.4 \mu\text{m}$ at >195 K), the peak detectivity is limited by the detector and not the photon shot noise. In these cases the maximum detector-limited peak D^* in the absence of excess bias-induced noise (both $1/f$ and shot noise) is

$$D_{\lambda_{\text{pk}}}^* = \frac{\eta q}{2} \sqrt{\frac{R_0 A}{kT}} \frac{\lambda}{hc} \quad (14)$$

The $R_0 A$ product of a detector thus describes detector quality even though other parameters may actually be more relevant for FPA operation.

Percentage of BLIP

Whereas D^* compares the performance of dissimilar detectors, FPA designers often need to quantify an FPA's performance relative to the theoretical limit at a specific operating background. Percentage of BLIP, η_{BLIP} , is one such parameter and is simply the ratio of photon noise to composite FPA noise

$$\eta_{\text{BLIP}} = \left(\frac{N_{\text{PHOTON}}^2}{N_{\text{PHOTON}}^2 + N_{\text{FPA}}^2} \right)^{1/2} \quad (15)$$

NE ΔT

The NE ΔT of a detector represents the temperature change, for incident radiation, that gives an output signal equal to the rms noise level. While normally thought of as a system parameter, detector NE ΔT and system NE ΔT are the same except for system losses (conservation of radiance). NE ΔT is defined:

$$\text{NE } \Delta T = v_n \left(\frac{\partial T}{\partial Q} \right) \bigg/ \left(\frac{\partial V_s}{\partial Q} \right) = v_n \frac{\Delta T}{\Delta V_s} \quad (16)$$

where v_n is the rms noise and ΔV_s is the signal measured for the temperature difference ΔT . It can be shown that

$$\text{NE } \Delta T = \left(\tau_o C_{T\lambda} \eta_{\text{BLIP}} \sqrt{N_c} \right)^{-1} \quad (16a)$$

where τ_o is the optics transmission, $C_{T\lambda}$ is the thermal contrast from Fig. 1, and N_c is the number of photogenerated carriers integrated for one integration time, t_{int} :

$$N_c = \eta A_{\text{det}} t_{\text{int}} Q_B \quad (16b)$$

The distinction between an integration time and the FPA's frame time must be noted. It is often impossible at high backgrounds to handle the large amount of carriers generated over frame times compatible with standard video frame rates. The impact on system D^* is often not included in the FPA specifications provided by FPA manufacturers. This practice is appropriate for the user to assess relative detector quality, but must be coupled with usable FPA duty cycle, read noise, and excess noise to give a clear picture of FPA utility. Off-FPA frame integration can be used to attain a level of sensor sensitivity that is commensurate with the detector-limited D^* and not the charge-handling-limited D^* .

The inability to handle a large amount of charge nevertheless is a reason why the debate as to whether LWIR or MWIR operation is superior is still heated. While the LWIR band should offer order-of-magnitude higher sensitivity, staring readout limitations often reduce LWIR imager sensitivity to below that of competing MWIR cameras. However, submicron photolithography, alternative dielectrics, and refinements in readout architectures are ameliorating this shortfall and LWIR FPAs having sensitivity superior to MWIR counterparts are available. Figure 11 shows the effect on high quantum efficiency FPA performance and compares the results to PtSi at TV-type frame rate. The figure illustrates BLIP and measured NE ΔT s versus background temperature for several spectral bands assuming a 640×480 DI readout multiplexer ($27\text{-}\mu\text{m}$ pixel pitch and $\approx 1\text{-}\mu\text{m}$ -minimum

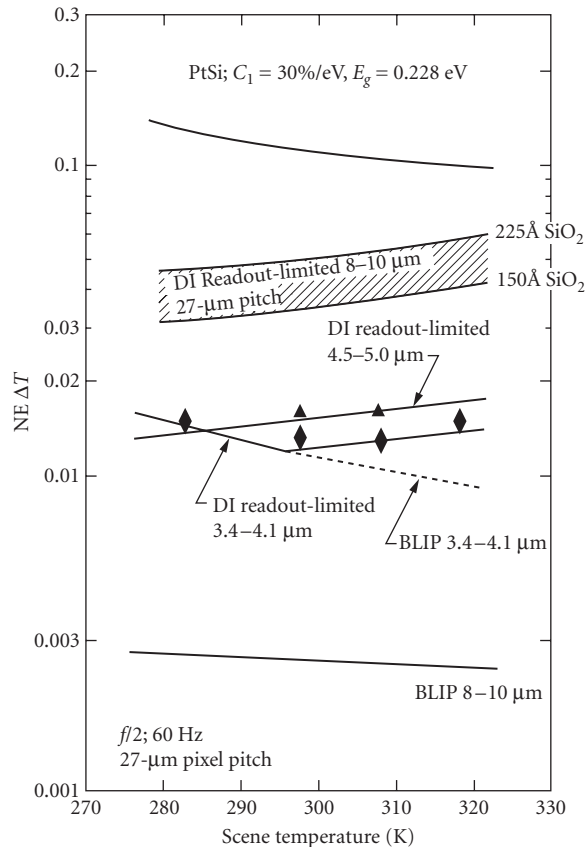


FIGURE 11 NE ΔT versus background temperature for several prominent spectral bands.

feature) is hybridized to high-quantum-efficiency MWIR and LWIR detectors. Though the device has large charge-handling capacity, there is large shortfall in the predicted LWIR FPA performance relative to the BLIP limit. The measured MWIR FPA performance values, as shown by the data points, are in good agreement with the predicted trends.

Spatial noise Estimation of IR sensor performance must include a treatment of spatial noise that occurs when FPA nonuniformities cannot be compensated correctly. This requires consideration of cell-to-cell response variations. Mooney et al.⁴³ comprehensively discussed the origin of spatial noise. The total noise determining the sensitivity of a staring array is the composite of the temporal noise and the spatial noise. The spatial noise is the residual nonuniformity U after application of nonuniformity compensation, multiplied by the signal electrons N . Photon noise, equal to \sqrt{N} , is the dominant temporal noise source for the high infrared background signals for which spatial noise is significant (except for TE-cooled or uncooled sensors). The total noise equivalent temperature difference is

$$\text{Total NE } \Delta T = \frac{\sqrt{N+U^2N^2}}{\frac{\partial N}{\partial T}} = \frac{\sqrt{1/N+U^2}}{\frac{1}{N} \frac{\partial N}{\partial T}} \quad (17)$$

where $\partial N/\partial T$ is the signal change for a 1 K source temperature change. The denominator, $(\partial N/\partial T)/N$, is the fractional signal change for a 1 K source temperature change. This is the relative scene contrast due to $C_{T\lambda}$ and the FPA's transimpedance.

The dependence of the total NE ΔT on residual nonuniformity is plotted in Fig. 12 for 300 K scene temperature, two sets of operating conditions, and three representative detectors: LWIR HgCdTe, MWIR HgCdTe, and PtSi. Operating case A maximizes the detected signal with $f/1.4$ optics, 30-Hz frame rate, and 3.4 to 5.0- μm passband. Operating case B minimizes the solar influence by shifting the passband to 4.2 to 5.0 μm and trades off signal for the advantages of lighter, less expensive optics ($f/2.0$) at 60-Hz frame rate. Implicit in the calculations are charge-handling capacities of 30 million e^- for MWIR HgCdTe, 100 million e^- for LWIR HgCdTe, and 1 million for PtSi. The sensitivity at the lowest nonuniformities is independent of nonuniformity and limited by the

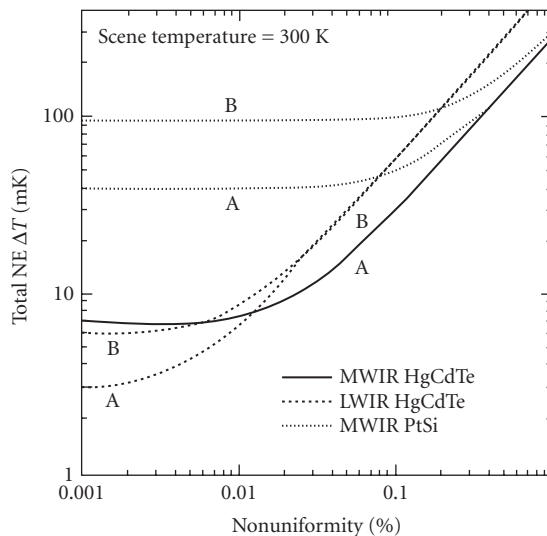


FIGURE 12 Total NE ΔT versus nonuniformity at 300 K scene temperature.

shot noise of the detected signal. The LWIR sensitivity advantage is achieved only at nonuniformities less than 0.01 percent, which is comparable to that achieved with buffered input circuits. At the reported direct-injection MWIR HgCdTe residual nonuniformity of 0.01 to 0.02 percent, the total $NE\Delta T$ is about 0.007 K, which is comparable to the MWIR BLIP limit. At the reported PtSi residual nonuniformity of 0.05 percent with direct detector integration, total $NE\Delta T$ is higher at 0.04K, but exceeds the BLIP limit for the lower quantum efficiency detectors.

Minimum Resolvable Temperature

The minimum resolvable temperature (MRT) is often the preferred figure of merit for imaging infrared sensors. MRT is a function of spatial resolution and is defined as the signal-to-noise ratio required for an observer to resolve a series of standard four-bar targets. While many models exist due to the influence of human psycho-optic response, a representative formula⁴⁴ is

$$MRT(f_s, T_{SCENE}) = \frac{2SNR_t NE \Delta T(T_{SCENE})}{MRT(f_s)} \left[\frac{f_s^2 \Delta x \Delta y}{L \tau_{eye} f_{frame} N_{OS} N_{SS}} \right]^{1/2} \quad (18)$$

where f_s is the spatial frequency in cycles/radian, a target signal-to-noise ratio (SNR_t) of five is usually assumed, the MTF describes the overall modulation transfer function including the optics, detector, readout, and the integration process, Δx and Δy are the respective detector subtenses in mRad, τ_{eye} is the eye integration time, f_{frame} is the display frame rate, N_{OS} is the overscan ratio, N_{SS} is the serial scan ratio, and L is the length-to-width ratio of a bar chart (always set to 7). While the MRT of systems with temporal noise-limited sensitivity can be adequately modeled using the temporal $NE\Delta T$, scan noise in scanned system and fixed pattern noise in staring cameras requires that the MRT formulation be appropriately modified.

Shown in Fig. 13 are BLIP (for 70 percent quantum efficiency) MRT curves at 300 K background temperature for narrow-field-of-view (high-resolution) sensors in the MWIR and LWIR spectral bands. Two LWIR curves are included to show the impact of matching the diffraction-limited blur

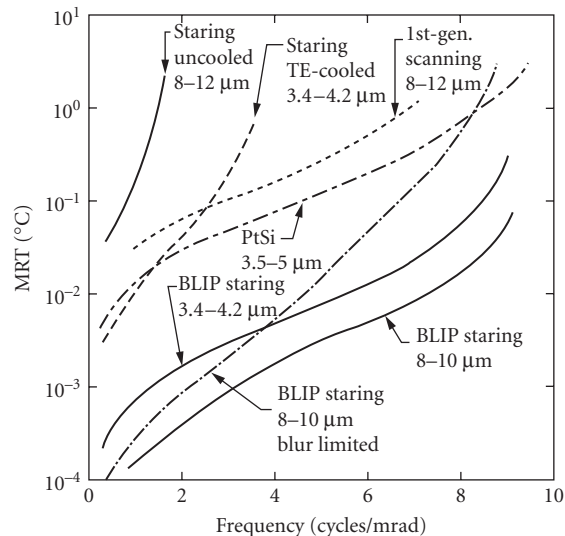


FIGURE 13 BLIP MRT for representative staring FPA configurations in the various bands.

to the pixel pitch versus $2\times$ oversampling of the blur. The latter case commonly arises, for example, when the LWIR FPA is miniaturized to minimize die size for enhancing yield and, for hybrid FPAs, alleviating thermal expansion mismatch. Also included for comparison are representative curves for first-generation scanning, staring uncooled, staring TE-cooled and staring PtSi sensors assuming 0.1, 0.1, 0.05, and 0.1 K NE Δ Ts, respectively. Theoretically, the staring MWIR sensors have order of magnitude better sensitivity while the staring LWIR bands have two orders of magnitude better sensitivity than the first-generation sensor. In practice, due to charge-handling limitations, an LWIR sensor has only slightly better MRT than the MWIR sensor. The uncooled sensor is useful for short-range applications such as a driver's aid in modern automobiles; the TE-cooled sensor provides longer range than the uncooled, but less than the high-density PtSi-based cameras.

33.6 CURRENT STATUS AND FUTURE TRENDS

Status

After over three decades of ongoing development, today's second-generation infrared focal plane arrays have typically 1000 times more pixels and up to 10 times higher sensitivity than first-generation devices. FPA format is consequently now set by application need, rather than a technological barrier. Nevertheless, there is ongoing motivation for fully achieving theoretically limited performance, especially at elevated operating temperatures. Development is hence continuing in the form of third-generation FPA technology. While specifications for third-generation FPAs encompass a broad range of needs, common objectives require overcoming recurring practical limits with respect to sensitivity, frame rate, power dissipation, multispectral capability, and cost. The common methodology going forward, regardless of specific mission requirements, is to dramatically increase on-FPA functionality via system-on-chip integration. While second-generation technology was largely driven by defense-oriented R&D funding, it is likely that third-generation technology will more directly leverage emerging commercial foundry capability for 3D-IC assembly; this emerging commercial interest in 3D-IC integration should dramatically lower infrared FPA cost while further improving FPA performance.

Figure 14 summarizes the typical performance of the most prominent detector technologies. The figure, a plot of the D_{th}^* (300 K, 0° field-of-view) versus operating temperature, clearly shows the performance advantage that the intrinsic photovoltaics have over the other technologies. Thermal detectivity is used here to compare the various technologies for equivalent NE Δ T irrespective of wavelength. While the extrinsic silicon detectors offer very high sensitivity, high producibility, and very long cutoff wavelengths, the very low operating temperature is often prohibitive. Also shown is the relatively low and slightly misleading detectivity of PtSi, which is offset by its

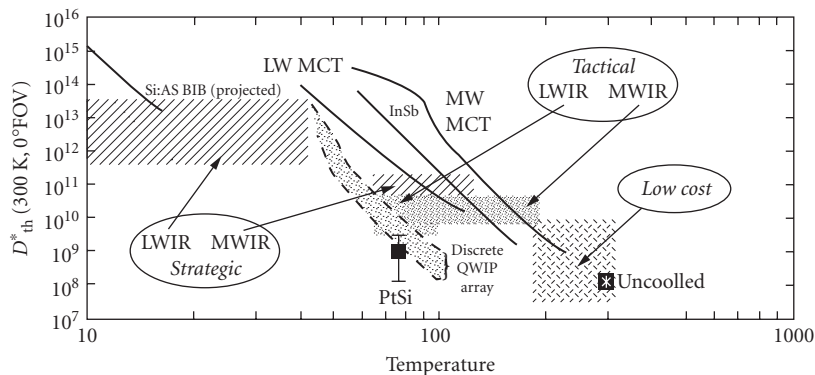


FIGURE 14 Comparison of photon detector D_{th}^* (300 K, 0°) for various IR technologies for equivalent NE Δ T.

Similar in performance at cryogenic temperatures, InSb and HgCdTe have comparable array size and pixel yield at MWIR cutoff wavelengths. Wavelength tunability and high radiative efficiency, however, have often made HgCdTe the preferred material because the widest possible bandgap semiconductor can be configured and thus the highest possible operating temperature achieved for a given set of operating conditions. The associated cooling and system power requirements can thus be optimally distributed.

FPA costs are currently very similar for all the second-generation FPA technologies. Though it is often argued that FPA cost for IR cameras will be irrelevant once it reaches a certain minimum production volume, nevertheless the key determinants as to which FPA technology becomes ubiquitous in the coming decade are availability and cost.

Future Trends and Technology Directions

The 1980s saw the maturation of PtSi and the emergence of HgCdTe and InSb as producible MWIR detector materials. Many indicators pointed to the 1990s being the decade that IR FPA technology would enter the consumer marketplace, but penetration did not actually occur until a decade later. Figure 16 shows the chronological development of IR FPAs including monolithic and hybrid technologies. Specifically compared is the development of various hybrids (primarily MWIR) to PtSi, the pace-setting technology with respect to array size. The hybrid FPA data includes Pb-salt, HgCdTe, InSb, and PtSi devices. This database suggests that monolithic PtSi led all other technologies with respect to array size by about 2 years and clearly shows the thermal mismatch barrier confronted by hybrid FPA developers in the mid-1980s.

In addition to further increases in pixel density to $>10^{16}$ pixels, several trends are clear. Future arrays will have much more on-chip signal processing, need less cooling, have higher sensitivity (particularly intrinsic LWIR FPAs), and offer multispectral capability. If the full performance potential of the uncooled technologies is realized, either the microbolometer arrays or, less likely, the pyroelectric arrays will capture the low-cost markets. It is not unreasonable that the uncooled arrays may obsolete “low-cost” PtSi, QWIP, and HIP FPAs, and render the intrinsic TE-cooled developments inconsequential. To improve hybrid reliability, alternative detector substrate materials including silicon along with alternative readout materials will become sufficiently mature to begin monolithic integration of the intrinsic materials with highest radiative performance. True optoelectronic FPAs consisting of IR sensors with optical output capability may be developed for reducing

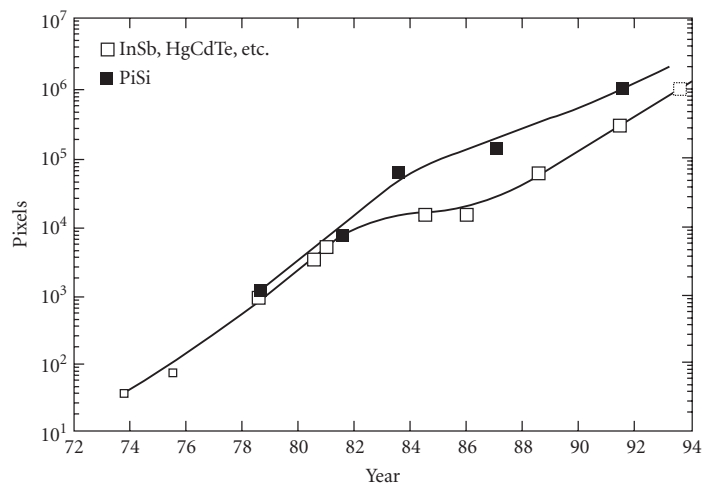


FIGURE 16 Chronological development of IR FPAs.

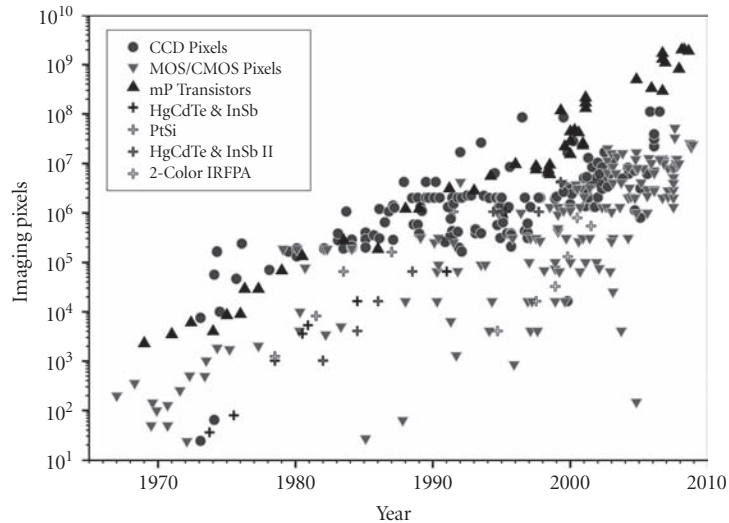


FIGURE 17 Integrated Circuit Chronology.

thermal loading and improving immunity to noise pickup. Availability of inexpensive, commercial devices is imminent along with the development of IR neural networks for additional signal processing capability.

In conclusion, it is likely instructive to compare the development of infrared sensors with both visible imaging sensors and other commercial integrated circuits, including semiconductor memory and microprocessors. All these products commonly share the benefits derived from Moore's law.⁵⁶ Fig. 17 is such a chronology showing notable devices of all types. Early infrared devices tended to develop at a pace only a few years behind the overlying trend known as Moore's law. Multicolor sensors initially developed at a rapid pace since the infrastructure was already in place. All infrared devices lag visible sensor development, which is driven by consumer demand. The largest visible sensors now boast about 100 million pixels; these foreshadow the size of upcoming infrared sensors once reliability issues and packaging costs are fully contained.

33.7 REFERENCES

1. W. P. McCracken, "CCDing in the Dark," *IEEE Spectrum*. (1992).
2. P. R. Norton, "Infrared Image Sensors," *Opt. Eng.* **30**:11 (1991).
3. A. Rogalski and J. Piotrowski, "Intrinsic Infrared Detectors," *Progress in Quantum Electronics* 12, Pergamon Press, 1988, pp. 2–3.
4. S. B. Stetson, D. B. Reynolds, M. G. Stapelbroek, and R. L. Stermer, "Design and Performance of Blocked-Impurity-Band Detector Focal Plane Arrays," *SPIE Proceedings*, vol. 686 (1986).
5. M. D. Petroff, G. Stapelbroek, and W. A. Kleinhans, "Solid State Photomultiplier," U.S. Patent No. 4,586,068 (1983).
6. B. F. Levine, K. K. Choi, C. G. Bethea, J. Walker, and R. J. Malik, *Appl. Phys. Lett.* **50**:1092 (1987).
7. C. S. Wu et al., "Novel GaAs/AlGaAs Multiquantum-Well Schottky Junction Device and Its Photovoltaic LWIR Detection," *IEEE Trans. Electron Devices*, **ED-39**:2 (1992).
8. J. Cooper, *Rev. Sci. Instrum.* **33**:92 (1962).

9. C. Hanson, H. Beratan, R. Owen, M. Corbin, and S. McKenney, "Uncooled Thermal Imaging at Texas Instruments," *SPIE Proceedings*, vol. 1735, 1992.
10. D. H. Pommerrenig, "Extrinsic Silicon Focal Plane Array," *SPIE Proceedings*, vol. 443, 1983.
11. (a) W. F. Kosonocky, "Review of Infrared Image Sensors with Schottky-Barrier Detectors," *Optoelectronics—Devices and Technologies*, vol. 6, no. 2, December 1991, pp. 173–203. (b) W. F. Kosonocky, "State-of-the-Art in Schottky-Barrier IR Image Sensors," *SPIE Proceedings*, vol. 1681, Orlando, Fla., 1992.
12. Private communications with N. A. Foss from Sensors and Systems Center, Honeywell, Inc., Bloomington, Minn.
13. C. G. Roberts, "HgCdTe Charge Transfer Focal Planes," *SPIE Proceedings*, vol. 443, 1983.
14. M. D. Gibbons et al., "Status of IrSb Charge Injection Device (CID) Detection Technology," *SPIE Proceedings*, vol. 443, 1983.
15. M. H. White et al., "Characterization of Surface CCD Image Arrays at Low Light Levels," *IEEE J. Solid-State Circuits*, vol. SC-9, 1974, pp. 1–12.
16. D. J. Sauer, F. V. Shallcross, F. L. Hsueh, G. M. Meray, P. A. Levine, H. R. Gilmartin, T. S. Villani, B. J. Esposito, and J. R. Tower, "640 × 480 MOS PtSI IR Sensor," *SPIE Proceedings*, vol. 1540, 1991, pp. 285–296.
17. L. J. Kozlowski et al., "Comparison of Passive and Active Pixel Schemes for CMOS Visible Imagers," *SPIE* 3360, 1998.
18. W. F. Kosonocky, T. S. Villani, F. V. Shallcross, G. M. Meray, and J. J. O'Neill, III, "A Schottky-Barrier Image Sensor with 100% Fill Factor," *SPIE Proceedings*, vol. 1308, 1990, pp. 70–80.
19. M. Denda, M. Kimata, S. Iwade, N. Yutani, T. Kondo, and N. Tsubouchi, "4 × 4096-Element SWIR Multispectral Focal Plane Array," *SPIE Proceedings*, vol. 819–824, 1987, pp. 279–286.
20. K. Chow, J. P. Rode, D. H. Seib, and J. D. Blackwell, "Hybrid Infrared Focal Plane Arrays," *IEEE Trans. Electron Devices* **ED-29**(1) (January, 1982).
21. E. R. Gertner, W. E. Tennant, J. D. Blackwell, and J. P. Rode, "HgCdTe on Sapphire—A New Approach to Infrared Detector Arrays," *J. Cryst. Growth* **72**:465 (1985).
22. E. R. Gertner, *Ann. Rev. Mater. Sci.* **15**:303–328 (1985).
23. E. R. Fossum, "Infrared Readout Electronics," *SPIE Proceedings*, vol. 1684 (1992).
24. D. F. Barbe, "Imaging Devices Using the Charge-Coupled Concept," *Proc. IEEE* **63**(1) (1975).
25. P. Felix, M. Moulin, B. Munier, J. Portmann, and J.-P. Reboul, "CCD Readout of Infrared Hybrid Focal Plane Arrays," *IEEE Trans. Electron Devices* **ED-27**(1) (1980).
26. S. T. Baier, "Complementary Heterostructure (CHFET) Readout Technology for Infrared Focal Plane Arrays," *SPIE Proceedings*, vol. 1684 (1992).
27. R. Sahai, R. L. Pierson, R. J. Anderson, E. H. Martin, E. A. Sover, and J. Higgins, "GaAs CCD's with Transparent (ITO) Gates for Imaging and Optical Signal Processing," *IEEE Electron Device Lett.* **EDL-4** (1983).
28. L. J. Kozlowski and R. E. Kezer, "2 × 64 GaAs Readout for IR FPA Application," *SPIE Proceedings*, vol. 1684 (1992).
29. J. S. Bugler and P. G. A. Jespers, "Charge Pumping in MOS Devices," *IEEE Trans. on Electron Devices* **ED-16**:3 (1969).
30. M. F. Tompsett, "Surface Potential Equilibration Method of Setting Charge in Charge-Coupled Devices," *IEEE Trans. Electron Devices* **ED-22**:6 (1975).
31. L. J. Kozlowski, K. Vural, W. E. Tennant, R. E. Kezer, and W. E. Kleinhans, "10 × 132 CMOS/CCD Readout with 25 μm Pitch and On-Chip Signal Processing Including CDS and TDI," *SPIE Proceedings*, vol. 1684 (1992).
32. M. F. Tompsett, "The Quantitative Effects of Interface States on the Performance of Charge-Coupled Devices," *IEEE Trans. Electron Devices* **ED-20**:1 (1973).
33. P. W. Bosshart, "A Multiplexed Switched Capacitor Filter Bank," *IEEE Journal Solid-State Circuits*, **SC-15**:6 (1980).
34. W. S. Chan, "Detector-Charge-Coupled Device (CCD) Interface Methods," *SPIE Proceedings*, vol. 244 (1980).
35. R. M. Swanson and J. D. Meindl, "Ion-Implanted Complementary MOS Transistors in Low-Voltage Circuits," *IEEE Journal of Solid-State Circuits* **SC-7**:2 (1972).
36. R. Troutman and S. N. Chakravarti, "Subthreshold Characteristics of Insulated-Gate Field-Effect Transistors," *IEEE Trans. Circuit Theory* **CT-20**:6 (1973).

37. N. Bluzer and R. Stehlik, "Buffered Direct Injection of Photocurrents into Charge-Coupled Devices," *IEEE Journal of Solid-State Circuits* **SC-13**:1 (1978).
38. C. C. Enz, E. A. Vittoz, and F. Krummenacher, "A CMOS Chopper Amplifier," *IEEE Journal of Solid-State Circuits* **SC-22**:3 (1987).
39. S. G. Chamberlain and J. P. Y. Lee, "A Novel Wide Dynamic Range Silicon Photodetector and Linear Imaging Array," *IEEE Trans. Electron Devices* **ED-31**:2 (1984).
40. F. Krummenacher, E. Vittoz, and M. DeGrauwe, "Class AB CMOS Amplifier for Micropower SC Filters," *Electron. Lett.* **17**:13 (1981).
41. E. Vittoz and J. Fellrath, "CMOS Analog Integrated Circuits Based on Weak Inversion Operation," *IEEE Journal Solid State Circuits* **SC-12**:3 (1977).
42. R. D. Hudson, *Infrared System Engineering*, John Wiley and Sons, 1969.
43. J. M. Mooney et al., "Responsivity Nonuniformity Limited Performance of Infrared Staring Cameras," *Opt. Eng.* **28**:1151 (1989).
44. D. L. Shumaker, J. T. Wood, and C. R. Thacker, *FLIR Performance Handbook*, DCS Corporation, Alexandria, Va. (1988).
45. N. Yutani, M. Kimata, H. Yagi, J. Nakanishi, S. Nagayoshi, and N. Tsubouchi, "1040 × 1040 Element PtSi Schottky-Barrier IR Image Sensor," IEDM, Washington, D.C., 1991.
46. T. S. Villani, W. F. Kosonocky, F. V. Shallcross, J. V. Groppe, G. M. Meray, J. J. O'Neill, III, and B. J. Esposito, "Construction and Performance of a 320 × 244-Element IR-CCD Imager with PtSi SBDs," *SPIE Proceedings*, vol. 1107-01, 1989, pp. 9-21.
47. D. J. Sauer, F. V. Shallcross, F. L. Hsueh, G. M. Meray, P. A. Levine, H. R. Gilmartin, T. S. Villani, B. J. Esposito, and J. R. Tower, "640 × 480 MOS PtSi IR Sensor," *SPIE Proceedings*, vol. 1540, 1991, pp. 285-296.
48. K. Konuma, N. Teranishi, S. Tohyama, K. Masubuchi, S. Yamagata, T. Tanaka, E. Oda, Y. Moriyama, N. Takada, and N. Yoshioka, "324 × 487 Schottky-Barrier Infrared Imager," *IEEE Trans. Electron Devices*, **37**(3):629-635 (1990).
49. K. Konuma, S. Tohyama, A. Tanabe, K. Masubuchi, N. Teranishi, T. Saito, and T. Muramatsu, "A 648 × 487 Pixel Schottky-Barrier Infrared CCD Image Sensor," *1991 ISSCC Digest of Tech. Papers*, 1991, pp. 156-157.
50. H. Elabd, Y. Abedini, J. Kim, M. Shih, J. Chin, K. Shah, J. Chen, F. Nicol, W. Petro, J. Lehan, M. Duron, M. Manderson, H. Balopole, P. Coyle, P. Cheng, and W. Shieh, "488 × 512 and 244 × 256-Element Monolithic PtSi Schottky IR Focal Plane Arrays," *SPIE Proceedings*, vol. 1107-29, presented at *SPIE Aerospace Sensor Symposium*, Orlando, Fla., March 1989.
51. E. T. Nelson, K. Y. Wong, S. Yoshizumi, D. Rockafellow, W. DesJardin, M. Elzinga, J. P. Lavine, T. J. Tredwell, R. P. Khosla, P. Sorlie, B. Howe, S. Brickman, and S. Refermat, "Wide Field of View PtSi Infrared Focal Plane Array," *SPIE Proceedings*, vol. 1308, 1990, pp. 36-44.
52. M. Kimata, M. Denda, N. Yutani, S. Iwade, and N. Tsubouchi, "A 512 × 512-Element PtSi Schottky-Barrier Infrared Image Sensor," *IEEE J. Solid-State Circuits*, vol. SC-22, no. 6, 1987, pp. 1124-1129.
53. H. Yagi, N. Yutani, S. Nagayoshi, J. Nakanishi, M. Kimata, and N. Tsubouchi, "Improved 512 × 512 IR CSD with Large Fill Factor and Large Saturation Level," *SPIE Proceedings*, vol. 1685-04, 1992.
54. J. Edwards, J. Gates, H. Altin-Mees, W. Connelly, and A. Thompson, "244 × 400 Element Hybrid Platinum Silicide Schottky Focal Plane Array," *SPIE Proceedings*, vol. 1308, 1990, pp. 99-100.
55. J. L. Gates, W. G. Connelly, T. D. Franklin, R. E. Mills, F. W. Price, and T. Y. Wittwer, "488 × 640-Element Hybrid Platinum Silicide Schottky Focal Plane Array," *SPIE Proceedings*, vol. 1540, *Infrared Technology XVII*, 1991, pp. 262-273.
56. G. E. Moore, "Cramming More Components onto Integrated Circuits," *Electronics Magazine* **38**:8, 1965.

This page intentionally left blank.

PART

7

**RADIOMETRY
AND
PHOTOMETRY**

This page intentionally left blank.

RADIOMETRY AND PHOTOMETRY

Edward F. Zalewski

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

34.1 GLOSSARY

A	area
A_1, A_2, A_s, A_d	area of surface 1, surface 2, a source, a detector, respectively
A_r	area of an image on the retina of a human eye
A_p	area of the pupil of a human eye
A_{in}, A_{out}, A_{sph}	area of an input port, output port, and sphere surface, respectively
b	distance from optic axis
c	the speed of light in a vacuum
C_e	photon-to-electron conversion efficiency, i.e., quantum efficiency of a photodetector
D	diameter
dA	infinitesimal element of area
dA_1, dA_2, dA_s, dA_d	infinitesimal element of area of surface 1, surface 2, a source, a detector, respectively
dL_λ	infinitesimal change in radiance per wavelength interval
dT	infinitesimal change in temperature
$d\Phi_{12}$	infinitesimal amount of radiant power transferred from point 1 to point 2
$d\lambda$	infinitesimal wavelength interval
$d\nu$	infinitesimal frequency interval
$d\Omega$	infinitesimal change in solid angle
E	irradiance, the incident radiant power per the projected area of a surface
E_v	illuminance, the photometric equivalent of irradiance
E_r	average illuminance in an image on the retina of a human eye
E_T	retinal illuminance in units of trolands
$E_T(\lambda)$	photopic retinal illuminance from a monochromatic source in trolands

$E'_T(\lambda)$	scotopic retinal illuminance from a monochromatic source in trolands
$E_{r\lambda}$	retinal spectral irradiance in absolute units: $\text{W nm}^{-1}\text{m}^{-2}$
f	focal length
$f\#$	F -number
g	fraction of light lost through the input and output ports of an averaging sphere
h	Planck's constant
h_s, h_d	object (source) height, image (detector) height
I	radiant intensity, the emitted or reflected radiant power per solid angle
i	photoinduced current from a radiation detector
I_v	luminous intensity, the photometric equivalent of radiant intensity
k	Boltzmann's constant
K_m	luminous efficacy (i.e., lumen-to-watt conversion factor) for photopic vision
K'_m	luminous efficacy for scotopic vision
K_{ab}	nonlinearity correction factor for a photodetector
L	radiance, the radiant power per projected area and solid angle
L_{12}	radiance from point 1 into the direction of point 2
L_a, L_b	radiance in medium a , in medium b
L_e	radiance within the human eye
L_λ	radiance per wavelength interval
L_ν	radiance per frequency interval
L_v	luminance, the photometric equivalent of radiance
$L_v(\lambda)$	luminance of a monochromatic light source
M	exitance, the emitted or reflected radiant power per the projected area of a source
m	mean value
N	photon flux, the number of photons per second
n	index of refraction
n_a, n_b	index of refraction in medium a , in medium b
n_e	index of refraction of the ocular medium of the human eye
n_s, n_d	index of refraction in the object (i.e., source) region, in the image (i.e., detector) region
N_λ	photon flux per wavelength interval
N_ν	photon flux per frequency interval
$N_{E\lambda}$	photon flux irradiance on the retina of a human eye
Q	radiant energy
Q_λ	radiant energy per wavelength interval
Q_ν	radiant energy per frequency interval
R	responsivity of a photodetector, i.e., electrical signal out per radiant signal in
r	radius
r_s, r_d, r_{sph}	radius of a source, detector, sphere, respectively
$R(\lambda)$	spectral (i.e., per wavelength interval) responsivity of a photodetector
s	distance
s_{12}	length of the light ray between points 1 and 2
s_{sd}	length of the light ray between points on the source and detector
s_{pr}	distance from the pupil to the retina in a human eye

T	absolute temperature
t	time
U	photon dose, the total number of photons
$V(\lambda)$	spectral luminous efficiency function (i.e., peak normalized human visual spectral responsivity) for photopic vision
$V'(\lambda)$	spectral luminous efficiency function for scotopic vision
w	width
x_i	the i th sample in a set of measurements
α	absorptance, fraction of light absorbed
β_a, β_b	angle of incidence or refraction
γ	absorption coefficient of a solute
δ	angle of rotation between crossed polarizers
ε	emittance of a blackbody simulator
E	étendue
η	total number of sample measurements
θ_s, θ_d	angle between the light ray and the normal to a point on the surface of a source, of a detector
θ_1, θ_2	angle between the light ray and the normal to a surface at point 1, at point 2
κ	concentration of a solute
λ	wavelength
ν	frequency
ρ	fraction of light scattered or reflected
σ	standard deviation
σ_m	standard deviation of the mean
τ	transmittance, radiant signal out per radiant signal into a material
$\tau(\lambda)$	spectral (i.e., per wavelength interval) transmittance
$\tau_c(\lambda)$	spectral transmittance of the ocular medium of the human eye
Φ	radiant power or equivalently radiant flux
ϕ	half angle subtended by a cone
Φ_{in}, Φ_{out}	incoming radiant power, outgoing radiant power
Φ_r	luminous flux at the retina of the human eye
Φ_λ	radiant power per wavelength interval
Φ_ν	radiant power per frequency interval
Φ_v	photopic luminous flux, radiant power by photopic detectable human vision
Φ'_v	scotopic luminous flux, radiant power detectable by scotopic human vision
Ω	solid angle, a portion of the area on the surface a sphere per of the square of the sphere radius
Ω_a, Ω_b	solid angle in medium a , in medium b

34.2 INTRODUCTION

Radiometry is the measurement of the energy content of electromagnetic radiation fields and the determination of how this energy is transferred from a source, through a medium, and to a detector. The results of a radiometric measurement are usually obtained in units of power, i.e., in watts. However, the result may also be expressed as photon flux (photons per second) or in units of energy

(joules) or dose (photons). The measurement of the effect of the medium on the transfer of radiation, i.e., the absorption, reflection, or scatter, is usually called *spectrophotometry* and will not be covered here. Rather, the assumption is made here that the radiant power is transferred through a lossless medium.

Traditional radiometry assumes that the propagation of the radiation field can be treated using the laws of geometrical optics. That is, the radiant energy is assumed to be transported along the direction of a ray and interference or diffraction effects can be ignored. In those situations where interference or diffraction effects are significant, the flow of energy will be in directions other than along those of the geometrical rays. In such cases, the effect of interference or diffraction can often be treated as a correction to the result obtained using geometrical optics. This assumption is equivalent to assuming that the energy flow is via an incoherent radiation field. This assumption is widely applicable since most radiation sources are to a large degree incoherent. For a completely rigorous treatment of radiant energy flow, the degree of coherence of the radiation must be considered via a formalism based on the theory of electromagnetism as derived from Maxwell's equations.^{1,2} This complexity is not necessary for most of the problems encountered in radiometry.

In common practice, radiometry is divided according to regions of the spectrum in which different measurement techniques are used. Thus, vacuum ultraviolet radiometry, intermediate-infrared radiometry, far-infrared radiometry, and microwave radiometry are considered separate fields, and all are distinguished from radiometry in the visible and near-visible optical spectral region.

The reader should note that there is considerable confusion regarding the nomenclatures of the various radiometries. The terminology for radiometry that we have inherited is dictated not only by its historical origin,³ but also by that of related fields of study. By the late 1700s, techniques were developed to measure light using the human eye as a null detector in comparisons of sources. At about the same time, radiant heating effects were studied with liquid-in-glass thermometers and actinic (i.e., chemical) effects of solar radiation were studied by the photoinduced decomposition of silver compounds into metallic silver. The discovery of infrared radiation in 1800 and ultraviolet radiation in 1801 stimulated a great deal of effort to study the properties of these radiations. However, the only practical detectors of ultraviolet radiation at that time were the actinic effects—for infrared radiation it was thermometers and for visible radiation it was human vision. Thus actinometry, radiometry, and photometry became synonymous with studies in the ultraviolet, infrared, and visible spectral regions. Seemingly independent fields of study evolved and even today there is confusion because the experimental methods and terminology developed for one field are often inappropriately applied to another. Vestiges of the confusion over what constitutes photometry and radiometry are to be found in many places. The problems encountered are not simply semantic, since the confusion can often lead to substantial measurement error.

As science progressed, radiometry was in the mainstream of physics for a short time at the end of the nineteenth century, contributing the absolute measurement base that led to Planck's radiation law and the discovery of the quantum nature of radiation. During this period, actinic effects, which were difficult to quantify, became part of the emerging field of photochemistry. In spite of the impossibility of performing an absolute physical measurement using the human eye, it was photometry, however, that grew to dominate the terminology and technology of radiant energy measurement practice in this period. At the beginning of the nineteenth century, the reason photometry was dominant was that the most precise (not absolute) studies of radiation transfer relied on the human eye. By the end of the nineteenth century, the growth of industries such as electric lighting and photography became the economic stimulus for technological developments in radiation transfer metrology and supported the dominance of photometry. Precise photometric measurements using instrumentation in which the human eye was the detector continued into the last half of the twentieth century. The fact that among the seven internationally accepted base units of physical measurement there remains one unit related to human physiology—the candela—is an indication of the continuing economic importance of photometry.

Presently, the recommended practice is to limit the term photometry to the measurement of the ability of electromagnetic radiation to produce a visual sensation in a physically realizable manner, that is, via a defined simulation of human vision.⁴⁻⁵ Radiometry, on the other hand, is used to describe the measurement of radiant energy independent of its effect on a particular detector.

Actinometry is used to denote measurement of photon flux (photons per second) or dose (total number of photons) independent of the subsequent photophysical, photochemical, or photobiological process. Actinometry is a term that is not extensively used, but there are current examples where measurement of the “actinic effect of radiation” is an occasion to produce a new terminology for a specific photoprocess, such as for the Caucasian human skin reddening effect commonly known as sunburn. We do not attempt here to catalog the many different terminologies used in photometry and radiometry, instead the most generally useful definitions are introduced where appropriate.

This chapter begins with a discussion of the basic concepts of the geometry of radiation transfer and photon flux measurement. This is followed by several approximate methods for solving simple radiation transfer problems. Next is a discussion of radiometric calibrations and the methods whereby an absolute radiant power or photon flux measurement is obtained. The discussion of photometry that follows is restricted to measurements employing physical detectors rather than those involving a human observer. Because many esoteric terms are still in use to describe photometric measurements, the ones most likely to be encountered are listed and defined in the section on photometry.

It is not the intention that this chapter be a comprehensive listing or a review of the extensive literature on radiometry and photometry; only selected literature citations are made where appropriate. Rather, it is hoped that the reader will be sufficiently introduced to the conceptual basis of these fields to enable an understanding of other available material. There are many texts on general radiometry. Some of the recent books on radiometry are listed in the reference section.⁶⁻¹⁰ In addition, the subject of radiometry or photometry is often presented as a subset of another field of study and can therefore be found in a variety of texts. Several of these texts are also listed in the reference section.¹¹⁻¹⁴ Finally, the reader will also find material related to radiometry, photometry, colorimetry, and spectrophotometry in Chaps. 34 to 40 in this volume and Chap. 10, “Colorimetry” in Vol. III.

34.3 RADIOMETRIC DEFINITIONS AND BASIC CONCEPTS

Radiant Power and Energy

For a steadily emitting source, that is a radiation source with a continuous and stable output, radiometric measurement usually implies measurement of the power of the source. For a flashing or single-pulse source, radiometric measurement implies a measurement of the energy of the source.

Radiometric measurements are traditionally measurements of thermal power or energy. However, because of the quantum nature of most photophysical, photochemical, and photobiological effects, in many applications it is not the measurement of the thermal power in the radiation beam but measurement of the number of photons that would provide the most physically meaningful result. The fact that most radiometric measurements are in terms of watts and joules is due to the history of the field. The reader should examine the particular application to determine if a measurement in terms of photon dose or photon flux would not be more meaningful and provide insight for the interpretation of the experiment. (See section on “Actinometry” later in this chapter.)

Radiant Energy Radiant energy is the energy emitted, transferred, or received in the form of electromagnetic radiation.

Symbol: Q *Unit:* joule (J)

Radiant Power Radiant power or radiant flux is the power (energy per unit time t) emitted, transferred, or received in the form of electromagnetic radiation.

Symbol: Φ *Unit:* watt (W)

$$\Phi = \frac{dQ}{dt} \quad (1)$$

Geometrical Concepts

The generally accepted terminology and basic definitions for describing the geometry of radiation transfer are presented below. More extensive discussions of each of these definitions and concepts can be found in the references.⁴⁻¹⁴

The concepts of irradiance, intensity, and radiance involve the density of the radiant power (or energy) over area, solid angle, and area times solid angle, respectively.

In situations where the density or distribution of the radiation on a surface is the required quantity, then it is the irradiance that must be measured. An example of where an irradiance measurement would be required is the exposure of a photosensitive surface such as the photoresists used in integrated circuit manufacture. The irradiance distribution over the surface determines the local degree of exposure of the photoresist. A nonuniform irradiance distribution will result in overexposure and/or underexposure of regions across the piece and results in a defect in manufacture.

In an optical system where the amount of radiation transfer through the system is important, then it is the radiance that must be measured. The amount of radiation passing through the optical system is determined by the area of the source from which the radiation was emitted and the field of view of the optic, also known as the solid angle or collection angle. Radiance is often thought of as a property of a source, but the radiance at a detector is also a useful concept.

Both irradiance and radiance are defined for infinitesimal areas and solid angles. However, in practice, measurements are performed with finite area detectors and optics with finite fields of view. Therefore all measurements are in fact measurement of average irradiance and average radiance.

Irradiance and radiance must be defined over a projected area in order to account for the effect of area change with angle of incidence. This is easily seen from the observation that the amount of a viewed area diminishes as it is tilted with respect to the viewer. Specifically, the view of the area falls off as the cosine of the angle between the normal to the surface and the line of sight. This effect is sometimes called the *cosine law of emission* or the *cosine law of irradiation*.

Intensity is a term that is part of our common language and often a point of confusion in radiometry. Strictly speaking, intensity is definable only for a source that is a point. An average intensity is not a measurable quantity since the source must by definition be an infinitesimal point. All intensity measurements are an approximation, since a true point source is physically impossible to produce. It is an extrapolation of a series of measurements that is the approximation of the intensity. An accurate intensity measurement is one that is made at a very large distance and, consequently, with a very small signal at the detector and an unfavorable signal-to-noise ratio. Historically speaking, however, intensity is an important concept in photometry and, to a much lesser extent, it has some application in radiometry. Intensity is a property of a source, not a detector.

Irradiance Irradiance is the ratio of the radiant power incident on an infinitesimal element of a surface to the projected area of that element, dA_d , whose normal is at an angle θ_d to the direction of the radiation.

Symbol: E Unit: watt/meter² (W m⁻²)

$$E = \frac{d\Phi}{\cos\theta_d dA_d} \quad (2)$$

Exitance The accepted convention makes a distinction between the irradiance, the surface density of the radiation incident on a radiation detector (denoted by the subscript d), and the exitance, the surface density of the radiation leaving the surface of a radiation source (denoted by the subscript s).

Exitance is the ratio of the radiant power leaving an infinitesimal element of a source to the projected area of that element of area dA_s , whose normal is at an angle θ_s to the direction of the radiation.

Symbol: M Unit: watt/meter² (W m⁻²)

$$M = \frac{d\Phi}{\cos\theta_s dA_s} \quad (3)$$

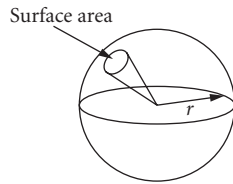


FIGURE 1 The solid angle at the center of the sphere is the surface area enclosed in the base of the cone divided by the square of the sphere radius.

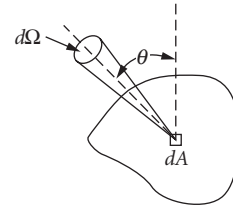


FIGURE 2 The radiance at the infinitesimal area dA is the radiant flux divided by the solid angle times the projection of the area dA onto the direction of the flux.

Intensity Radiant intensity (often simply “intensity”) is the ratio of the radiant power leaving a source to an element of solid angle $d\Omega$ propagated in the given direction.
Symbol: I Unit: watt/steradian (W sr^{-1})

$$I = \frac{d\Phi}{d\Omega} \quad (4)$$

Note that in the field of physical optics, the word *intensity* refers to the magnitude of the Poynting vector and thus more closely corresponds to irradiance in radiometric nomenclature.

Solid Angle The solid angle is the ratio of a portion of the area on the surface of a sphere to the square of the radius r of the sphere. This is illustrated in Fig. 1.
Symbol: Ω Unit: steradian (sr)

$$d\Omega = \frac{dA}{r^2} \quad (5)$$

It follows from the definition that the solid angle subtended by a cone of half angle ϕ , the apex of which is at the center of the sphere, is given by

$$\Omega = 2\pi(1 - \cos\phi) = 4\pi \sin^2 \frac{\phi}{2} \quad (6)$$

Radiance Radiance, shown in Fig. 2, is the ratio of the radiant power, at an angle θ_s to the normal of the surface element, to the infinitesimal elements of both projected area and solid angle. Radiance can be defined either at a point on the surface of either a source or a detector, or at any point on the path of a ray of radiation.

Symbol: L Unit: watt/steradian meter² ($\text{W sr}^{-1}\text{m}^{-2}$)

$$L = \frac{d\Phi}{\cos\theta_s dA_s d\Omega} \quad (7)$$

Radiance plays a special role in radiometry because it is the propagation of the radiance that is conserved in a lossless optical system; see “Radiance Conservation Theorem, Homogeneous Medium.” Radiance was often referred to as the brightness or the specific intensity, but this terminology is no longer recommended.

Spectral Dependence of Radiometric Quantities

Polychromatic Radiation Definitions For polychromatic radiation, the spectral distribution of radiant power (or radiant energy) is denoted as either radiant power (energy) per wavelength interval or radiant power (energy) per frequency interval.

Symbol: $\Phi_\lambda(Q_\lambda)$ Unit: watt/nanometer (W nm^{-1}); joule/nanometer (J nm^{-1});
or

Symbol: $\Phi_\nu(Q_\nu)$ Unit: watt/hertz (W Hz^{-1}); joule/hertz (J Hz^{-1})

It follows that $\Phi_\lambda d\lambda$ is the radiant power in the wavelength interval λ to $\lambda + d\lambda$, and $\Phi_\nu d\nu$ is the radiant power in the frequency interval ν to $\nu + d\nu$. The total radiant power over the entire spectrum is therefore

$$\Phi = \int_0^\infty \Phi_\lambda d\lambda \quad (8a)$$

or

$$\Phi = \int_0^\infty \Phi_\nu d\nu \quad (8b)$$

If λ is the wavelength in the medium corresponding to the frequency ν , and since $\nu = c/n\lambda$, where c is the speed of light in a vacuum and n is the index of refraction of the medium, then

$$d\nu = -\frac{c}{n\lambda^2} d\lambda \quad (9)$$

and

$$\lambda\Phi_\lambda = \nu\Phi_\nu \quad (10)$$

Since the wavelength changes with the index of refraction of the medium, it is becoming more common to use the vacuum wavelength, $\lambda = c/\nu$. It is particularly important in high-accuracy applications to state explicitly whether or not the vacuum wavelength is being used.

Spectral versions of the other radiometric quantities, i.e., radiant energy, radiance, etc., are defined similarly.

Polychromatic Radiation Calculations As an example of the application of the concept of the spectral dependence of a radiometric quantity, consider the calculation of the response of a radiometer consisting of a detector and a spectral filter. The spectral responsivity of a detector $R(\lambda)$ is the ratio of the output signal to the radiant input at each wavelength λ . The output is usually an electrical signal, such as a photocurrent i , and the input is a radiometric quantity, such as radiant power. The spectral transmittance of a filter $\tau(\lambda)$ is the ratio of the output radiant quantity to the input radiant quantity at each wavelength λ . For a spectral radiant power Φ_λ , the photocurrent i of the radiometer is

$$i = \int_0^\infty R(\lambda)\tau(\lambda)\Phi_\lambda d\lambda \quad (11)$$

In practice, either the responsivity of the detector or the transmittance of the filter are nonzero only within a limited spectral range. The integral need be evaluated only within the wavelength limits where the integrand is nonzero.

Photometry

The radiation transfer concepts, i.e., geometrical principles, of photometry are the same as those for radiometry. The exception is that the spectral responsivity of the detector, the human eye, is specifically defined. Photometric quantities are related to radiometric quantities via the spectral efficiency functions defined for the photopic and scotopic CIE Standard Observer. The generally accepted values of the photopic and scotopic human eye response function are represented in the "Photometry" section in Table 2.

Luminous Flux The photometric equivalent of radiant power is luminous flux, and the unit that is equivalent to the watt is the lumen. Luminous flux is spectral radiant flux weighted by the appropriate eye response function. The definition of luminous flux for the photopic CIE Standard Observer is

Symbol: Φ_v Unit: lumen (lm)

$$\Phi_v = K_m \int \Phi_\lambda V(\lambda) d\lambda \quad (12)$$

where $V(\lambda)$ is the spectral luminous efficiency function and K_m is the luminous efficacy for photopic vision. The spectral luminous efficacy is defined near the maximum, $\lambda_m = 555$ nm, of the photopic efficiency function to be approximately 683 lm W^{-1} .

Definitions of the Density of Luminous Flux

Illuminance Illuminance is the photometric equivalent of irradiance; that is, illuminance is the luminous flux per unit area.

Symbol: E_v Unit: lumen/meter² (1 m m^{-2})

$$E_v = \frac{d\Phi_v}{\cos\theta_d dA_d} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{\cos\theta_d dA_d} \quad (13)$$

Luminous intensity Luminous intensity is the photometric equivalent of radiant intensity. Luminous intensity is the luminous flux per solid angle. For historical reasons, the unit of luminous intensity, the candela—not the lumen—is defined as the base unit for photometry. However, the units for luminous intensity can either be presented as candelas or lumens/steradian.

Symbol: I_v Unit: candela or lumen/steradian (cd or lm sr^{-1})

$$I_v = \frac{d\Phi_v}{d\Omega} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{d\Omega} \quad (14)$$

Luminance Luminance is the photometric equivalent of radiance. Luminance is the luminous flux per unit area per unit solid angle.

Symbol: L_v Unit: candela/meter² (cdm^{-2})

$$L_v = \frac{d\Phi_v}{\cos\theta_s dA_s d\Omega} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{\cos\theta_s dA_s d\Omega} \quad (15)$$

Actinometry

Radiant Flux to Photon Flux Conversion Actinometric measurement practice closely follows that of general radiometry except that the quantum nature of light rather than its thermal effect is emphasized. In actinometry, the amount of electromagnetic radiation being transferred is measured in units of photons per second (photon flux). The energy of a single photon is

$$Q = h\nu \quad (16)$$

where ν is the frequency of the radiation and h is Planck's constant, $6.6261 \times 10^{-34} \text{ J s}$. For monochromatic radiant power Φ_λ , measured as watts and wavelength λ , measured as nanometers, the number of photons per second N_λ in the monochromatic radiant beam is

$$N_\lambda = 5.0341 \times 10^{15} n \lambda \Phi_\lambda \quad (17)$$

Photon Dose and the Einstein Dose is the total number of photons impinging on a sample. For a monochromatic beam of radiant power Φ_λ that irradiates a sample for a time t seconds, the dose U measured as Einsteins is

$$U = 8.3593 \times 10^{-9} n \lambda \Phi_\lambda t \quad (18)$$

The Einstein is a unit of energy used in photochemistry. An Einstein is the amount of energy in one mole (Avogadro's number, 6.0221×10^{23}) of photons.

TABLE 1 Radiation Transfer Terminology, Spectral Relationships

	Radiometric	Photometric	Actinometric
Base quantity:	Radiant power (also radiant flux)	Luminous flux	Photon flux
Units:	Watts/nanometer	Lumens	Photons/second
Conversion:	—	[W/nm] $K_m V(\lambda)$	[W/nm] $\lambda (hc)^{-1}$
Surface density:	Irradiance	Illuminance	Photon flux irradiance
Solid angle density:	Radiant intensity	Luminous intensity	Photon flux intensity
Solid angle and surface density:	Radiance	Luminance	Photon flux radiance

Conversions between Radiometry, Photometry, and Actinometry

Conversions between radiometric, photometric, and actinometric units is not simply one of determining the correct multiplicative constant to apply. As seen previously, the conversion between radiant power and photon flux requires that the spectral character of the radiation be known. It was also shown that, for radiometric to photometric conversions, the spectral distribution of the radiation must be known. Furthermore, there is an added complication for photometry where one must also specify the radiant power level in order to determine which CIE Standard Observer function is appropriate. Table 1, which summarizes the spectral radiation transfer terminology, may be helpful to guide the reader in determining the relationship between radiometric, photometric, and actinometric concepts. In Table 1, the power level is assumed to be high enough to restrict the photometric measurements to the range of the photopic eye response function.

Basic Concepts of Radiant Power Transfer

Radiance Conservation Theorem, Homogeneous Medium In a lossless, homogeneous isotropic medium, for a perfect optical system (i.e., having no aberrations) and ignoring interference and diffraction effects, the radiance is conserved along a ray through the optical system. In other words, the spectral radiance at the image always equals the spectral radiance at the source.

It follows from Eq. (7), the definition of radiance, that for a surface A_1 with radiance L_{12} in the direction of a second surface A_2 with radiance L_{21} in the direction to a first surface, and joined by a light ray of length s_{12} , the net radiant power exchange between elemental areas on each surface is given by

$$\Delta\Phi = d\Phi_{12} - d\Phi_{21} = \frac{(L_{12} - L_{21}) \cos\theta_1 \cos\theta_2 dA_1 dA_2}{s_{12}^2} \quad (19)$$

where θ_1 and θ_2 are the angles between the ray s_{12} and the normals to the surfaces A_1 and A_2 , respectively. The transfer of radiant power and the terminology used in this discussion is depicted in Fig. 3.

The total amount of radiation transferred between the two surfaces is given by the integral over both areas as follows:

$$\Phi = \iint \frac{(L_{12} - L_{21}) \cos\theta_1 \cos\theta_2}{s_{12}^2} dA_1 dA_2 \quad (20)$$

This is the generalized radiant power transfer equation for net exchange between two sources. In the specialized case of a source and receiver, the radiant power emitted by a receiver is zero by definition. In this case, the term L_{21} in Eq. (20) is zero.

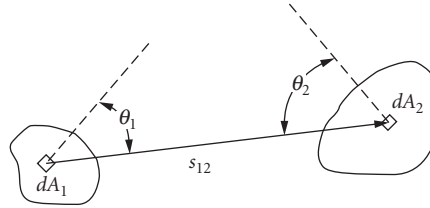


FIGURE 3 The radiant flux transferred between the infinitesimal areas dA_1 to dA_2 .

Refractive Index Changes In the case of a boundary between two homogeneous isotropic media having indices of refraction n_a and n_b , the angles of incidence and refraction at the interface β_a and β_b are related by Snell's law. If the direction of the light ray is oblique to the boundary between n_a and n_b , the solid angle change at the boundary will be

$$d\Omega_a = \frac{n_b^2 \cos \beta_b}{n_a^2 \cos \beta_a} d\Omega_b \quad (21)$$

Therefore the radiance change across the boundary will be

$$\frac{L_a}{n_a^2} = \frac{L_b}{n_b^2} \quad (22)$$

This result is obtained directly by substituting the optical path for the distance in Eq. (19) and considering that the radiance transferred across the boundary between the two media is unchanged. Optical path is the distance within the medium times the index of refraction of the medium.

In the case of an optical system having two or more indices of refraction, the radiance conservation theorem is more precisely stated as: In a lossless, homogeneous isotropic medium, for a perfect optical system (i.e., having no aberrations) and ignoring interference and diffraction effects, at a boundary between two media having different indices of refraction the radiance divided by the square of the refractive index is conserved along a ray through the optical system.

Radiative Transfer through Absorbing Media For radiation transmitted through an absorbing and/or scattering medium, the radiance is not conserved. This is not only because of the loss due to the absorption and/or scattering but the medium could also emit radiation. The emitted light will be due to thermal emission (see the discussion on blackbody radiation later in this chapter). In some cases, the medium may also be fluorescent. Fluorescence is the absorption of radiant energy at one wavelength with subsequent emission at a different wavelength.

Historically, the study of radiative transfer through absorbing and/or scattering media dealt with the properties of stellar atmospheres. Presently, there is considerable interest in radiative transfer measurements of the earth and its atmosphere using instruments on board satellites or aircraft. An accurate measure of the amount of reflected sunlight (approximately 400 to 2500 nm) or the thermally emitted infrared (wavelengths >2500 nm) requires correction for the absorption, scattering, and, in the infrared, the emission of radiation by the atmosphere. This specialized topic will not be considered here. Detailed discussion is available in the references.¹⁵⁻¹⁷

34.4 RADIANT TRANSFER APPROXIMATIONS

The solution to the generalized radiant power transfer equation is typically quite complex. However, there are several useful approximations that in some instances can be employed to obtain an estimate of the solution of Eq. (20). We shall consider the simpler case of a source and a detector rather

then the net radiant power exchange between two sources, since this is the situation commonly encountered in an optical system. In this case, Eq. (20) becomes

$$\Phi = \iint \frac{L \cos \theta_s \cos \theta_d}{s_{sd}^2} dA_s dA_d \quad (23)$$

where the subscripts s and d denote the source and detector, respectively. Here it is assumed the detector behaves as if it were a simple aperture. That is, it responds equally to radiation at any point across its surface and from any direction. Such a detector is often referred to as a cosine corrected detector. Of course, deviations from ideal detection behavior within the spatial and angular range of the calculation reduces the accuracy of the calculation.

Point-to-point Approximation: Inverse Square Law

The simplest approximations are obtained by assuming radiant flux transfer between a point source emitting uniformly in all directions and a point detector. The inverse square law is an approximation that follows directly from the definitions of intensity, solid angle, and irradiance, Eqs. (2), (4), and (5), respectively. The irradiance (at an infinitesimal area whose normal is along the direction of the light ray) times the square of the distance from a point source equals the intensity of the source

$$I = \frac{\Phi}{A} s^2 = E s^2 \quad (24)$$

The relationship between the uniformly emitted radiance and the intensity of a point source is obtained similarly from Eqs. (4) and (7):

$$L = \frac{I}{A_s} \quad (25)$$

These point-to-point relationships are perhaps most important as a test of the accuracy of a radiation transfer calculation at the limit as the areas approach zero.

Lambertian Approximation: Uniformly Radiant Areas

Lambertian Sources A very useful concept for the approximation of radiant power transfer is that of a source having a radiance that is uniform across its surface and uniformly emits in all directions from its surface. Such a uniform source is commonly referred to as a lambertian source.

For the case of a lambertian source, Eq. (23) becomes

$$\Phi = L \iint \frac{\cos \theta_s \cos \theta_d}{s_{sd}^2} dA_s dA_d \quad (26)$$

Configuration factor The double integral in Eq. (26) has been given a number of different names: configuration factor, radiation interaction factor, and projected solid angle. There is no generally accepted terminology for this concept, although configuration factor appears most frequently. Analytical solutions to the double integral have been found for a variety of different shapes of source and receiver. Tabulations of these exact solutions to the integral in Eq. (26) are usually found in texts on thermal engineering,^{18,19} under the heading of radiant heat transfer or configuration factor.

Radiation transfer between complex shapes can often be determined by using various combinations of configuration factors. This technique is often referred to as configuration factor algebra.¹⁸ The surfaces are treated as pieces, each with a calculable configuration factor, and the separate configuration factors are combined to obtain the effective configuration factor for the complete surface.

Étendue The double integral in Eq. (26) is often used as a means to characterize the flux-transmitting capability of an optical system in a way that is taken to be independent of the radiant properties of the source. Here the double integral is written as being over area and solid angle:

$$\Phi = L \iint \cos\theta_d dA_s d\Omega \quad (27)$$

In this case, the surface of the lambertian source is assumed perpendicular to the optic axis and to lie in the entrance window of the optical system. The solid angle is measured from a point on the source to the entrance pupil. The étendue E of an optical system of refractive index n is defined as

$$E = n^2 \iint \cos\theta_d dA_s d\Omega \quad (28)$$

Equation (28) is sometimes referred to as the throughput of an optical system.

Total flux into a hemisphere The total amount of radiation emitted from a lambertian source of area dA_s into the hemisphere centered at dA_s (or received by a hemispherical, uniform detector centered at dA_s) is obtained from integrating Eq. (26) over the area A_d . Note that the ray s_{sd} is everywhere normal to the surface of the hemisphere; i.e., $\cos\theta_d = 1$.

$$\Phi = L\pi \int dA_s \quad (29)$$

Using Eq. (3), the definition of the exitance, the radiance at each point on the surface of the source is

$$L = \frac{M}{\pi} \quad (30)$$

Because of the relationship expressed in Eq. (30), Eq. (26) is often written in terms of the exitance.

$$\Phi = M\pi \iint \frac{\cos\theta_s \cos\theta_d}{s_{sd}^2} dA_s dA_d \quad (31)$$

In this case, the factor π is considered to be part of the configuration factor. Note again that there is no generally accepted definition of the configuration factor.

Radiation transfer between a circular source and detector The particular case of radiation transfer between circular apertures, the centers of which are located along the same optical axis as shown in Fig. 4, is a configuration common to many optical systems and is therefore illustrated here. The

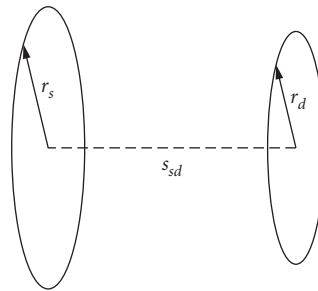


FIGURE 4 Radiant flux transfer between two circular apertures normal and concentric to the axis joining them.

radius of the source (or first aperture) is r_s , the detector (second aperture) radius is r_d , and the distance between the centers is s_{sd} . The exact solution of the integral in Eq. (26) yields

$$\Phi = \frac{2L(\pi r_s r_d)^2}{r_s^2 + r_d^2 + s_{sd}^2 + [(r_s^2 + r_d^2 + s_{sd}^2)^2 - 4r_s^2 r_d^2]^{1/2}} \quad (32)$$

This result can be approximated for the case where the sum of the squares of the distance and radii is large compared to the product of the radii, that is, $(r_s^2 + r_d^2 + s_{sd}^2) \gg 2r_s r_d$, so that Eq. (32) reduces to

$$\Phi \cong \frac{L(\pi r_s r_d)^2}{r_s^2 + r_d^2 + s_{sd}^2} \quad (33)$$

From this expression the irradiance at the detector can be obtained

$$E = \frac{\Phi}{A_d} \cong \frac{L A_s}{r_s^2 + r_d^2 + s_{sd}^2} \cong \frac{L A_s}{s_{sd}^2} \quad (34)$$

where A_s is the area of the lambertian disk and A_d is the detector area. The approximation at the extreme right is obtained by assuming that the radii are completely negligible with respect to the distance. This is the same result that would be obtained from a point-to-point approximation.

Off-axis irradiance: cosine-to-the-fourth approximation Equation (34) describes the irradiance from a small lambertian disk to a detector on the ray axis and where both surfaces are perpendicular to the ray. If the detector is moved off-axis by a distance b as depicted in Fig. 5, the ray from A_s to A_d will then be at an angle with respect to the normal at both surfaces as follows

$$\theta_s = \theta_d = \theta = \tan^{-1} \left(\frac{b}{s_{sd}} \right) \quad (35)$$

The projected areas are then $(A_s \cos \theta)$ and $(A_d \cos \theta)$. In addition, the distance from the source to the detector increases by the factor $(1/\cos \theta)$. The radiant power at a distance b away from the axis therefore decreases by the fourth power of the cosine of the angle formed between the normal to the surface and the ray.

$$\Phi \cong \frac{L A_s A_d}{s_{sd}^2} \cos^4 \theta \quad (36)$$

Since the radiance is conserved for propagation in a lossless optical system, Eq. (36) also approximates the radiant power from an off-axis region of a large lambertian source received at a small detector. The approximate total radiant power received at the detector would then be the sum of the radiant power contributed by each region of the source.

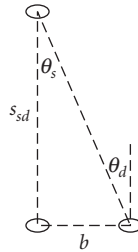


FIGURE 5 illustration of the cosine-fourth effect on irradiance, displacement of the receiving surface by a distance b .

Spherical lambertian source In order to compute the radiant power at a point at a distance s_{sd} from the center of a spherical lambertian source of radius r_{sph} , it is not necessary to explicitly solve the integrals in the radiation transfer equation. The solution is readily obtained from the symmetry of the lambertian sphere. Using the relationship between the exitance and radiance of a lambertian source [Eq. (30)], the total radiation power emitted by the source is obtained from the product of the surface area of the source times the exitance.

$$\Phi = 4\pi^2 r_{sph}^2 L \quad (37)$$

The radiant power is isotropically emitted. Therefore, the irradiance at any point on an enclosing sphere of radius s_{sd} is the total radiant power divided by the area of the enclosing sphere.

$$E = \frac{\pi r_{sph}^2 L}{s_{sd}^2} \quad (38)$$

Note that the irradiance from a spherical lambertian source follows the inverse square law at all distances from the surface of the sphere. The intensity of a spherical lambertian source is

$$I = \pi r_{sph}^2 L \quad (39)$$

Radiant Flux Transfer through a Lambertian Reflecting Sphere A lambertian reflector is a surface that uniformly scatters a fraction ρ of the radiation incident upon it.

$$L = \frac{\rho E}{\pi} \quad (40)$$

where E is the irradiance.

A spherical enclosure whose interior is coated with a material that approximates a lambertian reflector is a widely used tool in radiometry and photometry.²⁰ Such spheres are used either for averaging a nonuniform radiant power distribution (averaging sphere) or for measuring the total amount of radiant power emitted from a source (integrating sphere).

The sphere has the useful property whereby the solid angle subtended by any one section of the wall times the projected area is constant over all other points on the inside surface of the sphere. Therefore, if radiation falling on any point within the sphere is uniformly reflected, the reflected radiation will be uniformly distributed, i.e., produce uniform irradiance, throughout the interior. This result follows directly from the symmetry of the sphere.

Consider a sphere of radius r_{sph} and the radiant power transfer between two points on the inner surface. The normals to the two points are radii of the sphere and form an isosceles triangle when taken with the ray joining the points. Therefore, the angles between the ray and the normals to each point are equal. From Eq. (26)

$$\Phi = L \iint \frac{\cos^2 \theta}{s_{sd}^2} dA_s dA_d \quad (41)$$

The length of the ray joining the points is $2r_{sph} \cos \theta$. The irradiance is therefore

$$E = \frac{\Phi}{A_d} = \frac{LA_s}{4r_{sph}^2} \quad (42)$$

which is independent of the angle θ . If Φ_{in} is the radiant power entering the sphere, the irradiance at any point on the sphere after a single reflection will be

$$E = \frac{\rho \Phi_{in}}{4\pi r_{sph}^2} \quad (43)$$

A fraction ρ of the flux will be reflected and again uniformly distributed over the sphere. After multiple reflections the irradiance at any point on the wall of the sphere is

$$E = \frac{(\rho + \rho^2 + \rho^3 + \dots)\Phi_{\text{in}}}{4\pi r_{\text{sph}}^2} = \frac{\rho\Phi_{\text{in}}}{(1-\rho)A_{\text{sph}}} \quad (44)$$

where A_{sph} is the surface area of the sphere. The flux Φ_{out} exiting the sphere through a port of area A_{out} is

$$\Phi_{\text{out}} = \frac{\rho\Phi_{\text{in}}A_{\text{out}}}{(1-\rho)A_{\text{sph}}} \quad (45)$$

In Eq. (45) it is assumed that the loss of radiation at the entrance and exit ports is negligible and does not affect the symmetry of the radiation distribution.

The effect of the radiation lost through the entrance and exit ports is approximated as follows. After the first reflection, the fraction of radiation lost in each subsequent reflection is equal to the combined areas of the ports divided by the sphere area. Therefore the fraction reflected within the sphere is

$$g = 1 - \frac{A_{\text{in}} + A_{\text{out}}}{A_{\text{sph}}} \quad (46)$$

Using this in Eq. (44) yields

$$\Phi_{\text{out}} = \frac{\rho\Phi_{\text{in}}A_{\text{out}}}{(1-\rho g)A_{\text{sph}}} \quad (47)$$

Since the sphere is approximately a lambertian source, the radiance at the exit port is

$$L = \frac{\rho\Phi_{\text{in}}}{(1-\rho g)\pi A_{\text{sph}}} \quad (48)$$

Radiometric Effect of Stops and Vignetting

Refer to Fig. 6 for an illustration of these definitions. The *aperture stop* of an optical system is an aperture near the entrance to the optical system that determines the size of the bundle of rays leaving the source that can enter the optical system.

The *field stop* is an aperture within the optical system that determines the maximum angle of the rays that pass through the aperture stop that can reach the detector. The position and area of the field stop determines the field of view of the optical system. The field stop limits the extent of the source that is represented in its image at the detector.

The image of the aperture stop in object space, i.e., in the region of the source, is the *entrance pupil*. The image of the aperture stop in image space, i.e., in the region of the detector, is the *exit pupil*. Light rays that pass through the center of the aperture stop also pass through the centers of the images of the aperture stop at the entrance and exit pupils. Since all of the light entering the optical system must pass through the aperture stop, all of the light reaching the detector appears to pass through the exit pupil.

The field stop defines the solid angle within the optical system, the system field of view. When viewed from the image, the field stop of an optical system takes on the radiance of the object being imaged. This is a useful radiometric concept since a complex optical system can often be approximated as an exit pupil having the same radiance as the object being imaged (modified by the system transmission losses). The direction in which the radiation in the image appears to be emitted is,

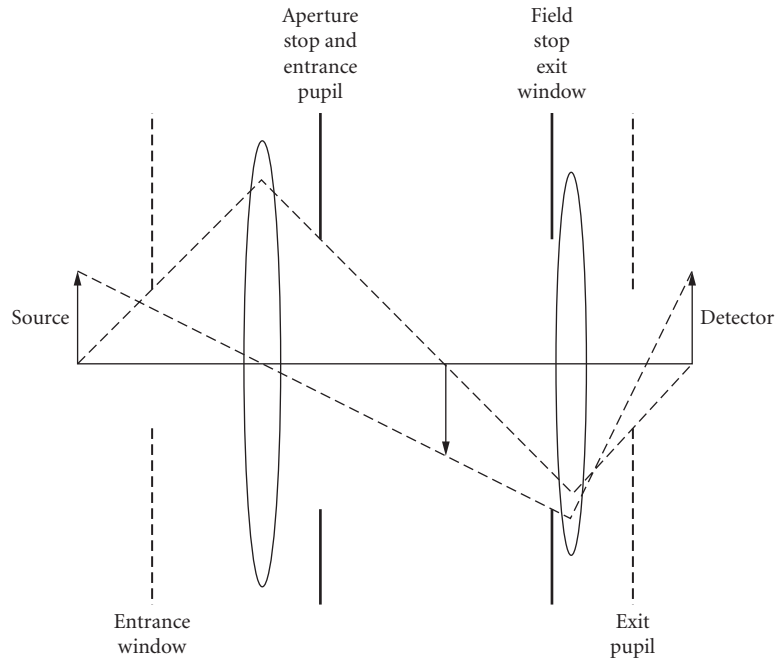


FIGURE 6 Schematic showing the relative positions of the stops, pupils, and windows in a simple optical system.

of course, limited by the aperture at the field stop. A word of caution: if the object is small, its image will be limited by diffraction effects and its radiance will depart even further from the extended source (large area) approximations used here.

The *entrance window* is the image of the field stop at the source and the *exit window* is the image of the field stop at the detector. If the field stop coincides with the detector, i.e., the detector is in the image plane of the optical system, then the entrance window will correspond with the object plane on the source. If the field stop does not coincide with the image plane at the detector, then because of parallax, different portions of the source will be visible from different points within the exit pupil. This condition, known as *vignetting*, causes a decrease in the irradiance at the off-axis points on the detector or image plane.

Approximate Radiance at an Image

Aplanatic Optical Systems Except for rays that lie on the optic axis, the radiance of an image must be based on a knowledge of the image quality since any aberrations introduced by the optical system divert some of the off-axis rays away from the image.

Consider a well-corrected optical system that is assumed to be aplanatic for the source and image points. That is, the optical system obeys Abbe's sine condition which is

$$n_s h_s \sin \theta_s = n_d h_d \sin \theta_d \quad (49)$$

where n_s and n_d are the refractive indices of the object (source) and image (detector) spaces, h_s and h_d are the object and image heights, and θ_s and θ_d are the angles between the off-axis rays and the

optic axis in object and image space. From Eq. (27) the flux radiated by a small lambertian source of area A_s into the solid angle of the optical system is

$$\Phi \cong 2\pi L A_s \int_0^{\theta_s} \cos\theta \sin\theta d\theta = \pi L A_s \sin^2 \theta_s \quad (50)$$

The differential of the solid angle is obtained from Eq. (6). Since Φ is the radiant flux at the image and A_d is the area of the image, the irradiance at the image is

$$E = \frac{\Phi}{A_d} = \frac{\pi L A_s}{A_d} \sin^2 \theta_s \quad (51)$$

If h_s and h_d are the radii of circular elements A_s and A_d , then according to Abbe's sine condition

$$\frac{A_s}{A_d} = \frac{n_d^2 \sin^2 \theta_d}{n_s^2 \sin^2 \theta_s} \quad (52)$$

The irradiance at the image is

$$E = \frac{\pi L n_d^2}{n_s^2} \sin^2 \theta_d \quad (53)$$

Numerical Aperture and F-number The quantity $n_d \sin \theta_d$ in Eq. (53) is called the *numerical aperture* of the imaging system. The irradiance of the image is proportional to the square of the numerical aperture. Geometrically speaking, the image irradiance increases with the angle of the cone of light converging on the image.

Another approximate measure of the image irradiance of an optical system is the *F-number*, $f\#$ (sometimes called the focal ratio) defined by the ratio of focal length (in the image space) f to the diameter D of the entrance pupil. For a source at a very large distance

$$f\# = \frac{f}{D} = \frac{1}{2 \tan \theta_d} \cong \frac{1}{2 \sin \theta_d} \quad (54)$$

The approximate image irradiance expressed in terms of the *F-number* is

$$E \cong \pi L \left(\frac{n_d}{2n_s f\#} \right)^2 \quad (55)$$

34.5 ABSOLUTE MEASUREMENTS

An absolute measurement, often referred to as an absolute calibration, is a measurement that is based upon, i.e., derived from, one of the internationally recognized units of physical measurement. These units are known as the SI units (Système International d'Unités²¹). The absolute SI base units are the meter, second, kilogram, kelvin, ampere, candela, and mole. The definitions of the SI units, the methods for their realization, or their physical embodiment are a matter of international agreement under the terms of the 1875 Treaty of the Meter. A convenient method (but often not a sufficient condition) for achieving absolute accuracy is to obtain traceability to one of the SI units via a calibration transfer standard issued by one of the national standards laboratories. The United States standards laboratory is the National Institute of Standards and Technology (NIST, formerly the National Bureau of Standards).

A relative measurement is one that is not required to be traceable to one of the SI units. Relative measurements are usually obtained as the ratio of two measurements. An example of a relative

measurement is the determination of the transmittance of an optical material wherein the ratio of the output radiant power to the input radiant power is measured; the measurement result is independent of SI units.

Absolute Accuracy and Traceability

Establishment of legal traceability to an SI unit requires that one obtain legally correct documentation, i.e., certification, of the device that serves as the calibration transfer standard and sometimes of the particular measurement process in which the device is to be used. Certification of legal traceability within each nation is obtained from the national standards laboratory of that nation. Often another nation's standards laboratory can be used to establish legal traceability, provided that there exists the legal framework for mutual recognition of the legality of each other's standards.

In order to establish accurate traceability to an SI unit, one needs to determine the total accumulated error arising from: (1) the realization of the base SI unit; (2) if applicable, the derivation of an associated measurement quantity; (3) if applicable, scaling to a higher or lower value; and (4,5...) transfer of the calibration from one device to another. The last entries must include the instability of the calibration transfer devices; the others may or may not involve a transfer device.

Legal traceability to SI units does not guarantee accurate traceability and vice versa. In order to obtain accurate traceability, it is not necessary to prove traceability to a national standards laboratory. Instead, the measurement must trace back to one of the SI units. However, it is usually convenient to establish accurate traceability via one of the national standards laboratories. The degree of convenience and accuracy will depend upon the accuracy of the measurement method and type of calibration transfer device available from the particular national standards laboratory.

Although relative measurements do not require traceability to one of the SI base units often to satisfy legal requirements traceability to a national standards laboratory may be necessary.

Types of Errors, Uncertainty Estimates, and Error Propagation

It is almost pointless to state a value for an absolute or relative measurement without an estimate of the uncertainty and the degree of confidence to be placed in the uncertainty estimate. Verification of the accuracy and the confidence limits is not only desirable but is often a legal requirement.

The accuracy and the uncertainty of a measurement are synonymous. The usual terminology is that a measurement is "accurate to within $\pm x$ " or "uncertain to within $\pm x$ ", where x is either a fraction (percent) of the measured value or an interval within which the true value is known to within some degree of confidence. The degree of confidence in the uncertainty estimate is the confidence interval or σ -level.

Errors are classified as type A errors, also known as random errors, and type B, or systematic errors. Type A errors are the variations due to the effects of uncontrolled variables. The magnitude of these effects is usually small and successive measurements form a random sequence. Type B errors are not detectable as variations since they do not change for successive measurements with a given apparatus and measurement method. Type B errors arise because of differences between the ideal behavior embodied in fundamental laws of physics and real behavior embodied in an experimental simulation of the ideal. A type B error could also be a function of the quantity being measured; for example, in a blackbody radiance standard using the freezing point of a metal and its defined temperature instead of the true absolute temperature.

Type A errors are estimated using standard statistical methods. If the distribution of the measurements is known (e.g., either Gaussian, which is often called a normal distribution, or Poisson), then one uses the formalism appropriate to the distribution. Unless enough data is obtained to establish that the distribution is not Gaussian, it is usual to assume a gaussian distribution. A brief discussion of Gaussian statistical concepts and terminology is given here to guide the reader in interpreting or determining the uncertainty in a radiometric or photometric calibration. A thorough discussion of these topics is available via the web from NIST.²²

The mean value m , the standard deviation σ , and the standard deviation of the mean σ_m , of a set of measurement x_i , are estimated for a small sample from a gaussian distribution of measurements as follows:

$$m = \sum_i \frac{x_i}{\eta} \quad (56)$$

$$\sigma^2 = \frac{\sum (x_i - m)^2}{\eta - 1} \quad (57)$$

$$\sigma_m = \frac{\sigma}{\sqrt{\eta}} \quad (58)$$

where $i = 1$ to η , and η is the total number of measurements.

The standard deviation is an estimate of the spread of the individual measurements within a sample, and it approaches a constant value as η is increased.

The standard deviation of the mean is an estimate of the spread of the values of the mean that would be obtained from several different sets of sample measurements. The standard deviation of the mean decreases as the number of samples in a set increases, since the estimate of the mean approaches the true mean for an infinite data set. The standard deviation of the mean is used in the estimate of the confidence interval assigned to the reported value of the mean.

The degree of confidence to which a reported value of the mean is valid is known as the *confidence interval* (CI). If it is assumed that a very large set of measurements has been sampled, then the CI is often given in terms of the number of standard deviations of the mean (one- σ level, two- σ level, etc.) within which the type A error of a reported value is known.

The CI is the probability that the mean from a normal distribution will be within the estimated uncertainty. That is, for a z -percent confidence interval, z -percent of the measurements will fall inside and $(100 - z)$ percent will fall outside of the uncertainty estimate. For small measurement samples from a gaussian distribution, Student's t -distribution is used to estimate the CI. Tables of Student's t -distribution along with discussions concerning its use are presented in most statistics textbooks. For large sets of measurements, a one- σ level corresponds approximately to a 68-percent CI, a two- σ level to a 95-percent CI, and a three- σ level to a 99.7-percent CI.

The reader is cautioned about using the σ level designation to describe the CI for a small sample of measurements. As an example of the small versus large sample difference, consider two data sets, one consisting of three samples and the other ten. Using a Student's t -distribution to estimate the CI for the three sample set, the one- σ , two- σ , and three- σ levels correspond to CIs of 61 percent, 86 percent, and 94 percent, respectively. For the ten-sample set, the respective CIs are 66 percent, 93 percent, and 99 percent. It can be quite misleading to state only the σ level of the uncertainty estimate without an indication of the size of the measurement set from which it was drawn. In order to avoid misleading accuracy statements, it is recommended that, instead of simply reporting the σ level, either the estimated CI be reported or the standard deviation of the mean be reported along with the number of measurement samples obtained.

Type B error estimates are either educated guesses of the magnitude of the difference between the real and the ideal or they are the result of an auxiliary measurement. If an appropriate auxiliary experiment can be devised to measure a systematic or type B error, then it need no longer be considered an error. The result obtained from the auxiliary measurement can usually be used as a correction factor. If a correction factor is applied, then the uncertainty is reduced to the uncertainty associated with the auxiliary experiment.

Most of the effort in high-accuracy radiometry and photometry is devoted to reducing type B errors. The first rule for reducing type B errors is to ensure that the experiment closely simulates the ideal. The second rule is that the differences between real and ideal should be investigated and that a correction be applied. Unlike type A errors for which an objective theory exists, the educated guess for a type B error is often subjective. For type B errors, neither a confidence interval nor a σ level is objectively quantifiable.

Error propagation, error accumulation, or a combined uncertainty analysis is the summation of all the type A and type B uncertainties that contribute to the final measurement in the chain. Because type A errors are truly random, they are uncorrelated and the accumulated type A error is obtained from the square-root of the sum of the squares (also known as root-sum-square, RSS) of the several type A error estimates. Type B uncertainties, however, may be either correlated or uncorrelated. If they are uncorrelated, the total uncertainty is the RSS of the several estimates. Type B uncertainties that are correlated must be arithmetically summed in a way that accounts for their correlation. Therefore, it is usually desirable to partition type B uncertainties so that they are uncorrelated.

Absolute Sources

Planckian or Blackbody Radiator A blackbody, or planckian, radiator is a thermal radiation source with a predictable absolute radiance output. An ideal blackbody is a uniform, i.e., lambertian, source of radiant power having a predictable distribution over area, solid angle, and wavelength. It is used as a standard radiance source from which the other radiometric quantities, e.g., irradiance, intensity, etc., can be derived.

Blackbody simulators are in widespread use not only at national standards laboratories but also in many other industrial, academic, and government laboratories. Blackbody simulators are commercially available from a number of manufacturers and cover a wide range of temperatures and levels of accuracy. Because they are in such widespread use as absolute standard sources for a variety of radiometric applications, particularly in the infrared, they are discussed here in some detail. Furthermore, since many practical sources of radiation can be approximated as a thermal radiation source, a blackbody function is often used in developing the radiometric model of an optical system.

An ideal blackbody is a completely enclosed volume containing a radiation field which is in thermal equilibrium with the isothermal walls of the enclosure that is at a known absolute temperature. The radiation in equilibrium with the walls does not depend upon the shape or constitution of the walls provided that the cavity dimensions are much larger than the wavelengths involved in the spectrum of the radiation.

Since the radiometric properties of a blackbody source are completely determined by its temperature, the SI base unit traceability for blackbody-based radiometry is to the kelvin. Because of recent improvements in the accuracy of absolute detector-based measurements thermodynamic temperatures are obtained by radiometric detector methods.²³

Since the radiation field and the walls are in equilibrium, the energy in the radiation field is determined by the temperature of the walls. The relationship between the absolute temperature T and the spectral radiance L_λ is given by Planck's law:

$$L_\lambda = \frac{2hc^2}{n^2\lambda^5} [e^{(hc/n\lambda kT)} - 1]^{-1} \quad (59)$$

Here h is Planck's constant, c is the speed of light in a vacuum, k is Boltzmann's constant, λ is the wavelength, and n is the index of refraction of the medium. Incorporating the values of the constants in this equation yields,

Spectral radiance units: $\text{W m}^{-2} \text{sr}^{-1} \mu\text{m}^{-1}$

$$L_\lambda = \frac{1.1910 \times 10^8}{n^2\lambda^5} [e^{(1.4388 \times 10^4/n\lambda T)} - 1]^{-1} \quad (60)$$

It follows that the peak of the spectrum of a blackbody is determined by its temperature (Wein displacement law).

$$n\lambda_{\max} T = 2898 \mu\text{mK} \quad (61)$$

It is often useful to measure blackbody spectral radiance in units of photons per second N_λ . The form of Planck's law in this case is

Spectral radiance units: photons $s^{-1} m^{-2} sr^{-1} \mu m^{-1}$

$$N_{\lambda} = \frac{2c}{n\lambda^4} [e^{(hc/n\lambda kT)} - 1]^{-1} \quad (62)$$

The peak of this curve is not at the same wavelength as in the case of radiance measured in units of power. Wein's displacement law for blackbody radiance measured in photons per second is

$$n\lambda_{\max} T = 3670 \mu m K \quad (63)$$

In other applications, the spectral distribution of the blackbody radiation may be required in units of photons per second per frequency interval (symbol: N_{ν}). This form of Planck's law is
Spectral radiance units: photons $s^{-1} m^{-2} sr^{-1} Hz^{-1}$

$$N_{\nu} = \frac{2n^2\nu^2}{c^2} [e^{(h\nu/kT)} - 1]^{-1} \quad (64)$$

and that of Wein's displacement law is

$$\frac{T}{\nu_{\max}} = 1.701 \times 10^{-11} K Hz^{-1} \quad (65)$$

Planck's law integrated over all wavelengths (or frequencies) leads to the Stefan-Boltzmann law which describes the temperature dependence of the total radiance of a blackbody. For blackbody radiance measured as radiant power, the Stefan-Boltzmann law is

Radiance units: $W m^{-2} sr^{-1}$

$$L = 1.8047 \times 10^{-8} n^2 T^4 \quad (66)$$

Equation (66) is the usual form of the Stefan-Boltzmann law; however, it can also be derived for blackbody radiance measured as photon flux.

Radiance units: photons $s^{-1} m^{-2} sr^{-1}$

$$N = 4.8390 \times 10^{14} n^2 T^3 \quad (67)$$

The preceding expressions are valid provided that the cavity dimensions are much larger than the wavelengths involved in the spectrum of the radiation. The restriction imposed by the cavity dimension may lead to significant errors in very high accuracy radiometry or very long wavelength radiometry. For example, in a cube 1 mm on a side and at a wavelength of 1 μm , the approximate correction to Planck's equation is only 3×10^{-7} ; however, if the measurement is made within a 1-nm bandwidth or less, the root mean square fluctuation of the signal is about 2×10^{-3} which may not be negligible. Recent work describes how well the Planck and Stefan-Boltzmann equations describe the radiation in small cavities and at long wavelengths.²⁴⁻²⁶

Blackbody Simulators An ideal blackbody, being completely enclosed, does not radiate into its surrounds and therefore cannot serve as an absolute radiometric source. A blackbody simulator is a device that does emit radiation but only approximates the conditions under which Planck's law is valid. In general, a blackbody simulator is an enclosure at some fixed temperature with a hole in it through which some of the radiation is emitted. Some low-accuracy blackbody simulators are fabricated as a flat surface held at a fixed temperature.

A blackbody simulator can be used as an absolute source provided that the type B errors introduced by the deviations from the ideal Planck's-law conditions are evaluated and the appropriate corrections are applied. In a blackbody simulator there are three sources of type B error: inaccurate surface temperature, nonequilibrium between the radiant surface and the radiation field due to openings in the enclosure, and nonuniformity in the temperature of the radiant surface.

Calculation of the effect of a temperature error on the spectral radiance is obtained from the derivative of Planck's law with respect to temperature.

$$\frac{dL_\lambda}{L_\lambda} = \frac{hc}{n\lambda kT} [1 - e^{-(hc/n\lambda kT)}]^{-1} \frac{dT}{T} \quad (68)$$

Since the radiation field is in equilibrium with the surface of the cavity, it is the absolute temperature of the surface that must be measured. It is usually impractical to have the thermometer located on the emitting surface and it is the temperature within the wall that is measured. The difference between the temperature within the wall and the surface must therefore be measured, or calculated from a thermal model, and the correction applied.

The error due to nonequilibrium occurs because a practical radiation source cannot be a completely closed cavity. The correction factor for the effect on the radiance due to the escaped radiation is obtained from application of Kirchhoff's law. Simply stated, Kirchhoff's law states that the absorptive power of a material is equal to its emissive power. According to the principle of detailed balancing, for a body to be in equilibrium in a radiation field, the absorption of radiation by a given element of the surface for a particular wavelength, state of polarization, and in a particular direction and solid angle must equal the emission of that same radiation. If this were not true, the body would either emit more than it absorbs or vice versa, it would not be in equilibrium with the radiation, and it would either heat up or cool off.

Radiation impinging upon a body is either reflected, transmitted, or absorbed. The fraction of the incident radiation that is reflected ρ (reflectance), plus the fraction absorbed α (absorptance), plus the fraction transmitted τ (transmittance), is equal to one.

$$1 = \rho + \alpha + \tau \quad (69)$$

From Kirchhoff's law for a surface in radiative equilibrium, the fraction of absorbed radiation equals the fraction emitted ε (emittance or emissivity). Therefore, the sum of the reflectance, transmittance, and emittance must also be equal to one. If the body is opaque, the transmittance is zero and the emittance is just equal to one minus the reflectance.

$$\varepsilon = 1 - \rho \quad (70)$$

For a body not in an enclosed volume to be in equilibrium with a radiation field, it must absorb all the radiation impinging upon it, because any radiation lost through reflection will upset the equilibrium. An emittance less than one is the measure of the departure from a perfect absorber and, therefore, it is a measure of the radiance change due to the departure from closed-cavity equilibrium. In general, cavities with an emittance nearly equal to unity are those for which the size of the hole is very small in comparison to the size of the cavity.

Temperature nonuniformity modifies the radiant flux over the whole cavity in much the same way as the presence of a hole in that it is a departure from equilibrium. Radiation loss from the region of the cavity near the hole is typically larger than from other regions and this loss produces a temperature change near the hole and a nonuniformity along the cavity wall. In addition, the temperature nonuniformity is another source of uncertainty in the absolute temperature. In practice, the limiting factor in the accuracy of a high-emittance blackbody simulator is typically the nonuniformity of the temperature.

Accurate calculation of the emittance of a cavity radiator requires a detailed knowledge of the geometry of the cavity and the viewing system. This is a radiance transfer calculation and, in order to perform it accurately, one must know the angular emitting or reflecting properties of the cavity surface. The regions that contribute most to the accuracy of the calculation are those that radiate directly out the hole into the direction of the solid angle of the optical detection system.

There are many methods of calculating the emittance. The most popular are based upon the assumption of uniform emission that is independent of direction, i.e., lambertian emission. One can calculate the spectral emittance and temperature of each element along the cavity wall and sum the contribution from each element to the cavity radiance. Extensive discussion of the diffuse emittance and temperature nonuniformity calculation methods can be found elsewhere.²⁷⁻³¹

Instead of calculating the emittance of a cavity directly, the problem may be transformed into one of calculating the absorptance for a ray incident from the direction in which the emittance is required.³²⁻³⁴ The quantity to be calculated in this case is that fraction of the radiation entering the hole from a particular direction which is subsequently reflected out of the hole into a hemisphere.

Real surfaces are not perfectly diffuse reflectors and often have a higher reflectance in the specular direction. A perfect specularly reflecting surface is at the other extreme for calculating the emittance of a blackbody simulator. In some applications, a specular black surface might perform better than a diffusely reflecting one, particularly if the viewing geometry is highly directional and well known. The calculation of the emittance of a cavity made from a perfectly specular reflector is obtained in terms of the number of reflections undergone by an incident ray before it leaves the cavity.³⁵

One can reduce the error due to temperature nonuniformity by reducing the emittance of those regions along the cavity wall that do not contribute radiation directly to that emitted from the cavity.³⁶ That is, by fabricating the “hidden” portions of the cavity wall from a specular, highly reflecting material and by proper orientation of these surfaces, the highly reflecting surfaces absorb almost none of the radiation but reflect it back to the highly absorbing surfaces. Since the highly reflective surfaces absorb and emit very little radiation, their temperature will have a minimal effect on the equilibrium within the cavity.

In high-accuracy applications, it is preferable to measure rather than calculate the emittance of the blackbody cavity. This can be done either by comparison of the radiance of the device under test to that of a higher-quality blackbody simulator (emittance closer to unity) or by a direct measurement of the reflectance of the cavity.³⁷ Accurate measurement of thermal nonuniformity by measurement of the variations in the radiance from different regions within the cavity is made difficult by the fact that radiance variations depend not only on the local temperature but also upon the emittance of the region.

Synchrotron Radiation A synchrotron is an electronic radiation source that if well-characterized has a predictable absolute radiance output. A synchrotron source is a very nonuniform, i.e., highly directional and highly polarized, radiance standard in contrast to a blackbody which is uniform and unpolarized. However, like a blackbody, a synchrotron has a predictable spectral output and it is useful as a standard radiance source from which the other radiometric quantities, e.g., irradiance, intensity, etc., can be derived.

Classical electrodynamic theory predicts that an accelerated charged particle will emit radiation. A synchrotron is a type of electron accelerator where the electron beam is accelerated in a closed loop and synchrotron radiation is the radiation emitted by the electrons undergoing acceleration. The development of these and other charged particle accelerators led to closer experimental and theoretical scrutiny of the radiation emitted by an accelerated charged particle. These studies culminated in Schwinger’s complete theoretical prediction, including relativistic effects, of the spectral and angular distribution of the radiation emitted by a beam in a particle accelerator.³⁸ The accuracy of Schwinger’s predictions have been verified in numerous experimental studies.³⁹⁻⁴¹

Schwinger’s theoretical model of the absolute amount of radiation emitted by an accelerated charged particle is analogous to the Planck equation for blackbody sources in that both predict the behavior of an idealized radiation source. Particle accelerators, when compared to even the most elaborate blackbodies, are, however, far more expensive. Furthermore, in order to accurately predict the spectral radiance of the beam in a particle accelerator, much detailed information is required of the type not found for most accelerators. Accurately predictable radiometric synchrotron sources are consequently found only in a few laboratories throughout the world.

The magnitude of the radiant power output from a synchrotron source is proportional to the number of electrons in the beam and their velocity, i.e., the number of electrons per second or current and their energy. Therefore, synchrotron radiometry is traceable to the SI unit of electricity, the ampere.

Because absolute synchrotron sources are so rare, a detailed discussion of Schwinger’s model of synchrotron radiation and the various sources of uncertainty will not be presented here. It is generally useful, however, to know some of the characteristics of synchrotron radiation. For example, the radiance from a synchrotron beam is highly polarized and very nonuniform: radiant power is almost entirely in the direction of the electron velocity vector and tangent to the electron beam. The peak of the synchrotron radiation spectrum varies from the vacuum ultraviolet to the soft x-ray region depending upon the energy in the beam. Higher-energy beams have a shorter wavelength peak: 1-GeV

peaks near 10 nm, 6-GeV peaks near 0.1 nm. Radiant power decreases to longer wavelengths by very roughly two decades for every decade increase of wavelength, so that for the typical radiometric-quality synchrotron source, there is usually sufficient energy to perform accurate radiometric measurements in the visible for intercomparison to other radiometric standards.⁴¹

Absolute Detectors

Electrical Substitution Radiometers An electrical substitution radiometer, often called an electrically calibrated detector, is a device for measuring absolute radiant power by comparison to electrical power.⁸ As a radiant power standard, an electrical substitution radiometer can be used as the basis for the derivation of the other radiometric quantities (irradiance, radiance, or intensity) by determining the geometrical distribution (either area and/or solid angle) of the radiation.^{42–48} Since an electrical substitution radiometer measures the spectrally total radiant power, it is used primarily for the measurement of monochromatic sources or those with a known relative spectral distribution.

An electrical substitution radiometer consists of a thermal detector (i.e., a thermometer) that has a radiation-absorbing surface and an electrical heater within the surface, or the heater is in good thermal contact with the surface. When the device is irradiated, the thermometer senses the temperature of the radiantly heated surface. The radiation source is then blocked and the power to the electrical heater adjusted to reproduce the temperature of the radiantly heated surface. The electrical power to the heater is measured and equated to the radiant power on the surface. The absolute base for this measurement is the electrical power measurement which is traceable to the SI ampere. In order for the measurement to be accurate, differences between the radiant and electrical heating modes must be evaluated and the appropriate corrections applied.

Electrical substitution radiometers predate the planckian radiator as an absolute radiometric standard.^{49–50} They were the devices used to quantify the radiant power output of the experimental blackbody simulators studied at the end of the nineteenth century. Electrical substitution radiometers are in widespread use today and are commercially available in a variety of forms that can be classified either as to the type of thermometer, the type of radiant power absorber, or the temperature at which the device operates.

Early electrical substitution radiometers operated at ambient temperature and used either a thermocouple, a thermopile, or a bolometer as the detector. Thermopile- and bolometer-based radiometers are presently used in a variety of applications. They have been refined over the years to produce devices of either greater accuracy, sensitivity, and/or faster response time. Thermopile-based, ambient temperature electrical substitution radiometers used for radiant power (and laser power, see later discussion) measurements at about the 1-mW level at several national standards laboratories have estimated uncertainties reported to be within ± 0.1 percent.^{51–53} Electrical substitution radiometers have also been used for very high-accuracy absolute radiant power measurements of the total solar irradiance both at the surface of the earth,^{51,53} and above its atmosphere.⁵² The type of high-accuracy radiometer used at various national standards laboratories is a custom-built device and is not commercially available in general. On the other hand, electrical substitution radiometers for solar and laser power measurements at a variety of accuracy levels are available commercially.

An ambient temperature electrical substitution radiometer based on a pyroelectric as the thermal detector was developed in the 1970s.^{54,55} A pyroelectric detector is a capacitor containing a dielectric with a temperature-sensitive spontaneous electrical polarization; a change in temperature results in a change in polarization. Small and rapid changes of polarization are readily detectable, making the pyroelectric a sensitive and fast thermal detector. It is most useful as a detector of a pulsed or chopped radiant power signal. During the period when the radiant power signal is blocked, electrical power can be introduced to a heater in the absorptive surface of the radiometer. As in the method for a thermopile- or bolometer-based electrical substitution radiometer, the electrical power is adjusted to equal the heating produced by the radiant power signal. Chopping can be done at a reasonable frequency, hence the electrical heating can be adjusted to achieve a balance in a comparatively short time. Because the radiant-to-electrical heating balance is more rapidly obtained, a pyroelectric radiometer is often more convenient to use than a thermopile radiometer. Pyroelectric

electrical substitution radiometers are generally more sensitive but are usually less accurate than the room temperature thermopile or bolometer electrical substitution radiometers.

Electrical substitution radiometers are further distinguished by two types of radiant power absorber configurations: a flat surface coated with a highly absorbing material or a cavity-shaped, light-trapping detector. Cavity-shaped radiometers are usually more accurate over a greater spectral range than flat-surface radiometers. However, a flat-surface receiver can usually be fabricated with less thermal mass than a cavity-shaped receiver and therefore may have greater sensitivity and/or a faster response time.

Electrical substitution radiometers are further distinguished by the temperature at which the electrical-to-radiant-power comparison is performed. In the last two decades there have been significant advancements^{56,57} made in instruments that perform the radiant-to-electrical comparison at a temperature near to that of liquid helium (4.2 K). Such devices are known as cryogenic electrical substitution radiometers or electrically calibrated cryogenic detectors and they are commercially available. Cryogenic electrical substitution radiometers are presently the most accurate absolute radiometric devices; the uncertainty of some measurements of radiant power has been estimated to be within ± 0.005 percent.⁴⁷

Sources of Error in Electrical Substitution Radiometers The relative significance of each of the possible sources of error and the derived correction factor depends upon the type of radiometer being used and the particular measurement application. It is possible to determine the total error occurring in equating radiant to electrical power and thence the accuracy of the traceability to the absolute electrical SI unit of measurement. Most manufacturers provide extensive characterization of their instruments. In such cases, the traceability to SI units is independent of radiometric standards such as blackbodies, hence electrical substitution radiometers are sometimes called absolute detectors. A commercially produced electrical substitution radiometer is capable of far greater accuracy (within ± 0.01 percent for the cryogenic instruments) than any of the typical radiometric transfer devices available from a national standards laboratory. Hence, establishing traceability through a radiometric standard from a national standards laboratory is almost pointless for a cryogenic electrical substitution radiometer.

The sources of error in an electrical substitution radiometer can be divided into three categories: errors in traceability to the absolute base unit, errors due to differences in the radiant-versus-electrical heating modes of operation, and errors arising in a particular application. The major error sources common to all electrical substitution radiometers as well as some of the less common are briefly described here. An extensive listing and description of all of these errors is given in Ref. 8, Chap. 1.

Electrical power measurement accuracy is first determined by the accuracy of the voltage and resistance standards (or voltmeter and resistance meter) used to measure voltage and current. Electrical power measurement accuracy within ± 0.01 percent is readily achievable and if needed it can be improved by an order of magnitude or better. Additional error is possible due to improper electrical measurement procedures such as those giving rise to ground-loops (improper connection to earth).

Differences between electrical-versus-radiant heating appear as differences in radiative, conductive, or convective losses. Most of these differences can be measured and a correction factor applied to optimize accuracy. The most obvious example is probably that of the radiative loss due to reflection from the receiver surface. Less obvious perhaps is the effect due to extraneous heating in the portion of the electrical conductors outside the region defined by the voltage sensing leads.

Differences between electrical heating and radiant heating may also arise due to spatial nonuniformity of the thermal sensor and/or differences in the heat conduction paths in the electrical-versus-radiant heating modes. These effects are specific to the materials and design of each radiometer. The electrical heater is typically buried within the device, whereas radiant heating occurs at the surface, so that the thermal conductivity paths to the sensor may be very different. Also, the distribution of the radiant power across the receiver is usually quite different compared to the distribution of the electrical heating. A detailed thermal analysis is required to create a design which minimizes these effects, but for optimum accuracy, the measurement of the magnitude of the nonuniformity effects is required to test the thermal model. Nonuniformity can be measured either by placing small auxiliary electrical heaters in various locations or by radiative heating of the receiver in several regions by moving a small spot of light across the device.

It should be noted that the thermal conduction path differences may also be dependent upon the environment in which the radiometer is to be operated. For example, atmospheric-pressure-dependent differences between the electrical-to-radiant power correction factor have been detected for many radiometers. These differences are, of course, greatest for a device for which the correction factors have been characterized in a normal atmosphere and which is then used in a vacuum.

Application-dependent errors arise from a variety of sources. Some examples are window transmission losses if a window is used, the accuracy of the aperture area and diffraction corrections are critical for measurements of irradiance; and, if a very intense source such as the sun is measured, heating of the instrument case and the body of the aperture could be an important correction factor. The last effect might also be very sensitive to atmospheric pressure changes.

Photoionization Devices Another type of absolute detector is a photoionization detector which can be used for absolute photon flux, i.e., radiant power, measurements of high-energy photon beams. Since a photoionization detector is a radiant power standard like the electrical substitution radiometer, it can in principle be used as the basis for the derivation of the other radiometric quantities (irradiance, radiance, or intensity) by determining the geometrical distribution (either area and/or solid angle) of the radiation.

A photoionization detector is a low-pressure gas-filled chamber through which a beam of high-energy (vacuum ultraviolet) photons is passed between electrically charged plates, the electrodes. The photons absorbed by the gas, if of sufficient energy, ionize the gas and enable a current to pass between the electrodes. The ion current is proportional to the number of photons absorbed times the photoionization yield of the gas and is, therefore, proportional to the photon flux.

The photoionization yield is the number of electrons produced per photon absorbed. If the photon is of sufficiently high energy, the photoionization yield is 100 percent for an atomic gas. The permanent atomic gases are the rare gases: helium, neon, argon, krypton, and xenon. Their photoionization yields have been measured relative to each other and shown to be 100 percent over specific wavelength ranges.^{58,59} If an ionization chamber is constructed properly and filled with the appropriate gas so that all of the radiation is absorbed, then the number of photons per second incident on the gas is simply equal to the ion current produced. If instead of measuring the ion current one were to measure each pulse produced by a photon absorption, then one would have a photon counter.

Carefully constructed ion current measurement devices have been used as absolute detectors from 25 to 102.2 nm and photon counters from 0.2 to 30 nm. Careful construction implies that all possible systematic error sources have either been eliminated or can be estimated, with an appropriate correction applied. Because of the difficulty in producing accurate and well-characterized devices, ion chambers and high-energy photon counters are not claimed to be high-accuracy radiometric devices. Furthermore, they are limited to applications in vacuum ultraviolet radiometry and are consequently of restricted interest.

Predictable Quantum Efficiency Devices A useful and quite economical type of absolute detector is a predictable quantum efficiency (PQE) device using high-quality silicon photodiodes. Quantum efficiency is the photon flux-to-photocurrent conversion efficiency. Because there have been many technological advancements made in the production of solid-state electronics, it is now possible to obtain very high quality silicon photodiodes whose performance is extremely close to that of the theoretical model.^{60,61} The technique for predicting the quantum efficiency of a silicon photodiode is also known as the self-calibration of a silicon photodiode.^{62,63} It is a relatively new absolute radiometric technique, quite simple to implement and of very high accuracy.^{64,65}

Conversion of a detector calibration from spectral responsivity $R(\lambda)$, in units of A/W, to quantum efficiency, i.e., photon-to-electron conversion efficiency, is as follows:

$$C_e = 1239.85 \frac{R(\lambda)}{\lambda} \quad (71)$$

where λ is the in-vacuum wavelength in nm and C_e is in units of electrons per photon.

As in the case of the other absolute detectors discussed previously, a PQE device is used for absolute photon flux, i.e., radiant power, measurements. It can also be used as the basis for the derivation

of the other radiometric quantities such as irradiance, radiance, or intensity. The extension to other radiometric measurements is by the determination of the geometrical distribution (area and/or solid angle) of the radiation. Also, like other absolute detectors, it measures spectrally total flux (within its spectral response range) and is therefore used primarily for the measurement of monochromatic sources or those with a known relative spectral distribution.

In a solid-state photodiode, the process for the conversion of a photon to an electronic charge is as follows. Photons not lost through reflection or by absorption in a coating at the front surface are absorbed in the semiconductor—if the photon is of high enough energy. To be absorbed, the photon energy must be greater than the band gap; the band gap for silicon is 1.11 eV (equivalent wavelength, 1.12 μm). In silicon, the absorption of a photon causes a promotion of a charge carried to the conduction band. Absorption of very high energy photons will create charge carriers with sufficient energy to promote a second, third, or possibly more charge carriers into the conduction band by collision processes. However, for silicon, the photon energy throughout the visible spectral range is insufficient for such impact ionization processes to occur. Therefore, in the visible to near-ir spectral region (about 400 to 950 nm), one absorbed photon produces one electron in the conduction band of silicon.

In a photodiode, impurity atoms diffused into a portion of the semiconductor material create an electric field. The internal electric field causes the newly created charge carriers to separate, eventually promoting the flow of an electron in an external measurement circuit. The efficiency with which the charge carriers are collected depends upon the region of the photodiode in which they are created. In the electric field region of a high-quality silicon photodiode, this collection efficiency has been demonstrated to approach 100 percent to within about 0.01 percent. Outside the field region, the collection efficiency can be determined by simple electrical bias measurements.

For the spectral regions in which the collection efficiency is 100 percent, the only loss in the photon-to-electron conversion process is due to reflection from the front surface of the detector. Several silicon photodiodes can be positioned to more effectively collect the radiation, acting as a light trap.^{66,67} If the radiation reflected from the first photodiode is directed to a second photodiode, then onto a third photodiode, etc., almost all the radiation will be collected in a small number of reflections. The photocurrents from all of the photodiodes are then summed and the total current (electrons per second) will be nearly equal, within 0.1 percent or less, to the photon flux (photons per second).

The more common type of silicon photodiode is the pn-type (positive charge impurity diffused into negative charge impurity starting material). High-quality pn-type detectors have their high collection efficiency in the long wavelength visible to near-ir spectral region. On the other hand, np-type silicon photodiodes have high collection efficiency in the short wavelength spectral region. At this time, the silicon photodiodes with the highest quantum efficiency (closest to ideal behavior) in the blue spectral region are the np-type devices, while nearly ideal red region performance is obtained with pn-type devices. The predictable quantum efficiency technique for silicon photodiodes has been demonstrated^{64–67} to be absolutely accurate to within ± 0.1 percent from about 400 nm to 900 nm.

A disadvantage of the light-trap geometry is the limited collection angle (field of view) of the device. Light-trap silicon photodiode devices are now commercially available using large area devices and a compact light-trap configuration that maximizes the field of view.

An np-type silicon photodiode trap detector optimized for short-wavelength performance and a pn-type silicon photodiode trap detector optimized for long-wavelength performance can be used as an almost ideal radiometric standard. The pair covers the 400- to 900-nm spectral range, has direct absolute SI base unit traceability via convenient electrical standards, and they are sufficiently independent to be meaningfully cross-checked to verify absolute accuracy and long-term stability. These detectors are not only useful radiometric standards by themselves but can be used with various source standards to either verify the absolute accuracy or to correct for the instabilities in the source standards.

The concept of a PQE light-trapping device is extendable to other high-quality photodiodes. Very recently, InGaAs devices with nearly 100-percent collection efficiency in the 1000- to 1600-nm spectral range have been developed. A light-trapping device employing these new detectors is now commercially available.

Calibration Transfer Devices

The discussion to this point focused on absolute radiometric measurements using methods that in themselves can be made traceable to absolute SI units. It is often more convenient (and sometimes required by contractual agreements) to obtain a device that has been calibrated in radiometric units at one of the national standards laboratories. Specific information as to the type and availability of various calibration transfer devices and calibration services may be obtained by directly contacting any of the national standards laboratories in the world. The products and services offered by the various standards laboratories cover a range of applications and accuracies, and differ from country to country.

Radiometric calibration transfer devices are either sources or detectors. The calibration transfer sources are either incandescent, tungsten filament lamps, deuterium lamps, or argon arc discharge sources.⁶⁸⁻⁷⁰ Generally, calibration transfer detectors are photodiodes of silicon, germanium, or indium gallium arsenide. The most prevalent calibration transfer sources are incandescent lamps and the typical calibration transfer detector is a silicon photodiode.⁷¹

The commonly available spectral radiance calibration transfer devices that span the 250- to 5000-nm region are typically tungsten strip filament lamps. Lamps calibrated in the 250- to 2500-nm region by a national standards laboratory are available. Lamps calibrated in the 2.5- to 5- μm region by comparison to a blackbody are commercially available. These devices are calibrated within specific geometrical constraints: the area on the filament, and the direction and solid angle of observation. The calibration is reported at discrete wavelengths, for a specified setting of the current through the filament and the ambient laboratory temperature. The optimum stability of spectral radiance is obtained with vacuum rather than gas-filled lamps, and with temperature controlled, i.e., water-cooled electrodes. Vacuum lamps cannot be operated at high filament temperatures and consequently do not have sufficient uv output. Gas-filled lamps cover a broader spectral and dynamic range and are the more commonly available calibration transfer device.

The commonly available spectral irradiance calibration transfer devices that span the 250- to 2500-nm region are tungsten coiled filament lamps. These are usually gas-filled lamps that have a halogen additive to prolong filament life and enable higher-temperature operation. Lamps calibrated in the 250- to 2500-nm region by a national standards laboratory are available. These devices are calibrated within the specific geometrical constraints of the distance and the direction with respect to a location on the lamp base or the filament. The calibration is reported at discrete wavelengths, for a specified setting of the current through the filament and the ambient laboratory temperature. Because the filament is operated at a higher temperature, the spectral irradiance lamps are usually less stable than the radiance lamps.

The drift of an incandescent lamp's radiance or irradiance output is not reliably predictable. It is for this reason that the calibration is most reliably maintained not by an individual lamp but by a group of lamps. The lamps are periodically intercompared and the average radiance (irradiance) of the group is considered to be the calibration value. The calculated differences between the group average and the individual lamps is used as a measure of the performance of the individual lamp. Lamps that have drifted too far from the mean are either recalibrated or replaced.

Spectral radiance and irradiance calibration transfer devices for the vacuum to near-uv (from about 160 to 400 nm) are typically available as deuterium lamps.

The commonly available calibrated transfer detectors for the 250- to 1100-nm spectral region are silicon photodiodes and for the 1000- to 1700-nm region, they are either germanium or indium gallium arsenide photodiodes. The calibration is reported at discrete wavelengths in absolute responsivity units (A/W) or irradiance response units ($A\text{ cm}^2/W$). In the first case, the calibration of the detector is performed with its active area underfilled, while in the second case, it is overfilled. If the detector is fitted with a precision aperture and if its spatial response is acceptably uniform, then the area of the aperture can be used to calculate the calibration in either units. The conditions under which the device was calibrated should be reported. The critical parameters are the location and size of the region within the active area in which it was calibrated, the radiant power in the calibration beam (alternately the photocurrent), and the temperature at which the calibration was performed. The direction in which the device was calibrated is usually assumed to be normal to its surface and

the irradiation geometry is usually that from a nearly collimated beam. Significant departures from normal incidence or near collimation should be noted.

Lasers

Power and Energy Measurement Lasers are highly coherent sources and the previous discussion of radiometry has been limited to the radiometry of incoherent sources. Nevertheless, the absolute power (or energy) in a laser beam can be determined to a very high degree of accuracy (within ± 0.01 percent in some cases) using some of the detector standards discussed here. The most accurate laser power measurements are made with cryogenic and room temperature electrical substitution radiometers and with predictable quantum efficiency devices. In order to measure the laser power (energy) it is necessary to ensure that all the radiation is impinging on the sensitive area of the detector and, if the absolute detector characterization was obtained at a different power (energy) level, that the detector is operating in a linear fashion. For pulsed lasers, the peak power may substantially exceed the dynamic range of the detector's linear performance. (A discussion of detector linearity is presented later in this chapter.) Furthermore, caution should be exercised to ensure that the detector not be damaged by the high photon flux levels achieved with many lasers.

In addition to ensuring that the detector intercept all of the laser beam, it is necessary to determine that all coherence effects have been eliminated (or minimized and corrected).^{72,73} The predominant effects of coherence are, first, interference effects at windows or beam splitters in the system optics and, second, diffraction effects at aperture edges. The use of wedged windows will minimize interference effects, and proper placement of apertures or the use of specially designed apertures⁷⁴ will minimize diffraction effects.

Lasers as a Radiometric Characterization Tool It should be noted that lasers, particularly the cw (continuous wave) variety, are particularly useful as characterization tools in a radiometric laboratory. Some of their applications are instrument response uniformity mapping, detector-to-detector spectral calibration transfer, polarization sensitivity, linearity verifications, and both diffuse and specular reflectance measurements.

Lasers are highly polarized and collimated sources of radiation. It is usually simple to construct an optical system as required for each measurement using mostly plane and spherical mirrors and to control scattered light with baffles and apertures. Lasers are high-power sources so that the signal-to-noise levels obtained are very good. If the power level is excessive it can usually be easily attenuated. Also, care must be taken to avoid local saturation of a detector at the peak of a laser's typical gaussian beam profile. They are highly monochromatic so that spectral purity, i.e., out-of-band radiation, is not usually a problem. However, in very high accuracy, within <0.1 percent, measurements, lasing from weaker lines may be significant and additional spectral blocking filters could be required.

Lasers are not particularly stable radiation sources. This problem is overcome by putting a beam splitter and stable detector into the optical system near the location of the measurement. The detector either serves to monitor the laser beam power and thereby supplies a correction factor to compensate for the instability, or its output is used to actively stabilize the laser.⁴³ In the latter case, an electronically controllable attenuator, such as an electro-optical, acousto-optical, or a liquid crystal system, is used to continuously adjust the power in the laser beam at the beam splitter. Feedback stabilization systems for cw lasers, both the electro-optical and liquid crystal type, are commercially available. For the highest-accuracy measurements, i.e., optimum signal-to-noise ratios, it is necessary both to actively stabilize the laser source and also to monitor the beam power close to the measurement in order to correct for the residual system drifts.

Various Type B Error Sources

Offset Subtraction One common error source, which is often simply an oversight, is the incorrect (or sometimes neglected) adjustment of an instrument reading for electronic and radiometric offsets.

This is often called the dark signal or dark current correction since it is obtained by shutting off the radiation source and reading the resulting signal. The shuttered condition needs to be close to radiant zero, at least within less than the expected accuracy of the measurement.

A dark signal reading is usually easy to achieve in the visible and near-visible spectral regions. However, in the long-wavelength infrared a zero radiance source is one that is at a temperature of ideally 0 K. Often an acceptably cold shutter is not easily obtained so that the radiance, i.e., temperature, of the “zero” reference source must be known in order to determine the true instrument offset.

Scattered Radiation and Size of Source Effect An error associated with the offset correction is that of scattered radiation from regions outside the intended optical path of the measurement system. Often, by judicious placement of the shutter, the principal light path can be blocked while the scattered light is not. In this case, the dark signal measurement includes the scattered light which is then subtracted from the measurement of the unshuttered signal. It is not possible to formulate a general scattered light elimination method so that each radiometric measurement system needs to be evaluated on an individual basis. The effects of scattered radiation can often be significantly reduced by using an optical chopper, properly placed, and lock-in amplifier system to read the output of the photodetector.

No optical element will produce a perfect image and there will be an error due to geometrically introduced stray light. Sharp edges between bright and dark regions will be blurred by aberrations, instrument fabrication errors, scattering due to roughness and contamination of the optical surfaces, and scattered light from baffles and stops within the instrument enclosure. Diffraction effects will also introduce stray light. Light originating from the source will be scattered out of the region of the image and light from the area surrounding the source will be scattered into the image. The error resulting from scattering at the objective lens or mirror is related to the size of the source since the scattering is proportional to the irradiance of the objective element. Thus, the error introduced by the lack of image quality is commonly referred to as the size-of-source effect.

The effect of the aberrations on the radiant power both into and out of an image can, in principle, be calculated. The diffraction-related error can also be calculated in some situations.⁷³⁻⁷⁵ However, the effect due to scattering is very difficult to model accurately and usually will have to be measured. In addition, the amount of scattering can be expected to change in time due to contamination of the optics, baffles, and stops. It is often more practical to measure the size of source effect and determine a correction factor for the elimination of this systematic error.

There are two different methods for measuring the size-of-source effect. The first method measures the response of the instrument as the size of the source is increased from the area imaged to the total area of the source. In the second method, a dark target of the same size as the image is placed at the imaged region on the source, and the surrounding area is illuminated. The second method has the advantage in that the effect being measured is the error signal above zero, whereas in the first method a small change in a large signal is being sought. In either case, the total error signal is measured; it includes aberrations, diffraction, and scattering effects.

Polarization Effects These are often significant perturbations of radiant power transfer due to properties of the radiation field other than its geometry. One such possible error is that due to the polarization state of the radiation field. The signal from a photodetector that is polarization-sensitive will be dependent upon the relative orientation of the polarization state of the radiation with respect to the detector orientation. Examples of polarization-dependent systems are grating monochromators and radiation transfer through a scattering medium or at a reflecting surface. In principle, the polarization state of the radiation field may be included in the geometrical transfer equation as a discrete transformation that occurs at each boundary or as a continuous transformation occurring as a function of position in the medium. Often it is sufficient to perform a calibration at two orthogonal rotational positions of the instrument or its polarization-sensitive components. However, it is recommended that other measurements at rotations intermediate between the two orthogonal measurements be included to test if the maximum and minimum polarization sensitivities have been sampled. The average of the maximum and minimum polarization measurements is then the calibration factor of the instrument for a nonpolarized radiation source.

Detector Nonlinearity

Nonlinearity measurement by superposition of sources Another possibly significant error source is photodetector and/or the electronic signal processing system nonlinearity. If the calibration and subsequent measurements are performed at the same radiant power level, then nonlinearity errors are avoided. Often conditions require that the measurements be performed over a range of power levels. In general, a separate measurement is required either to verify the linearity of the photodetector (and/or the electronics) or to deduce the form of the nonlinearity function in order to apply the appropriate correction.^{72,76–78}

The typical form of a nonlinearity appears as a saturation of the photoelectronic process at high irradiance levels. At low radiant flux levels what often appears to be a nonlinearity may be the result of failing to apply a dark signal or offset correction. There are, of course, other effects that will appear as a nonlinearity of the photodetector and/or electronics.⁷⁹

Either the linearity of the detector and electronics can be directly verified by experiment or it can be determined by comparison to a photodetector/electronics system of verified linear performance. It is useful to note that several types of silicon photodiodes using a transimpedance or current amplifier have been demonstrated to be linear within ± 0.1 percent over up to eight decades for most of its principal spectral range.⁷⁸

The fundamental experimental method for determining the dynamic range behavior of a photodetector is the superposition-of-sources method.^{76–78} The principle of the method is as follows. If a photodetector/electronics system is linear, then the arithmetic sum of the individual signals obtained from different radiant power sources should equal the signal obtained when all the sources irradiate the photodetector at the same time. There are many variations of the multiple source linearity measurement method using combinations of apertures or beam splitters. A note of caution: Interference effects must be avoided when combining beams split from the same source or when combining highly coherent sources such as lasers.

The difference between the arithmetic sum and the measured signal from the combined sources is used as the nonlinearity correction factor. Consider the superposition of two sources having approximately equal radiant powers ϕ_a and ϕ_b , which when combined have a radiant power of $\phi_{(a+b)}$. The signals from the photodetector when irradiated by the individual and combined sources is i_a , i_b , and $i_{(a+b)}$. The following equation would be equal to unity for a linear detector:

$$K_{ab} = \frac{i_{(a+b)}}{i_a + i_b} \quad (72)$$

For a calibration performed at the radiant power level ϕ_a (or ϕ_b), the detector responsivity is R and

$$i_a = R\phi_a \quad (73)$$

For a measurements at the higher radiant power level $\phi_{(a+b)}$,

$$i_{(a+b)} = K_{ab} R\phi_{(a+b)} \quad (74)$$

Scaling up to much higher radiant power levels (or down to lower levels) requires repeated application of the superposition-of-sources method. For example, in order to scale up to the next higher radiant power level, the source outputs from the first level are increased to match the second level (e.g., by using larger apertures). The increased source outputs are then combined to reach a third level and a new correction factor calculated. The process is repeated to cover the entire dynamic range of a photodetector/electronics system in factor-of-two steps.

Note that when type B errors, such as the interference effects noted above, are eliminated, the accumulated uncertainty in the source superposition method is the accumulated imprecision of the individual measurements.

Various nonlinearity measurement methods Other techniques for determining the dynamic range behavior of a photodetector are derivable from predictable attenuation techniques.⁷² One such method is based upon chopping the radiation signal using apertures of known area in a rotating disk. This is often referred to as Talbot's law: the average radiant power from a source viewed through the

apertures of a rotating disk is given by the product of the radiant power of the source and the transmittance of the disk. The transmittance of the disk is the ratio of the open area to the blocked area of the disk. The accuracy of this technique depends upon the accuracy with which the areas are known and may also be limited by the time dependence of the photodetector and/or electronics.

Another predictable attenuation technique is based upon the transmittance obtained when rotating, i.e., crossing two polarizers.

$$\tau = \tau_0 \cos^4 \delta \quad (75)$$

Here δ is the angle of rotation between the linear polarization directions of the two polarizers and τ_0 is the transmittance at $\delta = 0^\circ$. This technique, of course, assumes ideal polarizers that completely extinguish the transmitted beam at $\delta = 90^\circ$, and its accuracy is limited by polarization efficiency of the polarizers.

A third predictable attenuation technique is the application of Beer's law which states that the transmittance of a solution is proportional to the concentration κ of the solute

$$\tau = e^{-\gamma\kappa} \quad (76)$$

Here γ is the absorption coefficient of the solute. The accuracy of this technique depends upon the solubility of the solute and the absence of chemical interference, i.e., concentration-dependent chemical reactions.

Time-dependent Error For measurements of pulsed or repeatedly chopped sources of radiation, the temporal response of the detector could introduce a time-dependent error. A photodetector that has a response that is slow compared to the source's pulse width or the chopping frequency will not have reached its peak signal during the short time interval. Time-dependent error is avoided by determining if the detector's frequency response is suitable before undertaking the calibration and measurement of pulsed or chopped radiation sources.

Nonuniformity The nonuniformity of the distribution of radiation over an image or within the area sampled in an irradiance or radiance measurement may lead to an error if the response of the instrument is nonuniform over this area. The calibration factor for a nonuniform instrument will be different for differing distributions of radiation. The size of the error will depend upon the relative magnitudes of the source and instrument nonuniformities and it is a very difficult error to correct. This type of error is usually minimized either by measuring only sources that are uniform or by ensuring that the instrument response is uniform. It is usually easier to ensure that the instrument response is uniform.

Nonideal Aperture For very high-accuracy radiometric calibrations, the error due to the effect of the land on an aperture must be correctly taken into account. An ideal aperture is one that has an infinitesimally thin edge that intercepts the radiation beam. In practice an aperture will have a surface of finite thickness at its edge. This surface is referred to as the land; see Fig. 7. The effective

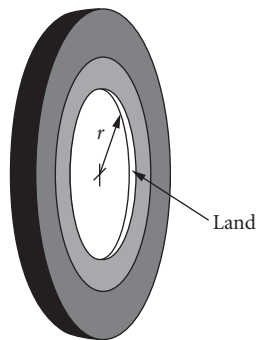


FIGURE 7 The nonideal aperture showing the location of the land.

radius of the aperture will be slightly reduced by vignetting caused by the land on its edge, assuming that the vignetting is small compared to the aperture radius and that all the radiation reflected by the land eventually falls on the detector (i.e., the land has a highly reflective surface). The effective radius of the aperture r' is

$$r' = r \left[1 - \frac{(1-\rho)w}{s} \right] \quad (77)$$

Here ρ is the reflectance, s is the distance between this aperture and the mirror or lens (or other aperture), and w is the width of the land.

Spectral Errors

Wavelength error In those cases where a spectrally selective element, such as a monochromator or a filter, is included in the radiometric instrument, spectral errors must be taken into account.⁴⁴ The first type of error is called the wavelength error and is due to misassignment of the wavelength of the spectrum of the filter or the monochromator in the instrument. That is, either the monochromator used to calibrate the filter transmission or the monochromator within the radiometric instrument has an error in its wavelength setting. This error is eliminated by calibration of the monochromator wavelength setting using one or more atomic emission lines from either a discharge lamp (usually a mercury and/or a rare gas lamp) or one of the many spectra of the elements available as a hollow cathode lamp. The wavelengths of most atomic emission lines are known with an accuracy that exceeds the requirements of radiometric calibrations.

A special note regarding the wavelength error and the use of interference filters in a radiometric instrument. The typical angular sensitivity of an interference filter is 0.1 nm per angular degree of rotation. If the transmission of the interference filter is measured in a collimated beam and then used in a convergent beam there will be an error due to the angularly dependent shift in the spectral shape of the transmission. For an accurate radiometric instrument it is important to measure interference filter transmission in nearly the same geometry as it will be used. Furthermore, the temperature coefficient of the transmission of an interference filter is about 0.2 nm K⁻¹. Therefore, it is also important in subsequent measurements to ensure that the filter remains at the same temperature as that used during the calibration.

It should also be noted that in order to accurately determine the spectral transmission of a monochromator it is necessary to completely fill the aperture of the dispersing element in the monochromator. This usually means that the field of view of the monochromator must be filled by the beam from the spectral calibration instrument.

Out-of-band radiation error The second type of spectral error is called the out-of-band or spectrally stray radiation error. This error is due to the radiation transmitted at both longer and shorter wavelengths that are beyond the edges of the principle transmission band of the filter or monochromator. This radiation is not taken into account if the limits in the integral in Eq. (11) are restricted to the edges of the principle transmission band. Although the relative amount of radiation transmitted at any wavelength beyond the edges of the principle transmission band may appear to be negligible with respect to the amount of radiation within the band, it is the spectrally total radiation that “leaks through” that is the error signal. It is therefore necessary to determine the transmission of the filter or monochromator over the entire spectrum of either the detector’s response and/or the light source’s output, whichever is greater. If the out-of-band radiation effect is small, it is possible to determine the correction factor from nominal values or limited accuracy measurements of the out-of-band spectra of the source, detector, and filter (monochromator).

Temperature Dependence The effect of temperature on the various elements in a radiometric instrument must not be overlooked. Unless the temperature of the instrument at the time of its calibration is maintained during subsequent applications, there may be substantial changes introduced in the instrument calibration factor that could well be above the uncertainty of its traceability to absolute SI units. The simple solution is to control the temperature of the system from the calibration

to the subsequent measurements. A more practical solution is often to measure the relative change in the instrument calibration factor as a function of temperature and then apply this as a correction factor to account for the temperature difference in the subsequent measurements.

34.6 PHOTOMETRY

Definition and Scope

Photometry is the measurement of radiation in a way that characterizes its effectiveness in stimulating the normal human visual system.^{4–5,80–82} Since visual sensation is a subjective experience, it is not directly quantifiable in absolute physical units. Attempts to quantify human visual sensation, therefore, were by comparison to various specified sources of light. The first sources used as standards were candles and later flames of prescribed construction. About the turn of the century, groups of incandescent lamps were selected as photometric standards and eventually a planckian radiator at a specific temperature was adopted by international agreement as the standard source. At present, the SI base unit for photometry, the candela, is no longer defined in terms of a given light source but is related to the radiant intensity by a multiplicative constant. Therefore, either an absolute source or detector can be used to establish an internationally recognized photometric calibration. Furthermore, there is no need for a human observer to effect a quantitative photometric measurement.

Photometry, as discussed here, is more precisely referred to as physical photometry to distinguish it from psychophysical photometry. Early photometric calibrations relied on human observers to compare an unknown light source to a standard. Presently, photometric calibrations are based on measurements using physical instruments. The instrument simulates human visual response either by having a detector with a spectral response that approximates that of the CIE standard observer or by using the CIE standard observer spectral response function in the data analysis.

Psychophysical photometry is the measurement of the effectiveness of light in individual observers and is more generally referred to as visual science. An individual's visual system may differ from that of the CIE standard observer defined for physical photometry, and these differences are sometimes important in experiments in visual science.

Photometry is restricted to the measurement of the magnitude of the visual sensation without regard to color, although it is well known that the perception of brightness is highly dependent on color in many circumstances. Measurement of the human response to color in terms of color matching is known as colorimetry. See Vol. III, Chap. 10, "Colorimetry."

Under reasonable light levels, the human eye can detect a difference of as little as 0.5 percent between two adjacent fields of illumination. For fields of illumination which are not adjacent, or are viewed at substantially different times, the eye can only detect differences of 10 to 20 percent. A discussion of the performance of the human visual system can be found in Vol. III, Chap. 2, "Visual Performance." Extensive treatments of photometry can be found in Walsh,⁸⁰ and Wyszecki and Stiles.⁸¹

Photopic, Scotopic, and Mesopic Vision

Electromagnetic radiation of sufficient power and in the wavelength range from about 360 to 830 nm, will stimulate the human visual system and elicit a response from an observer. The spectral range given here is the range over which measurements in physical photometry are defined. The range of reasonably perceptible radiation is usually given as about 400 to 700 nm. After light enters through the optical system of the eye—the cornea, iris or pupil, lens, and vitreous humor—the next stage of the visual response occurs in the retina. The retina contains two types of receptor cells: cones, which are the dominant sensors when the eye is adapted to higher radiance levels of irradiation (*photopic* vision), and rods, the dominant sensors at lower radiance levels (*scotopic* vision). Between the higher and lower levels of light adaptation is the region of *mesopic* vision, the range of radiance levels where both cones and rods contribute in varying degrees to the visual process.

Three types of cones having different spectral sensitivity functions exist in the normal human eye. The brain is able to distinguish colors by comparison of the signals from the three cone types. Of the three cone types, only the middle- and long-wavelength-sensitive cones contribute to the photopic sensation of the radiation entering the eye. The relative spectral sensitivity functions of the photopically and scotopically adapted human eye have been measured for a number of observers. From averages of these measurements, a set of values has been adopted by international agreement as the spectral efficiency for the CIE standard observer for photopic vision and another set for the CIE standard observer for scotopic vision (CIE, Commission Internationale de l'Eclairage). Because of the complexity of the spectral sensitivity of the eye at intermediate irradiation levels, there is no standard of spectral efficiency for mesopic vision. Values of the photopic and scotopic spectral efficiency functions are listed in Table 2.

TABLE 2 Photopic and Scotopic Spectral Luminous Efficiency Functions

Wavelength	Photopic	Scotopic	Wavelength	Photopic	Scotopic
375	0.00002	—	575	0.91540	0.1602
380	0.00004	0.00059	580	0.87000	0.1212
385	0.00006	0.00111	585	0.81630	0.0899
390	0.00012	0.00221	590	0.75700	0.0655
395	0.00022	0.00453	595	0.69490	0.0469
400	0.00040	0.00929	600	0.63100	0.03315
405	0.00064	0.01852	605	0.56680	0.02312
410	0.00121	0.03484	610	0.50300	0.01593
415	0.00218	0.0604	615	0.44120	0.01088
420	0.00400	0.0966	620	0.38100	0.00737
425	0.00730	0.1436	625	0.32100	0.00497
430	0.01160	0.1998	630	0.26500	0.00334
435	0.01684	0.2625	635	0.21700	0.00224
440	0.02300	0.3281	640	0.17500	0.00150
445	0.02980	0.3931	645	0.13820	0.00101
450	0.03800	0.455	650	0.10700	0.00068
455	0.04800	0.513	655	0.08160	0.00046
460	0.06000	0.567	660	0.06100	0.00031
465	0.07390	0.620	665	0.04458	0.00021
470	0.09098	0.676	670	0.03200	0.00015
475	0.11260	0.734	675	0.02320	0.00010
480	0.13902	0.793	680	0.01700	0.00007
485	0.16930	0.851	685	0.01192	0.00005
490	0.20802	0.904	690	0.00821	0.00004
495	0.25860	0.949	695	0.00572	0.00003
500	0.32300	0.982	700	0.00410	0.00002
505	0.40730	0.998	705	0.00293	0.00001
510	0.50300	0.997	710	0.00209	0.00001
515	0.60820	0.975	715	0.00148	0.00001
520	0.71000	0.935	720	0.00105	0.00000
525	0.79320	0.880	725	0.00074	0.00000
530	0.86200	0.811	730	0.00052	0.00000
535	0.91485	0.733	735	0.00036	0.00000
540	0.95400	0.650	740	0.00025	0.00000
545	0.98030	0.564	745	0.00017	0.00000
550	0.99495	0.481	750	0.00012	0.00000
555	1.00000	0.402	755	0.00008	0.00000
560	0.99500	0.3288	760	0.00006	0.00000
565	0.97860	0.2639	765	0.00004	0.00000
570	0.95200	0.2076	770	0.00003	0.00000

Photometric quantities can be calculated or measured as either photopic or scotopic quantities. Adaptation to luminance levels of $\geq 3 \text{ cd m}^{-2}$ (see further discussion) in the visual field usually leads to photopic vision, whereas adaptation to luminance levels of $\leq 3 \times 10^{-5} \text{ cd m}^{-2}$ usually leads to scotopic vision. Photopic vision is normally assumed in photometric measurements and photometric calculations unless explicitly stated to be otherwise.

Basic Concepts and Terminology

As noted in the earlier section on “Photometry,” the principles of photometry are the same as those for radiometry with the exception that the spectral responsivity of the detector is defined by general agreement to be specific approximations of the relative spectral response functions of the human eye. Photometric quantities are related to radiometric quantities via the spectral efficiency functions defined for the photopic and scotopic CIE standard observers.

Luminous Flux If physical photometry were to have been invented after the beginning of the twentieth century, then the physical basis of measurement might well have been the relationship between visual sensations and the energy of the photons and their flux density. It would follow naturally because vision is a photobiological process that is more closely related to the quantum nature of the radiation rather than its thermal heating effects. However, because of the weight of historical precedent, the basis of physical photometry is defined as the relationship between visual sensation and radiant power and its wavelength. The photometric equivalent of radiant power is luminous flux, and the unit that is equivalent to the watt is the lumen.

Luminous flux, Φ_v , is the quantity derived from spectral radiant power by evaluating the radiation according to its action upon the CIE standard observer.

$$\Phi_v = K_m \int \Phi_\lambda V(\lambda) d\lambda \quad (78)$$

where $V(\lambda)$ is the spectral efficiency function for photopic vision listed in Table 2, and K_m is the luminous efficacy for photopic vision. The spectral luminous efficacy is defined near the maximum, $\lambda_m = 555 \text{ nm}$, of the photopic efficiency function to be

$$K_m = 683 \frac{V(\lambda_m)}{V(555.016 \text{ nm})} \cong 683 \text{ lmW}^{-1} \quad (79)$$

The definitions are similar for scotopic vision

$$\Phi'_v = K'_m \int \Phi_\lambda V'(\lambda) d\lambda \quad (80)$$

where $V'(\lambda)$ is the spectral luminous efficiency function for scotopic vision listed in Table 2, and K'_m is the luminous efficacy for scotopic vision. The scotopic luminous efficiency function maximum occurs at $\lambda_m = 507 \text{ nm}$. The defining equation for K'_m is

$$K'_m = 683 \frac{V'(\lambda_m)}{V'(555.016 \text{ nm})} \cong 1700 \text{ lmW}^{-1} \quad (81)$$

The spectral shifts indicated in Eqs. (79) and (81) are required in order to obtain the precise values for the photopic and scotopic luminous efficacies. The magnitudes of the shifts follow from the specification of an integral value of frequency instead of wavelength in the definition of the SI base unit for photometry, the candela.

Luminous Intensity, Illuminance, and Luminance The candela, abbreviated cd, is defined by international agreement to be the luminous intensity in a given direction of a source that emits monochromatic radiation of frequency $540 \times 10^{12} \text{ Hz}$ (equal to 555.016 nm) and that has a radiant

intensity of $1/683 \text{ W sr}^{-1}$ in that direction. The spectral luminous efficacy of radiation at $540 \times 10^{12} \text{ Hz}$ equals 683 lm W^{-1} for all states of visual adaptation.

Because of the long history of using a unit of intensity as the basis for photometry, the candela was chosen as the SI base unit instead of the lumen, notwithstanding the fact that intensity is, strictly speaking, measurable only for point sources.

The functional form of the definitions of illuminance, luminous intensity, and luminance were presented in Eqs. (13), (14), and (15). The concepts are briefly reviewed here for the sake of convenience.

Luminous intensity is the photometric equivalent of radiant intensity, that is, luminous intensity is the luminous flux per solid angle. The symbol for luminous intensity is I_v . The unit for luminous intensity is the candela.

Illuminance is the photometric equivalent of irradiance, that is, illuminance is the luminous flux per unit area. The symbol for illuminance is E_v . The typical units for illuminance are lumens/meter².

Luminance is the photometric equivalent of radiance. Luminance is the luminous flux per unit area per unit solid angle. The symbol for luminance is L_v . The units for luminance are typically candelas/meter². In many older treatises on photometry, the term brightness is often taken to be equivalent to luminance, however, this is no longer the accepted usage.

In present usage, luminance and brightness have different meanings. In visual science (psychophysical photometry), two spectral distributions that have the same luminance typically do not have the same brightness. Operationally, spectral distributions of equal luminance are established with a psychophysical technique called heterochromatic flicker photometry. The observer views two spectral distributions that are rapidly alternated in time at the same spatial location, and the radiance of one is adjusted relative to the other to minimize the appearance of flicker. Spectral distributions of equal brightness are established with heterochromatic brightness matching, in which the two spectral distributions are viewed side-by-side and the radiance of one is adjusted relative to the other so that the fields appear equally bright. Though repeatable matches can easily be set with each technique, flicker photometric matches and brightness matches differ for many pairs of spectral distributions.

Photometric radiation transfer calculations and measurements are performed using the same methods and approximations that apply to the radiometric calculations discussed earlier. The exception, of course, is that the spectral sensitivity of the detector is specified.

Retinal Illuminance

In vision research it is frequently required to determine the effectiveness of a uniform, extended field of light (i.e., a large lambertian source that overfills the field of view of the eye) by estimating the illuminance on the retina. If it is assumed that the cornea, lens, and vitreous humor are lossless, then the luminous flux Φ_v in the image on the retina can be approximated from the conservation of the source luminance L_v as follows [see Eq. (22)],

$$L_v = \frac{L_e}{n_e^2} = \frac{\Phi_r S_{pr}^2}{n_e^2 A_r A_p} = E_r \frac{S_{pr}^2}{n_e^2 A_p} \quad (82)$$

where L_e is the radiance within the eye, n_e is the index of refraction of the ocular medium (the index of refraction of air is 1), Φ_r is the luminous flux at the retina, A_r is the area of the image of the retina, A_p is the area of the pupil, s_{pr} is the distance from the pupil to the retina, and E_r is the average illuminance in the image. Therefore, the average illuminance on the retina is

$$E_r = L_v \frac{n_e^2 A_p}{S_{pr}^2} \quad (83)$$

The luminance can, of course, be in units of either photopic or scotopic cd m^{-2} . The area of the pupil is measurable, but the distance between the pupil and retina is typically not available.

Therefore, a unit of retinal illuminance that avoids the necessity of determining this distance has been defined in terms of just the source luminance and pupil area. This unit is the troland, abbreviated td, and is defined as the retinal illumination for a pupil area of 1 mm^2 produced by a radiating surface having a luminance of 1 cd m^{-2} .

$$E_T = L_v A_p \quad (84)$$

Although it may be construed as an equivalent unit, one troland is *not equal* to one microcandela. The source is not a point but is infinite in extent. The troland is useful for relating several vision experiments where sources of differing luminance levels and pupil areas have been used.

The troland is, furthermore, not a measure of the actual illuminance level on the retina since the distance, index of refraction, and transmittance of the ocular medium are not included. For a schematic eye, which is designed to include many of the optical properties of the typical human eye, the effective distance between the pupil and the retina including the effect of the index of refraction is 16.7 mm^{83} (see also Vol. III, Chap. 10, "Colorimetry"). For the schematic eye with a 1-mm^2 pupil area, the effective solid angle at the retina is approximately 0.0036 sr . The retinal illuminance equivalent to one troland is therefore 0.0036 lm m^{-2} times the ocular transmittance.

Recall from the section on "Radiometric Effects of Stops and Vignetting" the effect of the aperture stop on the light entering an optical system; that is, all of the light entering the optical system appears to pass through the exit pupil, and the image of the aperture stop on the retina is the exit pupil. If the source is uniform and very large so that it overfills the field of view of the eye, the illuminance on the retina is independent of the distance between the source and the eye. If there is no intervening optic between the eye and the source, then the pupil is the aperture stop. However, if one uses an optical system to image the source into the eye, then the aperture stop need not be the pupil. An external optical system enables both the use of a more uniform, smaller source and the precise control of the retinal illumination by adjustment of an external aperture. The first configuration is called the newtonian view and the second is the maxwellian view of a source⁸⁰ (see also Vol. III, Chap. 5, "Optical Generation of the Visual Stimulus").

Though the troland is a very useful and a commonly used photometric unit among vision researchers, it should be interpreted with some caution in situations where one wishes to draw quantitative inferences about the effect of light falling on the retina. The troland is not a precise predictor because, besides not including transmission losses, no angular information is conveyed. The photoreceptors exhibit directional sensitivity where light entering through the center of the pupil is more effective than light entering through the pupil margin (the Stiles-Crawford effect, see Vol. III, Chap. 1, "Optics of the Eye"). Finally, specifying retinal illuminance in photometric units of trolands does not completely define the experimental conditions because the spectral distribution of the light on the retina is unspecified. Rather, the relative spectral responsivity of the eye (including the spectral dependence of the transmittance) is assumed by the inclusion of the $V(\lambda)$ or $V'(\lambda)$ functions. Experiments performed under mesopic conditions will be particularly prone to error.

If one measures the absolute spectral radiance of the light source, then Eq. (83) may be used in the radiometric form; that is, one substitutes $E_{r\lambda}$ and L_λ for E_r and L_v . The retinal spectral irradiance will then be in absolute units: $\text{W nm}^{-1}\text{m}^{-2}$.

Because the process of vision is a photobiological effect determined by the number and energy of the incident photons, the photon flux irradiance may be a more meaningful measure of the effect of the light on the retina. Using the radiometric form of Eq. (83) and the conversion to photon flux in Eq. (17), the photon flux irradiance $N_{E\lambda}$ on the retina is as follows.

$$N_{E\lambda} = 5.03 \times 10^{15} \lambda \tau_e(\lambda) L_\lambda \frac{n_e^2 A_p}{s_{pr}^2} = 1.80 \times 10^{13} \lambda \tau_e(\lambda) L_\lambda A_p \quad (85)$$

The ocular transmittance is included in this expression as $\tau_e(\lambda)$. The wavelength is in nm and, for radiance in units of $\text{W m}^{-2} \text{sr}^{-1} \text{nm}^{-1}$, the photon flux irradiance is in units of $\text{photons s}^{-1} \text{m}^{-2} \text{nm}^{-1}$.

If one uses a monochromatic light source, then a relationship between a monochromatic troland $E_T(\lambda)$ and the photon flux irradiance may be derived.

$$N_{E\lambda} = 1.80 \times 10^{13} \lambda A_p \tau_e(555) \frac{L_v(\lambda)}{K_m V(\lambda)} = 1.53 \times 10^{10} \lambda \frac{E_T(\lambda)}{V(\lambda)} \quad (86)$$

Only the transmittance at the peak of the $V(\lambda)$ curve, $\tau_e(555) = 0.58^{83}$ needs to be included since the spectral dependence of the ocular transmittance is already included in the $V(\lambda)$ function. The term $L_v(\lambda)$ is the luminance of a monochromatic light source.

An equivalent expression can be derived for the scotopic form of the monochromatic troland $E'_T(\lambda)$. Here, an ocular transmittance at 505 nm of 0.55^{83} has been used.

$$N_{E\lambda} = 5.82 \times 10^9 \lambda \frac{E'_T(\lambda)}{V'(\lambda)} \quad (87)$$

The reader is reminded that Eqs. (86) and (87) are valid only for a monochromatic source.

Absolute Photometric Calibrations

Photometric calibrations are in principle derived from the SI base unit for photometry, the candela. However, as one can see from the definitions of the candela and the other photometric quantities, photometric calibrations are in fact derived from absolute radiometric measurements using either a planckian radiator or an absolute detector. Typically, the relationship between illuminance and irradiance, Eq. (13), is used as the defining equation in deriving a photometric calibration.

The photometric calibration transfer devices available from national standards laboratories are usually incandescent lamps of various designs.⁸⁴ The photometric quantities commonly offered as calibrations are luminous intensity and total luminous flux.

The luminous intensity of a lamp, at a specified minimum distance and in a specified direction, is derived from a calibration of the spectral irradiance of the lamp, in the specified direction and at measured distance(s). The radiometric-to-photometric conversion [see Eq. (13)] is used to convert from spectral irradiance to illuminance. The inverse-square-law approximation, Eq. (24), is then used to derive luminous intensity.

Total luminous flux is a measure of all the flux emitted in every direction from a lamp. Total luminous flux is derived from illuminance (or luminous intensity) by measuring the flux emitted in all directions around the lamp. This procedure is known as goniophotometry. For an illuminance-based derivation, the total flux is the average of all the illuminance measurements times the surface area of the sphere described by the locus of the points at which the average illuminance was sampled. In the case of an intensity-based derivation, the total flux is the average of all the intensity measurements times 4π steradians. These are, in principle, the calculation methods for goniophotometry. In practice, the average illuminance (or intensity) is measured in a number of zones of fixed area (or solid angle) around the lamp. The product of the illuminance times the area of the zone (or the intensity times the solid angle of the zone) is the flux. The flux from each of the zones is then summed to obtain the total flux from the lamp.

A number of national standards laboratories provide luminance calibration transfer devices. These are typically in the form of a translucent glass plate that is placed at a specified distance and direction from a luminous intensity standard. One method of deriving the luminance calibration of the lamp/glass unit is to restrict the area of the glass plate with an aperture of known area. The intensity of the lamp/glass combination is then calibrated by comparison to an intensity standard lamp and the average luminance calculated by dividing the measured intensity by the area of the aperture.

Some national standards laboratories also offer calibrations of photometers,⁸⁵ also known as illuminance meters. A photometer is a photodetector that has been fitted with a filter to tailor its relative (peak normalized) spectral responsivity to match that of the CIE standard photopic observer.

Calibration of a photometer is usually obtained by reference to a luminous intensity standard positioned at a measured distance from the detector aperture. The inverse-square-law approximation is invoked to obtain the value of the illuminance at the measurement distance.

Other Photometric Terminology

Foot-candles, Foot-lamberts, Nits, etc The following units of illuminance are often used in photometry, particularly in older texts:

lux (abbreviation: lx) = lumen per square meter
 phot (abbreviation: ph) = lumen per square centimeter
 meter candle = lumen per square meter
 footcandle (abbreviation: fc) = lumen per square foot

One foot-candle = 0.0929 lux.

The following units of luminance are often used:

nit (abbreviation: nt) = candela per square meter
 stilb (abbreviation: sb) = candela per square centimeter

It is sometimes the practice, particularly in illuminating engineering, to express the luminance of an actual surface in any given direction in relation to the luminance of a uniform, diffuse, i.e., lambertian, source that emits one lumen per unit area into a solid angle of π steradians [see Eq. (30)]. This concept is one of relative luminance and its units (given following) are not equatable to the units of luminance. Furthermore, in spite of what may appear to be a similarity, this concept is not, strictly speaking, the photometric equivalent of the exitance of a source because, in its definition, the integral of the flux over the entire hemisphere is referenced. In other words, the equivalence to exitance is true only for a perfectly uniform radiance source. For all other sources it is the luminance in a particular direction divided by π . The units for luminance normalized to a lambertian source are:

1 apostilb (abbreviation: asb) = $(1/\pi)$ candela per square meter
 1 lambert (abbreviation: L) = $(1/\pi)$ candela per square centimeter
 1 foot-lambert (abbreviation: fL) = $(1/\pi)$ candela per square foot

Sometimes the total flux of a source is referred to as its *candlepower* or *spherical candlepower*. This term refers to a point source that uniformly emits in all directions, that is, into a solid angle of 4π steradians. Such a source does not exist, of course, so that the terminology is more precisely stated as the *mean spherical candle power*, which is the mean value of the intensity of the source averaged over the total solid angle subtended by a sphere surrounding the source.

Distribution Temperature, Color Temperature Distribution temperature is an approximate characterization of the spectral distribution of the visible radiation of a light source. Its use is restricted to sources having relative spectral outputs similar to that of a blackbody such as an incandescent lamp. The mathematical expression for evaluating distribution temperature is

$$\int_{\lambda_1}^{\lambda_2} \left[1 - \frac{\phi_x(\lambda)}{a\phi_b(\lambda, T)} \right]^2 d\lambda \Rightarrow \text{minimum} \quad (88)$$

where $\phi_x(\lambda)$ is the relative spectral radiant power distribution function of the test source, $\phi_b(\lambda, T)$ is the relative spectral radiant power distribution function of a blackbody at the temperature T , and a is an arbitrary constant. The limits of integration are the limits of visible radiation. Since distribution temperature is only an approximation, the exact values of the integration limits are arbitrary: typical limits are 400 and 750 nm. Values of a and T are adjusted simultaneously until the value of

the integral is minimized. The temperature of the best-fit blackbody function is the distribution temperature.

Color temperature and correlated color temperature are defined in terms of the perceived color of a source and are obtained by determining the chromaticity of the radiation rather than its relative spectral distribution. Because they are not related to physical photometry they are not defined in this section of the *Handbook*. These quantities do not provide information about the spectral distribution of the source except when the source has an output that closely approximates a blackbody. Although widely, and mistakenly, used in general applications to characterize the relative spectral radiant power distribution of light sources, color temperature and correlated color temperature relate only to the three types of cone cell receptors for photopic human vision and the approximate manner in which the human brain processes these three signals. As examples of incompatibility, there are the obvious differences between the spectral sensitivity of the human eye and physical receptors of visible optical radiation, e.g., photographic film, TV cameras. In addition, ambiguities occur in the ability of humans to distinguish the perceived color from different spectral distributions (metameric pairs). These ambiguities and the spectral sensitivity functions of the eye are not replicated by physical measurement systems. Caution must be exercised when using color temperature or correlated color temperature to predict the performance of a physical measurement system.

34.7 REFERENCES

1. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, Oxford, 1980.
2. E. Wolf, "Coherence and Radiometry," *J. Opt. Soc. Am.* **68**:6 (1978).
3. J. Geist, *McGraw-Hill Encyclopedia of Science and Technology*, McGraw-Hill, New York, 1987, p. 156.
4. W. Blevin et al., "Principles Governing Photometry," *Metrologia* **19**:97 (1983).
5. "The Basis of Physical Photometry," 2d ed., *Commission Internationale de L'Eclairage Publ. No. 18.2*, Central Bureau of the CIE, Vienna, 1983.
6. R. W. Boyd, *Radiometry and the Detection of Optical Radiation*, Wiley, New York, 1983.
7. F. Grum and R. J. Becherer, *Optical Radiation Measurements*, Academic Press, New York, 1979.
8. F. Hengstberger, ed. *Absolute Radiometry*, Academic Press, New York, 1989.
9. C. L. Wyatt, *Radiometric Calibration: Theory and Methods*, Academic Press, New York, 1978.
10. A. C. Parr, R. U. Datla, and J. L. Gardner, eds., *Optical Radiometry*, Elsevier, Amsterdam, 2005.
11. E. L. Dereniak and D. G. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984, pp. 1–14, and Apps. A and C.
12. M. V. Klein and T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986, pp. 203–222.
13. T. J. Quinn, *Temperature*, Academic Press, New York, 1983, pp. 284–363.
14. W. S. Smith, *Modern Optical Engineering*, 2d ed., McGraw-Hill, New York, 1990, pp. 135–136, 142–145, and 205–231.
15. M. E. Chahine, D. J. McCleese, P. W. Rosenkranz, and D. H. Staelin, in *Manual of Remote Sensing*, 2d ed., American Society of Photogrammetry, Falls Church, VA, 1983, pp. 172–179.
16. E. Hansen and L. D. Travis, "Light Scattering in Planetary Atmospheres," *Space Science Reviews* **16**:527 (1974).
17. Y. J. Kaufman, in Ghassem Asrar (ed.), *Theory and Applications of Optical Remote Sensing*, Wiley, New York, 1989, pp. 350–378.
18. H. Y. Wong, *Handbook of Essential Formulae and Data on Heat Transfer for Engineers*, Longman, London, 1977, pp. 89–128.
19. E. M. Sparrow and R. D. Cess, *Radiation Heat Transfer*, Brooks-Cole, Belmont, CA, 1966.
20. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**:125 (1967).
21. "Le Système International d'Unités," 3d ed., Bureau International des Poids et Mesures, Sèvres, France, 1977.

22. B. N. Taylor and C. E. Kuyatt, *Guidelines for Evaluating and Expressing the Uncertainty of NIST Measurement Results*, NIST Technical Note 1297, U.S. Government Printing Office, Washington, D.C., 1994.
23. N. P. Fox, J. E. Martin, and D. H. Nettleton, "Absolute Spectral Radiometric Determination of The Thermodynamic Temperatures of The Melting/Freezing Points of Gold, Silver and Aluminium," *Metrologia* **28**:357 (1991).
24. H. P. Baltes and E. R. Hilf, *Spectra of Finite Systems*, Bibliographisches Institut Mannheim, Vienna, Zurich, 1976.
25. H. P. Baltes, "Deviations from the Stefan-Boltzmann Law at Low Temperatures," *Appl. Phys.* **1**:39 (1973).
26. H. P. Baltes, "Planck's Law for Finite Cavities and Related Problems," *Infrared Phys.* **16**:1 (1976).
27. E. M. Sparrow, L. U. Ulbers, and E. R. Eckert, "Thermal Radiation Characteristics of Cylindrical Enclosures," *J. Heat Transfer* **C84**:188 (1962).
28. B. A. Peavy, "A Note on the Numerical Evaluation of Thermal Radiation Characteristics of Diffuse Cylindrical and Conical Cavities," *J. Res. Natl. Bur. Stds.* **70C**:139 (1966).
29. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities: Isothermal and Non-isothermal Cones and Cylinders," *J. Opt. Soc. Amer.* **64**:339 (1974).
30. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities, II: Isothermal and Non-isothermal Cylindrocones," *J. Opt. Soc. Amer.* **65**:565 (1975).
31. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities, III: Isothermal and Non-isothermal Double Cones," *J. Opt. Soc. Amer.* **66**:724 (1976).
32. J. C. de Vos, "Evaluation of the Quality of a Blackbody," *Physica* **20**:669 (1954).
33. T. J. Quinn, "The Calculation of the Emissivity of Cylindrical Cavities Giving near Blackbody Radiation," *Brit. J. Appl. Phys.* **18**:1105 (1967).
34. E. M. Sparrow and V. K. Johnson, "Absorption and Emission Characteristics of Diffuse Spherical Enclosures," *J. Heat Transfer* **C84**:188 (1962).
35. T. J. Quinn, "The Absorptivity of a Specularly Reflecting Cone for Oblique Angles of View," *Infrared Phys.* **21**:123 (1981).
36. T. J. Quinn and J. E. Martin, "Blackbody Source in the -50 to +200°C Range for the Calibration of Radiometers and Radiation Thermometers," *Appl. Opt.* **30**:4486 (1991).
37. E. F. Zalewski, J. Geist, and R. C. Willson, "Cavity Radiometer Reflectance," *Proc. of the SPIE* **196**:152 (1979).
38. J. Schwinger, "On the Classical Radiation of Accelerated Electrons," *Phys. Rev.* **75**:1912 (1949).
39. D. H. Tomboulian and P. L. Hartman, "Spectral and Angular Distribution of Ultraviolet Radiation from the 300-Mev Cornell Synchrotron," *Phys. Rev.* **102**:1423 (1956).
40. D. Lemke and D. Labs, "The Synchrotron Radiation of the 6-GeV DESY Machine as a Fundamental Radiometric Standard," *Appl. Opt.* **6**:1043 (1967).
41. N. P. Fox, P. J. Key, P. J. Riehle, and B. Wende, "Intercomparison between Two Independent Primary Radiometric Standards in the Visible and near Infrared: A Cryogenic Radiometer and the Electron Storage Ring BESSY," *Appl. Opt.* **25**:2409 (1986).
42. L. P. Boivin, "Calibration of Incandescent Lamps for Spectral Irradiance by Means of Absolute Radiometers," *Appl. Opt.* **19**:2771 (1980).
43. E. F. Zalewski and W. K. Gladden, "Absolute Spectral Irradiance Measurements Based on the Predicted Quantum Efficiency of a Silicon Photodiode," *Opt. Pura y Aplicada* **17**:133 (1984).
44. L. P. Boivin and A. A. Gaertner, "Realization of a Spectral Irradiance Scale in the Near Infrared at the National Research Council of Canada," *Appl. Opt.* **28**:6082 (1992).
45. T. J. Quinn and J. E. Martin, "A Radiometric Determination of the Stefan-Boltzmann Constant and Thermodynamic Temperatures between -40°C and +100°C," *Phil. Trans. Roy. Soc. London* **316**:85 (1985).
46. V. E. Anderson and N. P. Fox, "A New Detector-based Spectral Emission Scale," *Metrologia* **28**:135 (1991).
47. N. P. Fox, J. E. Martin, and D. H. Nettleton, "Absolute Spectral Radiometric Determination of the Melting/freezing Points of Gold, Silver and Aluminum," *Metrologia* **28**:357 (1991).
48. L. Jauniskis, P. Foukal, and H. Kochling, "Absolute Calibration of an Ultraviolet Spectrometer Using a Stabilized Laser and a Cryogenic Radiometer," *Appl. Opt.* **31**:5838 (1992).
49. F. Kurlbaum, "Über eine Methode zur Bestimmung der Strahlung in Absolutem Maass un die Strahlung des schwarzen Körpers zwischen 0 und 100 Grad," *Wied. Ann.* **65**:746 (1898).

50. K. Ångström, "The Absolute Determination of the Radiation of Heat with the Electrical Compensation Pyrheliometer, with Examples of the Application of this Instrument," *Astrophys. J.* **9**:332 (1899).
51. W. R. Blevin and W. J. Brown, "Development of a Scale of Optical Radiation," *Austr. J. Phys.* **20**:567 (1967).
52. R. C. Willson, "Active Cavity Radiometer Type V," *Appl. Opt.* **19**:3256 (1980).
53. L. P. Boivin and F. T. McNeely, "Electrically Calibrated Absolute Radiometer Suitable for Measurement Automation," *Appl. Opt.* **25**:554 (1986).
54. J. Geist and W. R. Blevin, "Chopper-Stabilized Radiometer Based on an Electrically Calibrated Pyroelectric Detector," *Appl. Opt.* **12**:2532 (1973).
55. R. J. Phelan and A. R. Cook, "Electrically Calibrated Pyroelectric Optical-radiation Detector," *Appl. Opt.* **12**:2494 (1973).
56. J. E. Martin, N. P. Fox, and P. J. Key, "A Cryogenic Radiometer for Absolute Radiometric Measurements," *Metrologia* **21**:147 (1985).
57. C. C. Hoyt and P. V. Foukal, "Cryogenic Radiometers and Their Application to Metrology," *Metrologia* **28**:163 (1991).
58. J. A. R. Samson, "Absolute Intensity Measurements in the Vacuum Ultraviolet," *J. Opt. Soc. Amer.* **54**:6 (1964).
59. F. M. Matsunaga, R. S. Jackson, and K. Watanabe, "Photoionization Yield and Absorption Coefficient of Xenon in the Region of 860–1022Å," *J. Quant. Spectrosc. Radiat. Transfer* **5**:329 (1965).
60. J. Geist, "Quantum Efficiency of the p–n Junction in Silicon as an Absolute Radiometric Standard," *Appl. Opt.* **18**:760 (1979).
61. J. Geist, W. K. Gladden, and E. F. Zalewski, "The Physics of Photon Flux Measurements with Silicon Photodiodes," *J. Opt. Soc. Amer.* **72**:1068 (1982).
62. E. F. Zalewski and J. Geist, "Silicon Photodiode Absolute Spectral Response Self-calibration," *Appl. Opt.* **19**:1214 (1980).
63. J. Geist, E. F. Zalewski, and A. R. Schaefer, "Spectral Response Self-calibration and Interpolation of Silicon Photodiodes," *Appl. Opt.* **19**:3795 (1980).
64. J. L. Gardner and W. J. Brown, "Silicon Radiometry Compared to the Australian Radiometric Scale," *Appl. Opt.* **26**:2341 (1987).
65. E. F. Zalewski and C. C. Hoyt, "Comparison Between Cryogenic Radiometry and the Predicted Quantum Efficiency of Silicon Photodiode Light Traps," *Metrologia* **28**:203 (1991).
66. E. F. Zalewski and C. R. Duda, "Silicon Photodiode Device with 100 Percent External Quantum Efficiency," *Appl. Opt.* **22**:2867 (1983).
67. N. P. Fox, "Trap Detectors and Their Properties," *Metrologia* **28**:197 (1991).
68. J. H. Walker, R. D. Saunders, J. K. Jackson, and D. A. McSparron, *Spectral Irradiance Calibrations*, National Bureau of Standards Special Publication No. 250-20, U. S. Government Printing Office, Washington, D.C., 1987.
69. J. H. Walker, R. D. Saunders, and A. T. Hattenburg, *Spectral Radiance Calibrations*, National Bureau of Standards Special Publication No. 250-1, U.S. Government Printing Office, Washington, D.C., 1987.
70. J. Z. Klose, J. M. Bridges, and W. R. Ott, *Radiometric Standards in the Vacuum Ultraviolet*, National Bureau of Standards Special Publication No. 250-3, U.S. Government Printing Office, Washington, D.C., 1987.
71. E. F. Zalewski, *The NBS Photodetector Spectral Response Calibration Transfer Program*, National Bureau of Standards Special Publication No. 250-17, U.S. Government Printing Office, Washington, D.C., 1987.
72. W. Budde, *Physical Detectors of Optical Radiation*, Academic Press, New York, 1983.
73. L. P. Boivin, "Some Aspects of Radiometric Measurements Involving Gaussian Laser Beams," *Metrologia* **17**:19 (1981).
74. L. P. Boivin, "Reduction of Diffraction Errors in Radiometry by Means of Toothed Apertures," *Appl. Opt.* **17**:3323 (1978).
75. W. R. Blevin, "Diffraction Losses in Photometry and Radiometry," *Metrologia* **6**:31 (1970).
76. C. L. Sanders, "A Photocell Linearity Tester," *Appl. Opt.* **1**:207 (1962).
77. C. L. Sanders, "Accurate Measurements of and Corrections for Non-linearities in Radiometers," *J. Res. Natl. Bur. Stand.* **A76**:437 (1972).
78. W. Budde, "Multidecade Linearity Measurements on Silicon Photodiodes," *Appl. Opt.* **18**:1555 (1979).

-
79. A. R. Schaefer, E. F. Zalewski, and J. Geist, "Silicon Detector Non-linearity and Related Effects," *Appl. Opt.* **22**:1232 (1983).
 80. J. W. T. Walsh, *Photometry*, Dover, New York, 1965.
 81. G. Wyszecki and W. S. Stiles, *Colour Science: Concepts and Methods*, Wiley, New York, 1967.
 82. "Light as a True Visual Quantity: Principles of Measurement," *Commission Internationale de L'Eclairage Publ. No. 41*, Central Bureau of the CIE, Vienna, 1978.
 83. E. N. Pugh, "Vision: Physics and Retinal Physiology," in R. C. Atkinson, R. J. Herrnstein, G. Lindsey, and R. D. Luce, (eds.), *Steven's Handbook of Experimental Psychology*, 2d ed., Wiley, New York, 1988, pp. 75–163.
 84. R. L. Booker and D. A. McSparron, *Photometric Calibrations*, National Bureau of Standards Special Publication No. 250–15, U.S. Government Printing Office, Washington, D.C., 1987.
 85. "Methods of Characterizing the Performance of Radiometers and Photometers," *Commission Internationale de L'Eclairage Publ. No. 53*, Central Bureau of the CIE, Vienna, 1982.

This page intentionally left blank.

MEASUREMENT OF TRANSMISSION, ABSORPTION, EMISSION, AND REFLECTION

James M. Palmer*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

35.1 GLOSSARY

A	Area
a	absorption
bb	blackbody
c_2	second radiation constant
E	irradiance
f	bidirectional scattering distribution function
i	internal
L	radiance
P	electrical power
R	reflectance factor
r	reflection
T	temperature
t	transmission
α	absorptance
α'	absorption coefficient
ε	emittance (emissivity)
θ, ϕ	angles
ρ	reflectance
σ	Stefan Boltzmann constant
τ	transmittance
Φ	power (flux)
Ω	projected solid angle

*Deceased.

35.2 INTRODUCTION AND TERMINOLOGY

When radiant flux is incident upon a surface or medium, three processes occur: transmission, absorption, and reflection. Figure 1 shows the ideal case, where the transmitted and reflected components are either specular or perfectly diffuse. Figure 2 shows the transmission and reflection for actual surfaces.

The symbols, units, and nomenclature employed in this chapter follow the established usage as defined in *ISO Standards Handbook 2*,¹ Cohen and Giacomo,² and Taylor.³ Additional general terminology applicable to this chapter is from ASTM,⁴ IES,⁵ IES,⁶ Drazil,⁷ and CIE.⁸ The prefix *spectral* is used

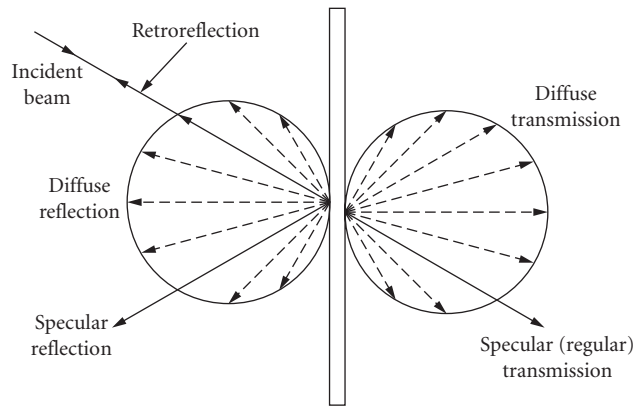


FIGURE 1 Idealized reflection and transmission.

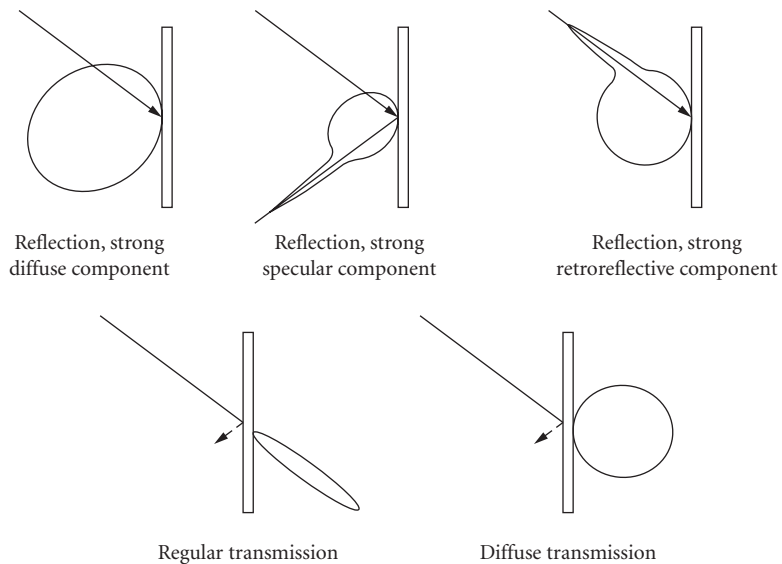


FIGURE 2 Actual reflection and transmission.

to denote a characteristic at a particular wavelength and is indicated by the symbol λ . The absence of the *spectral* prefix implies integration over all wavelengths with a source function included (omission of the source function is meaningless except where the characteristic is constant with wavelength).

There has been a continuing dialog over terminology, particularly between the suffixes *-ance* and *-ivity*.⁹⁻¹³ The suggested usage reserves terms ending with *-ivity* (such as transmissivity, absorptivity, and reflectivity) for properties of a pure material and employs the suffix *-ance* for the characteristics of a specimen or sample. For example, one can distinguish between the *reflectivity* of pure aluminum and the *reflectance* of a particular sample of 6061-T6 aluminum with a natural oxide layer. This distinction can be extended to differentiate between emissivity (of a pure substance) and emittance (of a sample). This usage of *emittance* should not be confused with the older term *radiant emittance*, now properly called *radiant exitance*. In this chapter, the suffix *-ance* will be used exclusively inasmuch as the measurement of radiometric properties of materials is under discussion.

35.3 TRANSMITTANCE

Transmission is the term used to describe the process by which incident radiant flux leaves a surface or medium from a side other than the incident side, usually the opposite side. The spectral transmittance $\tau(\lambda)$ of a medium is the ratio of the transmitted spectral flux $\Phi_{\lambda t}$ to the incident spectral flux $\Phi_{\lambda i}$, or

$$\tau(\lambda) = \frac{\Phi_{\lambda t}}{\Phi_{\lambda i}} \quad (1)$$

The transmittance τ is the ratio of the transmitted flux Φ_t , to the incident flux Φ_i , or

$$\tau = \frac{\int_0^\infty \tau(\lambda) \Phi_{\lambda i} d\lambda}{\int_0^\infty \Phi_{\lambda i} d\lambda} \neq \int_\lambda \tau(\lambda) d\lambda \quad (2)$$

Note that the integrated transmittance is *not* the integral over wavelength of the spectral transmittance, but must be weighted by a source function Φ_λ as shown.

The transmittance may also be described in terms of radiance as follows:

$$\tau = \frac{\int_0^\infty \int_\Omega L_{\lambda i} d\Omega_i d\lambda}{\int_0^\infty \int_\Omega L_{\lambda t} d\Omega_t d\lambda} \quad (3)$$

where $L_{\lambda i}$ represents the spectral radiance $L_{\lambda i}(\lambda; \theta_p, \phi_p)$ incident from direction (θ_p, ϕ_p) , $L_{\lambda t}$ represents the spectral radiance $L_{\lambda t}(\lambda; \theta_t, \phi_t)$ transmitted in direction (θ_t, ϕ_t) , and $d\Omega$ is the elemental projected solid angle $\sin \theta \cos \theta d\theta d\phi$.

The bidirectional transmittance distribution function (BTDF, symbol f_t) relates the transmitted radiance to the radiant incidence as

$$f_t(\lambda; \theta_i, \phi_i) \equiv \frac{dL_{\lambda t}}{dL_{\lambda i} d\Omega_i} = \frac{dL_{\lambda t}}{dE_{\lambda i}} (\text{sr}^{-1}) \quad (4)$$

Geometrically, transmittance can be classified as specular, diffuse, or total, depending upon whether the specular (regular) direction, all directions other than the specular, or all directions are considered.

35.4 ABSORPTANCE

Absorption is the process by which incident radiant flux is converted to another form of energy, usually heat. Absorptance is the fraction of incident flux that is absorbed. The absorptance α of an element is defined by $\alpha = \Phi_a / \Phi_i$. Similarly, the spectral absorptance $\alpha(\lambda)$ is the ratio of spectral power absorbed $\Phi_{\lambda a}$ to the incident spectral power $\Phi_{\lambda i}$,

$$\alpha = \frac{\int_0^\infty \alpha(\lambda) \Phi_{\lambda i} d\lambda}{\int_0^\infty \Phi_{\lambda i} d\lambda} \neq \int_\lambda \alpha(\lambda) d\lambda \tag{5}$$

An absorption coefficient α' (cm^{-1} or km^{-1}) is often used in the expression $\tau_i = e^{-\alpha' t}$, where τ_i is internal transmittance and t is pathlength (cm or km).

35.5 REFLECTANCE

Reflection is the process where a fraction of the radiant flux incident on a surface is returned into the same hemisphere whose base is the surface and which contains the incident radiation. The reflection can be specular (in the mirror direction), diffuse (scattered into the entire hemisphere), or a combination of both. Table 1¹⁴ shows a wide range of materials that have different goniometric (directional) reflectance characteristics.

TABLE 1 Goniometric Classification of Materials¹⁴

Material Classification	Scatter*	σ^\dagger	γ^\ddagger	Structure [§]	Example	
Exclusively reflecting materials	None	0	$\cong 0$	None	Mirror	
	Weak	≤ 0.4	$\leq 27^\circ$	Micro	Matte aluminum	
	$\tau = 0$	Strong	> 0.4	$> 27^\circ$	Macro	Retroreflectors
None					Laquer & enamel coatings	
Micro					Paint films, BaSO ₄ , Halon	
Macro					Rough tapestries, road surfaces	
Weakly transmitting, strongly reflecting Materials	None	0	$\cong 0$	None	Sunglasses, color filters cold mirrors	
	Weak	≤ 0.4	$\leq 27^\circ$	Micro	Matte-surface color filters	
	$\tau \leq 0.35$	Strong	> 0.4	$> 27^\circ$	Macro	Glossy textiles
					None	Highly turbid glass
Micro					Paper	
Strongly transmitting materials	None	0	$\cong 0$	Macro	Textiles	
				None	Window glass	
				None	Plastic film	
	Weak	≤ 0.4	$\leq 27^\circ$	Micro	Ground glass	
				Macro	Ornamental glass	
				Macro	prismatic glass	
$\tau > 0.35$	Strong	> 0.4	$> 27^\circ$	None	Opal glass	
				Micro	Ground opal glass	
				Macro	Translucent acrylic plastic with patterned surface	

*It is suggested that the diffusion factor is appropriate for strongly diffusing materials and that the half-angle is better for weakly diffusing materials.

† γ is a half-value angle, the angle from the normal where the radiance has dropped to one-half the value at normal.

‡ σ is a diffusion factor, the ratio of the mean of radiance measured at 20° and 70° to the radiance measured at 5° from the normal, when the incoming radiation is normal. $\sigma = [L(20) + L(70)]/[2L(5)]$. It gives an indication of the spatial distribution of the radiance, and is unity for a perfect (Lambertian) diffuser.

§Structure refers to the nature of the surface. In a microscattering structure, the scatterers cannot be resolved with the unaided eye. The macrostructure scatterers can be readily seen.

The most general definition for reflectance ρ is the ratio of the radiant flux reflected Φ_r to the incident radiant flux Φ_i , or

$$\rho = \frac{\Phi_r}{\Phi_i} \quad (6)$$

Spectral reflectance is similarly defined at a specified wavelength λ as

$$\rho(\lambda) = \frac{\Phi_{\lambda r}}{\Phi_{\lambda i}} \quad (7)$$

(Spectral) reflectance factor (symbol R) is the ratio of (spectral) flux reflected from a sample to the (spectral) flux which would be reflected by a perfect diffuse (Lambertian) reflector.

No single descriptor of reflectance will suffice for the wide range of possible geometries. The fundamental geometric descriptor of reflectance is the bidirectional reflectance distribution function (BRDF, symbol f_r). It is defined as the differential element of reflected radiance dL_r in a specified direction per unit differential element of radiant incidence dE_i , also in a specified direction,¹⁵ and carries unit of sr^{-1} :

$$f_r(\theta_i, \phi_i, \theta_r, \phi_r) = \frac{dL_r(\theta_i, \phi_i; \theta_r, \phi_r; E_i)}{dE_i(\theta_i, \phi_i)} \quad [\text{sr}^{-1}] \quad (8)$$

The polar angle θ is measured from the surface normal and the azimuth angle ϕ is measured from an arbitrary reference in the surface plane, most often the plane containing the incident beam. The subscripts i and r refer to the incident and reflected beams, respectively.

By integrating over varying solid angles, Nicodemus et al.,¹⁵ based upon earlier work by Judd,¹⁶ defined nine goniometric reflectances, and by extension, nine goniometric reflectance factors. These are shown in Tables 2 and 3 and Fig. 3. In these tables, the term *directional* refers to a differential solid angle $d\omega$ in the direction specified by (θ, ϕ) . *Conical* refers to a cone of finite extent centered in direction (θ, ϕ) ; the solid angle ω of the cone must also be specified.

Details on these definitions and further discussion can be found in ASTM STP475,⁴ ASTM E808,¹⁷ Judd,¹⁶ Nicodemus,¹⁸ Nicodemus,¹⁹ and Nicodemus et al.¹⁵

TABLE 2 Nomenclature for Nine Types of Reflectance¹⁵

1. Bidirectional reflectance	$d\rho(\theta_i, \phi_i; \theta_r, \phi_r) = f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
2. Directional-conical reflectance	$\rho(\theta_i, \phi_i; \omega_r) = \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
3. Directional-hemispherical reflectance	$\rho(\theta_i, \phi_i; 2\pi) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
4. Conical-directional reflectance	$d\rho(\omega_i; \theta_r, \phi_r) = (d\Omega_r / \Omega_i) \int_{\omega_i} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
5. Biconical reflectance	$\rho(\omega_i; \omega_r) = (1/\Omega_i) \int_{\omega_i} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
6. Conical-hemispherical reflectance	$\rho(\omega_i; 2\pi) = (1/\Omega_i) \int_{\omega_i} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
7. Hemispherical-directional reflectance	$d\rho(2\pi; \theta_r, \phi_r) = (d\Omega_r / \pi) \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
8. Hemispherical-conical reflectance	$\rho(2\pi; \omega_r) = (1/\pi) \int_{2\pi} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
9. Bihemispherical reflectance	$\rho(2\pi; 2\pi) = (1/\pi) \int_{2\pi} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$

TABLE 3 Nomenclature for Nine Types of Reflectance Factor¹⁵

1. Bidirectional reflectance factor	$R(\theta_i, \phi_i; \theta_r, \phi_r) = \pi f_r(\theta_i, \phi_i; \theta_r, \phi_r)$
2. Directional-conical reflectance factor	$R(\theta_i, \phi_i; \omega_r) = (\pi/\Omega_r) \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
3. Directional-hemispherical reflectance factor	$R(\theta_i, \phi_i; 2\pi) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
4. Conical-directional reflectance factor	$R(\omega_i; \theta_r, \phi_r) = (\pi/\Omega_i) \int_{\omega_i} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
5. Biconical reflectance factor	$R(\omega_i; \omega_r) = [\pi/(\Omega_i \Omega_r)] \int_{\omega_i} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
6. Conical-hemispherical reflectance factor	$R(\omega_i; 2\pi) = (1/\Omega_i) \int_{\omega_i} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
7. Hemispherical-directional reflectance factor	$R(2\pi; \theta_r, \phi_r) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
8. Hemispherical-conical reflectance factor	$R(2\pi; \omega_r) = (1/\Omega_r) \int_{2\pi} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
9. Bihemispherical reflectance factor	$R(2\pi; 2\pi) = (1/\pi) \int_{2\pi} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$

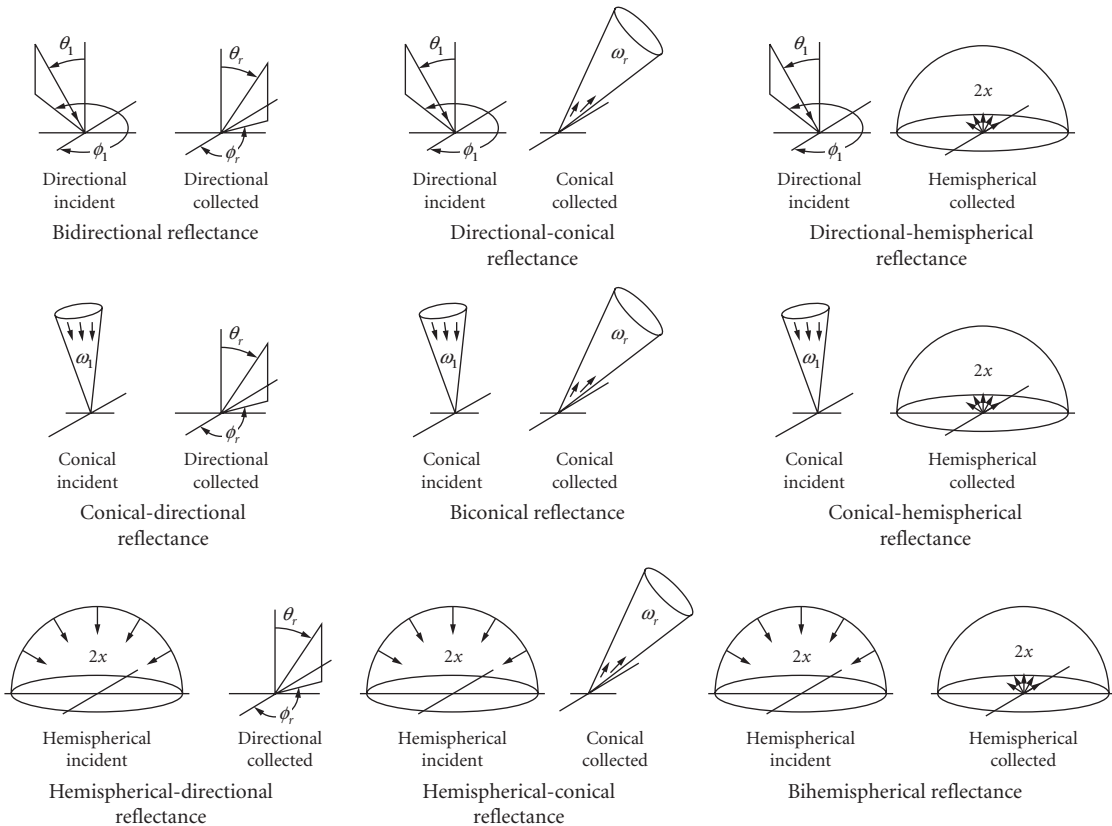


FIGURE 3 Nine geometrical definitions of reflectance.

35.6 EMITTANCE

Emittance (ε) is the ratio of the radiance of an object or surface to the radiance of a blackbody (planckian radiator) at the same temperature. It is therefore dimensionless and can assume values between 0 and 1 for thermal radiators at equilibrium. Spectral emittance $\varepsilon(\lambda)$ is the emittance at a given wavelength. If a radiator is neutral with respect to wavelength, with a constant spectral emittance less than unity, it is called a graybody.

$$\varepsilon = \frac{L}{L^{\text{bb}}}, \quad \varepsilon(\lambda) = \frac{L_\lambda}{L_\lambda^{\text{bb}}} \quad (9)$$

Directional emittance $\varepsilon(\theta, \phi)$ is defined by

$$\varepsilon(\theta, \phi) = \frac{L(\theta, \phi)}{L^{\text{bb}}} \quad (10)$$

Note that if the body is nongray, its emittance is dependent upon temperature inasmuch as the integral must be weighted by the source (Planck) function.

$$\varepsilon = \frac{\int_0^\infty \varepsilon(\lambda) L_\lambda^{\text{bb}} d\lambda}{\int_0^\infty L_\lambda^{\text{bb}} d\lambda} = \frac{1}{\pi} \frac{\int_0^\infty \varepsilon(\lambda) L_\lambda^{\text{bb}} d\lambda}{\sigma T^4} \quad (11)$$

35.7 KIRCHHOFF'S LAW

In a closed system at thermal equilibrium, conservation of energy necessitates that emitted and absorbed fluxes be equal. Since the radiation field in such a system is isotropic (the same in all directions), the directional spectral emittance and the directional spectral absorptance must be equal, i.e.,

$$\varepsilon(\lambda; \theta, \phi) = \alpha(\lambda; \theta, \phi) \quad (12)$$

This statement was first made by Kirchhoff.²⁰ Strictly, this equation holds for each orthogonal polarization component, and for it to be valid as written, the total radiation must have equal orthogonal polarization components. Kirchhoff's law is often simplified to the declaration $\alpha = \varepsilon$; however, this is not a universal truth; it may only be applied under a limited set of conditions. The geometrical and spectral averaging (integration) is governed by a specific set of rules as demonstrated by Siegel and Howell.²¹ Table 4, adapted from Siegel and Howell²¹ and Grum and Becherer,²² shows the various geometrical and spectral conditions under which the absorptance may be related to the emittance.

35.8 RELATIONSHIP BETWEEN TRANSMITTANCE, REFLECTANCE, AND ABSORPTANCE

Radiant flux incident upon a surface or medium undergoes transmission, reflection, and absorption. Application of conservation of energy leads to the statement that the sum of the transmission, reflection, and absorption of the incident flux is equal to unity, or

$$\alpha + \tau + \rho = 1 \quad (13)$$

TABLE 4 Summary of Absorptance-Emittance Relations²¹

Quantity	Equality	Required Conditions
Directional spectral	$\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$	None other than thermal equilibrium
Directional total	$\alpha(\theta, \phi, T_a) = \varepsilon(\theta, \phi, T_a)$	(1) Spectral distribution of incident energy proportional to blackbody at T_a , or (2) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of wavelength
Hemispherical spectral	$\alpha(\lambda, T_a) = \varepsilon(\lambda, T_a)$	(1) Incident radiation independent of angle, or (2) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle
Hemispherical total	$\alpha(T_a) = \varepsilon(T_a)$	(1) Incident energy independent of angle and spectral distribution proportional to blackbody at T_a , or (2) Incident energy independent of angle and $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of wavelength, or (3) Incident energy at each angle has spectral distribution proportional to blackbody at T_a and $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle, or (4) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle and wavelength

In the absence of nonlinear effects (i.e., the Raman effect, etc.),

$$\alpha(\lambda) + \tau(\lambda) + \rho(\lambda) = 1 \quad (14)$$

If the situation is such that one of the above Kirchhoff-type relations is applicable, then emittance ε may be substituted for absorptance α in the previous equations, or

$$\varepsilon = 1 - \tau - \rho \quad \varepsilon(\lambda) = 1 - \tau(\lambda) - \rho(\lambda) \quad (15)$$

35.9 MEASUREMENT OF TRANSMITTANCE

A knowledge of the transmission of optical materials and elements, gaseous atmospheres, and various liquids is necessary throughout the realm of optics. Most of these measurements are made with commercial spectrophotometers. It is beyond the scope of this chapter to discuss the design and operation of spectrophotometric equipment except sample-handling practices. For further discussion, see Gram and Becherer,²² ASTM E275,²³ and ASTM E409.²⁴

Conventional spectrophotometers are of the double-beam configuration, where the output is the ratio of the signal in the sample beam to the signal in the reference beam plotted as a function of wavelength. It is incumbent upon the experimenter to ensure that the only difference between the two beams is the unknown. Therefore, if liquid or gas cells are employed, one should be placed in each beam. For gas cells, an equal amount of carrier gas should be injected into each cell, with the unknown to be sample placed in only one cell, destined for the sample beam. For liquids, an equal amount of solute should be placed in each cell. A critical issue with liquid and solid samples is the beam geometry. Most spectrophotometers feature converging beams in the sample space. If the

optical path length (the product of index of refraction and actual distance) for each beam is not identical, a systematic difference is presented to either the entrance slit or the detector. In addition, some specimens (e.g., interference filters) are susceptible to errors when measured in a converging beam. Most instruments also have a single monochromator which is susceptible to stray radiation, the limiting factor when trying to make measurements of samples that are highly absorbing in one spectral region and transmitting in another. Some recent instruments feature linear detector arrays along with single monochromators to allow the acquisition of the entire spectrum in several milliseconds; these are particularly applicable to reaction rate studies.

Conventional double-beam instruments are limited by these factors to uncertainties on the order of 0.1 percent. For lower uncertainties, the performance deficiencies found in double-beam instruments can largely be overcome by the use of a single-beam architecture. The mode of operation is sample-in–sample-out. If the source is sufficiently stable with time, the desired spectral range can be scanned without the sample, then rescanned with the sample in place. Otherwise, the spectrometer can be set at a fixed wavelength and alternate readings with and without the sample in place must be made. Care should be taken to ensure that the beam geometry is not altered between sequential readings.

To achieve the ultimate in performance from conventional spectrophotometry, several design characteristics should be included. A double monochromator is essential to minimize stray light. The beam geometry in the sample compartment should be highly collimated to avoid focus shifts with optically thick samples. Some form of beam integration, such as an integrating sphere or other diffuser, should be employed to negate the effects of nonuniform detectors and beam shifts. An exemplary instrument is the high-accuracy spectrophotometer developed by the National Institute for Standards and Technology (NIST), described in Mielenz and Eckerle,²⁵ Mielenz et al.,²⁶ Venable et al.,²⁷ Eckerle,²⁸ and Eckerle et al.²⁹ A particularly useful review is Eckerle et al.³⁰ Similar laboratory instruments have also been built elsewhere by Clarke,³¹ Freeman,³² and Zwinkles and Gignac.³³

Numerous other instruments have been described in the literature; some have been designed for singular or limited purposes while others have a more universal appeal. Use of integrating spheres is common, both for the averaging effects and for the isolation of the specular and diffuse components, as shown in Fig. 4.^{14,34} Several useful instruments are described by Karras,³⁵ Taylor,³⁶ Zerlaut and Anderson,³⁷ Clarke and Larkin,³⁸ and Kessell.³⁹

Conventional instruments lack a wide dynamic range because there are simply not enough photons available in a narrow bandpass in a reasonable time. Solutions include Fourier transform spectrometers with a large multiplex advantage, the use of tunable lasers, and heterodyne spectrometry.⁴⁰

Simple instruments can be purchased or constructed for specific purposes. For example, solar transmittance can be determined using either the natural sun (if available) or simulated solar radiation as the source. A limited degree of spectral isolation can be achieved with an abridged spectrophotometer using narrow bandpass interference or glass absorption filters.

Several publications have suggested methods for making accurate and repeatable measurements including Hughes,⁴¹ Mielenz,⁴² Venable and Hsia,⁴³ Burke and Mavrodineanu,⁴⁴ ASTM F768,⁴⁵ ASTM E971,⁴⁶ ASTM E903,⁴⁷ and ASTM E179.⁴⁸ Calibration and performance assessment of spectrophotometers includes photometric accuracy, linearity, stray light analysis, and wavelength calibration.

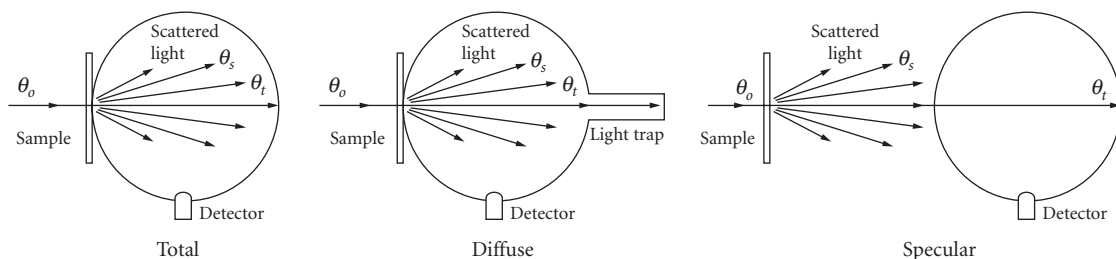


FIGURE 4 Measurement of total, diffuse, and specular transmittance using an integrating sphere.

Particular attention should be paid to luminescent samples that absorb radiant energy in one spectral region and re-emit it at longer wavelengths. Pertinent references include Hawes,⁴⁹ Bennett and Ashley,⁵⁰ ASTM E387,⁵¹ ASTM E275,²³ and ASTM, E409.²⁴

Standards of spectral transmittance are available as Standard Reference Materials from NIST. These take the form of metal-on-glass, metal-on-quartz, and solid glass filters. Most are used for verifying the photometric scale or for checking the wavelength calibration of a recording spectrophotometer. Descriptions of their development and use are given in Mavrodineanu and Baldwin,⁵² Mavrodineanu and Baldwin,⁵³ Eckerle et al.,³⁰ and Hsia.⁵⁴ The standardization laboratories of several countries occasionally conduct international intercomparisons of traveling standards. Recent intercomparisons have been reported in Eckerle et al.⁵⁵ and Fillinger and Andor.⁵⁶

35.10 MEASUREMENT OF ABSORPTANCE

In most cases, absorptance is not directly measured, but is inferred from transmission measurements, with appropriate corrections for reflection losses. These corrections can be calculated from the Fresnel equations if the surfaces are polished and the index of refraction is known. For materials where the absorption is extremely small, this method is unsatisfactory, as the uncertainties are dominated by the reflection contribution. In this case, direct measurements (such as laser calorimetry) must be made as discussed by Lipson et al.⁵⁷ and Hordvik.⁵⁸

35.11 MEASUREMENT OF REFLECTANCE

Instrumentation for the measurement of reflectance takes many forms. Only a few of the definitions for reflectance (Table 2) and reflectance factor (Table 3) have been adopted as standard configurations. The biconical configuration with small solid angles is most suited to a measurement of specular (regular, in the mirror direction) reflection. A simple reflectometer for the absolute measurement of specular reflectance was devised by Strong^{59,60} and is shown schematically in Fig. 5. Numerous detail improvements have been made on this fundamental design, including the use of averaging spheres. Designs range from simple⁶¹⁻⁶⁵ to complex.⁶⁶ Some reflectometers have been built specifically to measure at normal incidence.⁶⁷⁻⁶⁹ Measurement methods and data interpretation are also given in ASTM F768,⁴⁵ ASTM D523,⁷⁰ ASTM F1252,⁷¹ Hernicz and DeWitt,⁷² and Snail et al.⁷³

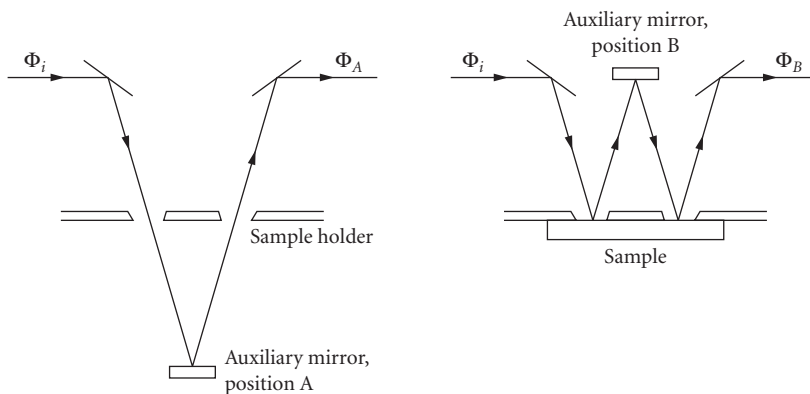


FIGURE 5 The Strong "VW" reflectometer.

The characterization of appearance of materials involves measurements of reflectance, both diffuse and specular. Numerous procedures and instruments have been devised for goniophotometry, the measurement of specular gloss with biconical geometry. Measurements are made at several angles from normal (20° , 30° , 45° , 60° , 75° , and 85°) depending upon the material under scrutiny. Further details can be found in ASTM C347,⁷⁴ ASTM E167,⁷⁵ ASTM D523,⁷⁰ ASTM E1349,⁷⁶ ASTM E179,⁷⁷ ASTM E430,⁷⁸ Erb,⁷⁹ and Hunter.⁸⁰

The measurement of diffuse reflectance can be accomplished using any one of the nine definitions from Table 2 and integrating where necessary. One could, for example, choose to measure the bidirectional reflectance distribution function (BRDF) as a function of incident beam parameters and to integrate over the hemisphere, but this would be a tedious process, and the large amount of data generated would be useful only to those involved with detailed materials properties research. Most practical measurements of diffuse reflectance involve the use of an integrating sphere. Several papers have discussed the general theory of the integrating sphere.^{81,82}

In the visible and near-IR spectral regions, the integrating sphere is the instrument of choice for both specular and diffuse specimens. Many papers have been written detailing instruments, methods, and procedures, some of which are shown in Fig. 6. The specular component of the reflected flux can be included to determine the total reflectance (Fig. 6a) or excluded to measure just the diffuse component (Fig. 6b). The angle of incidence can be varied by placing the sample at the center of the sphere (Fig. 6c), Edwards et al.⁸³ Others making contributions include Clarke and Compton,⁸⁴ Clarke and Larkin,⁸⁵ Dunkle,⁸⁶ Egan and Hilgeman,⁸⁷ Goebel et al.,⁸⁸ Hisdal,^{89,90} Karras,³⁵ McNicholas,⁹¹ Richter and Erb,⁹² Sheffer et al.,^{93,94} Taylor,⁹⁵ and Venable et al.²⁷ Some of these methods have been incorporated into standard methods and practices, such as ASTM C523,⁹⁶

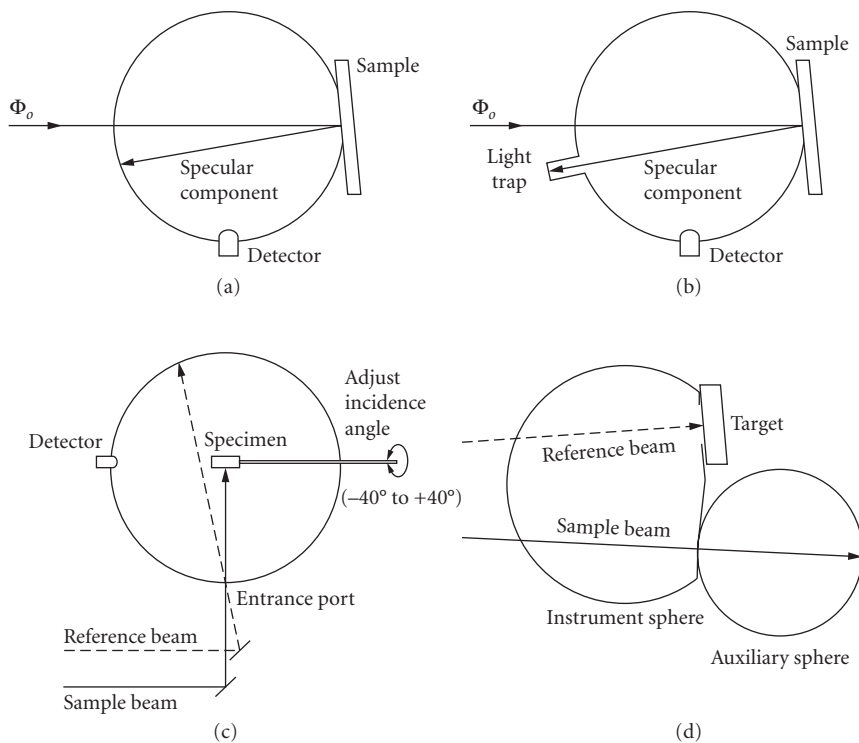


FIGURE 6 Measurement of diffuse reflectance using an integrating sphere.

ASTM E429,⁹⁷ ASTM E903,⁴⁷ CIE,⁹⁸ and IES.⁹⁹ Most integrating sphere measurements require reference to some form of an artifact standard, but the double-sphere method (Fig. 6d) produces absolute diffuse reflectance.^{37,38,100,101} Lindberg¹⁰² has demonstrated a method to scale relative measurements to absolute.

Alternative forms of hemispherical irradiation and/or collection have been described, several of which are shown in Fig. 7. Specular hemispherical (Fig. 7a), paraboloidal (Fig. 7b), and ellipsoidal (Fig. 7c) collectors have been used, particularly in those spectral regions where integrating sphere coatings are difficult to obtain.^{86,103–109} The Helmholtz reciprocity principle has been invoked to demonstrate the reversibility of the source and the collector.¹¹⁰ Hemispherical irradiation has also been employed by placing a cooled sample coplanar with the wall of a furnace as shown in Fig. 7d.^{86,111–113}

The procedures and instrumentation for the measurement of reflectance factor are identical with those for diffuse reflectance using the $0^\circ/45^\circ$ or $45^\circ/0^\circ$ geometry with annular, circumferential, or uniplanar illumination or viewing. Reference is made to a white reflectance standard characterized for reflectance factor, which must be compared with a perfect diffuse reflector. Pertinent references are ASTM E1349,⁷⁶ ASTM E97,¹¹⁴ ASTM E1348,¹¹⁵ Hsia and Weidner,¹¹⁶ and Taylor.^{36,117}

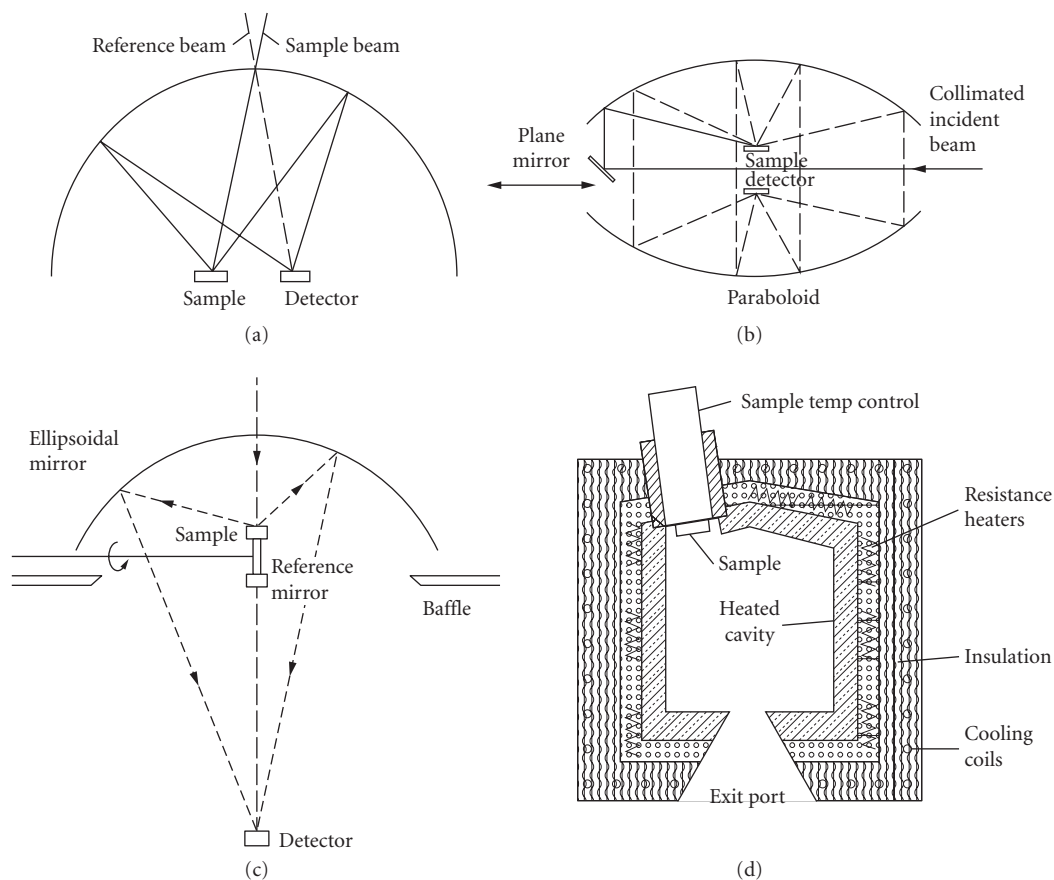


FIGURE 7 Measurement of diffuse reflectance using alternate methods.

Orbiting sensors measure the radiance of the earth-atmosphere system with some known geometry (generally nadir) and in well-defined wavelength bands. The quantity of interest is reflectance, as it is related to factors such as crop assessment, mineralization, etc. Corrections must be made for the atmospheric absorption, emission, and scattering and for the BRDF of the target. BRDF has been characterized in the field using the sun as the source, as described by Duggin.^{118,119}

Laboratory measurements of BRDF are made using goniometers where the sample-source angle and sample-receiver angle are independently adjustable. Coherent sources (lasers) are employed for the characterization of smooth specimens where a large source power is necessary for adequate SNR for off-specular angles and where speckle is not a concern. Incoherent sources (xenon arcs, black-body simulators, or tungsten-halogen) sources are often employed with spectral filters for more diffuse specimens. Similar measurements and techniques are employed to characterize bidirectional transmittance distribution function (BTDF) and bidirectional scattering distribution function (BSDF). For further information, see Asmail,¹²⁰ ASTM E1392,¹²¹ and Bartell et al.¹²²

The measurement of retroreflection poses the situation that the return beam coincides with the incident beam. The usual solution is to employ a beam splitter in the system, allowing the incident beam to pass and the return beam to be reflected. This immediately imposes both a significant loss in flux and the situation where the beam reaching the sample is partially polarized, unless a non-polarizing beamsplitter is employed. In addition, it is imperative that the reflected component of the incident beam be well-trapped, as the radiometer is looking in the same direction. The special vocabulary for retroreflection is given in CIE¹²³ and ASTM E808.¹⁷ Test methods are given in ASTM E810,¹²⁴ ASTM E809,¹²⁵ and Venable and Johnson.¹²⁶ An instrument specifically designed to measure retroreflection is detailed in Eckerle et al.¹²⁷

Most measurements of reflectance are relative and require artifact standards. Specular reflectance measurements are occasionally made using the absolute technique shown in Fig. 5, but are more commonly done using a simple reflectance attachment for a commercial spectrophotometer, requiring a calibrated reference. Freshly deposited metallic films have been used, with the assumption that an individual coating is the same as accepted data as shown in Table 5 for Al and Au.¹²⁸⁻¹³³ Standards for specular reflectance are also available.¹³⁵⁻¹³⁷

Diffuse reflectance standards exist in several forms. Ideally, a perfect ($\rho = 1$, lambertian) diffuser would be used, particularly for measuring reflectance factor.¹³⁸⁻¹⁴⁰ Certain materials approach the ideal over a limited angular and wavelength range. Historically, MgO was used in the visible spectrum.^{141,111} It was replaced first by BaSO₄¹⁴² and more recently by PTFE^{143,144} and ASTM E259.¹⁴⁵ Table 5 also shows a typical 6°/hemispherical reflectance factor for PTFE.¹⁴⁶ This fine, white powder, when pressed to a density of 1 g/cm³, is close to ideal over a wide spectral range. It is not quite lambertian, showing a falloff of BRDF at angles far removed from the specular direction¹⁴⁷ and exhibiting a slight amount of retroreflection. It may also be slightly luminescent when excited by far-ultraviolet.¹⁴⁸

In the infrared, two materials have proven useful. Flowers of sulfur¹⁴⁹⁻¹⁵¹ is suitable over the spectral range 1 to 15 μm . Gold is highly reflective and very stable. To be useful as a diffuse reflectance standard, it must be placed on top of a lambertian surface. Several substrates for gold have been suggested, including sandpaper¹⁵² and flame-sprayed aluminum.

PTFE is a satisfactory laboratory standard but is not well-suited for field use as it is not particularly rugged and is highly adsorbant and therefore subject to contamination. Several solutions have been proposed for working standards, including Eastman integrating sphere paint (BaSO₄), Vitriolite tile, and the Russian MS20 and MS14 opal glasses. These materials are, in general, more rugged, stable, and washable than PTFE. Further details can be found in CIE¹⁵³ and Clarke et al.¹⁵⁴

Discussion on the fabrication, calibration, and properties of various diffuse reflectance standards can be found in ASTM E259,¹⁵⁵ Budde,^{156,157} Egan and Hilgeman,¹⁵⁸ Fairchild and Daoust,¹⁴⁷ Morren et al.,¹⁵⁹ TAPPI,¹⁶⁰ and Weidner.^{161,162} International intercomparisons of laboratory standards of diffuse reflectance have been reported in Budde et al.,¹⁶³ and Weidner and Hsia,¹⁴⁶ and IES.⁹⁹

There are no standards available for retroreflection. However, NIST offers a Measurements Assurance Program (MAP) to enable laboratories to make measurements consistent with other national standards.

TABLE 5 Reflectance Standards

Wavelength (nm)	Aluminum	Gold	PTFE (6°/hemi)
250		0.295	0.973
300	0.921	0.346	0.984
350	0.921	0.330	0.990
400	0.919	0.360	0.993
450	0.918	0.358	0.993
500	0.916	0.453	0.994
550	0.916	0.800	0.994
600	0.912	0.906	0.994
650	0.906	0.947	0.994
700	0.898	0.963	0.994
750	0.886	0.970	0.994
800	0.868	0.973	0.994
850	0.868	0.973	0.994
900	0.891	0.974	0.994
950	0.924	0.974	0.994
1000	0.940	0.974	0.994
1100		0.975	0.994
1200	0.964	0.975	0.993
1300		0.975	0.992
1400		0.975	0.991
1500	0.974	0.975	0.992
1600		0.975	0.992
1700		0.976	0.990
1800		0.976	0.990
1900		0.976	0.985
2000	0.978	0.976	0.981
2100		0.976	0.968
2200		0.976	0.977
2300		0.976	0.972
2400		0.976	0.962
2500	0.979	0.977	0.960

35.12 MEASUREMENT OF EMITTANCE

Measurements of emittance can be done in several ways. The most direct method involves forming a material into the shape of a cavity in such a way that near-blackbody radiation is emitted. A measurement then compares the radiation from a location within the formed cavity to radiation from a flat, outside surface of the material, presumably at the same temperature.¹⁶⁴ The cavity can take the form of a cylinder, cone, or sphere. Similarly, a small-diameter, deep hole can be drilled into a specimen and radiation from the surface compared to radiation from the hole. Care must be taken that the specimen is isothermal and that the reflected radiation is considered. The definitive measurements of several materials, such as tungsten,¹⁶⁵ were determined in this fashion. The significant advantage in this direct method is that it is relative, depending on neither absolute radiometry or thermometry, but only requiring that the radiometer or spectroradiometer be linear over the dynamic range of the measurement. This linearity is also determinable by relative measurements.

If a variable-temperature blackbody simulator and a suitable thermometer are available, the specimen can be heated to the desired temperature T_s and the blackbody simulator temperature T_{bb} can be adjusted such that its (spectral) radiance matches that of the specimen. Then the (spectral)

emittance is calculable using the following equations for spectral emittance $\varepsilon(\lambda)$ and emittance ε (for a graybody only).

$$\varepsilon(\lambda) = \frac{e^{c_2/\lambda T_s} - 1}{e^{c_2/\lambda T_{bb}} - 1}, \quad \varepsilon = \frac{T_{bb}^4}{T_s^4} \quad (16)$$

If an absolutely calibrated radiometer and a satisfactory thermometer are available, a direct measurement can be made, as L_b is calculable if the temperature is known. Again, the reflected radiation must be considered.

Simple “inspection meter” techniques have been developed, and instrumentation is commercially available to determine the hemispherical emittance over a limited range of temperatures surrounding ambient. These instruments provide a single number, as they integrate both spatially and spectrally. A description of the technique can be found in ASTM E408.¹⁶⁶

Measurements of spectral emittance are most often made using spectral reflectance techniques, invoking Kirchhoff’s law along with the assumption that the transmittance is zero. A review of early work is found in Dunn et al.,¹⁶⁷ and Millard and Streed.¹⁶⁸ The usual geometry of interest is the directional-hemispherical. This can be achieved by either hemispherical irradiation-directional collection, or, using Helmholtz reciprocity,¹¹⁰ directional irradiation-hemispherical collection. Any standard reflectometry technique is satisfactory.

A direct method for the measurement of total (integrated over all wavelengths) hemispherical emittance is to use a calorimeter as shown in Fig. 8. A heated specimen is suspended in the center of a large, cold, evacuated chamber. The vacuum minimizes gaseous conduction and convection. If the sample suspension is properly designed, the predominant means of heat transfer is radiation. The chamber must be large to minimize the configuration factor between the chamber and the specimen. The chamber is cooled to T_c to reduce radiation from the chamber to the specimen. The equation used to determine emittance ε is

$$\varepsilon = \frac{P}{\sigma A(T_s^4 - T_c^4)} \quad (17)$$

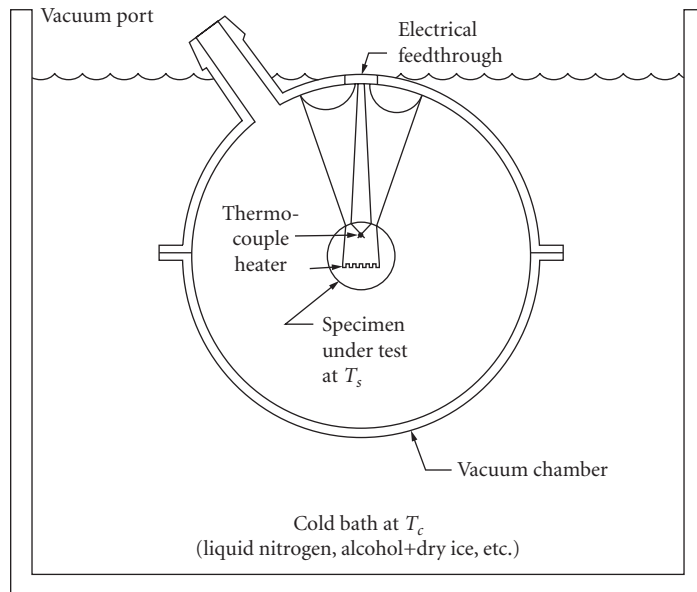


FIGURE 8 Calorimetric measurement of total hemispherical emittance.

where P is the power input to the specimen heater necessary to maintain an equilibrium specimen temperature T_s and A is the specimen area. The equation has been simplified with the aid of the following assumptions: (1) no thermal conduction from the specimen to the chamber, (2) no convective losses, (3) equilibrium has been achieved, and (4) the specimen area is much less than the chamber area. The power can be supplied electrically by means of a known heater or optically via a window in the chamber. In the latter case, a direct measurement of the ratio of solar absorptance α_s to thermal emittance ϵ_T can be directly obtained if the optical source simulates solar radiation. By varying the input power, the emittance can be determined as a function of temperature. There are numerous small corrections to account for geometry, lead conduction, etc. Details can be found in ASTM C835,¹⁶⁹ ASTM E434,¹⁷⁰ Edwards,¹⁷¹ and Richmond and Harrison.¹⁷²

Several attempts have been made to define and characterize artifact standards of spectral emittance for direct measurements.^{10,11} These were specimens of a thermally stable metal (i.e., Inconel) which were calibrated for emittance as a function of wavelength at several temperatures. No such standards are currently available. Interlaboratory comparisons have been made and reported.¹⁷³

Special problems include measurements at cryogenic temperatures¹⁷⁴ and effects of partially transparent materials.¹⁷⁵ Some additional references relating to emittance and its measurement are ASTM E307,¹⁷⁶ ASTM E423,¹⁷⁷ Clarke and Larkin,⁸⁵ DeWitt,¹⁷⁸ DeWitt and Richmond,¹⁷⁹ Hornbeck,¹⁸⁰ Millard and Streed,¹⁶⁸ Redgrove,¹⁸¹ Sparrow et al.,¹⁶⁴ Stierwalt,¹⁸² and Wittenberg.¹⁸³

35.13 REFERENCES

1. ISO, *Units of Measurement*, ISO Standards Handbook 2, International Organization for Standardization, Geneva, 1982.
2. E. R. Cohen and P. Giacomo, *Symbols, Units, Nomenclature and Physical Constants in Physics*, Document IUPAP-25, International Union of Pure and Applied Physics, 1987.
3. B. N. Taylor, "The International System of Units (SI)," *NIST Special Publication 330*, National Institute of Standards and Technology, Washington, D.C., 1991.
4. ASTM, "Nomenclature and Definitions Applicable to Radiometric and Photometric Characteristics of Matter," *ASTM Special Technical Publication 475*, ASTM, Philadelphia (1971).
5. IES Nomenclature Committee, "Proposed American National Standard Nomenclature and Definitions for Illuminating Engineering (proposed revision of Z7.1-R-1973)," *J. Illum. Eng. Soc.* **8**:2 (1979).
6. IES Nomenclature Committee, *American National Standard Nomenclature and Definitions for Illuminating Engineering*, ANSI/IES RP-16-1986. Illuminating Engineering Society of North America, N.Y., 1986.
7. J. V. Drazil, *Quantities and Units of Measurement: A Dictionary and Handbook*, Mansell, London, (1983).
8. CIE, "International Lighting Vocabulary," *CIE Publ. 17.4*, CIE, Paris (1987).
9. A. G. Worthing, "Temperature Radiation Emissivities and Emittances," in *Temperature, Its Measurement and Control in Science and Industry*, Reinhold, N.Y., p. 1164 (1941).
10. J. C. Richmond, "Physical Standards of Emittance and Reflectance," in H. H. Blau and H. Fischer (eds), *Radiative Transfer from Solid Material*, Macmillan, N.Y., 1962.
11. J. C. Richmond, W. N. Harrison, and F. J. Shorten, "An Approach to Thermal Emittance Standards, in J. C. Richmond (ed.), *Measurement of Thermal Radiation Properties of Solids*, NASA SP-31, NASA, Washington, D.C., 1963.
12. W. L. Wolfe, "Proclivity for Emissivity," *Appl. Opt.* **21**:1 (1982).
13. J. C. Richmond, J. J. Hsia, V. R. Weidner, and D. B. Wilmering, *Second-Surface Mirror Standards of Spectral Specular Reflectance*, NBS Special Publication SP260-79, U.S. National Bureau of Standards, Washington, D.C., 1982.
14. CIE, "Radiometric and Photometric Characteristics of Materials and their Measurement," *CIE Publication 38*, CIE, Paris (1977).
15. F. E. Nicodemus, J. C. Richmond, and J. J. Hsia, *Geometrical Considerations and Nomenclature for Reflectance*, NBS Monograph 160, U.S. National Bureau of Standards, Washington, D.C., 1977.
16. D. B. Judd, "Terms, Definitions and Symbols in Reflectometry," *J. Opt. Soc. Am.* **57**:445 (1967).

17. ASTM, "Terminology for Retroreflection and Retroreflectors," *ASTM E808*, ASTM, Philadelphia (1981).
18. F. E. Nicodemus, "Directional Reflectance and Emissivity of an Opaque Surface," *Appl. Opt.* **4**:767 (1965).
19. F. E. Nicodemus, "Reflectance Nomenclature and Directional Reflectance and Emissivity," *Appl. Opt.* **9**:1474 (1970).
20. G. Kirchhoff, "On the Relation between the Radiating and Absorbing Powers of Different Bodies for Light and Heat," *Phil. Mag.* **20**:1 (1860).
21. R. Siegel, and J. R. Howell, *Thermal Radiation Heat Transfer*, 2d ed., Hemisphere, N.Y., p. 63 (1981).
22. F. Grum, and R. J. Becherer, "Radiometry," in *Optical Radiation Measurements*, vol. 1, Academic, N.Y., p. 115 (1979).
23. ASTM, "Practice for Describing and Measuring Performance of UV/VIS/Near-IR Spectrophotometers," *ASTM E275*, ASTM, Philadelphia (1989).
24. ASTM, "Procedure for Description and Performance in the Spectrophotometer," *ASTM E409*, ASTM, Philadelphia (1990).
25. K. D. Mielenz, and K. L. Eckerle, *Design, Construction, and Testing of a New High Accuracy Spectrophotometer*, NBS Tech Note 729, U.S. National Bureau of Standards, Washington, D.C., 1972.
26. K. D. Mielenz, K. L. Eckerle, R. P. Madden, and J. Reader, "New Reference Spectrophotometer," *Appl. Opt.* **12**:1630 (1973).
27. W. H. Venable, J. J. Hsia, and V. R. Weidner, *Development of an NBS Reference Spectrophotometer for Diffuse Transmittance and Reflectance*, NBS Tech Note TN594-11, U.S. National Bureau of Standards, Washington, D.C., 1976.
28. K. L. Eckerle, *Modification of an NBS Reference Spectrophotometer*, NBS Technical Note TN913, U.S. National Bureau of Standards, Washington, D.C., 1976.
29. K. L. Eckerle, V. R. Weidner, J. J. Hsia, and Z. W. Chao, *Extension of a Reference Spectrophotometer into the Near Infrared*, NBS Technical Note TN1175, U.S. National Bureau of Standards, Washington, D.C., 1983.
30. K. L. Eckerle, J. J. Hsia, K. D. Mielenz, and V. R. Weidner, *Regular Spectral Transmittance*, NBS Special Publication SP250-6, U.S. National Bureau of Standards, Washington, D.C., 1987.
31. F. J. J. Clarke, "High-Accuracy Spectrophotometry at the National Physical Laboratory," *J. Res. Natl. Bur. Stand.* **A76**:375 (1972).
32. G. H. C. Freeman, "The New Automated Reference Spectrophotometer at NPL," in C. Burgess and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1986, p. 69.
33. J. C. Zwinkles and D. S. Gignac, "Design and Testing of a New High-Accuracy Ultraviolet-Visible- Near-Infrared Spectrophotometer," *Appl. Opt.* **31**:1557 (1992).
34. A. Roos, "Interpretation of Integrating Sphere Signal Output for Nonideal Transmitting Samples," *Appl. Opt.* **30**:468 (1991).
35. E. Karras, "The Use of the Ulbricht Sphere in Measuring Reflection and Transmission Factors," *J Opt. Soc. Am.* **11**:96 (1921).
36. A. H. Taylor, "A Simple Portable Instrument for the Absolute Measurement of Reflection and Transmission Factors," *Sci. Papers Bur. Standards* **17**:1 (1922).
37. G. A. Zerlaut and T. E. Anderson, "Multiple-Integrating Sphere Spectrophotometer for Measuring Absolute Spectral Reflectance and Transmittance," *Appl. Opt.* **20**:3797 (1981).
38. F. J. J. Clarke and J. A. Larkin, "Measurement of Total Reflectance, Transmittance and Emissivity over the Thermal IR Spectrum," *Infrared Phys.* **25**:359 (1985).
39. J. Kessell, "Transmittance Measurements in the Integrating Sphere," *Appl. Opt.* **25**:2752 (1986).
40. A. L. Migdall, B. Roop, Y. C. Zheng, J. E. Hardis, and G. J. Xia, "Use of Heterodyne Detection to Measure Optical Transmittance over a Wide Range," *Appl. Opt.* **29**:5136 (1990).
41. H. K. Hughes, "Beer's Law and the Optimum Transmittance in Absorption Measurements," *Appl. Opt.* **2**:937 (1963).
42. K. D. Mielenz, "Physical Parameters in High-Accuracy Spectrophotometry," in R. Mavrodineanu, J. I. Schultz, and O. Menis (eds), *Accuracy in Spectrophotometry and Luminescence Measurements*, NBS SP378, U.S. National Bureau of Standards, Washington, D.C., 1973.
43. W. H. Venable and J. J. Hsia, *Describing Spectrophotometric Measurements*, NBS Technical Note TN594-9, U.S. National Bureau of Standards, Washington, D.C., 1974.

44. R. W. Burke and R. Mavrodineanu, *Standard Reference Material: Accuracy in Analytical Spectrophotometry*, NBS Special Publication SP260-81, U.S. National Bureau of Standards, Washington, D.C., 1983.
45. ASTM, "Specular Reflectance/Transmittance of Optically Flat Coated/Non-coated Specimens," *ASTM F768*, ASTM, Philadelphia (1987).
46. ASTM, "Test for Photometric Transmittance/Reflectance of Materials to Solar Radiation," *ASTM E971*, ASTM, Philadelphia (1988).
47. ASTM, "Standard Test Method for Solar Absorptance, Reflectance and Transmittance of Materials Using Spectrophotometers with Integrating Spheres," *ASTM E903*, ASTM, Philadelphia (1988).
48. ASTM, "Standard Guide for Selection of Geometric Conditions for Measurement of Reflectance and Transmittance Properties of Materials," *ASTM E179*, ASTM, Philadelphia (1990).
49. R. C. Hawes, "Technique for Measuring Photometric Accuracy," *Appl. Opt.* **10**:1246 (1971).
50. J. M. Bennett and E. J. Ashley, "Calibration of Instruments Measuring Reflectance and Transmittance," *Appl. Opt.* **11**:1749 (1972).
51. ASTM, "Estimating Stray Radiant Energy of Spectrophotometers," *ASTM E387*, ASTM, Philadelphia (1984).
52. R. Mavrodineanu and J. R. Baldwin, *Standard Reference Materials: Glass Filters as a Standard Reference Material for Spectrophotometry—Selection, Preparation, Certification, Use*, NBS Special Publication SP260-51, U.S. National Bureau of Standards, Washington, D.C., 1975.
53. R. Mavrodineanu and J. R. Baldwin, *Standard Reference Materials: Metal-on-Quartz Filters as a Standard Reference Material for Spectrophotometry*, NBS Special Publication SP260-68, U.S. National Bureau of Standards, Washington, D.C., 1980.
54. J. J. Hsia, "National Scales of Spectrometry in the U.S.," in C. Burgess and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
55. K. E. Eckerle, E. Sutter, G. H. Freeman, G. Andor, and L. Fillinger, "International Intercomparison for Transmittance," *Metrologia* **27**:33 (1990).
56. L. Fillinger and G. Andor, "International Intercomparison of Transmittance Measurement," *CIE Journal* **7**:21 (1990).
57. H. G. Lipson, L. H. Skolnik, and D. L. Stierwalt, "Small Absorption Coefficient Measurement by Calorimetric and Spectral Emittance Techniques," *Appl. Opt.* **13**:1741 (1974).
58. A. Hordvik, "Measurement Techniques for Small Absorption Coefficients: Recent Advances," *Appl. Opt.* **16**:2827 (1977).
59. J. Strong, *Procedures in Experimental Physics*, Prentice-Hall, N.Y., 1938, p. 376.
60. J. Strong, *Procedures in Applied Optics*, Dekker, N.Y., 1989, p. 162.
61. C. Castellini, G. Emiliani, E. Masetti, P. Poggi, and P. P. Polato, "Characterization and Calibration of a Variable-Angle Absolute Reflectometer," *Appl. Opt.* **29**:538 (1990).
62. R. S. Ram, O. Prakash, J. Singh, and S. P. Varma, "Simple Design for a Reflectometer," *Opt. Eng.* **30**:467 (1991).
63. R. F. Weeks, "Simple Wide Range Specular Reflectometer," *J. Opt. Soc. Am.* **48**:775 (1958).
64. V. R. Weidner and J. J. Hsia, "NBS Specular Reflectometer-Spectrophotometer," *Appl. Opt.* **19**:1268 (1980).
65. D. K. Zhuang and T. L. Yang, "Spectral Reflectance Measurements Using a Precision Multiple Reflectometer in the UV and VUV Range," *Appl. Opt.* **28**:5024 (1989).
66. H. E. Bennett and W. F. Koehler, "Precision Measurements of Absolute Specular Reflectance with Minimized Systematic Errors," *J. Opt. Soc. Am.* **50**:1 (1960).
67. A. Bittar and J. D. Hamlin, "High-Accuracy True Normal-Incidence Absolute Reflectometer," *Appl. Opt.* **23**:4054 (1984).
68. G. Boivin and J. M. Theriault, "Reflectometer for Precise Measurement of Absolute Specular Reflectance at Normal Incidence," *Rev. Sci. Instrum.* **52**:1001 (1981).
69. J. E. Shaw and W. R. Blevin, "Instrument for the Absolute Measurement of Direct Spectral Reflectances at Normal Incidence," *J. Opt. Soc. Am.* **54**:334 (1964).
70. ASTM, "Test Method for Specular Gloss," *ASTM D523*, ASTM, Philadelphia (1988).

71. ASTM, "Test Method for Measuring Optical Reflectivity of Transparent Materials," *ASTM F1252*, ASTM, Philadelphia (1990).
72. R. S. Hernicz and D. P. DeWitt, "Evaluation of a High Accuracy Reflectometer for Specular Materials," *Appl. Opt.* **12**:2454 (1973).
73. K. A. Snail, A. A. Morrish, and L. M. Hanssen, "Absolute Specular Reflectance Measurements in the Infrared," *Proc. SPIE* **692**:143 (1986).
74. ASTM, "Test for Reflectivity and Coefficient of Scatter of White Porcelain Enamels," *ASTM C347*, ASTM, Philadelphia (1983).
75. ASTM, "Recommended Practice for Goniophotometry of Objects and Materials," *ASTM E167*, ASTM, Philadelphia (1987).
76. ASTM, "Standard Test Method for Reflectance Factor and Color by Spectrophotometry Using Bidirectional Geometry," *ASTM E1349*, ASTM, Philadelphia (1990).
77. ASTM, "Measurement of Gloss of High-Gloss Surfaces by Goniophotometry," *ASTM E430*, ASTM, Philadelphia (1983).
78. W. Erb, "Computer-Controlled Gonioreflectometer for the Measurement of Spectral Reflection Characteristics," *Appl. Opt.* **19**:3789 (1980).
79. W. Erb, "High Accuracy Gonioreflectance Spectrometry," Chap. 2.2 in C. Burgess and K. D. Mielenz (eds), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
80. R. S. Hunter, *The Measurement of Appearance*, Wiley, N.Y., 1975.
81. J. A. J. Jacquez and H. F. Kuppenheim, "Theory of the Integrating Sphere," *J. Opt. Soc. Am.* **45**:460 (1955).
82. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**:125 (1967).
83. D. K. Edwards, J. T. Gier, K. E. Nelson, and R. D. Roddick, "Integrating Sphere for Imperfectly Diffuse Samples," *J. Opt. Soc. Am.* **51**:1279 (1961).
84. F. J. J. Clarke and J. A. Compton, "Correction Methods for Integrating Sphere Measurement of Hemispherical Reflectance," *Color. Res. Appl.* **11**:253 (1986).
85. F. J. J. Clarke and J. A. Larkin, "Improved Techniques for the NPL Hemispherical Reflectometer," *Proc. SPIE* **917**:7 (1988).
86. R. V. Dunkle, "Spectral Reflectance Measurements," in F. J. Clauss (ed.), *Surface Effects on Spacecraft Materials*, Wiley, N.Y., 1960, p. 117.
87. W. G. Egan and T. Hilgeman, "Integrating Spheres for Measurements between 0.185 micrometers and 12 micrometers," *Appl. Opt.* **14**:1137 (1975).
88. D. G. Goebel, B. P. Caldwell, and H. K. Hammond III, "Use of an Auxiliary Sphere with a Spectroreflectometer to Obtain Absolute Reflectance," *J. Opt. Soc. Am.* **56**:783 (1966).
89. B. J. Hisdal, "Reflectance of Perfect Diffuse and Specular Samples in the Integrating Sphere," *J. Opt. Soc. Am.* **55**:1122 (1965a).
90. B. J. Hisdal, "Reflectance of Nonperfect Surfaces in the Integrating Sphere," *J. Opt. Soc. Am.* **55**:1255 (1965b).
91. H. J. McNicholas, "Absolute Methods in Reflectometry," *J. Res. Natl. Bur. Stand.* **1**:29 (1928).
92. W. Richter and W. Erb, "Accurate Diffuse Reflectance Measurements in the IR Spectral Range," *Appl. Opt.* **26**:4620 (1987).
93. D. Sheffer, U. P. Oppenheim, D. Clement, and A. D. Devir, "Absolute Reflectometer for the 0.8–2.5 micrometer Region," *Appl. Opt.* **26**:583 (1987).
94. D. Sheffer, U. P. Oppenheim, and A. D. Devir, "Absolute Measurements of Diffuse Reflectances in the x°/d Configuration," *Appl. Opt.* **30**:3181 (1991).
95. A. H. Taylor, "Errors in Reflectometry," *J. Opt. Soc. Am.* **25**:51 (1935).
96. ASTM, "Test for Light Reflectance of Acoustical Materials by the Integrating Sphere Method," *ASTM C523*, ASTM, Philadelphia (1984).
97. ASTM, "Measurement and Calculation of Reflecting Characteristics of Metallic Surfaces Using Integrating Sphere Instruments," *ASTM E429*, ASTM, Philadelphia (1987).
98. CIE, "Absolute Methods for Reflection Measurements," *CIE Publication 44*, CIE, Paris (1979).

99. IES Testing Procedures Committee, "IES Approved Method for Total and Diffuse Reflectometry," IES LM-44-1985, *J. Illum. Eng. Soc.* **19**:195 (1985).
100. J. A. Van den Akker, L. R. Dearth, and W. M. Shilcox, "Evaluation of Absolute Reflectance for Standardization Purposes," *J. Opt. Soc. Am.* **56**:250 (1966).
101. W. H. Venable, J. J. Hsia, and V. R. Weidner, "Establishing a Scale of Directional-Hemispherical Reflectance Factor I: The Van den Akker Method," *J. Res. Natl. Bur. Stand.* **82**:29 (1977).
102. J. D. Lindberg, "Absolute Diffuse Reflectance from Relative Reflectance Measurements," *Appl. Opt.* **26**:2900 (1987).
103. P. Y. Barnes and J. J. Hsia, "45°/0° Bidirectional Reflectance Distribution Function Standard Development," *Proc. SPIE* **1165**:165 (1989).
104. W. R. Blevin and W. J. Brown, "An Infrared Reflectometer with a Spheroidal Mirror," *J. Sci. Instrum.* **42**:1 (1965).
105. S. T. Dunn, J. C. Richmond, and J. C. Weibel, "Ellipsoidal Mirror Reflectometer," *J. Res. Nat. Bur. Stand.* **70C**:75 (1966b).
106. L. M. Hanssen and K. A. Snail, "Infrared Diffuse Reflectometer for Spectral, Angular and Temperature Resolved Measurements," *Proc. SPIE* **807**:148 (1987).
107. P. L. Hartman and E. Logothetis, "An Absolute Reflectometer for Use at Low Temperatures," *Appl. Opt.* **3**:255 (1964).
108. J. T. Neu, R. S. Dummer, and O. E. Myers, "Hemispherical Directional Ellipsoidal Infrared Spectroreflectometer," *Proc. SPIE* **807**:165 (1987).
109. B. E. Wood, J. G. Pipes, A. M. Smith, and J. A. Roux, "Hemi-Ellipsoidal Mirror Infrared Reflectometer," *Appl. Opt.* **15**:940 (1976).
110. F. J. J. Clarke and D. J. Parry, "Helmholtz Reciprocity: Its Validity and Application to Reflectometry," *Light. Res. Technol.* **17**:1 (1985).
111. J. T. Agnew and R. B. McQuistan, "Experiments Concerning Infrared Diffuse-Reflectance Standards in the Range 0.8 to 20.0 Micrometers," *J. Opt. Soc. Am.* **43**:999 (1953).
112. J. T. Gier, R. V. Dunkle, and J. T. Bevans, "Measurement of Absolute Spectral Reflectivity from 1.0 to 15 microns," *J. Opt. Soc. Am.* **44**:558 (1954).
113. D. C. Reid and E. D. McAlister, "Measurement of Spectral Emissivity from 3 μ to 15 μ ," *J. Opt. Soc. Am.* **49**:78 (1959).
114. ASTM, "Test Method for (45-0) Directional Reflectance Factor of Opaque Specimens by Broad-Band Filter Reflectometry," *ASTM E97*, ASTM, Philadelphia (1987).
115. ASTM, "Standard Test Method for Reflectance Factor and Color by Spectrophotometry Using Hemispherical Geometry," *ASTM E1348*, ASTM, Philadelphia (1990).
116. J. J. Hsia and V. R. Weidner, "NBS 45-degree/Normal Reflectometer for Absolute Reflectance Factors," *Metrologia* **17**:97 (1981).
117. A. H. Taylor, "The Measurement of Diffuse Reflection Factors and a New Absolute Reflectometer," *J. Opt. Soc. Am.* **4**:9 (1920).
118. M. J. Duggin, "The Field Measurement of Reflectance Factors," *Photogram. Eng. Rem. Sens.* **46**:643 (1980).
119. M. J. Duggin and W. R. Philipson, "Field Measurement of Reflectance: Some Major Considerations," *Appl. Opt.* **21**:2833 (1982).
120. C. Asmail, "Bidirectional Scattering Distribution Function (BSDF): A Systematized Bibliography," *J. Res. Natl. Inst. Stand. Technol.* **96**:215 (1991).
121. ASTM, "Standard Practice for Angle Resolved Optical Scatter Measurements on Specular or Diffuse Surfaces," *ASTM E1392*, ASTM, Philadelphia (1990).
122. F. O. Bartell, E. L. Dereniak, and W. L. Wolfe, "Theory and Measurement of Bidirectional Reflectance Distribution Function (BRDF) and Bidirectional Transmittance Distribution Function (BTDF)," *Proc. SPIE* **257**:154 (1980).
123. CIE, "Retroreflection: Definition and Measurement," *CIE Publ. 54*, CIE, Paris (1982).
124. ASTM, "Test Method for Coefficient of Retroreflection on Retroreflective Sheeting," *ASTM E810*, ASTM, Philadelphia (1981).
125. ASTM, "Standard Practice for Measuring Photometric Characteristics of Retroreflectors," *ASTM E809*, ASTM, Philadelphia (1991).

126. W. H. Venable and N. L. Johnson, "Unified Coordinate System for Retroreflectance Measurements," *Appl. Opt.* **19**:1236 (1980).
127. K. L. Eckerle, J. J. Hsia, V. R. Weidner, and W. H. Venable, "NBS Reference Retroreflectometer," *Appl. Opt.* **19**:1253 (1980).
128. G. Hass and J. E. Waylonis, "Optical Constants and Reflectance and Transmittance of Evaporated Aluminum in the Visible and Ultraviolet," *J. Opt. Soc. Am.* **51**:719 (1961).
129. H. E. Bennett, J. M. Bennett, and E. J. Ashley, "Infrared Reflectance of Evaporated Aluminum Films," *J. Opt. Soc. Am.* **52**:1245 (1962).
130. H. E. Bennett, M. Silver, and E. J. Ashley, "Infrared Reflectance of Aluminum Evaporated in Ultra-High Vacuum," *J. Opt. Soc. Am.* **53**:1089 (1963).
131. G. Hass and R. E. Thun, *Physics of Thin Films*, vol. 2, Academic, N.Y., 1964, p. 337.
132. J. M. Bennett and E. J. Ashley, "Infrared Reflectance and Emittance of Silver and Gold Evaporated in Ultra-High Vacuum," *Appl. Opt.* **4**:221 (1965).
133. G. Hass, "Reflectance and Preparation of Front-Surface Mirrors for Use at Various Angles of Incidence from the Ultraviolet to the Far Infrared," *J. Opt. Soc. Am.* **72**:27 (1982).
134. J. C. Richmond and J. J. Hsia, *Preparation and Calibration of Standards of Spectral Specular Reflectance*, NBS Special Publication SP260-38, U.S. National Bureau of Standards, Washington, D.C., 1972.
135. J. C. Richmond, "Rationale for Emittance and Reflectivity," *Appl. Opt.* **21**:1 (1982).
136. J. F. Verrill, "Physical Standards in Absorption and Reflection Spectrometry," Chap. 3.1 in C. Burgess and K. D. Mielenz (eds), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
137. V. R. Weidner and J. J. Hsia, *Standard Reference Materials: Preparation and Calibration of First Surface Aluminum Mirror Spectral Reflectance Standard*, NBS Special Publication SP260-75, U.S. National Bureau of Standards, Washington, D.C., 1982.
138. W. Erb, "Requirements for Reflection Standards and the Measurement of their Reflection Value," *Appl. Opt.* **14**:493 (1975).
139. W. Erb and W. Budde, "Properties of Standard Materials for Reflection," *Color Res. Appl.* **4**:113 (1979).
140. D. Scheffer, U. P. Oppenheim, and A. D. Devir, "Absolute Reflectometer for the Mid-Infrared Region," *Appl. Opt.* **29**:129 (1990).
141. W. E. K. Middleton and C. L. Sanders, "The Absolute Spectral Diffuse Reflectance of Magnesium Oxide," *J. Opt. Soc. Am.* **41**:419 (1951).
142. F. Grum and G. W. Luckey, "Optical Sphere Paint and a Working Standard of Reflectance," *Appl. Opt.* **7**:2289 (1968).
143. F. Grum and M. Saltzman, "New White Standard of Reflectance," *CIE Publication 36*, CIE, Paris (1976).
144. V. R. Weidner and J. J. Hsia, "Reflection Properties of Pressed Polytetrafluorethylene Powder," *J. Opt. Soc. Am.* **71**:856 (1981).
145. ASTM, "Preparation of Pressed Power White Reflectance Factor Transfer Standards for Hemispherical Geometry," *ASTM E259*, ASTM, Philadelphia (1992).
146. V. R. Weidner and J. J. Hsia, *Spectral Reflectance*, NBS Special Publication SP250-8, U.S. National Bureau of Standards, Washington, D.C., 1987.
147. M. D. Fairchild and D. J. O. Daoust, "Goniospectrophotometric Analysis of Pressed PTFE Powder for Use as a Primary Transfer Standard," *Appl. Opt.* **27**:3392 (1988).
148. R. D. Saunders and W. R. Ott, "Spectral Irradiance Measurements: Effect of UV-Produced Luminescence in Integrating Spheres," *Appl. Opt.* **15**:827 (1976).
149. M. Kronstein, R. J. Kraushaar, and R. E. Deacle, "Sulfur as a Standard of Reflectance in the Infrared," *J. Opt. Soc. Am.* **53**:458 (1963).
150. S. T. Dunn, "Application of Sulfur Coatings to Integrating Spheres," *Appl. Opt.* **4**:877 (1965).
151. R. Tkachuk and F. D. Kuzina, "Sulfur as a Proposed Near Infrared Reflectance Standard," *Appl. Opt.* **17**:2817 (1978).
152. T. W. Stuhlinger, E. L. Dereniak, and F. O. Bartell, "Bidirectional Distribution Function of Gold-Plated Sandpaper," *Appl. Opt.* **20**:2648 (1981).
153. CIE, "A Review of Publications on Properties and Reflection Values of Material Reflection Standards," *CIE Publication 46*, CIE, Paris (1979b).

154. F. J. J. Clarke, F. A. Garforth, and D. J. Parr, "Goniophotometric and Polarization Properties of White Reflection Standard Materials," *Light. Res. Technol.* **15**:133 (1983).
155. ASTM, "Practice for Preparation of Reference White Reflectance Standards," *ASTM E259*, ASTM, Philadelphia (1987).
156. W. Budde, "Standards of Reflectance," *J. Opt. Soc. Am.* **50**:217 (1960).
157. W. Budde, "Calibration of Reflectance Standards," *J. Res. Natl. Bur. Stand.* **A80**:585 (1976).
158. W. G. Egan and T. Hilgeman, "Retroreflectance Measurements of Photometric Standards and Coatings," *Appl. Opt.* **15**:1845 (1976).
159. L. Morren, G. Vandermeersch, and P. Antoine, "A Study of the Reflection Factor of Usual Photometric Standards in the Near Infrared," *Light. Res. Technol.* **4**:243 (1972).
160. TAPPI, "Calibration of Reflectance Standards for Hemispherical Geometry," TAPPI Standard TIS 0804-07 in *1990 TAPPI Test Methods*, TAPPI, Atlanta, 1990.
161. V. R. Weidner, *Standard Reference Materials: White, White Opal Glass Diffuse Spectral Reflectance Standards for the Visible Spectrum*, NBS Special Publication SP260-82, U.S. National Bureau of Standards, Washington, D.C., 1983.
162. V. R. Weidner, "Gray Scale of Diffuse Reflectance for the 250–2500 nm Wavelength Range," *Appl. Opt.* **25**:1265 (1986).
163. W. Budde, W. Erb, and J. J. Hsia, "International Intercomparison of Absolute Reflectance Scales," *Color Res. Appl.* **7**:24 (1982).
164. E. M. Sparrow, P. D. Kruger, and R. P. Heinisch, "Cavity Methods for Determining the Emittance of Solids," *Appl. Opt.* **12**:2466 (1973).
165. J. C. DeVos, "A New Determination of the Emissivity of Tungsten Ribbon," *Physica* **20**:690 (1954).
166. ASTM, "Test for Total Normal Emittances of Surfaces Using Inspection Meter Techniques," *ASTM E408*, ASTM, Philadelphia (1990).
167. S. T. Dunn, J. C. Richmond, and J. F. Parmer, "Survey of Infrared Measurement Techniques and Computational Methods in Radiant Heat Transfer," *J. Spacecraft Rockets* **3**:961 (1966a).
168. J. P. Millard and E. R. Streed, "A Comparison of Infrared Emittance Measurements and Measurement Techniques," *Appl. Opt.* **8**:1485 (1969).
169. ASTM, "Total Hemispherical Emittance of Surfaces from 20 to 1400C," *ASTM C835*, ASTM, Philadelphia (1988).
170. ASTM, "Test for Calorimetric Determination of Hemispherical Emittance and the Ratio of Solar Absorbance to Hemispherical Emittance Using Solar Simulation," *ASTM E434*, ASTM, Philadelphia (1990).
171. D. K. Edwards, "Thermal Radiation Measurements," Chap. 9 in E. R. G. Eckert and R. J. Goldstein (eds), *Measurement Techniques in Heat Transfer*, AGARD 130, Technivision, Slough, England, 1970, p. 353.
172. J. C. Richmond and W. N. Harrison, "Equipment and Procedures for Evaluation of Total Hemispherical Emittance," *Am. Ceram. Soc. Bull.* **39**:668 (1960).
173. R. R. Willey, "Results of a Round-Robin Measurement of Spectral Emittance in the Mid-Infrared," *Proc. SPIE* **807**:140 (1987).
174. D. Weber, "Spectral Emissivity of Solids in the Infrared at Low Temperatures," *J. Opt. Soc. Am.* **49**:815 (1959).
175. R. Gardon, "The Emissivity of Transparent Materials," *J. Am. Ceram. Soc.* **39**:278 (1956).
176. ASTM, "Test for Normal Spectral Emittance at Elevated Temperatures," *ASTM E307*, ASTM, Philadelphia (1990).
177. ASTM, "Test for Normal Spectral Emittance at Elevated Temperatures of Non-Conducting Specimens," *ASTM E423*, ASTM, Philadelphia (1990).
178. D. P. DeWitt, "Inferring Temperature from Optical Radiation Measurements," *Opt. Eng.* **25**:596 (1986).
179. D. P. DeWitt and J. C. Richmond, "Theory and Measurement of the Thermal Radiative Properties of Metals," in *Techniques of Metals Research*, vol. 6, Wiley, N.Y., 1970.
180. G. A. Hornbeck, "Optical Methods of Temperature Measurement," *Appl. Opt.* **5**:179 (1966).
181. J. S. Redgrove, "Measurement of the Spectral Emissivity of Solid Materials," *Measurement (UK)* **8**:90 (1990).
182. D. L. Stierwalt, "Infrared Spectral Emittance Measurements on Optical Materials," *Appl. Opt.* **5**:1911 (1966).
183. A. M., Wittenberg, "Determination of Total Emittance of a Nongray Surface," *J. Appl. Phys.* **39**:1936 (1968).

35.14 FURTHER READING

- ASTM, *ASTM Standards on Color and Appearance Measurement*, 3d ed., ASTM, Philadelphia, 1991. Blau, Jr, H. H. and H. Fischer (eds.), *Radiative Transfer from Solid Materials*, Macmillan, N.Y., 1962. Burgess, C. and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
- Clauss, F. J. (ed.), *First Symposium, Surface Effects on Spacecraft Materials*, Wiley, N.Y., 1960.
- Frei, R. W. and J. D. MacNeil, *Diffuse Reflectance Spectroscopy in Environmental Problem-Solving*, CRC Press, Cleveland, 1973.
- Grum, F. and R. J. Becherer, "Radiometry," in *Optical Radiation Measurements*, vol. 1, Academic, N.Y., 1979.
- Hammond III, H. K., and H. L. Mason (eds.), "Selected NBS Papers on Radiometry and Photometry," NBS Special Publication SP300-7, *Precision Measurement and Calibration*, U.S. National Bureau of Standards, Washington, D.C., 1971.
- Hunter, R. S., *The Measurement of Appearance*, Wiley, N.Y., 1975.
- Kortum, G., *Reflectance Spectroscopy*, Springer-Verlag, N.Y., 1969.
- Nimeroff, L. (ed.), "Selected NBS Papers on Colorimetry," NBS Special Publication SP300-9, *Precision Measurement and Calibration*, U.S. National Bureau of Standards, Washington, D.C., 1972.
- Richmond, J. C. (ed.), "Measurement of Thermal Radiation Properties of Solids," *NASA Special Publication SP-31*, National Aeronautics and Space Administration, Washington, D.C., 1963.
- Walsh, J. W. T., *Photometry*, 3d ed., Dover, N.Y., 1958.
- Wendlandt, W. W. and H. G. Hecht, *Reflectance Spectroscopy*, Interscience, N.Y., 1966.

This page intentionally left blank.

RADIOMETRY AND PHOTOMETRY: UNITS AND CONVERSIONS

James M. Palmer*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

36.1 GLOSSARY[†]

A	area, m ²
A_{proj}	projected area, m ²
E, E_e	radiant incidence (irradiance), W m ⁻²
E_p, E_q	photon incidence, s ⁻¹ m ⁻²
E_v	illuminance, lm m ⁻² [lux (lx)]
I, I_e	radiant intensity, W sr ⁻¹
I_p, I_q	photon intensity, s ⁻¹ sr ⁻¹
I_v	luminous intensity, candela (cd)
K_m	absolute luminous efficiency at λ_p for photopic vision, 683 lm/W
K'_m	absolute luminous efficiency at λ_p for scotopic vision, 1700 lm/W
$K(\lambda)$	absolute spectral luminous efficiency, photopic vision, lm/W
$K'(\lambda)$	absolute spectral luminous efficiency, scotopic vision, lm/W
L, L_e	radiance, W m ⁻² sr ⁻¹
L_p, L_q	photon radiance (photonance), s ⁻¹ m ⁻² sr ⁻¹
L_v	luminance, cd sr ⁻¹
M, M_e	radiant exitance, W m ⁻²
M_p, M_q	photon exitance, s ⁻¹ m ⁻²
Q, Q_e	radiant energy, joule (J)
Q_p, Q_q	photon energy, J or eV
Q_v	luminous energy, lm s ⁻¹

*Deceased.

[†]Note. The subscripts are used as follows: e (energy) for radiometric, v (visual) for photometric, and q (or p) for photonic. The subscript e is usually omitted; the other subscripts may also be omitted if the context is unambiguous.

$\mathfrak{R}(\lambda)$	spectral responsivity, A/W or V/W
$V(\lambda)$	relative spectral luminous efficiency, photopic vision
$V'(\lambda)$	relative spectral luminous efficiency, scotopic vision
$V_q(\lambda)$	relative spectral luminous efficiency for photons
λ	wavelength, nm or μm
λ_p	wavelength at peak of function, nm or μm
Φ, Φ_e	radiant power (radiant flux), watt (W)
Φ_p, Φ_q	photon flux, s^{-1}
Φ_v	luminous power (luminous flux), lumen (lm)
ω	solid angle, steradian (sr)
Ω	projected solid angle, sr

36.2 INTRODUCTION AND BACKGROUND

After more than a century of turmoil, the symbols, units, and nomenclature (SUN) for radiometry and photometry seem somewhat stable at present. There are still small, isolated pockets of debate and even resistance; by and large, though, the community has settled on the International System of Units (SI) and the recommendations of the International Organization for Standardization (ISO)¹ and the International Union of Pure and Applied Physics (IUPAP).²

The seed of the SI was planted in 1799 with the deposition of two prototype standards, the meter and the kilogram, into the Archives de la République in Paris. Noted physicists Gauss, Weber, Maxwell, and Thompson made significant contributions to measurement science over the next 75 years. Their efforts culminated in the Convention of the Meter, a diplomatic treaty originally signed by representatives of 17 nations in 1875 (currently there are 48 member nations). The Convention grants authority to the General Conference on Weights and Measures (CGPM), the International Committee for Weights and Measures (CIPM), and the International Bureau of Weights and Measures (BIPM). The CIPM, along with a number of subcommittees, suggests modifications to the CGPM. In our arena, the subcommittee is the Consultative Committee on Photometry and Radiometry (CCPR). The BIPM, the international metrology institute, is the physical facility that is responsible for realization, maintenance, and dissemination of standards.

The SI was adopted by the CGPM in 1960 and is the official system in the 48 member states. It currently consists of seven base units and a much larger number of derived units. The base units are a choice of seven well-defined units that, by convention, are regarded as independent. The seven base units are as follows:

1. Meter
2. Kilogram
3. Second
4. Ampere
5. Kelvin
6. Mole
7. Candela

The derived units are those that are formed by various combinations of the base units.

International organizations involved in the promulgation of SUN include the International Commission on Illumination (CIE), the IUPAP, and the ISO. In the United States, the American National Standards Institute (ANSI) is the primary documentary (protocol) standards organization. Many other scientific and technical organizations publish recommendations concerning the use of SUN for their scholarly journals. Several examples are the Illuminating Engineering Society

TABLE 1 Projected Areas of Common Shapes

Shape	Area	Projected area
Flat rectangle	$A = L \times W$	$A_{\text{proj}} = L \times W \cos \beta$
Circular disc	$A = \pi r^2 = \pi d^2/4$	$A_{\text{proj}} = \pi r^2 \cos \beta = (\pi d^2 \cos \beta)/4$
Sphere	$A = 4\pi r^2 = \pi d^2$	$A_{\text{proj}} = A/4 = \pi r^2$

(IESNA), the International Astronomical Union (IAU), the Institute for Electrical and Electronic Engineering (IEEE), and the American Institute of Physics (AIP).

The terminology employed in radiometry and photometry consists principally of two parts: (1) an adjective distinguishing between a radiometric, photonic, or photometric entity, and (2) a noun describing the underlying geometric or spatial concept. In some instances, the two parts are replaced by a single term (e.g., radiance).

There are some background concepts and terminology that are needed before proceeding further.

Projected area is defined as the rectilinear projection of a surface of any shape onto a plane normal to the unit vector. The differential form is $dA_{\text{proj}} = \cos(\beta)dA$, where β is the angle between the local surface normal and the line of sight. Integrate over the (observable) surface area to get

$$A_{\text{proj}} = \int_A \cos \beta dA \quad (1)$$

Some common examples are shown in Table 1.

Plane angle and solid angle are both derived units in the SI system. The following definitions are from National Institute of Standards and Technology (NIST) SP811.³

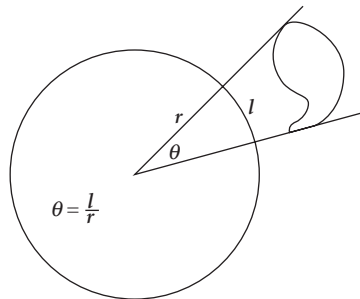
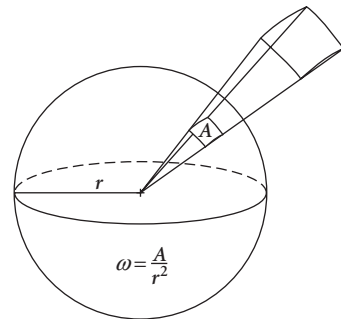
The radian is the plane angle between two radii of a circle that cuts off on the circumference an arc equal in length to the radius.

The abbreviation for the radian is *rad*. Since there are 2π rad in a circle, the conversion between degrees and radians is $1 \text{ rad} = (180/\pi)$ degrees. (See Fig. 1.)

A solid angle is the same concept that is extended to three dimensions.

One steradian (sr) is the solid angle that, having its vertex in the center of a sphere, cuts off an area on the surface of the sphere equal to that of a square with sides of length equal to the radius of the sphere.

The solid angle is the ratio of the spherical area (A_{proj}) to the square of the radius, r . The spherical area is a projection of the object of interest onto a unit sphere, and the solid angle is the surface area of that projection, as shown in Fig. 2. Divide the surface area of a sphere by the square of its radius to find that there are 4π sr of solid angle in a sphere. One hemisphere has 2π sr. The accepted symbols


FIGURE 1 Plane angle.

FIGURE 2 Solid angle.

for solid angle are either lowercase Greek omega (ω) or uppercase Greek omega (Ω). I recommend using ω exclusively for solid angle, and reserving Ω for the advanced concept of projected solid angle ($\omega \cos \theta$). The equation for solid angle is $d\omega = dA_{\text{proj}}/r^2$. For a right circular cone, $\omega = 2\pi(1 - \cos \theta)$, where θ is the half-angle of the cone.

Both plane angles and solid angles are dimensionless quantities, and they can lead to confusion when attempting dimensional analysis. For example, the simple inverse square law, $E = I/d^2$ appears dimensionally inconsistent. The left side has units W m^{-2} while the right side has $\text{W sr}^{-1} \text{m}^{-2}$. It has been suggested that this equation be written $E = I\Omega_0/d^2$, where Ω_0 is the unit solid angle, 1 sr. Inclusion of the term Ω_0 will render the equation dimensionally correct, but will far too often be considered a free variable rather than a constant equal to 1, leading to erroneous results.

Isotropic and lambertian both carry the meaning “the same in all directions” and, regrettably, are used interchangeably.

Isotropic implies a spherical source that radiates the same amount in all directions [i.e., the intensity (W/sr or cd) is independent of direction]. The term *isotropic point source* is often heard. No such entity can exist, for the energy density would necessarily be infinite. A small, uniform sphere comes very close. Small in this context means that the size of the sphere is much less than the distance from the sphere to the plane of observation, such that the inverse square law is applicable. An example is a globular tungsten lamp with a milky white diffuse envelope some 10 cm in diameter, as viewed from a distance greater than 1 m. From our vantage point, a distant star is considered an isotropic point source.

Lambertian implies a flat radiating surface. It can be an active emitter or a passive, reflective surface. The intensity decreases with the cosine of the observation angle with respect to the surface normal (Lambert’s law). The radiance ($\text{W m}^2 \text{sr}^{-1}$) or luminance (cd m^{-2}) is independent of direction. A good approximation is a surface painted with a quality matte, or flat, white paint. The intensity is the product of the radiance, L , or luminance, L_v , and the projected area A_{proj} . If the surface is uniformly illuminated, it appears equally bright from whatever direction it is viewed. Note that the flat radiating surface can be used as an elemental area of a curved surface.

The ratio of the radiant exitance (power per unit area, W m^{-2}) to the radiance (power per unit projected area per unit solid angle, $\text{W m}^2 \text{sr}^{-1}$) of a lambertian surface is a factor of π and not 2π . This result is not intuitive, as, by definition, there are 2π sr in a hemisphere. The factor of π comes from the influence of the $\cos\theta$ term while integrating over a hemisphere.

A sphere with a lambertian surface illuminated by a distant point source will display a radiance that is maximum at the point where the local normal coincides with the incoming beam. The radiance will fall off with a cosine dependence to zero at the terminator. If the intensity (integrated radiance over area) is unity when viewing from the direction of the source, then the intensity when viewing from the side is $1/\pi$. Think about this and ponder whether our Moon has a lambertian surface.

36.3 SYMBOLS, UNITS, AND NOMENCLATURE IN RADIOMETRY

Radiometry is the measurement of optical radiation, defined as electromagnetic radiation within the frequency range from 3×10^{11} to 3×10^{16} hertz (Hz). This range corresponds to wavelengths between 0.01 and 1000 micrometers (μm) and includes the regions commonly called *ultraviolet*, *visible*, and *infrared*.

Radiometric units can be divided into two conceptual areas: (1) those having to do with power or energy, and (2) those that are geometric in nature. In the first category are:

Energy is an SI-derived unit, measured in joules (J). The recommended symbol for energy is Q . An acceptable alternate is W .

Power (also called *radiant flux*) is another SI-derived unit. It is the derivative of energy with respect to time, dQ/dt , and the unit is the watt (W). The recommended symbol for power is uppercase Greek phi (Φ). An accepted alternate is P .

Energy is the integral of power over time, and it is commonly used with integrating detectors and pulsed sources. Power is used for continuous sources and nonintegrating detectors. The radiometric quantity power can now be combined with the geometric spatial quantities area and solid angle.

Irradiance (also referred to as *flux density* or *radiant incidence*) is an SI-derived unit and is measured in W m^{-2} . Irradiance is power per unit area *incident* from all directions within a hemisphere onto a surface that coincides with the base of that hemisphere. A related quantity is *radiant exitance*, which is power per unit area *leaving* a surface into a hemisphere whose base is that surface. The symbol for irradiance is E , and the symbol for radiant exitance is M . Irradiance (or radiant exitance) is the derivative of power with respect to area, $d\Phi/dA$. The integral of irradiance or radiant exitance over area is power. There is no compelling reason to have two quantities carrying the same units, but it is convenient.

Radiant intensity is an SI-derived unit and is measured in W sr^{-1} . Intensity is power per unit of solid angle. The symbol is I . *Intensity* is the derivative of power with respect to solid angle, $d\Phi/d\Omega$. The integral of radiant intensity over solid angle is power.

A great deal of confusion surrounds the use and misuse of the term *intensity*. Some use it for W sr^{-1} ; some use it for W m^{-2} ; others use it for $\text{W m}^{-2} \text{sr}^{-1}$. It is quite clearly defined in the SI system, in the definition of the base unit of luminous intensity—the candela. Attempts are often made to justify these different uses of intensity by adding adjectives like *optical* or *field* (used for W m^{-2}) or *specific* (used for $\text{W m}^{-2} \text{sr}^{-1}$). In the SI system, the underlying geometric concept for intensity is quantity per unit solid angle. For more discussion, see Palmer.⁴

Radiance is an SI-derived unit and is measured in $\text{W m}^{-2} \text{sr}^{-1}$. Radiance is a directional quantity, power per unit projected area per unit solid angle. The symbol is L . Radiance is the derivative of power with respect to solid angle and projected area, $d\Phi/d\omega dA \cos\theta$, where θ is the angle between the surface normal and the specified direction. The integral of radiance over area and solid angle is power.

Photon quantities are also common. They are related to the radiometric quantities by the relationship $Q_p = hc/\lambda$, where Q_p is the energy of a photon at wavelength λ , h is Planck's constant, and c is the velocity of light. At a wavelength of $1 \mu\text{m}$, there are approximately 5×10^{18} photons per second in a watt. Conversely, one photon has an energy of about $2 \times 10^{-19} \text{ J (W/s)}$ at $1 \mu\text{m}$.

36.4 SYMBOLS, UNITS, AND NOMENCLATURE IN PHOTOMETRY

Photometry is the measurement of light, electromagnetic radiation detectable by the human eye. It is thus restricted to the wavelength range from about 360 to 830 nanometers (nm; $1000 \text{ nm} = 1 \mu\text{m}$). Photometry is identical to radiometry *except* that everything is weighted by the spectral response of the nominal human eye. *Visual photometry* uses the eye as a comparison detector, while *physical photometry* uses either optical radiation detectors constructed to mimic the spectral response of the nominal eye, or spectroradiometry coupled with appropriate calculations to do the eye response weighting.

Photometric units are basically the same as the radiometric units, except that they are weighted for the spectral response of the human eye and have strange names. A few additional units have been introduced to deal with the amount of light that is reflected from diffuse (matte) surfaces. The symbols used are identical to the geometrically equivalent radiometric symbols, except that a subscript v is added to denote *visual*. Table 2 compares radiometric and photometric units.

The SI unit for light is the *candela* (unit of luminous intensity). It is one of the seven base units of the SI system. The candela is defined as follows:⁵

The candela is the luminous intensity, in a given direction, of a source that emits monochromatic radiation of frequency 540×10^{12} hertz and that has a radiant intensity in that direction of $1/683$ watt per steradian.

The candela is abbreviated as *cd*, and its symbol is I_v . This definition was adopted by the 16th CGPM in 1979. The candela was formerly defined as the luminous intensity, in the perpendicular direction,

TABLE 2 Comparison of Radiometric and Photometric Units

Quantity	Radiometric	Photometric
Power	Φ : watt (W)	Φ_v : lumen (lm)
Power per area	E, M : $W\ m^{-2}$	E_v : $lm\ m^{-2} = lux$ (lx)
Power per solid angle	I : $W\ sr^{-1}$	I_v : $lm\ sr^{-1} = candela$ (cd)
Power per area per solid angle	L : $W\ m^{-2}\ sr^{-1}$	L_v : $lm\ m^{-2}\ sr^{-1} = cd\ m^{-2} = nit$

of a surface of $1/600,000\ m^2$ of a blackbody at the temperature of freezing platinum under a pressure of 101,325 newtons per square meter ($N\ m^{-2}$). This earlier definition was initially adopted in 1948 and later modified by the 13th CGPM in 1968. It was abrogated in 1979 and replaced by the current definition.

The 1979 definition was adopted for several reasons.^{6–10} First, the realization of the candela using a platinum blackbody was extraordinarily difficult—only several were ever built, and there were large variations between the units realized by different national laboratories based upon the state of platinum at its freezing point. The difficulty in fabricating and operating the platinum point blackbody created an unacceptable uncertainty in the value of the candela. For example, if the platinum blackbody temperature is slightly off, possibly because of temperature gradients in the ceramic crucible or contamination of the platinum, the freezing point may change or the temperature of the cavity may differ. The sensitivity of the candela to a slight change in temperature is significant. At a wavelength of 555 nm, a change in temperature of only 1 K results in a luminance change approaching 1 percent. Second, the unit of the candela was realized on the specific broadband radiation, whose spectral power distribution was not known with satisfactory accuracy (because the platinum fix point temperature was not precisely known), thus there were large uncertainties in determining photometric quantities of various other practical light sources from their spectral power distributions. Third, recent advances in radiometry based on absolute radiometers offered new possibilities for realization of the candela using a much simpler device with much lower uncertainties if the candela is defined in relation to watt. In 1977, through an international comparison among several national laboratories, the Comité International des Poids et Mesures (CIPM) determined the numerical relationship (683 lm/W at 555 nm) to be recommended for the new standard for candela so that the magnitude of the unit was kept consistent with the previous unit of the candela.

The value of 683 lm/W was selected based upon the best measurements with existing platinum freezing point blackbodies at several national standards laboratories. It has varied over time from 620 to nearly 700 lm/W, depending largely upon the assigned value of the freezing point of platinum. The value of $1/600,000\ m^2$ was chosen to maintain consistency with prior standards. Note that neither the old nor the new definition of the candela say anything about the spectral responsivity of the human eye. There are additional definitions that include the characteristics of the eye, but the base unit (candela) and those SI units derived from it are “eyeless.”

Note also that in the definition of the candela, there is no specification for the spatial distribution of intensity. Luminous intensity, while often associated with an isotropic point (i.e., small) source, is a valid specification for characterizing any highly directional light source, such as a spotlight or an LED.

One other issue: since the candela is no longer independent but is now defined in terms of other SI-derived quantities, there is really no need to retain it as an SI base quantity. It remains so for reasons of history and continuity and perhaps some politics.

The *lumen* is an SI-derived unit for luminous flux (power). The abbreviation is *lm*, and the symbol is Φ_v . The lumen is derived from the candela and is the luminous flux that is emitted into unit solid angle (1 sr) by an isotropic point source having a luminous intensity of 1 cd. The lumen is the product of luminous intensity and solid angle (cd sr). It is analogous to the unit of radiant flux (watt), differing only in the eye response weighting. If a light source is isotropic, the relationship between lumens and candelas is $1\ cd = 4\pi\ lm$. In other words, an isotropic source that has a luminous intensity of 1 cd emits $4\pi\ lm$ into space, which is $4\pi\ sr$. Also, $1\ cd = 1\ lm\ sr^{-1}$, which is analogous to the equivalent radiometric definition.

If a source is not isotropic, the relationship between candelas and lumens is empirical. A fundamental method used to determine the total flux (lumens) is to measure the luminous intensity (candelas) in many directions using a goniophotometer, and then numerically integrate over the

entire sphere. Later on, this “calibrated” lamp can be used as a reference in an integrating sphere for routine measurements of luminous flux.

The SI-derived unit of luminous flux density, or illuminance, has a special name: *lux*. It is lumens per square meter (lm m^{-2}), and the symbol is E_v . Most light meters measure this quantity, as it is of great importance in illuminating engineering. The IESNA’s *Lighting Handbook*¹¹ has some 16 pages of recommended illuminances for various activities and locales, ranging from morgues to museums. Typical values range from 100,000 lx for direct sunlight to between 20 and 50 lx for hospital corridors at night.

Luminance should probably be included on the list of SI-derived units, but it is not. Luminance is analogous to radiance, giving the spatial and directional dependences. It also has a special name, *nit*, and is candelas per square meter (cd m^{-2}) or lumens per square meter per steradian ($\text{lm m}^{-2} \text{sr}^{-1}$). The symbol is L_v . Luminance is most often used to characterize the “brightness” of flat-emitting or -reflecting surfaces. A common use is the luminance of a laptop computer screen. They typically have between 100 and 250 nits, and the sunlight-readable ones have more than 1000 nits. Typical CRT monitors have luminances between 50 and 125 nits.

Other Photometric Units

There are other photometric units, largely historical. The literature is filled with now obsolete terminology, and it is important to be able to properly interpret these terms. Here are several terms for illuminance that have been used in the past.

$$\begin{aligned} 1 \text{ meter-candle} &= 1 \text{ lx} \\ 1 \text{ phot (ph)} &= 1 \text{ lm cm}^{-2} = 10^4 \text{ lx} \\ 1 \text{ footcandle (fc)} &= 1 \text{ lm ft}^{-2} = 10.76 \text{ lx} \\ 1 \text{ milliphot} &= 10 \text{ lx} \end{aligned}$$

Table 3 is useful to convert from one unit to another. Start with the unit in the leftmost column and multiply it by the factor in the table to arrive at the unit in the top row.

There are two classes of units that are used for luminance. The first is conventional, directly related to the SI unit, the cd m^{-2} (nit).

$$\begin{aligned} 1 \text{ stilb} &= 1 \text{ cd cm}^{-2} = 10^4 \text{ cd m}^{-2} = 10^4 \text{ nit} \\ 1 \text{ cd ft}^{-2} &= 10.76 \text{ cd m}^{-2} = 10.76 \text{ nit} \end{aligned}$$

The second class was designed to “simplify” characterization of light that is reflected from diffuse surfaces by incorporating within the definition the concept of a perfect diffuse reflector (lambertian, reflectance $\rho = 1$). If 1 unit of illuminance falls upon this ideal reflector, then 1 unit of luminance is reflected. The perfect diffuse reflector emits $1/\pi$ units of luminance per unit of illuminance. If the reflectance is ρ , then the luminance is ρ/π times the illuminance. Consequently, these units all incorporate a factor of $1/\pi$.

$$\begin{aligned} 1 \text{ lambert (L)} &= (1/\pi) \text{ cd cm}^{-2} = (10^4/\pi) \text{ cd m}^{-2} = (10^4/\pi) \text{ nit} \\ 1 \text{ apostilb} &= (1/\pi) \text{ cd m}^{-2} = (1/\pi) \text{ nit} \\ 1 \text{ foot-lambert (ft-lambert)} &= (1/\pi) \text{ cd ft}^{-2} = 3.426 \text{ cd m}^{-2} = 3.426 \text{ nit} \\ 1 \text{ millilambert} &= (10/\pi) \text{ cd m}^{-2} = (10/\pi) \text{ nit} \\ 1 \text{ skot} &= 1 \text{ milliblondel} = (10^{-3}/\pi) \text{ cd m}^{-2} = 10^{-3}/\pi \text{ nit} \end{aligned}$$

TABLE 3 Illuminance Unit Conversions

	fc	lx	phot	milliphot
1 fc (lm/ft^2) =	1	10.764	0.0010764	1.0764
1 lx (lm/m^2) =	0.0929	1	0.0001	0.1
1 phot (lm/cm^2) =	929	10,000	1	0.001
1 milliphot =	0.929	10	0.1	1

TABLE 4 Illuminance Unit Conversions*

	nit	stilb	cd/ft ²	apostilb	lambert	ft-lambert
1 nit (cd/m ²) =	1	10 ⁻⁴	0.0929	π	$\pi/10000$	0.0929 π
1 stilb (cd/cm ²) =	10,000	1	929	10 ⁴ π	π	929 π
1 cd/ft ² =	10.764	1.0764 × 10 ⁻³	1	10.764 π	$\pi/929$	π
1 apostilb =	1/ π	10 ⁴ / π	0.0929/ π	1	10 ⁻⁴	0.0929
1 lambert =	10 ⁴ / π	1/ π	929/ π	10 ⁴	1	929
1 ft · lambert =	10.76/ π	1/(929 π)	1/ π	10.764	1.076 × 10 ⁴	1

*Note: Photometric quantities are the result of an integration over wavelength. It therefore makes no sense to speak of spectral luminance or the like.

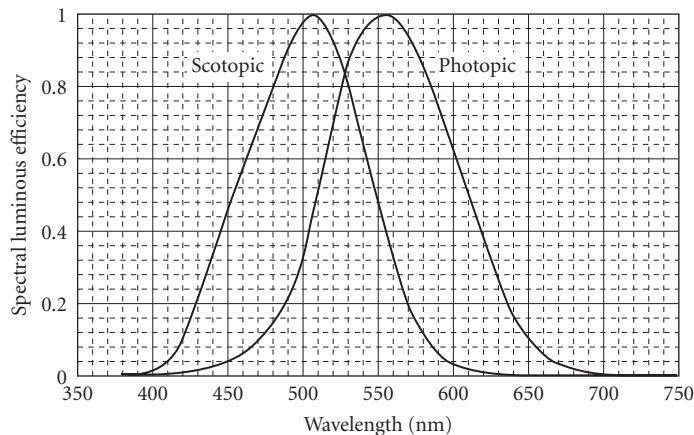
Table 4 is useful to convert from one unit to another. Start with the unit in the leftmost column and multiply it by the factor in the table to arrive at the unit in the top row.

Human Eye The SI base unit and units derived therefrom have a strictly physical basis; they have been defined monochromatically at a wavelength of 555 nm. But the eye does not see all wavelengths equally. For other wavelengths or for band or continuous-source spectral distributions, the spectral properties of the human eye must be considered. The eye has two general classes of photosensors: *cones* and *rods*.

Cones The cones are responsible for light-adapted vision; they respond to color and have high resolution in the central foveal region. The light-adapted relative spectral response of the eye is called the *spectral luminous efficiency function for photopic vision*, $V(\lambda)$, and is published in tabular form.¹² This empirical curve, shown in Fig. 3, was first adopted by the CIE in 1924. It has a peak that is normalized to unity at 555 nm, and it decreases to levels below 10⁻⁵ at about 370 and 785 nm. The 50 percent points are near 510 and 610 nm, indicating that the curve is slightly skewed. A logarithmic representation is shown in Fig. 4.

More recent measurements have shown that the 1924 curve may not best represent typical human vision. It appears to underestimate the response at wavelengths shorter than 460 nm. Judd,¹³ Vos,¹⁴ and Stockman and Sharpe¹⁵ have made incremental advances in our knowledge of the photopic response.

Rods The rods are responsible for dark-adapted vision, with no color information and poor resolution when compared with the foveal cones. The dark-adapted relative spectral response of the

**FIGURE 3** Spectral luminous efficiency for photopic and scotopic vision.

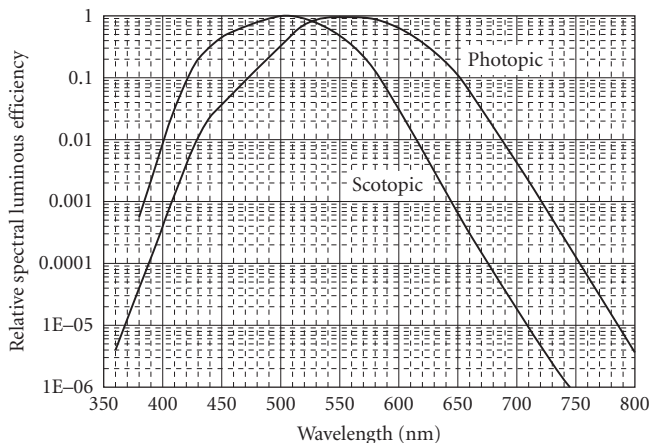


FIGURE 4 Spectral luminous efficiency for photopic and scotopic vision (log scale).

eye is called the *spectral luminous efficiency function for scotopic vision*, $V(\lambda)$, also published in tabular form.¹² Figures 3 and 4 also show this empirical curve, which was adopted by the CIE in 1951. It is defined between 380 and 780 nm. The $V(\lambda)$ curve has a peak of unity at 507 nm, and it decreases to levels below 10^{-3} at about 380 and 645 nm. The 50 percent points are near 455 and 550 nm.

Photopic (light-adapted cone) vision is active for luminances that are greater than 3 cd m^{-2} . Scotopic (dark-adapted rod) vision is active for luminances that are lower than 0.01 cd m^{-2} . In between, both rods and cones contribute in varying amounts, and in this range the vision is called *mesopic*. There have been efforts to characterize the composite spectral response in the mesopic range for vision research at intermediate luminance levels. Definitive values at 1-nm intervals for both photopic and scotopic spectral luminous efficiency functions may be found in CIE.¹² Values at 5-nm intervals are given by Zalewski.¹⁶

The relative spectral luminous efficiency curves can be converted for use with photon flux (s^{-1}) by multiplying by the spectrally dependent conversion from watts to photons per second. The results are shown in Fig. 5. The curves are similar to the spectral luminous efficiency curves, with the peaks shifted to slightly shorter wavelengths, and the skewness of the curves is different. This function

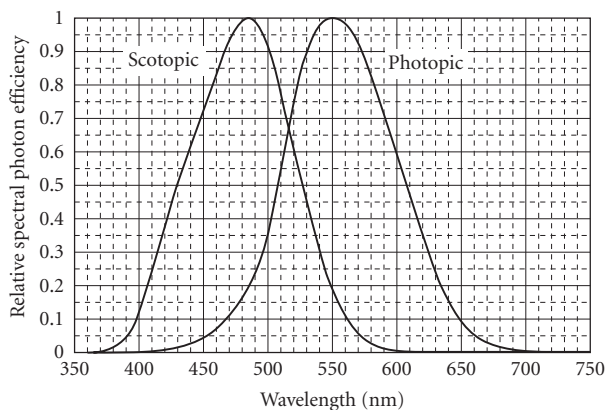


FIGURE 5 Spectral photon efficiency for scotopic and photopic vision.

can be called $V_q(\lambda)$ for photopic response or $V'_q(\lambda)$ for scotopic response. The conversion to an absolute curve is made by multiplying by the response at the peak wavelength. For photopic vision ($\lambda = 550$ nm), $K_{mp} = 2.45 \times 10^{-16}$ lm/photon s^{-1} . There are, therefore, 4.082×10^{15} photon s^{-1} lm^{-1} at 550 nm, and more at all other wavelengths. For scotopic vision ($\lambda_p = 504$ nm), $K'_{mp} = 6.68 \times 10^{-16}$ lm/photon s^{-1} . There are 1.497×10^{15} photon s^{-1} lm^{-1} at 504 nm, and more at all other wavelengths.

Approximations The $V(\lambda)$ curve appears similar to a Gaussian (normal) function. A nonlinear regression technique was used to fit the Gaussian shown in Eq. (2) to the $V(\lambda)$ data

$$V(\lambda) \cong 1.019e^{-285.4(\lambda - 0.559)^2} \quad (2)$$

The scotopic curve can also be fit with a Gaussian, although the fit is not quite as good as the photopic curve. My best fit is

$$V'(\lambda) \cong 0.992e^{-321.0(\lambda - 0.503)^2} \quad (3)$$

The results of the curve fitting are shown in Figs. 6 and 7. These approximations are satisfactory for application with continuous spectral distributions, such as sunlight, daylight, and incandescent sources. Calculations have demonstrated errors of less than 1 percent with blackbody sources from 1500 K to more than 20,000 K. The equations must be used with caution for narrow-band or line sources, particularly in those spectral regions where the response is low and the fit is poor.

Usage The SI definition of the candela was chosen in strictly physical terms at a single wavelength. The intent of photometry, however, is to correlate a photometric observation to the visual perception of a human observer. The CIE introduced the two standard spectral luminous efficiency functions $V(\lambda)$ (photopic) and $V'(\lambda)$ (scotopic) as spectral weighting functions, and they have been approved by the CIPM for use with light sources at other wavelengths. Another useful function is the CIE $V_M(\lambda)$ Judd-Vos modified $V(\lambda)$ function,¹⁴ which has increased response at wavelengths that are shorter than 460 nm. It is identical to the $V(\lambda)$ function for wavelengths that are longer than 460 nm. This function, while not approved by CIPM, represents more realistically the spectral responsivity of the eye. More recently, studies on cone responses have led to the proposal of a new, improved luminous spectral efficiency curve, with the suggested designation $V_2^*(\lambda)$.¹⁵

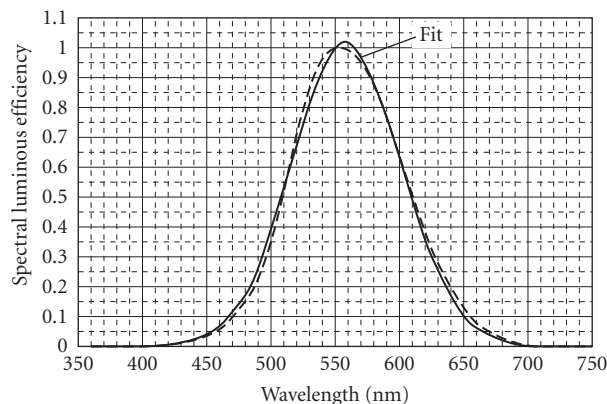


FIGURE 6 Gaussian fit to photopic relative spectral efficiency curve.

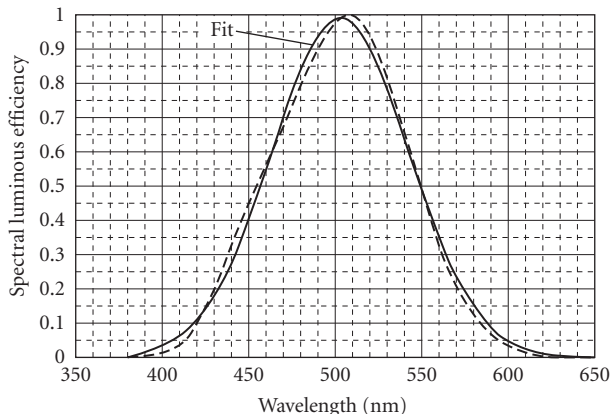


FIGURE 7 Gaussian fit to scotopic relative spectral efficiency curve.

36.5 CONVERSION OF RADIOMETRIC QUANTITIES TO PHOTOMETRIC QUANTITIES

The definition of the candela states that there are 683 lm W^{-1} at a frequency of 540 terahertz (THz), which is very nearly 555 nm (vacuum or air), the wavelength that corresponds to the maximum spectral responsivity of the photopic (light-adapted) human eye. The value 683 lm W^{-1} is K_m , the absolute luminous efficiency at λ_p for photopic vision. The conversion from watts to lumens at any other wavelength involves the product of the power (watts), K_m , and the $V(\lambda)$ value at the wavelength of interest. For example, a 5-mW laser pointer has $0.005 \text{ W} \times 0.032 \times 683 \text{ lm W}^{-1} = 0.11 \text{ lm}$. $V(\lambda)$ is 0.032 at 670 nm. At 635 nm, $V(\lambda)$ is 0.217, and a 5-mW laser pointer has $0.005 \text{ W} \times 0.217 \times 683 \text{ lm W}^{-1} = 0.74 \text{ lm}$. The shorter-wavelength laser pointer will create a spot that has nearly seven times the luminous power as the longer-wavelength laser.

Similar calculations can be done in terms of photon flux at a single wavelength. As was shown previously, there are $2.45 \times 10^{16} \text{ lm}$ in 1 photon s^{-1} at 550 nm, the wavelength that corresponds to the maximum spectral responsivity of the light-adapted human eye to photon flux. The conversion from lumens to photons per second at any other wavelength involves the product of the photon flux (s^{-1}) and the $V_p(\lambda)$ value at the wavelength of interest. For example, again compare laser pointers at 670 and 635 nm. As shown before, a 5-mW laser at 670 nm [$V_p(\lambda) = 0.0264$] has a luminous power of 0.11 lm. The conversion is $0.11 \times 4.082 \times 10^{15} / 0.0264 = 1.68 \times 10^{16} \text{ photon s}^{-1}$. At 635 nm [$V_p(\lambda) = 0.189$], the 5-mW laser has a luminous power of 0.74 lm. The conversion is $0.74 \times 4.082 \times 10^{15} / 0.189 = 1.6 \times 10^{16} \text{ photon s}^{-1}$. The 635-nm laser delivers just 5 percent more photons per second.

In order to convert a source with nonmonochromatic spectral distribution to a luminous quantity, the situation is decidedly more complex. The spectral nature of the source must be known, as it is used in an equation of the form

$$X_v = K_m \int_0^{\infty} X_\lambda V(\lambda) d\lambda \quad (4)$$

where X_v is a luminous term, X_λ is the corresponding spectral radiant term, and $V(\lambda)$ is the photopic spectral luminous efficiency function. For X_v , luminous flux (lm) may be paired with spectral power (W nm^{-1}), luminous intensity (cd) with spectral radiant intensity ($\text{W sr}^{-1} \text{ nm}^{-1}$), illuminance (lx) with spectral irradiance ($\text{W m}^{-2} \text{ nm}^{-1}$), or luminance (cd m^{-2}) with spectral radiance ($\text{W m}^{-2} \text{ sr}^{-1} \text{ nm}^{-1}$). This equation represents a weighting, wavelength by wavelength, of the radiant spectral term

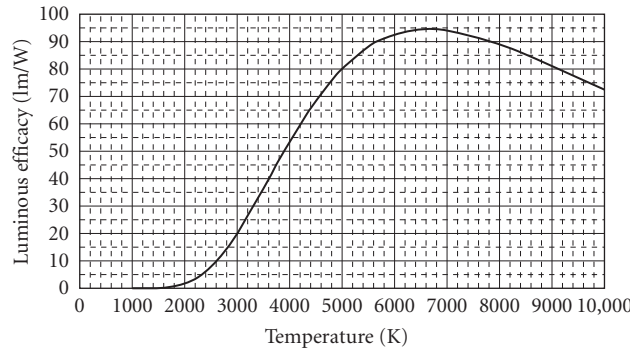


FIGURE 8 Luminous efficacy of blackbody radiation versus temperature (K).

by the visual response at that wavelength. The constant K_m is the maximum spectral luminous efficiency for photopic vision, 683 lm W^{-1} . The wavelength limits can be set to restrict the integration to only those wavelengths where the product of the spectral term X_λ and $V(\lambda)$ is nonzero. Practically, the limits of integration need only extend from 360 to 830 nm, limits specified by the CIE $V(\lambda)$ function. Since this $V(\lambda)$ function is defined by a table of empirical values,¹² it is best to do the integration numerically. Use of the Gaussian equation [Eq. (2)] is only an approximation.

For source spectral distributions that are blackbody-like (thermal source, spectral emissivity constant between 360 and 830 nm) and of known source temperature, it is straightforward to convert from power to luminous flux and vice versa. Equation (4) is used to determine a scale factor for the source term X_λ . Figure 8 shows the relationship between total power and luminous flux for blackbody (and graybody) radiation as a function of blackbody temperature. The most efficient temperature for the production of luminous flux is near 6630 K.

There is nothing in the SI definitions of the base or derived units concerning the eye response, so there is some flexibility in the choice of the weighting function. The choice can be made to use a different spectral luminous efficiency curve, perhaps one of the newer ones. The equivalent curve for scotopic (dark-adapted) vision can also be used for work at lower light levels. The $V'(\lambda)$ curve has its own constant, K'_m , the maximum spectral luminous efficiency for scotopic vision. K'_m is defined as 1700 lm/W at the peak wavelength for scotopic vision (507 nm). This value was deliberately chosen such that the absolute value of the scotopic curve at 555 nm coincides with the photopic curve, 683 lm/W at 555 nm. Some researchers are referring to “scotopic lumens,” a term that should be discouraged because of the potential for misunderstanding. In the future, expect to see spectral weighting to represent the mesopic region as well.

The CGPM has approved the use of the CIE $V(\lambda)$ and $V'(\lambda)$ curves for determination of the value of photometric quantities of luminous sources.

36.6 CONVERSION OF PHOTOMETRIC QUANTITIES TO RADIOMETRIC QUANTITIES

The conversion from watts to lumens in the previous section required only that the spectral function, X_λ , of the radiation be known over the spectral range from 360 to 830 nm, where $V(\lambda)$ is nonzero. Attempts to go in the other direction, from lumens to watts, are far more difficult. Since the desired quantity was inside of an integral, weighted by a sensor spectral responsivity function, the spectral function, X_λ , of the radiation must be known over the entire spectral range where the source emits, not just the visible.

For a monochromatic source in the visible spectrum (between the wavelengths of 380 and 860 nm), if the photometric quantity (e.g., lux) is known, apply the conversion $K_m \times V(\lambda)$ and determine the

radiometric quantity (e.g., $W\ m^{-2}$). In practice, the results that one obtains are governed by the quality of the $V(\lambda)$ correction of the photometer and the knowledge of the wavelength of the source. Both of these factors are of extreme importance at wavelengths where the $V(\lambda)$ curve is steep (i.e., other than very close to the peak of the $V(\lambda)$ curve).

Narrowband sources, such as LEDs, cause major problems. Typical LEDs have spectral bandwidths ranging from 10- to 40-nm full width at half-maximum (FWHM). It is intuitive that in those spectral regions where the $V(\lambda)$ curve is steep, the luminous output will be greater than that predicted using the $V(\lambda)$ curve at the peak LED wavelength. This expected result increases with wider-bandwidth LEDs. Similarly, it is also intuitive that the luminous output is less than that predicted by using the $V(\lambda)$ curve when the peak LED wavelength is in the vicinity of the peak of the $V(\lambda)$ curve. Therefore, there must be two wavelengths where the conversion ratio (lm/W) is largely independent of LED bandwidth. An analysis of this conversion ratio was done using a Gaussian equation to represent the spectral power distribution of an LED and applying Eq. (4). Indeed, two null wavelengths were identified (513 and 604 nm) where the conversion between radiometric and photometric quantities is constant (independent of LED bandwidth) to within 0.2 percent up to an LED bandwidth of 40 nm. These wavelengths correspond (approximately) to the wavelengths where the two maxima of the first derivative of the $V(\lambda)$ curve are located.

At wavelengths between these two null wavelengths (513 and 604 nm), the conversion ratio (lm/W) decreases slightly with increasing bandwidth. The worst case occurs when the peak wavelength of the LED corresponds with the peak of $V(\lambda)$. It is about 5 percent lower for bandwidths up to 30 nm, increasing to nearly 10 percent for 40-nm bandwidth. At wavelengths outside of the null wavelengths, the conversion (lm/W) increases with increasing bandwidth, and the increase is greater when the wavelength approaches the limits of the $V(\lambda)$ curve. Figure 9 shows that factor by which a conversion ratio (lm/W) should be multiplied as a function of LED bandwidth, with the peak LED wavelength as the parameter. Note that the peak wavelength of the LED is specified as the radiometric peak and not as the dominant wavelength (a color specification). The dominant wavelength shifts with respect to the radiometric peak, the difference increasing with bandwidth.

Most often, LEDs are specified in luminous intensity [cd or millicandela (mcd)]. The corresponding radiometric unit is watts per steradian (W/sr). In order to determine the radiometric power (watts), the details of the spatial distribution of the radiant intensity must be known prior to integration.

For broadband sources, a photometric quantity cannot in general be converted to a radiometric quantity unless the radiometric source function, Φ_λ , is known over all wavelengths. However, if the source spectral distribution is blackbody-like (thermal source, spectral emissivity constant between 360 and 830 nm), and the source temperature is also known, then an illuminance can be converted to a spectral irradiance curve over that wavelength range. Again, use Eq. (4) to determine a scale factor

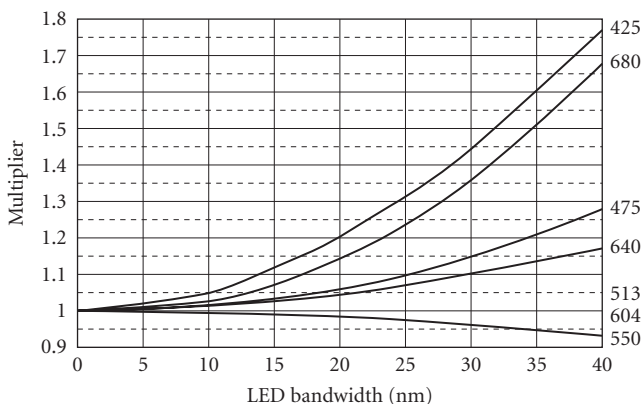


FIGURE 9 Multiplier for converting LED luminous intensity to radiometric intensity.

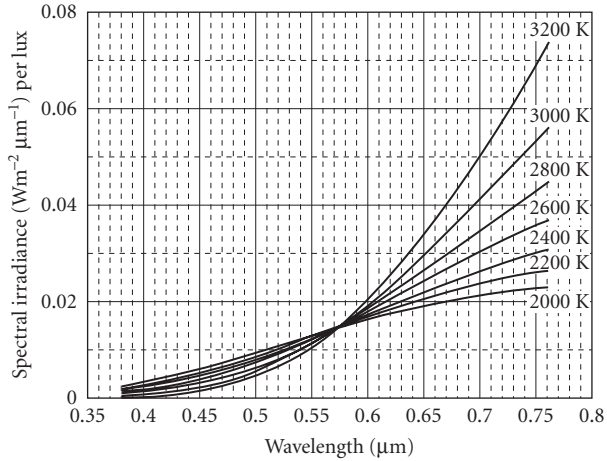


FIGURE 10 Spectral irradiance versus wavelength of blackbody radiation versus temperature.

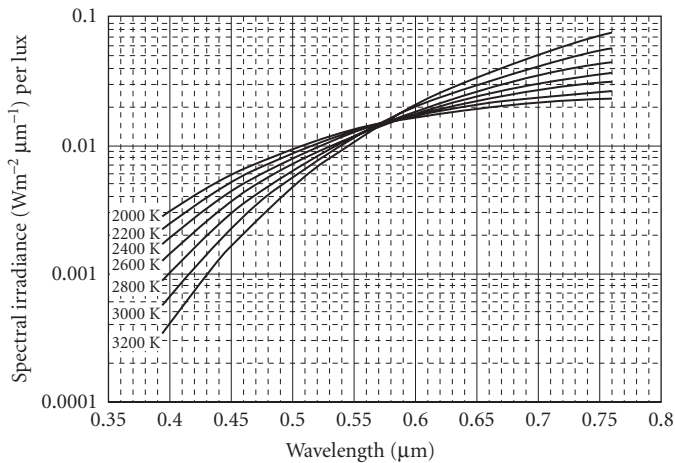


FIGURE 11 Spectral irradiance versus wavelength of blackbody radiation versus temperature.

for the source term, X_λ . Figures 10 and 11 show the calculated spectral irradiance versus wavelength for blackbody radiation with source temperature as the parameter.

36.7 RADIOMETRIC/PHOTOMETRIC NORMALIZATION

In radiometry and photometry, there are two mathematical activities that fall into the category called *normalization*. The first is bandwidth normalization. The second doesn't have an official title but involves a conversion between photometric and radiometric quantities, used to convert detector responsivities from amperes per lumen to amperes per watt.

The measurement equation relates the output signal from a sensor to its spectral responsivity and the spectral power that is incident upon it:

$$S = \int_0^{\infty} \Phi_{\lambda} \mathfrak{R}(\lambda) d\lambda \quad (5)$$

An extended discourse on bandwidth normalization¹⁷ showed that the spectral responsivity of a sensor (or detector) can be manipulated to yield more source information than is immediately apparent from the measurement equation. The sensor weights the input spectral power according to its spectral responsivity, $\mathfrak{R}(\lambda)$, such that only a limited amount of information about the source can be deduced. If either the source or the sensor has a sufficiently small bandwidth such that the spectral function of the other term does not change significantly over the passband, the equation simplifies to

$$S = \Phi_{\lambda} \cdot \mathfrak{R}(\lambda) \cdot \Delta\lambda \quad (6)$$

where $\Delta\lambda$ is the passband. Spectroradiometry and multifilter radiometry, using narrow-bandpass filters, take advantage of this simplified equation. For those cases where the passband is larger, the techniques of bandwidth normalization can be used. The idea is to substitute for $\mathfrak{R}(\lambda)$ an equivalent response that has a uniform spectral responsivity, \mathfrak{R}_n , between wavelength limits λ_1 and λ_2 and zero response elsewhere. Then, the signal is given by

$$S = \mathfrak{R}_n \int_{\lambda_1}^{\lambda_2} \Phi_{\lambda} d\lambda \quad (7)$$

and now the integrated power between wavelengths λ_1 and λ_2 is determined. There are many ways of assigning values for λ_1 , λ_2 , and \mathfrak{R}_n for a sensor. Some of the more popular methods were described by Nicodemus¹⁷ and Palmer.¹⁸ An effective choice is known as the *moments method*,¹⁹ an analysis of the zeroth, first, and second moments of the sensor spectral responsivity curve. The derivation of this normalization scheme involves the assumption that the source function is exactly represented by a second-degree polynomial. If this condition is met, the moments method of determining sensor parameters yields exact results for the source integral. In addition, the results are completely independent of the source function. The errors encountered are related to deviation of the source function from the said second-degree polynomial.

Moments normalization has been applied to the photopic spectral luminous efficiency function, $V(\lambda)$, and the results are given in Table 5 and shown in Fig. 12. These values indicate the skewed nature of the photopic and scotopic curves as the deviation from the centroid and the peak wavelengths. The results can be applied to most continuous sources, like blackbody and tungsten radiation, which are both continuous across the visible spectrum. To demonstrate the effectiveness of moments normalization, the blackbody curve was multiplied by the $V(\lambda)$ curve for temperatures ranging from 1000 to 20,000 K to determine a photometric function [e.g., lumens per square meter (or lux)]. Then, the blackbody curve was integrated over the wavelength interval between λ_1 and λ_2 to determine the equivalent (integrated between λ_1 and λ_2) radiometric

TABLE 5 Bandwidth Normalization on Spectral Luminous Efficiency

	Photopic	Scotopic
Peak wavelength (λ_p)	555 nm	507 nm
Centroid wavelength (λ_c)	560.19 nm	502.40 nm
Short wavelength (λ_1)	487.57 nm	436.88 nm
Long wavelength (λ_2)	632.81 nm	567.93 nm
Moments bandwidth	145.24 nm	131.05 nm
Normalized \mathfrak{R}_n	0.7357	0.7407
Absolute \mathfrak{R}_n	502.4 lm/W	1260 lm/W

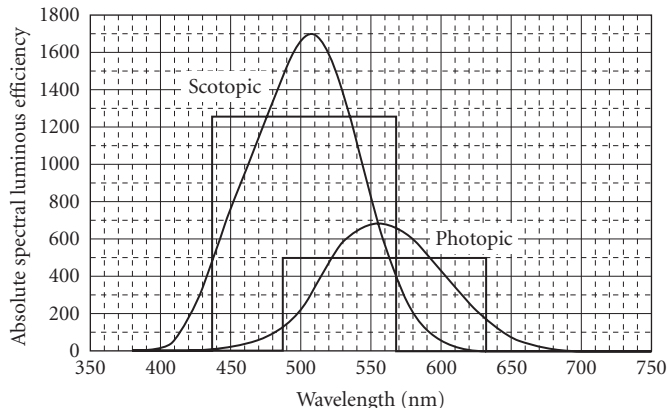


FIGURE 12 Absolute spectral luminous efficiency functions and equivalent normalized bandwidths

function (e.g., in-band watts per square meter). The ratio of lux to watt per square meter is 502.4 ± 1.0 (3σ) over the temperature range from 1600 to more than 20,000 K. This means that the in-band (487.6 to 632.8 nm) irradiance for a continuous blackbody-like source can be determined using a photometer that is properly calibrated in lux and that is well corrected for $V(\lambda)$. Simply divide the reading in lux by 502.4 to get the in-band irradiance in watts per square meter between 487.6 and 632.8 nm.

If a photometer is available with a $V'(\lambda)$ correction, calibrated for lux (scotopic) with K'_m of 1700 lm/W, a similar procedure is effective. The integration takes place over the wavelength range from 436.9 to 567.9 nm. The ratio of lux (scotopic) to watts per square meter is 1260 ± 2 (3σ) over the temperature range from 1800 K to more than 20,000 K. This means that the in-band (436.9 to 567.9 nm) irradiance for a continuous blackbody-like source can be determined using a $V'(\lambda)$ -corrected photometer that is properly calibrated in lux (scotopic). Simply divide the reading in lux (scotopic) by 1260; the result is the in-band irradiance in watts per square meter between 436.9 and 567.9 nm.

A common problem is the interpretation of specifications for photodetectors, which are given in photometric units. An example is a photomultiplier with an S-20 photocathode, which has a typical responsivity of 200 $\mu\text{A}/\text{lm}$. Given this specification and a curve of the relative spectral responsivity, the problem is to determine the output when exposed to a known power from an arbitrary source.

Photosensitive devices, in particular vacuum photodiodes and photomultiplier tubes, are characterized using CIE Illuminant A, a tungsten source at 2854 K color temperature. The illuminance is measured using a photooptically corrected photometer, and this illuminance is applied to the device under scrutiny. This technique is satisfactory only if the source being used is spectrally comparable to Illuminant A. If a source with a different spectral distribution is used, a photometric normalization must be done. Eberhart²⁰ generated a series of conversion factors for various sources and standardized detector spectral responsivities (S-1, S-11, S-20, etc.).

The luminous flux from any source is given by

$$\Phi_v = K_m \int_{360}^{830} \Phi_\lambda V(\lambda) d\lambda \quad (8)$$

and the output of a detector when exposed to the said luminous flux is

$$S = \int_0^\infty \Phi_\lambda \mathfrak{R}(\lambda) d\lambda \quad (9)$$

where Φ_λ is spectral radiant flux, $V(\lambda)$ is the spectral luminous efficiency of the photopic eye, $\mathfrak{R}(\lambda)$ is the absolute spectral responsivity of the photodetector, and S is the photodetector signal. The luminous responsivity of the detector when exposed to this source is

$$\mathfrak{R} = \frac{\int_0^\infty \Phi_\lambda \mathfrak{R}(\lambda) d\lambda}{K_m \int_{360}^{830} \Phi_\lambda V(\lambda) d\lambda} \text{ A/lm} \quad (10)$$

This luminous responsivity is specific to the source that is used to make the measurement, and it cannot be applied to other sources with differing spectral distributions.

36.8 OTHER WEIGHTING FUNCTIONS AND CONVERSIONS

The general principles that are outlined here can be applied to action spectra other than those already defined for the human eye (photopic and scotopic). Some action spectra take the form of defined spectral regions, such as UVA (315–400 nm), UVB (280–315 nm), UVC (100–280 nm), IR-A (770–1400 nm), IR-B (1400–3000 nm), and IR-C (3000–10⁶ nm). Others are more specific. $A(\lambda)$ is for aphakic hazard, $B(\lambda)$ is for photochemical blue-light hazard, $R(\lambda)$ is for retinal thermal hazard, and $S(\lambda)$ is an actinic ultraviolet action spectrum.²¹ PPF (a.k.a. PhAR) is a general action spectrum for plant growth. Many others have been defined, including those for erythema (sunburn), skin cancer, psoriasis treatment, mammalian and insect vision, and other generalized chemical and biological photoeffects. Various conversion factors from one action spectrum to another are scattered throughout the popular and archival literature. Ideally, they have been derived via integration over the appropriate spectral regions.

36.9 REFERENCES

1. ISO, "Quantities and Units—Part 6. Light and Related Electromagnetic Radiations," *ISO Standards Handbook, Quantities and Units*, 389.15 (1993).
2. IUPAP, *Symbols, Units, Nomenclature and Fundamental Constants in Physics*, prepared by E. Richard Cohan and Pierre Giacomo, Document IUPAP-25 (SUNAMCO 87-1), International Union of Pure and Applied Physics, 1987.
3. NIST, *Guide for the Use of the International System of Units (SI)*, prepared by Barry N. Taylor, NIST Special Publication SP811 (1995). (Available in PDF format from <http://physics.nist.gov/cuu/>.)
4. J. M. Palmer, "Getting Intense on Intensity," *Metrologia* **30**:371 (1993).
5. BIPM (Bureau International des Poids et Mesures), *The International System of Units (SI)*, 7th ed., 1998.
6. <http://www.bipm.org/en/CGPM/db/16/3/>.
7. W. R. Blevin and B. Steiner, "The Redefinition of the Candela and the Lumen," *Metrologia* **11**:97–104 (1975).
8. O. C. Jones, "Proposed Changes to the SI System of Photometric Units," *Lighting Research and Technology* **10**:37–40 (1978).
9. W. R. Blevin, "The Candela and the Watt," *CIE Proc.* P-79-02 (1979).
10. CGPM, *Comptes Rendus des Séances de la 16e Conférence Générale des Poids et Mesures*, Paris 1979, BIPM, Sèvres, France (1979).
11. IESNA, *Lighting Handbook: Reference and Application*, M. S. Rea, ed., Illuminating Engineering Society of North America, New York, 1993.
12. CIE, *The Basis of Physical Photometry*, CIE 18.2, Vienna, 1983.

13. D. B. Judd, "Report of U.S. Secretariat Committee on Colorimetry and Artificial Daylight," *Proceedings of the Twelfth Session of the CIE (Stockholm)*, CIE, Paris, 1951. (The tabular data is given in G. Wyszecki and W. S. Stiles, *Color Science*, 2nd ed., Wiley, New York, 1982.)
14. J. J. Vos, "Colorimetric and Photometric Properties of a 2-Degree Fundamental Observer," *Color Research and Application* 3:125 (1978).
15. A. Stockman and L. T. Sharpe, "Cone Spectral Sensitivities and Color Matching," *Color Vision: From Genes to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge, 1999. (This information is summarized at the Color Vision Lab at UCSD Web site located at cvision.uscd.edu/.)
16. E. F. Zalewski, "Radiometry and Photometry," chap. 24, *Handbook of Optics*, vol. II, McGraw-Hill, New York, 1995.
17. F. E. Nicodemus, "Normalization in Radiometry," *Appl. Opt.* 12:2960 (1973).
18. J. M. Palmer, "Radiometric Bandwidth Normalization Using r.m.s. Methods," *Proc. SPIE* 256:99 (1980).
19. J. M. Palmer and M. G. Tomasko, "Broadband Radiometry with Spectrally Selective Detectors," *Opt. Lett.* 5:208 (1980).
20. E. H. Eberhart, "Source-Detector Spectral Matching Factors," *Appl. Opt.* 7:2037 (1968).
21. ANSI, *Photobiological Safety of Lamps*, American National Standards Institute RP27.3-96 (1996).

36.10 FURTHER READING

Books, Documentary Standards, Significant Journal Articles

- American National Standard Nomenclature and Definitions for Illuminating Engineering*, ANSI Standard ANSI/IESNA RP-16 96.
- W R. Blevin and B. Steiner, "Redefinition of the Candela and the Lumen," *Metrologia* 11:97 (1975).
- C. DeCusatis, *Handbook of Applied Photometry*, AIP Press, 1997. (Authoritative, with pertinent chapters written by technical experts at BIPM, CIE, and NIST. Skip chapter 4!)
- J. W. T. Walsh, *Photometry*, Constable, London, 1958. (The classic!)

Publications Available on the World Wide Web

- All you ever wanted to know about the SI is contained at BIPM and at NIST. Available publications (highly recommended) include the following:
- "The International System of Units (SI)," 7th ed. (1998), direct from BIPM. This is the English translation of the official document, which is in French. Available in PDF format at www.bipm.fr/.
- NIST Special Publication SP330, "The International System of Units (SI)." The U.S. edition (meter rather than metre) of the above BIPM publication. Available in PDF format from <http://physics.nist.gov/cuu/>.
- NIST Special Publication SP811, "Guide for the Use of the International System of Units (SI)," Available in PDF format from <http://physics.nist.gov/cuu/>.
- Papers published in recent issues of the NIST Journal of Research are also available on the Web in PDF format from mvl.nist.gov/pub/nistpubs/jres/jres.htm. Of particular relevance is "The NIST Detector-Based Luminous Intensity Scale," vol. 101, p. 109 (1996).

Useful Web Sites

- AIP (American Institute of Physics): www.aip.org
- ANSI (American National Standards Institute): www.ansi.org/
- BIPM (International Bureau of Weights and Measures): www.bipm.fr/

CIE (International Commission on Illumination): www.de.co.at/cie/

Color Vision Lab at UCSD: cvision.uscd.edu/

CORM (Council for Optical Radiation Measurements): www.corm.org

IESNA (Illuminating Engineering Society of North America): www.iesna.org/

ISO (International Standards Organization): www.iso.ch/

IUPAP (International Union of Pure and Applied Physics): www.physics.umanitoba.ca/iupap/

NIST (National Institute of Standards and Technology): physics.nist.gov/

OSA (Optical Society of America): www.osa.org

SPIE (International Society for Optical Engineering): www.spie.org

This page intentionally left blank.

RADIOMETRY AND PHOTOMETRY FOR VISION OPTICS*

Yoshi Ohno

*Optical Technology Division
National Institute of Standards and Technology
Gaithersburg, Maryland*

37.1 INTRODUCTION

Radiometry is the measurement of optical radiation, which is electromagnetic radiation in the frequency range between 3×10^{11} Hz and 3×10^{16} Hz. This range corresponds to wavelengths between 10 nm and 1000 μm , and includes the regions commonly called the ultraviolet, the visible, and the infrared. Typical radiometric units include watt (radiant flux), watt per steradian (radiant intensity), watt per square meter (irradiance), and watt per square meter per steradian (radiance).

Photometry is the measurement of light, which is defined as electromagnetic radiation detectable by the human eye. It is thus restricted to the visible region of the spectrum (wavelength range from 360 nm to 830 nm), and all the quantities are weighted by the spectral response of the eye. Photometry uses either optical radiation detectors constructed to mimic the spectral response of the eye, or spectroradiometry coupled with appropriate calculations for weighting by the spectral response of the eye. Typical photometric units include lumen (luminous flux), candela (luminous intensity), lux (illuminance), and candela per square meter (luminance).

The difference between radiometry and photometry is that radiometry includes the entire optical radiation spectrum (and often involves spectrally resolved measurements), while photometry deals with the visible spectrum weighted by the response of the eye. This chapter provides some guidance in photometry and radiometry (Refs. 1 through 6 are available for further details). The terminology used in this chapter follows international standards and recommendations.⁷⁻⁹

37.2 BASIS OF PHYSICAL PHOTOMETRY

The primary aim of photometry is to measure visible optical radiation, light, in such a way that the results correlate with the visual sensation of a normal human observer exposed to that radiation. Until about 1940, visual comparison measurement techniques were predominant in photometry.

*Chapter 1 in Vol. III gives further treatment of issues in radiometry, and Chap. 35 in this volume describes measurement of optical properties of materials.

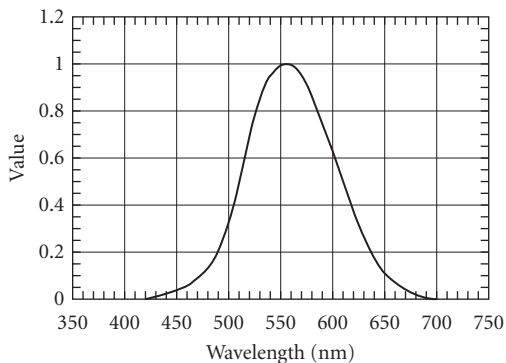


FIGURE 1 CIE $V(\lambda)$ function.

In modern photometric practice, measurements are made with photodetectors. This is referred to as *physical photometry*. In order to achieve the aim of photometry, one must take into account the characteristics of human vision. The relative spectral responsivity of the human eye was first defined by the CIE (Commission Internationale de l'Éclairage) in 1924,¹⁰ and redefined as part of colorimetric standard observers in 1931.¹¹ Called the *spectral luminous efficiency for photopic vision*, or $V(\lambda)$, it is defined in the domain from 360 nm to 830 nm, and is normalized at its peak, 555 nm (Fig. 1). This model has gained wide acceptance. The values were republished by CIE in 1983,¹² and published by CIPM (Comité International des Poids et Mesures) in 1982¹³ to supplement the 1979 definition of the candela. (The tabulated values of the function at 1-nm increments are available in Refs. 12 through 15.) In most cases, the region from 380 nm to 780 nm suffices for calculation with negligible errors because the value of the $V(\lambda)$ function falls below 10^{-4} outside this region. Thus, a photodetector having a spectral responsivity matched to the $V(\lambda)$ function replaced the role of human eyes in photometry.

Radiometry concerns physical measurement of optical radiation in terms of optical power, and in many cases, as a function of its wavelength. As specified in the definition of the candela by CGPM (Conférence Générale des Poids et Mesures) in 1979¹⁶ and by CIPM in 1982,¹³ a photometric quantity X_v is defined in relation to the corresponding radiometric quantity $X_{e,\lambda}$ by the equation:

$$X_v = K_m \int_{360\text{nm}}^{830\text{nm}} X_{e,\lambda} V(\lambda) d\lambda \quad (1)$$

The constant, K_m , relating the photometric quantities and radiometric quantities, is called the *maximum spectral luminous efficacy (of radiation) for photopic vision*. The value of K_m is given by the 1979 definition of candela that defines the spectral luminous efficacy of light at the frequency 540×10^{12} Hz (at the wavelength 555.016 nm in standard air) to be 683 lm/W. The value of K_m is calculated as $683 \times V(555.000 \text{ nm})/V(555.016 \text{ nm}) = 683.002 \text{ lm/W}$.¹² K_m is normally rounded to 683 lm/W with negligible errors.

It should be noted that $V(\lambda)$ is defined for the *CIE standard photometric observer for photopic vision*, which assumes additivity of sensation and a 2° field of view at relatively high luminance levels (higher than approximately 1 cd/m^2). The human vision in this level is called photopic vision. The spectral responsivity of human vision deviates significantly at very low levels of luminance (less than approximately 10^{-3} cd/m^2). This type of vision is called scotopic vision. Its spectral responsivity, peaking at 507 nm, is designated by $V'(\lambda)$, which was defined by CIE in 1951,¹⁷ recognized by CIPM in 1976,¹⁸ and republished by CIPM in 1982.¹³ Human vision in the region between photopic vision and scotopic vision is called mesopic vision. While there have been active researches in this area,¹⁹ there is no internationally accepted spectral luminous efficiency function for the mesopic region yet. In current practice, almost all photometric quantities are given in terms of photopic vision, even at low light levels. Quantities in scotopic vision are seldom used except for special calculations for research purposes. (Further details of the contents in this section are given in Ref. 12.)

37.3 PHOTOMETRIC BASE UNIT—THE CANDELA

The history of photometric standards dates back to the early nineteenth century, when the intensity of light sources was measured in comparison with a standard candle using visual bar photometers. At that time, the flame of a candle was used as a unit of luminous intensity that was called the *candle*. The old name for luminous intensity *candle power* came from this origin. Standard candles were gradually superseded by flame standards of oil lamps, and in the early twentieth century investigations on platinum point blackbodies began at some national laboratories. An agreement was first established in 1909 among several national laboratories to use such a blackbody to define the unit of luminous intensity, and the unit was recognized as the *international candle*. This standard was adopted by the CIE in 1921. In 1948, it was adopted by the CGPM¹⁶ with a new Latin name *candela* with the following definition:

The candela is the luminous intensity, in the perpendicular direction, of a surface of 1/600000 square meter of a blackbody (full radiator) at the temperature of freezing platinum under a pressure of 101325 newton per square meter.

Although the 1948 definition served to establish the uniformity of photometric measurements in the world, difficulties in fabricating the blackbodies and in improving accuracy were addressed. Beginning in the mid-1950s, suggestions were made to define the candela in relation to the unit of optical power, watt, so that complicated source standards would not be necessary. Finally, in 1979, a new definition of the candela was adopted by the CGPM¹⁶ as follows:

The candela is the luminous intensity, in a given direction, of a source that emits monochromatic radiation of frequency 540×10^{12} hertz and that has a radiant intensity in that direction of (1/683) watt per steradian.

The value of K_m (683 lm/W) was determined in such a way that the consistency from the prior unit was maintained, and was determined based on the measurements by several national laboratories. (Technical details on this redefinition of the candela are reported in Refs. 20 and 21.) This 1979 redefinition of the candela has enabled photometric units to be derived from radiometric units using a variety of techniques. (The early history of photometric standards is described in greater detail in Ref. 22.)

37.4 QUANTITIES AND UNITS IN PHOTOMETRY AND RADIOMETRY

In 1960, the SI (Système International) was established, and the candela became one of the seven SI base units.²³ (For further details on the SI, Refs. 23 through 26 may be consulted.) Several quantities and units, defined in different geometries, are used in photometry and radiometry. Table 1 lists the photometric quantities and units, along with corresponding quantities and units for radiometry.

While the candela is the SI base unit, the luminous flux (lumen) is perhaps the most fundamental photometric quantity, as the other photometric quantities are defined in terms of lumen with an appropriate geometric factor. The definitions of these photometric quantities are given in the following sections. (The descriptions given here are somewhat simplified from the rigorous definitions for ease of understanding. Refer to Refs. 7 through 9 for official rigorous definitions.)

Radiant Flux and Luminous Flux

Radiant flux (also called *optical power* or *radiant power*) is the energy Q (in joules) radiated by a source per unit of time, expressed as

$$\Phi = \frac{dQ}{dt} \quad (2)$$

The unit of radiant flux is the *watt* ($W = J/s$).

TABLE 1 Quantities and Units Used in Photometry and Radiometry

Photometric Quantity	Unit	Relationship with Lumen	Radiometric Quantity	Unit
Luminous flux	lm (lumen)		Radiant flux	W (watt)
Luminous intensity	cd (candela)	lm sr ⁻¹	Radiant intensity	W sr ⁻¹
Illuminance	lx (lux)	lm m ⁻²	Irradiance	W m ⁻²
Luminance	cd m ⁻²	lm sr ⁻¹ m ⁻²	Radiance	W sr ⁻¹ m ⁻²
Luminous exitance	lm m ⁻²		Radiant exitance	W m ⁻²
Luminous exposure	lx · s		Radiant exposure	W m ⁻² · s
Luminous energy	lm · s		Radiant energy	J (joule)
Total luminous flux	lm (lumen)		Total radiant flux	W (watt)
Color temperature	K (kelvin)		Radiance temperature	K (kelvin)

Luminous flux (Φ_v) is the time rate of flow of light as weighted by $V(\lambda)$. The unit of luminous flux is the *lumen* (lm). It is defined as

$$\Phi_v = K_m \int_{\lambda} \Phi_{e,\lambda} V(\lambda) d\lambda \tag{3}$$

where ($\Phi_{e,\lambda}$) is the spectral concentration of radiant flux as a function of wavelength λ . The term luminous flux is often used in the meaning of total luminous flux in photometry (see the following subsection entitled “Total Radiant Flux and Total Luminous Flux”).

Radiant Intensity and Luminous Intensity

Radiant intensity (I_e) or luminous intensity (I_v) is the radiant flux (luminous flux) from a point source emitted per unit solid angle in a given direction, as defined by

$$I = \frac{d\Phi}{d\Omega} \tag{4}$$

where $d\Phi$ is the radiant flux (luminous flux) leaving the source and propagating in an element of solid angle $d\Omega$ containing the given direction. The unit of radiant intensity is W/sr, and that of luminous intensity is the *candela* (cd = lm/sr). (See Fig. 2.)

Solid Angle The solid angle (Ω) of a cone is defined as the ratio of the area (A) cut out on a spherical surface (with its center at the apex of that cone) to the square of the radius (r) of the sphere, as given by

$$\Omega = \frac{A}{r^2} \tag{5}$$

The unit of solid angle is *steradian* (sr), which is a dimensionless unit. (See Fig. 3.)

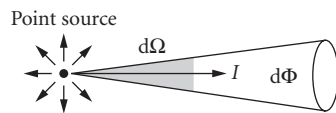


FIGURE 2 Radiant intensity and luminous intensity.

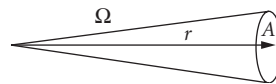


FIGURE 3 Solid angle.

Irradiance and Illuminance

Irradiance (E_e) or *illuminance* (E_v) is the density of incident radiant flux or luminous flux at a point on a surface, and is defined as radiant flux (luminous flux) per unit area, as given by

$$E = \frac{d\Phi}{dA} \quad (6)$$

where $d\Phi$ is the radiant flux (luminous flux) incident on an element dA of the surface containing the point. The unit of irradiance is W/m^2 , and that of illuminance is *lux* ($\text{lx} = \text{lm}/\text{m}^2$). (See Fig. 4.)

Radiance and Luminance

Radiance (L_e) or *luminance* (L_v) is the radiant flux (luminous flux) per unit solid angle emitted from a surface element in a given direction, per unit projected area of the surface element perpendicular to the direction. The unit of radiance is $\text{W sr}^{-1} \text{m}^{-2}$, and that of luminance is cd/m^2 . These quantities are defined by

$$L = \frac{d^2\Phi}{d\Omega \cdot A \cdot \cos\theta} \quad (7)$$

where $d\Phi$ is the radiant flux (luminous flux) emitted (reflected or transmitted) from the surface element and propagating in the solid angle $d\Omega$ containing the given direction. dA is the area of the surface element, and θ is the angle between the normal to the surface element and the direction of the beam. The term $dA \cos \theta$ gives the projected area of the surface element perpendicular to the direction of measurement. (See Fig. 5.)

Radiant Exitance and Luminous Exitance

Radiant exitance (M_e) or *luminous exitance* (M_v) is defined to be the density of radiant flux (luminous flux) leaving a surface at a point. The unit of radiant exitance is W/m^2 and that of luminous exitance is lm/m^2 (but it is not lux). These quantities are defined by

$$E = \frac{d\Phi}{dA} \quad (8)$$

where $d\Phi$ is the radiant flux (luminous flux) leaving the surface element. Luminous exitance is rarely used in the general practice of photometry. (See Fig. 6.)

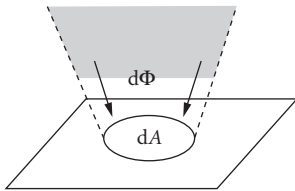


FIGURE 4 Irradiance and illuminance.

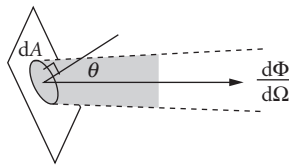


FIGURE 5 Radiance and luminance.

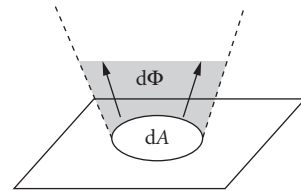


FIGURE 6 Radiant exitance and luminous exitance.

Radiant Exposure and Luminous Exposure

Radiant exposure (H_c) or *luminous exposure* (H_v) is the time integral of irradiance $E_c(t)$ or illuminance $E_v(t)$ over a given duration Δt , as defined by

$$H = \int_{\Delta t} E(t) dt \quad (9)$$

The unit of radiant exposure is J m^{-2} , and that of luminous exposure is *lux · second* ($\text{lx} \cdot \text{s}$).

Radiant Energy and Luminous Energy

Radiant energy (Q_c) or *luminous energy* (Q_v) is the time integral of the radiant flux or luminous flux (Φ) over a given duration Δt , as defined by

$$Q = \int_{\Delta t} \Phi(t) dt \quad (10)$$

The unit of radiant energy is *joule* (J), and that of luminous energy is *lumen · second* ($\text{lm} \cdot \text{s}$).

Total Radiant Flux and Total Luminous Flux

Total radiant flux or *total luminous flux* (Φ_v) is the geometrically total radiant (luminous) flux of a light source. It is defined as

$$\Phi = \int_{\Omega} I d\Omega \quad (11)$$

or

$$\Phi = \int_A E dA \quad (12)$$

where I is the radiant (luminous) intensity distribution of the light source and E is the irradiance (illuminance) distribution over a given closed surface surrounding the light source. If the radiant (luminous) intensity distribution or the irradiance (illuminance) distribution is given in polar coordinates (θ, ϕ) , the total radiant (luminous) flux of the source Φ is given by

$$\Phi = \int_{\phi=0}^{2\pi} \int_{\theta=0}^{\pi} I(\theta, \phi) \sin\theta d\theta d\phi \quad (13)$$

or

$$\Phi = r^2 \int_{\phi=0}^{2\pi} \int_{\theta=0}^{\pi} E(\theta, \phi) \sin\theta d\theta d\phi \quad (14)$$

For example, the total luminous flux of an isotropic point source having luminous intensity of 1 cd would be 4π lm.

Radiance Temperature and Color Temperature

Radiance temperature (unit: kelvin) is the temperature of the Planckian radiator for which the radiance at the specified wavelength has the same spectral concentration as for the thermal radiator considered.

Color temperature (unit: kelvin) is the temperature of a Planckian radiator with radiation of the same chromaticity as that of the light source in question. This term is commonly used to specify the

colors of incandescent lamps whose chromaticity coordinates are practically on the blackbody locus. The next two terms are also important in photometry.

Distribution temperature (unit: kelvin) is the temperature of a blackbody with a spectral power distribution closest to that of the light source in question, and is used for quasi-Planckian sources such as incandescent lamps (refer to Ref. 27 for details).

Correlated color temperature (unit: kelvin) is the temperature of the Planckian radiator whose perceived color most closely resembles that of the light source in question. Correlated color temperature is used for sources with a spectral power distribution significantly different from that of Planckian radiation (e.g., discharge lamps; refer to Ref. 28 for details).

Relationship Between SI Units and English Units

The SI units as described previously should be used in all radiometric and photometric measurements according to international standards and recommendations on SI units. However, some English units are still rather widely used in some countries, including the United States. The use of these non-SI units is discouraged. The definitions of these English units are given in Table 2 for conversion purposes only.

The definition of footlambert is such that the luminance of a perfect diffuser is 1 fL when illuminated at 1 fc. In the SI unit, the luminance of a perfect diffuser would be $1/\pi$ (cd/m²) when illuminated at 1 lx. For convenience of changing from English units to SI units, the conversion factors are listed in Table 3. For example, 1000 lx is the same illuminance as 92.9 fc, and 1000 cd/m² is the same luminance as 291.9 fL. (Conversion factors to and from some more units are given in Ref. 5.)

Troland

This unit is not an SI unit, not used in metrology, and is totally different from all other photometric units mentioned previously. It is introduced here because this unit is commonly used by vision scientists. Troland is defined as the retinal illuminance when a surface of luminance one candela per square meter is viewed through a pupil at the eye (natural or artificial) of area one square millimeter. Thus, the troland value, T , for the luminance, L (cd/m²), of an external field and the pupil size, p (mm²), is given by

$$T = L \cdot P \tag{15}$$

TABLE 2 English Units and Definition

Unit	Quantity	Definition
Footcandle (fc)	Illuminance	Lumen per square foot (lm ft ⁻²)
Footlambert (fL)	Luminance	$1/\pi$ candela per square foot (π^{-1} cd ft ⁻²)

TABLE 3 Conversion between English Units and SI Units

To Obtain the Value in	Multiply the Value in	By
lx from fc	fc	10.764
fc from lx	lx	0.09290
cd/m ² from fL	fL	3.4263
fL from cd/m ²	cd/m ²	0.29186
m (meter) from feet	feet	0.30480
mm (millimeter) from inch	inch	25.400

or, for pupil size p (m²),

$$T = L \cdot 10^6 \times p \quad (16)$$

The troland value is not the real illuminance in lux on the retina, but is a quantity proportional to it. Since the natural pupil size changes with luminance level, luminance changes do not have a proportional visual effect. Thus, troland value rather than luminance is often useful in visual experiments. There is no simple or defined conversion between troland value (for natural pupil) and luminance (cd/m²), without knowing the actual pupil size. (Further details of this unit can be found in Ref. 29.)

37.5 PRINCIPLES IN PHOTOMETRY AND RADIOMETRY

Several important theories in practical photometry and radiometry are introduced in this section.

Inverse Square Law

Illuminance E (lx) at a distance d (m) from a point source having luminous intensity I (cd) is given by

$$E = \frac{I}{d^2} \quad (17)$$

For example, if the luminous intensity of a lamp in a given direction is 1000 cd, the illuminance at 2 m from the lamp in this direction is 250 lx. Note that the inverse square law is valid only when the light source is regarded as a point source. Sufficient distances relative to the size of the source are needed to assume this relationship.

Lambert's Cosine Law

The luminous intensity of a Lambertian surface element is given by

$$I(\theta) = I_n \cos \theta \quad (18)$$

(See Fig. 7.)

Lambertian Surface A surface whose luminance is the same in all directions of the hemisphere above the surface.

Perfect (Reflecting/Transmitting) Diffuser A Lambertian diffuser with a reflectance (transmittance) equal to 1.

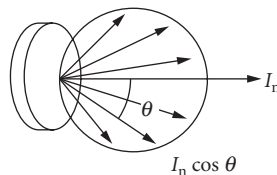


FIGURE 7 Lambert's cosine law.

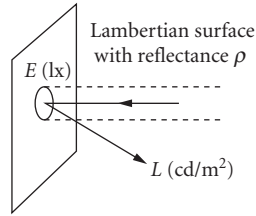


FIGURE 8 Relationship between illuminance and luminance.

Relationship between Illuminance and Luminance

The luminance $L \text{ (cd/m}^2\text{)}$ of a Lambertian surface of reflectance ρ , illuminated by $E \text{ (lx)}$ is given by

$$L = \frac{\rho \cdot E}{\pi} \quad (19)$$

(See Fig. 8.)

Reflectance (ρ) The ratio of the reflected flux to the incident flux in a given condition. The value of ρ can be between 0 and 1.

In the real world, there is no existing perfect diffuser nor perfectly Lambertian surfaces, and Eq. 19 does not apply. For real object surfaces, the following terms apply.

Luminance Factor (β) Ratio of the luminance of a surface element in a given direction to that of a perfect reflecting or transmitting diffuser, under specified conditions of illumination. The value of β can be larger than 1. For a Lambertian surface, reflectance is equal to the luminance factor. Equation (19) for real object is restated using β as

$$L = \frac{\beta \cdot E}{\pi} \quad (20)$$

Luminance Coefficient (q) Quotient of the luminance of a surface element in a given direction by the illuminance on the surface element, under specified conditions of illumination,

$$q = \frac{L}{E} \quad (21)$$

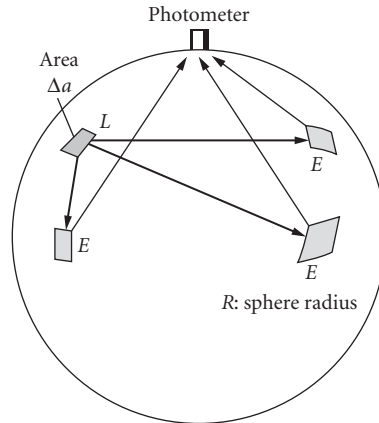
Using q , the relationship between luminance and illuminance is thus given by

$$L = q \cdot E \quad (22)$$

Luminance factor corresponds to radiance factor, and luminance coefficient corresponds to radiance coefficient in radiometry. BRDF (bidirectional reflectance distribution function) is also used for the same concept as radiance coefficient.

Integrating Sphere

An integrating sphere is a device to make a spatial integration of luminous flux (or radiant flux) generated (or introduced) in the sphere and to detect it with a single photodetector. In the case of measurement of light sources, the spatial integration is made over the entire solid angle (4π).


FIGURE 9 Flux transfer in a sphere.

In Fig. 1, assuming that the integrating sphere wall surfaces are perfectly Lambertian, the illuminance E on any part of the sphere wall created by luminance L of an element Δa is given by

$$E = \frac{L\Delta a}{4R^2} \quad (23)$$

where R is the radius of the sphere. This equation holds no matter where the two surface elements are. In other words, the same amount of flux incident anywhere on the sphere wall will create an equal illuminance on the detector port. In the case of actual integrating spheres, the surface is not perfectly Lambertian, but due to interreflections of light in the sphere, the distribution of reflected light will be uniform enough to assume the relationship of Eq. (23). (See Fig. 9.)

The direct light from an actual light source is normally not uniform; thus, it must be shielded from the detector. When a light source with luminous flux Φ is operated in a sphere having reflectance ρ , the flux created by interreflections is given by

$$\Phi(\rho + \rho^2 + \rho^3 + \dots) = \Phi \cdot \frac{\rho}{1 - \rho} \quad (24)$$

Then, the illuminance E_d created by all the interreflections is given by

$$E_d = \frac{\Phi \cdot \rho}{1 - \rho} \cdot \frac{1}{4\pi \cdot R^2} \quad (25)$$

The sphere efficiency (E_d/Φ) is strongly dependent on reflectance ρ due to the term $1 - \rho$ in the denominator. For example, the detector signal at $\rho = 0.98$ is 10 times larger than at $\rho = 0.8$.

Planck's Law

The spectral radiance of a blackbody at a temperature T (K) is given by

$$I_c(\lambda, T) = c_1 n^{-2} \pi^{-1} \lambda^{-5} \left[\exp\left(\frac{c_2}{n\lambda T}\right) - 1 \right]^{-1} \quad (26)$$

where $c_1 = 2\pi hc^2 = 3.7417749 \times 10^{-16} \text{ W} \cdot \text{m}^2$, $c_2 = hc/k = 1.438769 \times 10^{-2} \text{ m} \cdot \text{K}^2$ (1986 CODATA from Ref. 9), h is Planck's constant, c is the speed of light in vacuum, k is the Boltzmann constant, n ($= 1.00028$) is the refractive index of standard air (12),³⁰ and λ is the wavelength.

Wien's Displacement Law

Taking the partial derivative of the Planck's equation with respect to temperature T , and setting the result equal to zero, the solution yields the relationship between the peak wavelength λ_m for Planck's radiation and temperature T (K), as given by

$$\lambda_m T = 2897.8 \mu\text{m} \cdot \text{K} \quad (27)$$

This shows that the peak wavelength of blackbody radiation shifts to shorter wavelengths as the temperature of the blackbody increases.

Stefan-Boltzmann's Law

The (spectrally total) radiant exitance M_e from a blackbody in a vacuum is expressed in relation to the temperature T (K) of the blackbody, in the form

$$M_e(T) = \int_0^\infty M_e(\lambda, T) d\lambda = \sigma T^4 \quad (28)$$

where $M_e(\lambda, T)$ is the spectral radiant exitance of the blackbody, and σ is the Stefan-Boltzmann constant, equal to $5.67051 \times 10^{-8} \text{ W} \cdot \text{m}^{-2} \cdot \text{K}^{-4}$ (1986 CODATA from Ref. 9). Using this value, the unit for M_e is $\text{W} \cdot \text{m}^{-2}$.

37.6 PRACTICE IN PHOTOMETRY AND RADIOMETRY

Photometry and radiometry are practiced in many different areas and applications, dealing with various light sources and detectors, and cannot be covered in this chapter. Various references are available on practical measurements in photometry and radiometry.

Further references in practical radiometry include books on absolute radiometry,³¹ optical detectors,³² spectroradiometry,³³ photoluminescence,³⁴ radiometric calibration,³⁵ etc. There are a number of publications from CIE that are regarded as international recommendations or standards. CIE publications in radiometry include reports on absolute radiometry,³⁶ reflectance,³⁷ spectroradiometry,³⁸ detector spectral response,³⁹ photobiology and photochemistry,⁴⁰ etc. There are also a number of publications from the National Institute of Standards and Technology (NIST) in radiometry, on spectral radiance,⁴¹ spectral irradiance,⁴² spectral reflectance,⁴³ spectral responsivity,⁴⁴ and so on (Ref. 45 provides greater depths of knowledge in radiometry).

For practical photometry, Ref. 4 provides the latest information on standards and practical measurements of photometry in many aspects. A recent publication from NIST⁴⁶ is also available. CIE publications are also available on many subjects in photometry, including characterization of illuminance meters and luminance meters,⁴⁷ luminous flux measurement,⁴⁸ measurements of LEDs,⁴⁹ characteristics of displays,⁵⁰ and many others. A series of measurement guide documents are published from the Illuminating Engineering Society of North America (IESNA) for operation and measurement of particular types of lamps⁵¹⁻⁵³ and luminaires. The American Society for Testing and Materials (ASTM) provides many useful standards and recommendations on optical properties of materials and color measurements.⁵⁴ Colorimetry is a branch of radiometry and is becoming increasingly important among color imaging industry and multimedia applications. The basis of colorimetry is provided by CIE publications^{28,55,56} and many other authoritative references are available.^{29,57}

37.7 REFERENCES

1. F. Grum and R. J. Becherer, *Optical Radiation Measurements, Vol. 1 Radiometry*, Academic Press, San Diego, CA, 1979.
2. R. McCluney, *Introduction to Radiometry and Photometry*, Artech House, Norwood, MA, 1994.
3. W. L. Wolfe, *Introduction to Radiometry*, SPIE—The International Society for Optical Engineering, P.O. Box 10, Bellingham, WA 98227-0010, 1998.
4. Casimer DeCusatis (ed.), *OSA/AIP Handbook of Applied Photometry*, AIP Press, Woodbury, NY, 1997.
5. *IES Lighting Handbook, 8th edition, Reference and Application*, Illuminating Engineering Society of North America, New York, 1993.
6. J. M. Palmer, Radiometry and Photometry FAQ, <http://www.optics.Arizona.EDU/Palmer/rpfaq/rpfaq.htm>.
7. *International Lighting Vocabulary*, CIE Publication 17.4 (1987).
8. *International Vocabulary of Basic and General Terms in Metrology*, BIPM, IEC, IFCC, ISO, IUPAC, IUPAP, and OIML, 1994.
9. *Quantities and Units*, ISO Standards Handbook, 3rd edition, 1993.
10. CIE Compte Rendu, p. 67 (1924).
11. CIE Compte Rendu, Table II, pp. 25–26 (1931).
12. *The Basis of Physical Photometry*, CIE Publication 18.2 (1983).
13. CIPM, *Comité Consultatif de Photométrie et Radiométrie 10e Session—1982*, BIPM, Pavillon de Breteuil, F-92310 Sèvres, France (1982).
14. *Principles Governing Photometry*, Bureau International Des Poids et Mesures (BIPM) Monograph, BIPM, F-92310 Sèvres, France (1983).
15. CIE Disk D001, Photometric and Colorimetric Tables (1988).
16. CGPM, *Comptes Rendus des Séances de la 16e Conférence Générale des Poids et Mesures*, Paris 1979, BIPM, F-92310 Sèvres, France (1979).
17. CIE Compte Rendu, Vol. 3, Table II, pp. 37–39 (1951).
18. CIPM Procès-Verbaux 44, 4 (1976).
19. *Mesopic Photometry: History, Special Problems and Practical Solutions*, CIE Publication 81 (1989).
20. W. R. Blevin and B. Steiner, *Metrologia* 11:97 (1975).
21. W. R. Blevin, “The Candela and the Watt,” *CIE Proc.* P-79-02 (1979).
22. J. W. T. Walsh, *Photometry*, Constable, London, 1953.
23. *Le Système International d’Unité (SI), The International System of Units (SI)*, 6th edition, Bur. Intl. Poids et Mesures, Sèvres, France, 1991.
24. B. N. Taylor, *Guide for the Use of the International System of Units (SI)*, Natl. Inst. Stand. Technol. Spec. Publ. 811, 1995.
25. B. N. Taylor (ed.), *Interpretation of the SI for the United States and Metric Conversion Policy for Federal Agencies*, Natl. Inst. Stand. Technol. Spec. Publ. 814, 1991.
26. *SI Units and Recommendations for the Use of Their Multiples and of Certain Other Units*, ISO 1000: 1992, International Organization for Standardization, Geneva, Switzerland, 1992.
27. *CIE Collection in Photometry and Radiometry*, Publication No. 114/4, 1994.
28. *Colorimetry* 3rd edition, CIE Publication 15:2004 (2004).
29. G. Wyszecki and W. S. Stiles, *Color Science*, John Wiley and Sons, Inc., New York, 1982.
30. W. R. Blevin, “Corrections in Optical Pyrometry and Photometry for the Refractive Index of Air,” *Metrologia* 8:146 (1972).
31. F. Hengstberger (ed.), *Absolute Radiometry*, Academic Press, San Diego, CA, 1989.
32. W. Budde, *Optical Radiation Measurements, Vol. 4—Physical Detectors of Optical Radiation*, Academic Press, Orlando, FL, 1983.
33. H. Kostkowski, *Reliable Spectroradiometry*, Spectroradiometry Consulting, P.O. Box 2747, La Plata, MD 20646-2747, USA.

34. K. D. Mielenz (ed.), *Optical Radiation Measurements, Vol. 3, Measurement of Photoluminescence*, Academic Press, New York, 1982.
35. C. L. Wyatt, *Radiometrie Calibration: Theory and Models*, Academic Press, New York, 1978.
36. *Electrically Calibrated Thermal Detectors of Optical Radiation (Absolute Radiometers)*, CIE Publication 65 (1985).
37. *Absolute Methods for Reflection Measurements*, CIE Publication 44 (1979).
38. *The Spectroradiometric Measurement of Light Sources*, CIE Publication 63 (1984).
39. *Determination of the Spectral Responsivity of Optical Radiation Detectors*, CIE Publication 64 (1984).
40. *CIE Collection in Photobiology and Photochemistry*, CIE Publication 106 (1993).
41. J. H. Walker, R. D. Saunders, and A. T. Hattenburg, *Spectral Radiance Calibrations*, NBS Special Publication 250-1, 1987.
42. J. H. Walker, R. D. Saunders, J. K. Jackson, and D. A. McSparron, *Spectral Irradiance Calibrations*, NBS Special Publication 250-20, 1987.
43. P. Y. Barnes, E. A. Early, and A. C. Parr, *Spectral Reflectance*, NIST Special Publication 250-48, 1998.
44. T. C. Larason, S. S. Bruce, and A. C. Parr, *Spectroradiometric Detector Measurements*, NIST Special Publication 250-41, 2008.
45. F. Nicodemos (ed.), *Self-Study Manual on Optical Radiation Measurements*, NBS Technical Note 910 Series, Parts 1-12, 1978-1985.
46. Y. Ohno, *Photometric Calibrations*, NIST Special Publication 250-37 (1997).
47. *Methods of Characterizing Illuminance Meters and Luminance Meters*, CIE Publication 69 (1987).
48. *Measurement of Luminous Flux*, CIE Publication 84 (1989).
49. *Measurement of LEDs*, 2nd edition, CIE Publication 127: 2007 (2007).
50. *The Relationship between Digital and Colorimetric Data for Computer-Controlled CRT Displays*, CIE Publication 122, 1996.
51. *IES Approved Method for the Electric and Photometric Measurement of Fluorescent Lamps*, IESNA LM-9.
52. *Electrical and Photometric Measurements of General Service Incandescent Filament Lamps*, IESNA LM-45.
53. *Electrical and Photometric Measurements of Compact and Fluorescent Lamps*, IESNA LM-66.
54. *ASTM Standards on Color and Appearance Measurement*, 5th edition, 1996.
55. ISO 11664-2:2008(E)/CIE S 014-2/E:2006, CIE Colorimetry Part 2: Standard Illuminants for Colorimetry.
56. ISO 11664-1:2008(E)/CIE S 014-1/E:2006, CIE Colorimetry Part 1: Standard Colorimetric Observers.
57. F. Grum and C. J. Bartleson (eds.), *Optical Radiation Measurements, Vol. 2, Color Measurement*, Academic Press, New York, 1980.

This page intentionally left blank.

Carolyn J. Sher DeCusatis

*Pace University
White Plains, New York*

38.1 INTRODUCTION

Spectroradiometry is the measurement of the spectral content of optical radiation. This has many important applications. The measure of terrestrial, direct, solar spectral irradiance between 295 and 305 nm can be used to calculate atmospheric ozone thickness.¹ More close to home, irradiance measurements are used to characterize light fixtures, and solar UV spectroradiometry methods also apply to the measurement of artificial sources that mimic the sun for applications like phototherapy to treat seasonal affective disorder (SAD) and tanning booths.² Transmission spectra are used to analyze the chemical composition of samples, such as the concentration of chlorophyll in a solution. Spectral reflectance quantifies the color of surfaces, with many practical applications to building, lighting, and design. Spectral responsivity is a necessary part of calibrating photodetectors.

38.2 DEFINITIONS, CALCULATIONS, AND FIGURES OF MERIT

Defining Quantities

There is a relationship between radiometric, photometric, and spectroradiometric quantities. Radiometric and photometric quantities, such as irradiance and luminous flux have been defined in other chapters of this book. Photometric quantities that are similar to radiometric quantities, such as radiant energy versus luminous energy, have the same symbol with a subscript of γ . In general, the spectroradiometric quantity that is defined by the similar radiometric quantity is preceded by the term “spectral,” and designated by the symbol λ , either in parenthesis or with a subscript.

Spectral Irradiance is the quantity most frequently measured in spectroradiometry.¹ Irradiance E is the total radiant flux incident on an element of surface divided by the surface area of that element

$d\Phi/dA$, in watts per meter squared. The average spectral irradiance \bar{E}_λ is the irradiance for a wavelength interval, or

$$\bar{E}_\lambda = \frac{\Delta\Phi}{\Delta\lambda \cdot \Delta A} \quad (1)$$

where $\Delta\Phi$ is the radiant flux within a wavelength interval $\Delta\lambda$ incident on a surface area ΔA . As the area and wavelength are made smaller, \bar{E}_λ becomes the spectral irradiance for that wavelength,

$$E_\lambda(\lambda, P) = \frac{d^2\Phi}{d\lambda \cdot dA} \quad (2)$$

where P is the position and λ is the wavelength. The SI unit for spectral irradiance is the watt per meter cube. However, more intuitive units, such as microwatt per square centimeter of area and per nanometer of wavelength ($\mu\text{W} \cdot \text{cm}^{-2} \cdot \text{nm}^{-1}$) are also used.

The spectral irradiance is dependant on the position, the size of the solid angle subtended, and the orientation of the surface. This leads to the cosine dependence, where the spectral irradiance of a point source is proportional to cosine θ , where θ is the angle between the normal to the surface and the direction of the source, and the inverse square law, relating to the distance between the point source and the detector.

When the radiation from a point source is emitted into a solid angle $d\Omega_s$ the *spectral radiance* L_λ is the flux per unit solid angle and per unit projected area perpendicular to the specified direction, per wavelength.

$$L_\lambda = \frac{d^3\Phi}{dA \cdot \cos\theta \cdot d\Omega \cdot d\lambda} \quad (3)$$

The unit for spectral radiance is watt · centimeter⁻² · steradian⁻¹ · nanometer⁻¹.

Spectral radiance represents the flux density at a point for a particular direction through that point. While it is not usually the desired quantity an experiment is designed to measure, it is required for quantitatively analyzing data, and is also very useful for flux transfer calculations.

If you know the wavelength dependence of the radiant flux, you can also calculate *radiant flux* Φ , and *luminous flux* Φ_v .

$$\phi = \int_\lambda \Phi(\lambda) d\lambda \quad \text{in watts} \quad (4)$$

$$\phi_v = 683 \int_\lambda \Phi(\lambda) V(\lambda) d\lambda \quad \text{in lumens} \quad (5)$$

where $V(\lambda)$ is the relative photopic luminous efficiency curve (normalized at 555 nm). The absolute luminous efficacy at 555 nm is 683 lumens per watt.

Spectral transmittance $\tau(\lambda)$, which is widely measured during spectrometry, is

$$\tau(\lambda) = \phi^t(\lambda) / \phi^i(\lambda) \quad (6)$$

where t and i refer to transmitted and incident flux. The transmittance has two parts: the regular and diffuse transmittance. The regular transmittance follows Snell's law, while the diffuse transmittance is scattered by the roughness of the surface.

Similarly, *spectral reflectance* $\rho(\lambda)$ is also a flux ratio, but in this case, it is the ratio of the reflected radiant flux to the incident flux.

$$\rho(\lambda) = \phi^r(\lambda) / \phi^i(\lambda) \quad (7)$$

where r refers to reflected flux. The total reflectance is also composed of specular and diffuse components.

The *spectral responsivity* $R(\lambda)$ is the current from a detector divided by the incident flux for a specific wavelength.

$$R(\lambda) = r(\lambda) / \phi(\lambda) \quad (8)$$

where $r(\lambda)$ is the electric signal generated by the photodetector. The units of responsivity are amps per watt. Spectral responsivity is an important part of calibrating photodetectors. The spectroradiometer itself will be calibrated using a standard lamp.

Calculations

Since spectroradiometers provide the spectral content of light, it is natural to use their data to calculate color values. The *tristimulus values* for sources can be calculated,³

$$X = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{x}(\lambda) \Delta\lambda \quad (9)$$

$$Y = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{y}(\lambda) \Delta\lambda \quad (10)$$

$$Z = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{z}(\lambda) \Delta\lambda \quad (11)$$

where $E(\lambda)$ is the irradiance in watt · meters⁻², and \bar{x} , \bar{y} , and \bar{z} are the CIE 1931 spectral tristimulus values.

When analyzing the reflection off an object, the values assume the color as seen under a standard light source.

$$X = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{x}(\lambda) \Delta\lambda \quad (12)$$

$$Y = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{y}(\lambda) \Delta\lambda \quad (13)$$

$$Z = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{z}(\lambda) \Delta\lambda \quad (14)$$

where $\Gamma(\lambda)$ is the *spectral reflectance* or *transmittance data*, and

$$k = \frac{100}{\sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{y}(\lambda) \Delta\lambda} \quad (15)$$

Since $\bar{y}(\lambda)$ is the photopic curve, k is a constant that can be used to couple the colorimetric (photometric) quantities with the radiometric ones. This can be expressed in equation form as

$$E_y [\text{lm cm}^{-2}] = 683 Y [\text{W cm}^{-2}] \quad (16)$$

because absolute luminous efficacy of the photopic curve at 555 nm is 683 lumens per watt.

The CIE 1931 chromaticity x , y , z coordinates are

$$x = \frac{X}{X + Y + Z} \quad (17)$$

$$y = \frac{Y}{X + Y + Z} \quad (18)$$

$$z = \frac{Z}{X + Y + Z} \quad (19)$$

Similarly, the R , G , B tristimulus values are⁴

$$R = k \sum_{\lambda=380}^{780} P(\lambda) \bar{r}(\lambda) \Delta\lambda \quad (20)$$

$$G = k \sum_{\lambda=380}^{780} P(\lambda) \bar{g}(\lambda) \Delta\lambda \quad (21)$$

$$B = k \sum_{\lambda=380}^{780} P(\lambda) \bar{b}(\lambda) \Delta\lambda \quad (22)$$

where $P(\lambda)$ is the spectral power distribution in watts, and \bar{r} , \bar{g} , and \bar{b} are the color-matching functions of the CIE 1931 Colorimetric Observer.

The UCS 1960 u and v coordinates, and the UCS 1976 u' and v' coordinates are³

$$v = \frac{6y}{12y - 12x + 3} = \frac{2}{3}v' \quad (23)$$

$$u = \frac{4x}{12y - 12x + 3} = u' \quad (24)$$

The CIE LAB/LUV color space calculations (1976) are calculated using the tristimulus values normalized equally to $Y = 100$. X_n , Y_n , and Z_n are the tristimulus values of the reference white. The coordinates can be defined in either the $L^*a^*b^*$ color space, or the $L^*u^*v^*$ color space. When X/X_n , Y/Y_n , and Z/Z_n are all greater than 0.01,

$$L^* = 116 \left(\frac{Y}{Y_n} \right)^{1/3} - 16 \quad (25)$$

$$a^* = 500 \left[\left(\frac{X}{X_n} \right)^{1/3} - \left(\frac{Y}{Y_n} \right)^{1/3} \right] \quad (26)$$

$$b^* = 200 \left[\left(\frac{Y}{Y_n} \right)^{1/3} - \left(\frac{Z}{Z_n} \right)^{1/3} \right] \quad (27)$$

Otherwise

$$L^* = 116 \left[f \left(\frac{Y}{Y_n} \right) - \left(\frac{16}{116} \right) \right] \quad (28)$$

$$a^* = 500 \left[f \left(\frac{X}{X_n} \right)^{1/3} - f \left(\frac{Y}{Y_n} \right)^{1/3} \right] \quad (29)$$

$$b^* = 200 \left[f \left(\frac{Y}{Y_n} \right)^{1/3} - f \left(\frac{Z}{Z_n} \right)^{1/3} \right] \quad (30)$$

where

$$f(Y/Y_n) = \left(\frac{Y}{Y_n} \right)^{1/3} \quad \text{for } Y/Y_n > 0.008856 \quad (31)$$

$$f(Y/Y_n) = 7.787 \left(\frac{Y}{Y_n} \right) + 16/116 \quad \text{for } Y/Y_n \leq 0.008856$$

and $f(X/X_n)$ and $f(Z/Z_n)$ are defined in the same way, and

$$u^* = 13L^*(u' - u'_n) \quad (32)$$

$$v^* = 13L^*(v' - v'_n) \quad (33)$$

These values represent a comparison to a standard illuminant for sources and an ideal white object illuminated by a standard illuminant for objects.³

The *correlated color temperature* is calculated using interpolation from a table of 30 isothermperature lines. Robertson's method, which uses successive approximation, should be accurate to within 0.1 μ rad, with a maximum error from 1600 to 3000 K of less than 0.2 K plus the measurement uncertainty. This technique should only be used for sources with chromaticities farther than 0.01 from the Planckian locus.³

When using a spectroradiometer or spectrophotometer to measure transmission through a sample, the *Beer-Lambert law*, also known as *Beer's law* is used to calculate the *concentration* of a sample solution. The absorbance A is defined as

$$A = -\log_{10}(\tau) \quad (34)$$

where τ is the transmittance. The Beer-Lambert Law states:

$$A = cl\alpha \quad (35)$$

where c is the concentration, l is the path length, and α is the absorption coefficient.

The absorption coefficient is related to the wavelength by

$$\alpha = \frac{4\pi k}{\lambda} \quad (36)$$

where k is the extinction coefficient.

Figures of Merit

Spectroradiometry measurements have large errors compared to other physical measurements. During an intercomparison of solar ultraviolet monitoring between 14 instruments the measured solar irradiances agreed to within 3 percent when the instruments remained outdoors, but the spectral irradiance responsivities changed upon moving the instruments.⁵ In a 2002 intercomparison study by the project Quality Assurance of Ultraviolet Measurements in Europe (QASUME), the spread of absolute irradiance between spectroradiometers was 12 percent (± 6 percent).⁶ There are two major reasons for the large uncertainties:

- The measurement has many dimensions—it is dependent on the magnitude of the flux, its position on the entrance aperture, its direction, its wavelength distribution, and its polarization.
- The instability of measuring instruments and standards, which are very dependent on room conditions such as temperature, and are frequently off by 1 percent or more.

Potential errors in spectroradiometer measurements include measurement noise, detector instability, wavelength instability, nonlinearity, directional and positional effects, spectral scattering, spectral distortion, polarization effects, and size of source effect.¹

These errors can be characterized as³

- Random noise from the detector, electronics, and light source
- Systematic errors from
 - The measurement of the geometry
 - The calibration, including uncertainty from the calibration standard

Noncosine collection of light
 Stray light
 Nonlinearity of the detector and its electronics
 Dark noise subtraction errors

- Periodic errors from drifts due to temperature, humidity, air movement, electronics, beating of AC sources, and changes in stray light

The way of calculating uncertainty was standardized in 1992 by the International Committee for Weights and Measures (CIPM), and the *Guide to the Expression of Uncertainty in Measurement* was published in 1993.⁷ There are two ways to determine the uncertainty of a component: A. statistically or B. “usually based on scientific judgment using all the relevant information available.”⁸

In spectroradiometry, Type B evaluation finds the upper and lower limits of the value (or correction), and then assumes a probability distribution between these values to obtain the standard uncertainty. If you have nothing to base the probability distribution on, you are instructed to assume that it is rectangular (uniform).

For example, for value a where $a = (a_+ - a_-)/2$, where a_+ is the upper limit and a_- is the lower limit, the standard uncertainty is $a/\sqrt{3}$, but assuming a triangular distribution makes it $a/\sqrt{6}$, and a Gaussian distribution makes it $a/3$.

The collected uncertainties for each identified potential error combines as the square root of the sum of the squares, called the *suggested* or *overall uncertainty*. Often this value is multiplied by a constant, under current international practice of value 2, to form the expanded uncertainty.

CIPM requires that all the standard uncertainties and their derivation are included in the uncertainty report.^{7,8}

If your value for spectral irradiance can be expressed of the form

$$E_{\lambda}^{\text{report}} = E_{\lambda}^{\text{observed}} + c_1 + c_2 + c_3 + \dots + c_n \quad (37)$$

where $E_{\lambda}^{\text{report}}$ is the reported spectral irradiance, $E_{\lambda}^{\text{observed}}$ is the measured value of the spectral irradiance, and the c_i 's are the corrections mentioned earlier in this section, then by the CIPM method, the uncertainty U is calculated from the uncertainties u of the parts as follows

$$U = 2\sqrt{u^2(E_{\lambda}^{\text{observed}}) + u^2(c_1) + u^2(c_2) + u^2(c_3) + \dots + u^2(c_n)} \quad (38)$$

Spectroradiometers are calibrated by use of a standard, which has a known value with a reported error. Assuming your system is linear,

$$E_{\lambda}^{\text{obs}} = \frac{S}{S^s} E_{\lambda}^s \quad (39)$$

where S is the measurement and the superscript s refers to the standard.

Then, the uncertainty of the observed spectral irradiance can be calculated using

$$u(E_{\lambda}^{\text{obs}}) = E_{\lambda}^{\text{obs}} \sqrt{\left(\frac{u(S)}{S}\right)^2 + \left(\frac{u(S^s)}{S^s}\right)^2 + \left(\frac{u(E_{\lambda}^s)}{E_{\lambda}^s}\right)^2} \quad (40)$$

The values for $u(S)$ and $u(S^s)$ are typically calculated by a Type B evaluation, while the standard lamp's uncertainty $u(E_{\lambda}^s)$ is calculated from the uncertainty U reported by the standard lamp's supplier.¹

38.3 GENERAL FEATURES OF SPECTRORADIOMETRY SYSTEMS

There are four parts in every spectroradiometer system

- Input or fore-optics
- A monochromator
- A detector
- Electronics and software to analyze data

There is a fifth aspect to every spectroradiometer system, although it is not usually included on these types of lists, because it is not a “part” of the spectroradiometer. However, I believe it is important to consider it while considering other fundamental parts of the system, because it is essential for the accurate measurement of optical radiation:

- Calibration, usually using standard lamps, reflectance standards, or a standard detector

The Input or Fore-Optics

The input optics gathers light from a specified field of view. The layout determines the quantity which is being measured. For example, when measuring spectral irradiance, the light must be diffused, so an integrating sphere or diffusing plate is used, but when measuring spectral radiance, imaging optics control the solid angle and source area, so a focusing mirror is typically used. Transmittance can be measured by placing a light source at the entrance of the system, and measuring the signal twice: with and without the object to be measured. However, an instrument dedicated to measuring transmittance may have a double beam optical design where the light source is split and recombines at the photodetector.

For measurements in the ultraviolet, or light below 190 nm, the radiation is absorbed by the oxygen in the air, so the whole system will be designed to be enclosed and under vacuum.

Telescoping input optics are used when the sources are large distances from the measurement system, turning a spectroradiometer into a telespectroradiometer. Mounting a microscope to the entrance port of the monochromator can make it possible to measure small radiating sources. Fiber-optic probes can be coupled directly to the monochromator, or in combination with any of the previously mentioned input optical devices.³

The Monochromator

The monochromator is the heart of the spectroradiometric system, because separating the radiation into its component wavelengths is the fundamental aspect of the system. While monochromators used to be made with prisms, they are now always made with diffraction gratings. Monochromators come in different sizes; a large monochromator will be more accurate, but a smaller monochromator can be easier to evacuate to measure the ultraviolet, or place in a dry carbon dioxide-free enclosure to measure the infrared.

The monochromator is designed to collimate and focus light. After the entrance slit, light hits a collimating element. Since light is often diverging when it reaches the slit, a concave mirror can form it into a collimated beam directed at the grating.

Generally, the grating will rotate so this beam hits it at different angles. There are also monochromators with curved gratings, but they are limited in wavelength range.

The grating equation, which defines the wavelength of the diffracted flux to the angle of diffraction, is

$$m\lambda = d(\sin\theta \pm \sin\beta) \quad (41)$$

where m is an integer known as the *order of diffraction*, λ is the wavelength, d is the distance between grooves, θ is the angle of incidence, and β is the angle of diffraction. To remove higher orders ($m > 1$), *blocking filters* that absorb short wavelengths while transmitting long wavelengths are used.

The maximum theoretical groove density is $2/\lambda$, although a practical limit is usually 0.85 of the maximum.

The efficiency of a monochromator is directly proportional to large grating area, short focal length, high groove density, long slit length, and high transmittance.¹ However, there is a trade-off between this efficiency and accuracy, which requires large focal lengths. The f -number is the focal distance divided by the entrance slit. Large f -numbers (> 3 or 4) are required because the mirrors are spherical, not parabolic, and introduce errors.³

The *dispersion* is the width of the band of wavelengths per unit of slit width, in nm/mm. The *band pass* is the spectral interval that may be isolated. If the dispersion at a groove density k is known, then the band pass can be calculated

$$B = (n_k D(n_k) S) / n \quad (42)$$

where B is the bandpass (in nm), n_k (in grooves/mm) is the known groove density where $D(n_k)$ is its dispersion (in nm/mm), S is the slit width (in mm), and n (in grooves/mm) is the groove density of the grating used.

Bandpass should be small for the best precision. However, there is a trade-off between bandpass and *geometrical entendue* G the light gathering power of an optical system.

$$G = \frac{hnmG_A B}{F 10^6} \quad (43)$$

where h is the slit height (mm), n groove density (groove/mm), m is order of the grating, G_A area of the grating (mm²), B is the bandpass (nm), and F is the focal length (mm). The ratio h/F implies that the etendue may be increased by making the height of the entrance slit larger. However, this does not work as well in practice; increasing the height of the slit will increase stray light and may also increase the system aberrations, reducing the bandpass.

Geometrical entendue is a limiting function of system *throughput*.⁹

High signal, for which high throughput is necessary, is limited by bandpass. Sometimes the slit size is determined by other factors, like the field of view. But the slit size might be chosen to optimize other factors. In this case, monochromatic sources behave differently from broadband sources, and mixed sources are a combination of the two. For a fluorescent lamp, which is a mixed source, as you decrease the bandpass, the peaks due to the monochromatic lines become much higher in proportion to the broad emission spectra of the phosphors.³

In night vision systems, stray light becomes an important factor which affects system design. To limit stray light, a double monochromator system is used, where the output of the first monochromator is the input of the second one. This can reduce typical stray light levels from 10^{-4} to 10^{-8} .³

Some spectroradiometers are designed with multichannel detectors within the monochromator, or so multichannel detectors can be easily installed. Multichannel detectors can also eliminate the need to scan, reducing moving parts and allowing for longer integration times or quick measurements of unstable or short-lived sources. However, they are not suitable for all types of spectral irradiance and radiance measurements, and spectral transmission, reflection and responsivity require the monochromatic light to exit the monochromator and interact with samples.

The Detector

The desired wavelength range will strongly influence the type of detector used. In Table 1 the approximate wavelength ranges of different spectroradiometry detectors are shown. Other important factors in choosing a detector include the dynamic range, sensitivity, and response time required for the data, as well as environmental factors determining how rugged a detector is needed.

TABLE 1 Approximate Wavelength Ranges of Different Spectroradiometry Detector Types

Detector	Wavelength Range
PMT	200–850 nm
Si photodiode	200–1100 nm
Ge photodiode	1100–1800 nm
InGaAs photodiode	850–1700 nm
PbS photoconductor	1–4.5 μm
PbSe photodiode	1–4.5 μm
InSb photodiode	1.5–5 μm
Pyroelectric	500 nm–50 μm
CCD	200–1100 nm
InGas PDA	800–1700 nm

The two types of detectors most commonly used in spectroradiometry are photomultiplier tubes and semiconductor devices, although thermal detectors have some very limited applications. Detector sensitivity is measured in noise equivalent power (NEP) or equivalent noise input (ENI), which basically mean the same thing, the minimum detectable signal, whose units can be taken in watts.¹ The detectivity is the inverse of the NEP.

Photomultiplier tubes are the most sensitive detectors when used in their wavelength region, with ENIs ranging from 10^{-15} W at 1100 nm to 10^{-17} W from 850 nm to 200 nm and 5×10^{-16} W at 110 nm. They are usually used from 200–850 nm, their range of greatest sensitivity. Silicon photodiodes have NEPs reported of 2×10^{-14} W at 1100 nm, 10^{-15} W at 850 nm, 5×10^{-15} W at 350 nm, and 10^{-14} W between 300 and 200 nm. Silicon photodiodes are used in these wavelength ranges when the signal is sufficiently large, because they are more temperature stable and more rugged.¹

From 1100 to 1800 nm, germanium photodiodes are the most sensitive, with NEPs in the 10^{-13} W range, and lead sulfide photoconductors the most sensitive from 1800–3800 nm, with NEPs ranging from 4×10^{-13} W to 4×10^{-12} W, although InAs may also be used in this range.¹

Photomultipliers, germanium photodiodes, and PbS photoconductors would be cooled, but silicon photodiodes can operate at room temperature (25°C).

Thermal detectors are less sensitive, having NEPs of about 6×10^{-10} W, and have a flat response.¹

When comparing detectors, it is also common to compare the normalized detectivity. There is some confusion in the nomenclature, in that sometimes this is D^* , and sometimes D^* refers to the specific detectivity, which has a different definition. Therefore, I will refer to the normalized detectivity as D_N .

$$D_N = D\sqrt{AB_w} \quad (44)$$

where D is the detectivity, A is the detector area, and B_w is the detector bandwidth.

For example, the D_N of an InAs photodiode is roughly the same as that of a PbS photoconductor at 3000 nm, and the rise time for InAs is one thousandth of PbS. At that specific wavelength, the two detector types have equivalent sensitivity, but InAs has a faster response. However, the D_N drops off more rapidly for InAs at longer and shorter wavelengths than for PbS.¹

Multichannel detectors are sometimes used for spectroradiometry. They have the advantage in spectral irradiance and radiance measurements that they can eliminate the need to scan, reducing moving parts and allowing for longer integration times, as well as making measurements of nonstable sources or short-lived, such as flashbulbs, explosions, and solar measurements during changing weather conditions, possible. Silicon photodiode linear arrays (PDAs) and charge-coupled detectors (CCDs) both singly and with microchannel plates (intensified arrays) are used, as well as a combination of microchannel plates with resistive film called resistive anode (MCP-RA's).¹ These detectors are temperature dependent and therefore require cooling. Noise is also a factor.

The NEP of array detectors is measured with respect to integration time. This is because array detectors operate in the capacitive-discharge mode, which collects charge for a period of time before

TABLE 2 A Comparison of the NEP in Watts of Five Different Multichannel Detector Types with 5-Second Integration Times.

Detector Type	250 nm	550 nm	850 nm	1100 nm
Si PDA at -40°C	1.4×10^{-15}	2.7×10^{-16}	2.2×10^{-16}	1.3×10^{-15}
CCD at -110°C	6.0×10^{-18}	1.6×10^{-18}	1.0×10^{-18}	2.5×10^{-17}
Intensified Si PDA at -40°C	1.3×10^{-18}	8.2×10^{-19}	1.1×10^{-17}	
Intensified CCD at -110°C	1.3×10^{-18}	8.3×10^{-19}	1.1×10^{-17}	
MCP-RA at -30°C	2.0×10^{-18}	1.3×10^{-18}	7.3×10^{-18}	

This data is compiled from Tables 15.4, 15.6, 15.7, 15.9, and 15.11 of Ref. 1.

discharging, which will be discussed in more detail in the next section on electronics and software. However, in Table 2 you can see a comparison of the NEP of the five different types of detectors at integration times of 5 seconds. It can be noted that CCD detectors have NEPs comparable with photomultiplier tubes, even without the advantage of multichannel detection.

While the detectors compared above were all cooled, there exist many commercial spectroradiometers with multichannel detectors which are not cooled for applications where sensitivity is less important, such as projector calibration and film and video post production.

Signal to noise ratio is also a major concern when comparing detector types, but that is also dependent on the integration time. As a general rule, the signal to noise ratio is much lower for PDAs than other multichannel detectors, but at 0.5 second integration times, CCDs are superior.¹

Spectral scattering is a more important concern with multichannel detectors than with single detectors. Interference filters are used to block the short wavelength scattered flux.

The Electronics and Software

In “The Detectors” section we have discussed photomultiplier tubes, semiconductor detectors (photodiodes and photoconductors), PDA-based detectors, and CCD-based detectors.

Photomultiplier tubes use a 250- to 2500-V power supply, and the gain of the detector is adjusted by changing the voltage of the power supply. The signal is the anode current, and the simplest way of measuring it is measuring the voltage drop over a load resistor. However, many problems can be eliminated by using an operational amplifier with a feedback resistor instead, creating a transimpedance current to voltage converter. For even lower signal to noise ratios, integrated current is used as the signal, and the amplifier is chosen to have a short enough time constant that the entire anode current is integrated during the entire sampling time. It is important to make sure that the integrating time is properly matched with the scanning time.

Photodiodes are usually used in the photovoltaic or unbiased mode to minimize noise. In this mode a voltage amplifier with a high-input impedance measures the voltage generated across the photodiode.

Photoconductors decrease in resistance when absorbing radiant flux. The detector is placed in series with a load resistor and a bias voltage. As the resistance of the detector decreases there is an increase in the voltage over the load resistor.

Lock-in amplifiers are often used with semiconductor detectors to decrease $1/f$ noise.

PDA-based detectors operate in the capacitive-discharge mode where the photodiodes store charge as well as sense light. The diodes are reverse biased for 5 V and then electrically isolated while being exposed to light. After an integration time, the charge required to bring the diode back to 5 V is measured.

CCD-based detectors have arrays of pixels that convert radiation to charge. Parallel gates between the pixels can release the charge at will. The charge is then moved along the channels of the arrays by changing the voltage at the gates, until it reaches an output shift register at the end of the channel. The charge is then amplified and measured. Sometimes neighboring pixel's values are intentionally combined, which is known as *binning* or *pixel summation*.

While the readout of a single detector's signal can be as simple as a digital voltmeter, in most cases spectroradiometric data will be recorded by a computer. Often it will be part of an automated system that also controls the drive mechanism. After the information is recorded by the computer, it is analyzed by the computer software. Examples of typical calculations performed on spectroradiometric data are in Sec. 38.2 "Definitions, Calculations, and Figures of Merit" under Calculations.

Calibration

To calibrate a spectroradiometer for irradiance and radiance, standard lamps are usually used, although standard detectors are available. While spectroradiometric standards available for use over the visible are based on the spectral radiance of a blackbody as defined by Planck's radiation law, most commercially available blackbodies are used primarily in the infrared at wavelengths above 1000 nm. Blackbodies suitable for the visible are very expensive because they must operate at temperatures of 2500 K or higher, and are not practical for normal laboratory calibrations.³

Standard lamps introduce a certain degree of error; in the ultraviolet intercomparison in 2002 where the error was 12 percent (± 6 percent), 6 percent of that error was attributed to the calibration lamps.⁶ In earlier intercomparisons, it was more like 1.4 percent from 250 to 2400 nm, but 3 to 4 percent in the infrared.¹

Different lamps are used to calibrate over different wavelength regions. From wavelengths of 250 to 2400 nm, tungsten lamps are used; below 250 nm, a deuterium lamp, argon miniarc, or synchrotron radiation are used.¹

In Table 3 the wavelength range of frequently used spectral irradiance and spectral radiance source standards is listed.

For a discussion of the history of calibration standards, and some additional standards see reference.³

The FEL 1000-W tungsten lamp is the most widely used source standard for spectral irradiance. It is a commercially available 1000-W clear quartz envelope tungsten-halogen, coiled coil filament lamp that is modified to a medium bipost base with 1/4-inch diameter stainless steel posts. It is 5 inches high with a filament about 1 inch long and 1/4 inch in diameter. It is operated at 8-A DC and about 120 V. It is mounted base down with the steel post vertical and the optic axis of the spectroradiometer is horizontal. For more details of the calibration procedure, see reference.¹

It is important to recognize that any individual standard lamp may develop problems. It is recommended to have three standard lamps, so if one changes significantly, you still will have two that agree. It is also standard procedure to transfer the calibration of your newly acquired standards to working standards, lamps that you calibrate using your detector and the NIST traceable standards that have been shipped to you. These lamps can be commercially acquired, but then you must age them for 40 hours at 120-V DV and check them for stability; they must have a drift of less than 0.5 percent at 650 nm in 24 hours. The working standard should be compared to your three purchased standards after 50 to 100 hours of use.

The responsivity of the spectroradiometer can be modeled from a measurement equation, although it is necessary to actually calibrate any system.

TABLE 3 The Wavelength Range of Frequently Used Spectral Irradiance and Spectral Radiance Source Standards

Irradiance	
Tungsten FEL lamp	250–2400 nm
Deuterium lamp	165–350 nm
Argon mini-arc	90–350 nm
Radiance	
Tungsten strip lamp	225–2400 nm
Blackbody below 1000 K	1000–4000 nm
Argon miniarc	90–350 nm

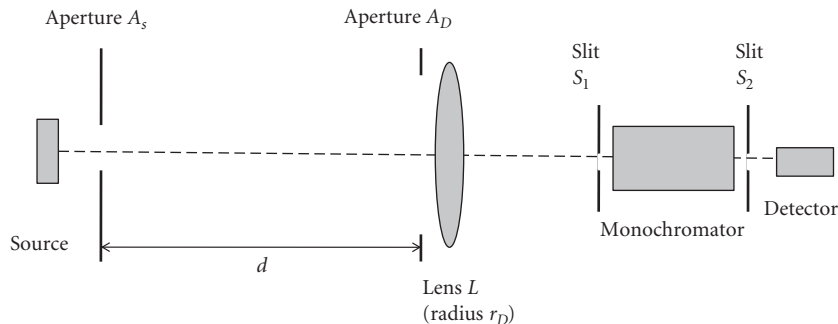


FIGURE 1 A schematic of the experimental setup to measure the spectroradiometer's responsivity.

If the monochromator is set to a wavelength λ_0 , the spectral responsivity $R(\lambda_0, \lambda)$ can be separated into two parts: the slit scattering function $\rho(\lambda_0, \lambda)$ which is an envelope around wavelength λ_0 and the overall responsivity as a function of wavelength $R^f(\lambda)$

$$R(\lambda_0, \lambda) = \rho(\lambda_0, \lambda) R^f(\lambda) \quad (45)$$

When put in an experimental setup shown in Fig. 1, the signal response for a monochromator, ignoring small corrections, is

$$r(\lambda_0) = \frac{A_s A_D}{(r_s^2 + r_D^2 + d^2)} \int_{\lambda} L_{\lambda}(\lambda) \rho(\lambda_0, \lambda) R^f(\lambda) d\lambda \quad (46)$$

where A_s is a circular aperture in front of a standard radiometric source of radius r_s , A_D is an aperture a distance d away from the source in front of a focusing lens L of radius r_D , and $L_{\lambda}(\lambda)$ is of the Lambertian calibration source.¹⁰

If the constants are collected in a term C , and a suitably averaged luminance is used, the signal response can be calculated by

$$r(\lambda_0) = CL_{\lambda}(\lambda_0) \int_{\lambda} \rho(\lambda_0, \lambda) R^f(\lambda) d\lambda \quad (47)$$

Detector standards are also available. In this case, the assumption is that any drift is not very wavelength specific. It is possible to interpolate between calibration laser wavelengths using a blackbody,¹⁰ but that is not usually done, because it is very time consuming and probably an insignificant difference within the larger sources of error.¹ Detectors that presently are suitable for use in transfer standards include silicon, germanium, and InGas photodiodes, and certain types of thermal detectors. The basic approach is to have the two detectors measure the same sources. A typical source is a continuous wave (CW) laser directed at the entrance aperture of an integrating sphere. Types of lasers which can be used are helium-neon, argon, krypton, helium-cadmium, Nd:Yag, and Ti:sapphire.¹¹ When using a detector standard to calibrate a spectroradiometer's spectral irradiance, the area of the entrance aperture and spectral slit width become important and must be taken into account, so it is not a simple measurement.¹

Spectroradiometers are used for reflectance measurements as well as irradiance and radiance measurements. For these purposes, reflectance standards are used for calibration. A number of reflectance standards are available which have been developed for spectroscopy applications, as well as to calibrate colorimeters and spectrophotometers. Specular reflectance standards are calibrated mirrors, and diffuse reflectance standards are made of material similar to the inside of integrating spheres.

Labsphere makes Spectralon into a diffuse white standard and a selection of diffuse gray and color standards, which are calibrated and NIST traceable. Halon, a trade name for polytetrafluoroethylene

(PTFE) powder (which is also used to coat integrating spheres) is also used to make NIST traceable reflectance standards. In a round-robin intercomparison of bidirection diffuse reflectance (BRDR) four types of diffuse reflectors (spectralon, halon, sintered halon, and vacuum deposited aluminum on a ground aluminum surface) were measured at five laboratories. These four types of standards were chosen because of their different scattering mechanisms; Spectralon and pressed PTFE scatter from the bulk, aluminum scatters from the surface, and sintered PTFE scatters from both the bulk and the surface. The purpose of this experiment was to test the laboratories, not the standards, but there was general agreement with the NIST specifications to 2 percent.¹²

38.4 TYPICAL SPECTRORADIOMETRY SYSTEM DESIGNS

Spectral Irradiance and Radiance

The input optics for spectral irradiance measurements require a diffuser in order to eliminate or reduce directional, positional, and polarization effects. This diffuser can be an integrating sphere coated with Halon, which is the best choice for sunlight and large or irregularly shaped sources, a plane reflector diffuser coated with BaSO_4 , or a transmitting diffuser made of teflon.

In Fig. 2 we see a block diagram of a spectroradiometer designed to measure spectral irradiance. The basic steps that have to be taken are the signal must be diffused (cosine corrected), then wavelength selected by the monochromator, after which it is detected by the detector, amplified and analyzed.

In certain special cases, such as when measuring the spectral irradiance of point sources or columnated sources, and if the spectroradiometer responds uniformly over the angular field viewed, no input optics are necessary. However, most of the time input optics are necessary. In Fig. 3 we see a typical setup of the input optics for measuring the irradiance of a large source using a small integrating sphere and a spherical mirror. A typical small sphere may have a 2.5-cm outer diameter with a 3-mm-thick PTFE coating (or BaSO_4 in the 310 to 350 nm wavelength region, where PTFE fluoresces weakly), and a circular entrance port of 1 cm^2 and a rectangular exit port with dimensions

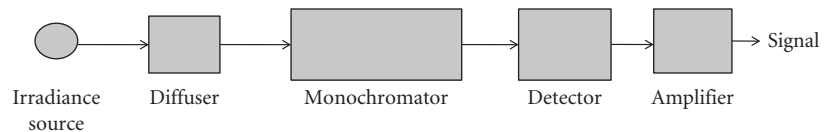


FIGURE 2 Block Diagram to measure spectral irradiance.

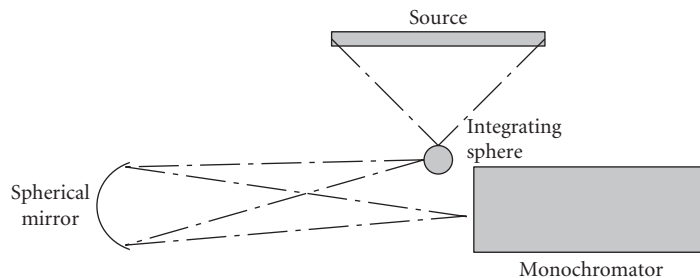


FIGURE 3 Typical setup of the input optics for measuring the irradiance of a large source using a small integrating sphere and a spherical mirror.

3 mm × 12 mm. The radiation which has one reflection does not get into the monochromator because it is reflected away from the spherical mirror. The radiation that does reach the monochromator is independent of the direction and position of the incident flux. However, it is worth noting that the attenuation of a system of this type is large.

In Fig. 4 we see more compact input optics to measure irradiance. The disadvantages are that the sphere is harder to reach, position, and orient, and more stray flux reaches the monochromator.

In Fig. 5 we see typical input optics using a plane diffuser instead of an integrating sphere. In Fig. 6 we see typical input optics using a transmitting diffuser.

In Fig. 7 we see the schematic of a single monochromator and Fig. 8 we see the schematic of a double monochromator. Both of these designs assume an external detector outside the exit slit. In Fig. 9 we see the layout of a single grating monochromator with a built-in multichannel detector.

The input optics for measuring spectral radiance form an image of the source on the entrance port of the monochromator. Possible geometries include a plane and spherical mirror which focus

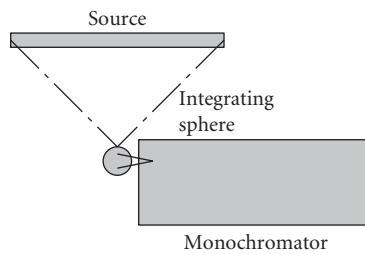


FIGURE 4 Measuring irradiance directly using an integrating sphere.

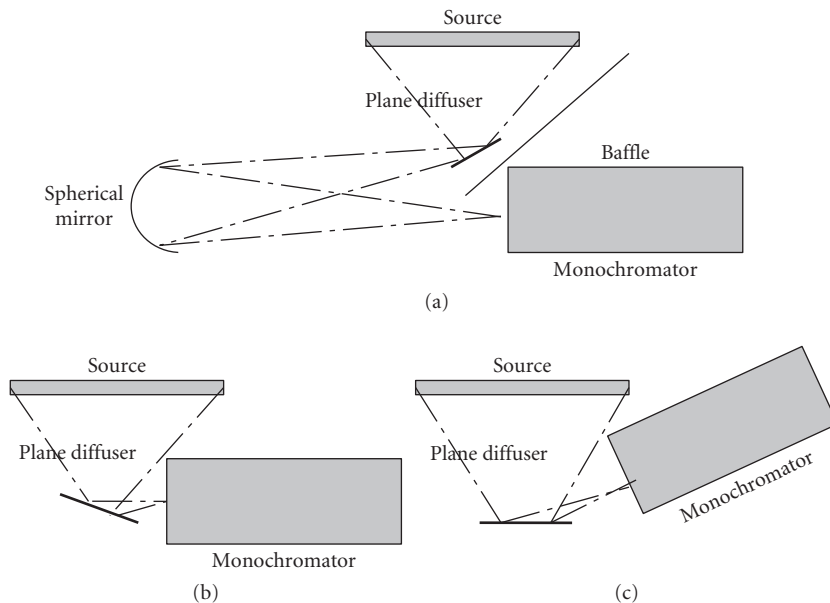


FIGURE 5 Typical setups using a plane diffuser to measure irradiance: (a) with a spherical mirror, (b) and (c) directly reflected off the diffuser into the sphere.

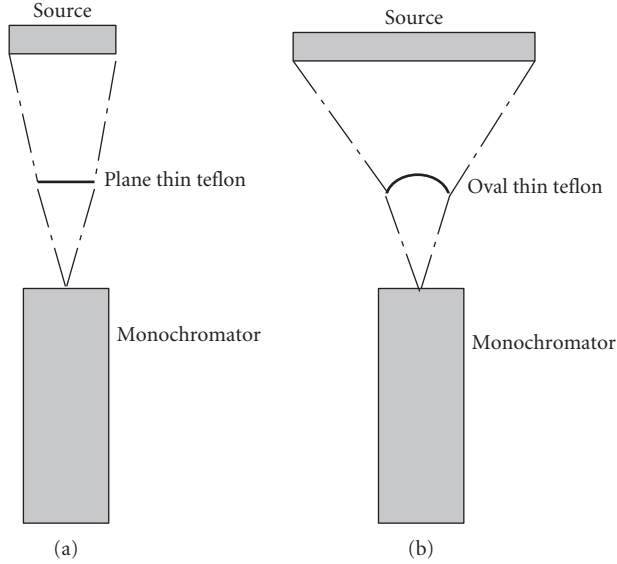


FIGURE 6 Typical input optics for irradiance measurements using a (a) plane thin teflon diffuser and (b) oval shaped thin teflon diffuser.

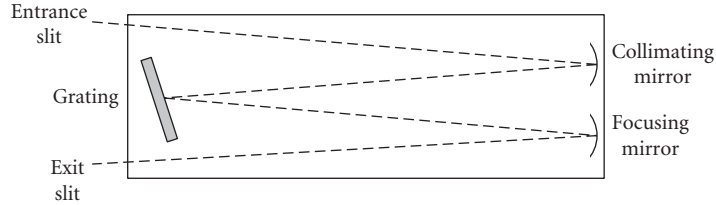


FIGURE 7 Schematic of a single monochromator.

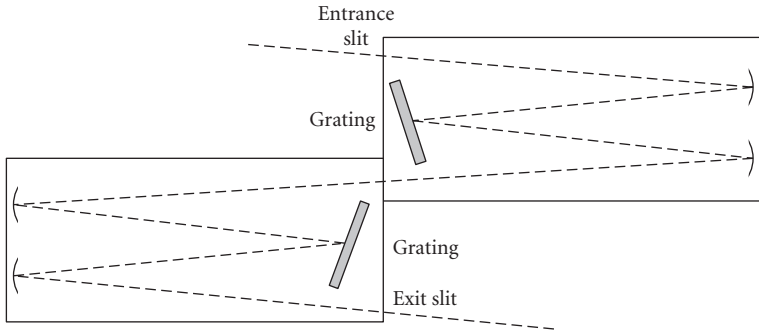


FIGURE 8 Schematic of a double monochromator.

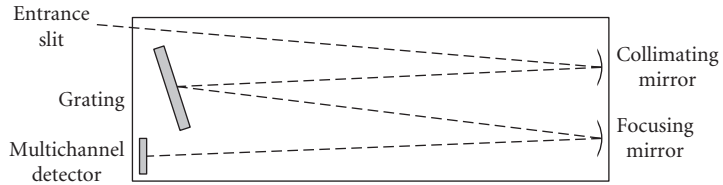


FIGURE 9 Schematic of a single monochromator with a built-in multichannel detector.

the radiation on the entrance slit, and a field of view baffle attachment to limit the acceptance angle of the monochromator.

In Fig. 10 we see the simplest input optics for measuring spectral radiance. However, it is much more likely you are using a system that measures both irradiance and radiance. In that case, your input optics are more likely to look like Fig. 11, where a mirror will reflect the light away from the integrating sphere, and toward the spherical mirror.

The same monochromator and detection electronics would be used for irradiance and radiance measurements.

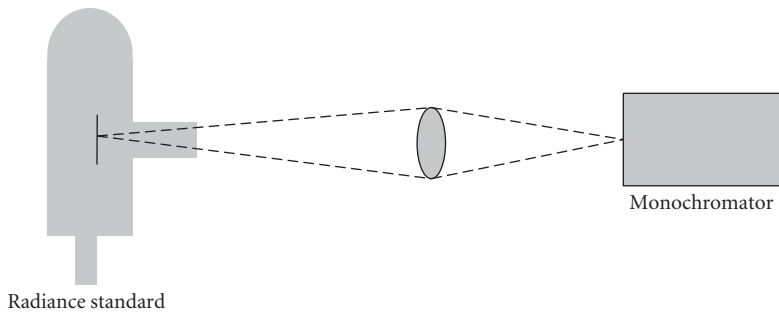


FIGURE 10 Measuring spectral radiance.

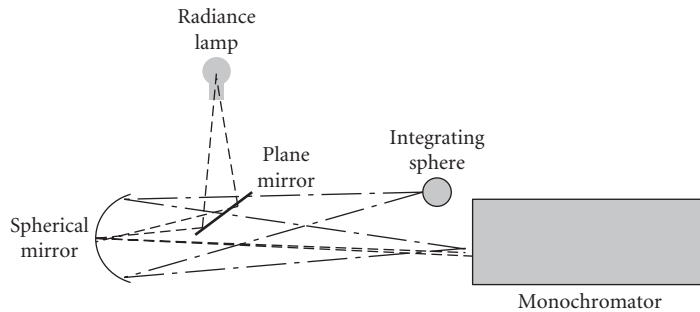


FIGURE 11 System for measuring both irradiance and radiance. The radiance setup adds a plane mirror and measures the radiance lamp. The irradiance setup removes the plane mirror and measures light after it passes the integrating sphere.

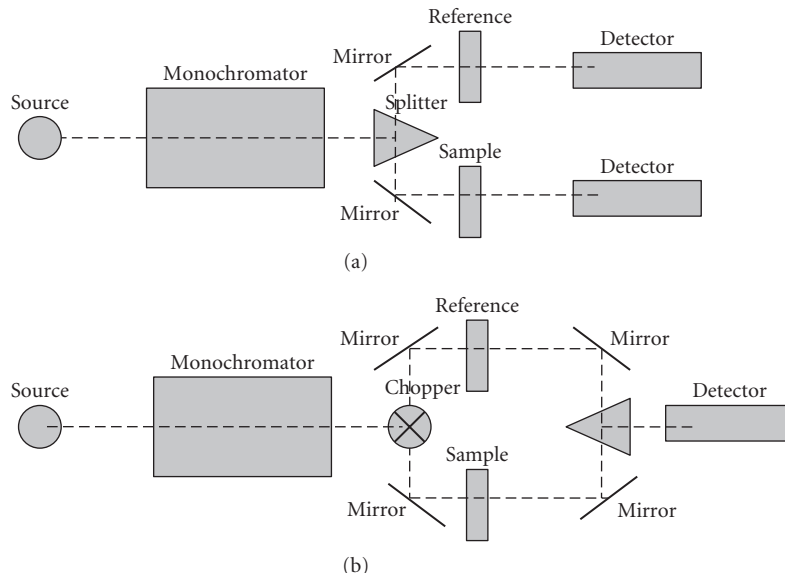


FIGURE 12 Spectrophotometer measuring spectral transmission: (a) dual-beam design and (b) double beam design.

Spectral Transmittance

Spectral transmittance measurements use the source and monochromator to create monochromatic radiation, rather than to measure it. The goal is to pass the radiation through a sample, to measure the sample's properties. This application is of vast importance to biology and chemistry, and there are many commercially available spectroradiometers dedicated to measuring transmittance and reflectance, known as *spectrophotometers*.

Chemists use absorbance spectroscopy to obtain qualitative and quantitative information about samples, using the Beer-Lambert law, also known as Beer's law, as was discussed in Sec. 38.2. Most spectrophotometers use a dual or double beam configuration, as is shown in Fig. 12. This experimental design measures regular transmittance, the signal that passes directly through the sample without being scattered and follows Snell's law. The output is the ratio of the signal in the sample beam to the signal in the reference beam with respect to wavelength. It is necessary to ensure that the only difference between the sample beam and the reference beam is the quantity to be measured, which implies that liquid cells with equal amounts of solute, or gas cells with equal amount of carrier gas, should be placed in the reference beam.¹³

Because of fluorescence, broadband illumination may have different results than monochromatic illumination. This should be considered when measuring transmittance and the setup should approximate the same manner in which the material will be used. Total spectral transmission, which is a combination of regular and diffuse transmission requires the addition of an integrating sphere after the beam transmits through the sample.³

Spectral Reflectance

Manufacturers use spectral reflectance information to provide color information about inks and textiles. There exist several spectral libraries (USGS,¹⁴ Johns Hopkins,¹⁵ JPL¹⁶) which contain almost 2000 spectra of powdered materials for use in spectroscopy, measured by using spectrophotometers and spectrometers in the diffuse reflectance mode.

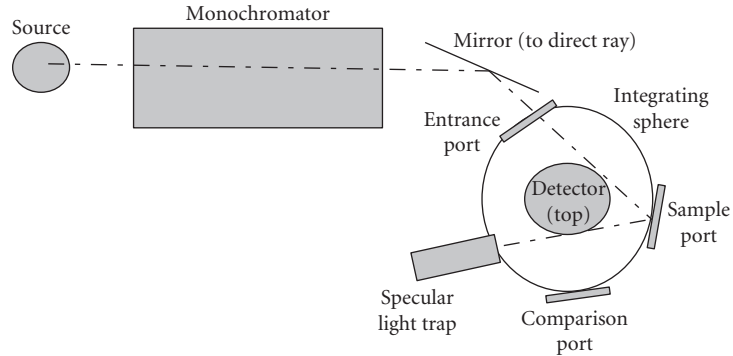


FIGURE 13 Schematic of a spectroradiometer measuring diffuse spectral reflectance.

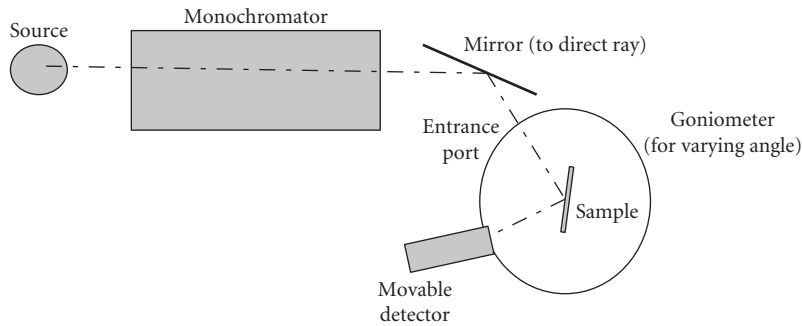


FIGURE 14 A schematic of a spectroradiometer measuring specular spectral reflectance.

Figure 13 is a spectroradiometer configured for measuring diffuse spectral reflectance. There is a double beam design, allowing for the comparison method of using a diffuse reflectance standard. An integrating sphere collects the diffuse radiation, with a removable light trap allows the user to block specular reflectance. If the light trap is not in place, total reflectance, rather than diffuse reflectance, is measured. The sample to be measured is placed in the sample port, and the standard is placed in the comparison port. The detector is perpendicular to the samples and the incident radiation.

Figure 14 is a spectroradiometer configured for measuring specular spectral reflectance. Specular reflectance can be measured at various angles of incidence, including 0° for a 100 percent reading. A calibrated mirror (specular reflectance standard) is not necessary in this design.

Spectral Responsivity

In 38.3 General Features of Spectroradiometer Systems under “Calibration” we discussed how to calculate the spectral responsivity (see Fig. 1) of the spectroradiometer. Once that is known, spectroradiometers can be used to find the spectral responsivity of other detectors.

A spectroradiometer can be configured to measure detector spectral responsivity. The first step uses the spectroradiometer’s standard detector to measure the monochromatic flux or irradiance of the source, generating a function $r^s(\lambda)$ for the standard detector, which has a known responsivity of $R^s(\lambda)$.

These values should not be confused with the responsivity of the total spectroradiometer system. This responsivity of the spectroradiometer system includes both the scatter from the slits and the detector responsivity. In this case we are singling out the detector responsivity. A NIST traceable standard silicon detector is usually used for measurements over the visible spectrum.³

Then the detector is replaced with the detector to be tested, and a signal response of $r^t(\lambda)$ is generated. The responsivity of the test detector $R^t(\lambda)$ will be

$$R^t(\lambda) = r^t(\lambda) / \Phi(\lambda) = r^t(\lambda) R^s(\lambda) / r^s(\lambda) \quad (48)$$

38.5 REFERENCES

1. H. J. Kostkowski (1997) *Reliable Spectroradiometry*, La Plata, MD: Spectroradiometry Consulting.
2. *A Guide to Spectroradiometry: Instruments & Applications for the Ultraviolet*, Reading: Bentham Instruments (1997) (p. 24—tanning booth reference).
3. W. E. Schneider and R. Young (1997) “Spectroradiometry Methods,” *Handbook of Applied Photometry*, Casimer DeCusatis (ed.), Chap. 8, pp. 239–287. New York: AIP Press.
4. J. D. Schanda (1997) “Colorimetry,” *Handbook of Applied Photometry*, C. DeCusatis (ed.), Chap. 10, p. 347. New York: AIP Press.
5. E. Early, A. Thompson, C. Johnson, J. DeLuisi, P. Disterhoft, D. Wardle, and E. Wu, et al., “The 1995 North American Interagency Intercomparison of Ultraviolet Monitoring Spectroradiometers,” *Journal of Research of the National Institute of Standards and Technology* **103**: 15 (1997).
6. A. R. Webb and D. Cotton (2002) “Report of Ispra Intercomparison, May 2002” *Assurance of Ultraviolet Measurements in Europe (QASUME)*. lap.physics.auth.gr/qasume/Files/EvalPdfs/QASUMEREPORT.pdf. Accessed May 19, 2009.
7. ISO (1993) *Guide to the Expression of Uncertainty in Measurement*, International Organization for Standardization, 1, rue de Varembe, Case postale 56, CH-1211 Geneve 20, Switzerland.
8. B. N. Taylor and C. E. Kuyatt (1993) *Guidelines for Evaluating and Expressing the Uncertainty of NIST Measurement Results*, NIST Technical note 1297, Washington, D.C.: U.S. Government Printing Office.
9. J. M. Lerner and A. Thevenon (1988) *The Optics of Spectroscopy: A Tutorial*, vol. 2.0, pp. 1–56. Jobin-Yvon: Instruments SA Inc.
10. R. U. Datla and A. C. Parr (2005) “Introduction to Optical Radiometry,” *Optical Radiometry*, A. C. Parr and R. U. Datla (eds.), Chap. 3, pp. 97–154. New York: Elsevier/AIP Press.
11. L. P. Boivin (2005) “Realization of Spectral Responsivity Scales,” *Optical Radiometry*, A. C. Parr and R. U. Datla (eds.), Chap. 3, pp. 97–154. New York: Elsevier/AIP Press.
12. E. A. Early, P. Y. Barnes, B. C. Johnson, J. J. Butler, C. J. Bruegge, S. F. Biggar, P. R. Spyak, and M. M. Pavlov, “Bidirectional Reflectance Round-Robin in Support of the Earth Observing System Program,” *Journal of Atmospheric and Oceanic Technology* **17**: 1077 (2000).
13. J. M. Palmer (2001) “The Measurement of Transmission, Absorption, Emission and Reflection,” *Handbook of Optics*, M. Bass (ed.), I.25.5. New York: McGraw Hill.
14. R. N. Clark, G. A. Swayze, A. J. Gallagher, T. V. V. King, and W. M. Calvin (1993) “The U. S. Geological Survey, Digital Spectral Library: Version 1: 0.2 to 3.0 microns,” *U.S. Geological Survey Open File Report 93–592*, 1340 pages, <http://speclab.cr.usgs.gov>. Accessed May 19, 2009.
15. J. W. Salisbury, L. S. Walter, N. Vergo, and D. M. D’Aria (1992) *Infrared (2.1–25 μm) Spectra of Minerals*, Baltimore, MD: The Johns Hopkins University Press.
16. C. I. Grove, S. J. Hook, and E. D. Paylor (1992) *Laboratory Reflectance Spectra of 160 Minerals, 0.4 to 2.5 Micrometers*. JPL-Publication 92-2. Pilot Land Data System. Pasadena, California: Jet Propulsion Laboratory.

This page intentionally left blank.

NONIMAGING OPTICS: CONCENTRATION AND ILLUMINATION

William Cassarly

*Optical Research Associates
Pasadena, California*

39.1 INTRODUCTION

Nonimaging optics is primarily concerned with efficient and controlled transfer of radiation. *Nonimaging* refers to the fact that image formation is not a fundamental requirement for efficient radiation transfer; however, image formation is not excluded and is useful in many cases. The two main aspects of radiation transfer that nonimaging optics attempts to solve are maximizing radiation transfer and creating a controlled illuminance distribution. These two problem areas are often described as *concentration* and *illumination*. Many systems require a blending of the concentration and illumination aspects of nonimaging optics, with an important example being the creation of a uniform illuminance distribution with maximum collection efficiency.

Solar collection is an area in which one of the dominant requirements is to maximize the concentration of flux collected by a receiver of a given size (e.g., the irradiance). Any small patch of area on the receiver surface can collect radiation from a hemispherical solid angle. A related problem exists in the detection of low levels of radiation (e.g., Cherenkov counters,¹ infrared detection,² and scintillators³). Nonimaging optics provides techniques to simultaneously maximize the concentration and collection efficiency. The compound parabolic concentrator⁴ is the most well known nonimaging device and has been extensively investigated over the past 40 years.

In addition, there are situations in which collected solid angle must be less than a hemisphere, but the flux density must still be maximized. For a given solid angle and receiver size, the problem is then to maximize the average radiance rather than the irradiance. For example, many fibers only propagate flux that enters the lightpipe within the fiber's numerical aperture (NA), and minimizing the size of the fiber is a key aspect to making practical fiber optic lighting systems.⁵ Projector systems often require maximum flux density at the "film gate," but only the flux reimaged by the projection lens is of importance.^{6,7}

The collected flux must often satisfy certain uniformity criteria. This can be especially important when the luminance of the source of radiation varies or if the flux is collected in a way that introduces nonuniformities. Nonimaging optics provides techniques to control the uniformity. Uniformity control can sometimes be built into the optics used to collect the radiation. In other cases, a separate component is added to the optical system to improve the uniformity.

Uniformity control tends to be created using two main approaches: tailoring and superposition. Tailoring transfers the flux in a controlled manner and uses knowledge of the source's radiance distribution.⁸ Superposition takes subregions of a starting distribution and superimposes them to provide an averaging effect that is often less dependent on the details of the source than tailoring. In many systems, a combination of tailoring and superposition is used. Examples of nonimaging approaches to uniformity control include mixing rods, lens arrays, integrating spheres, faceted reflectors, and tailored surfaces.

The emphasis of this chapter is on the transfer of incoherent radiation. Many of the techniques can be applied to coherent radiation, but interference and diffraction are only mentioned in passing.

39.2 BASIC CALCULATIONS

Some of the terms used in nonimaging optics have yet to be standardized. This section is an effort to collate some of the most common terms.

Photometric and Radiometric Terminology

Understanding nonimaging optics can sometimes be confusing if one is not familiar with photometric and radiometric terminology. This is especially true when reading papers that cross disciplines because terminology usage has not been consistent over time. Reference 9 discusses some of these terminology issues.

This chapter is written using photometric terminology. Photometry is similar to radiometry except that photometry weights the power by the response of the human eye. For readers more familiar with radiometric terminology, a chart showing the relationship between some of the most common photometric terms and radiometric terms is shown in Table 1. Reference 10 describes radiometry and photometry terminology in more detail.

Exitance and *emittance* are similar terms to *irradiance*; however, they denote the case of flux at the surface of a source, whereas irradiance applies to any surface. There are numerous textbooks and handbooks on radiometry and photometry.^{11–16}

Etendue

Etendue describes the integral of the area and the angular extents over which a radiation transfer problem is defined. Etendue is used to determine the trade-off between the required area and angular extents in nonimaging optic designs. Reference 17 provides a brief review of various etendue descriptions with copious references to other work.

TABLE 1 Photometric and Radiometric Terminology

Quantity	Radiometric	Photometric
Power or flux	Watt (W)	Lumen (lm)
Power per unit area	Irradiance, W/m ²	Illuminance, lm/m ² = lux (lx)
Power per unit solid angle	Radiant intensity, W/sr	Luminous intensity, lm/sr = candela (cd)
Power per unit solid angle per unit projected area or Power per unit projected solid angle per unit area	Radiance, W/m ² -sr	Luminance, cd/m ²

One definition of etendue is

$$\text{etendue} = n^2 \iint \cos(\theta) dA d\Omega \quad (1)$$

where n is the index of refraction and θ is the angle between the normal to the differential area dA and the centroid of the differential solid angle $d\Omega$.

In phase space nomenclature (e.g., Ref. 4, Appendix A), etendue is described by

$$\text{etendue} = \iint dx dy dL dM = \iint dx dy dp dq \quad (2)$$

where dL and dM are differential changes in the direction cosines (L, M, N), $dx dy$ is the differential area, and $dp dq$ is the differential projected solid angle within the material of index n . The term *phase space* is often used to describe the area and solid angle over which the etendue integral is performed.

Luminance

Luminance divided by the index of refraction squared is the ratio of the differential flux $d\Phi$ to the differential etendue:

$$L/n^2 = d\Phi/d \text{ etendue} = d\Phi/[n^2 \cos(\theta) dA d\Omega] \quad (3)$$

The L/n^2 over a small area and a small angular extent is constant for a blackbody source (e.g., Ref. 18, p. 189). A consequence of constant L/n^2 is that if optical elements are added to modify the apparent area of the small region, then the angular extent of this small area must also change.

If the source of radiation is not a blackbody source, then flux that passes back to the source can either pass through the source or be reflected by the source. In this case, L/n^2 can increase. One example occurs with discharge sources where a spherical mirror is used to reflect flux back through the emitting region. Another example is observed by evaluating the luminance of tungsten in a coiled filament. The luminance is higher at the filament interior surface than the exterior surface because the interior filament surface emits radiation and also reflects radiation that is emitted from other coils.

Lambertian

In many situations the luminance of a radiation source does not vary as a function of angle or position. Such a source is often called a *Lambertian radiator*. In some nonimaging literature, the term *isotropic* is used. The context in which *isotropic* is used should be evaluated because *isotropic* is sometimes used to describe isotropic intensity instead of isotropic luminance.

If constant L/n^2 can be assumed for a given system, then etendue is an important tool for understanding the trade-offs between angular and spatial distributions. Etendue has been used to understand the limits to concentration,⁴ projection display illumination,¹⁹ and backlit display illumination.²⁰ Reference 21 investigates the situation where there is a spectral frequency shift.

In the imaging community, etendue conservation arises in many situations and is often described using the Lagrange invariant. Because imaging systems can often make paraxial assumptions, the approximation $\tan(\theta) = \sin(\theta)$ is often used; however, $\sin(\theta)$ should be used when the collection angles become large.

Clipped Lambertian

An aperture with a clipped Lambertian distribution is one where the source of flux appears to be Lambertian, but only over a finite range of angles. Outside of that range of angles, there is no flux and the range of angles is assumed to be constant across the aperture. The most common example is

when the source is at infinity and the flux at a planar aperture is considered. A clipped Lambertian distribution is also often found at an intermediate surface within an optical system. The terms *limited Lambertian* and *restricted Lambertian* are also used.

Generally, a clipped Lambertian distribution is defined across an aperture and the illuminance across the aperture is constant (e.g., a spatially uniform clipped Lambertian distribution). Another common distribution is a spatially uniform apodized Lambertian. In this case, the angular distribution is spatially uniform but the luminance is not Lambertian. A clipped Lambertian is a special case of an apodized Lambertian. The output of a fiber-optic cable is often assumed to have a spatially uniform apodized Lambertian distribution.

The etendue for clipped Lambertian situations is straightforward to compute. Consider the case of an infinite strip of width $2R$ with a clipped Lambertian distribution defined between $\pm\theta_{\max}$ relative to the surface normal. The etendue per unit length for this 2D clipped Lambertian case is

$$\text{etendue}_{2D} = n(2R)(2\sin\theta_{\max}) \tag{4}$$

2D refers to a trough or extruded geometry and typically assumes that the distribution is infinite in the third dimension. Mirrors can often be placed at the ends of a finite length source so that the source appears to be infinitely long.

For a Lambertian disk with a half cone angle of θ_{\max} , the etendue is

$$\text{etendue}_{3D} = n^2 \text{Area} \pi \sin^2 \theta_{\max} = n^2 \pi R^2 \pi \sin^2 \theta_{\max} \tag{5}$$

In a system where the etendue is preserved, these etendue relationships highlight that increasing either θ_{\max} or R requires a reduction in the other, as depicted in Fig. 1.

Hotell Strings

The etendue relationships described by Eqs. (4) and (5) are primarily defined for the case of a clipped Lambertian source. Such situations arise when the source is located at infinity (e.g., solar collection) or when considering the output of a fiber-optic cable. When the angular distributions vary spatially, Ref. 22 provides a method to compute the etendue that is very simple to use with 2D systems. The method is straightforward for the case of a symmetric system or an off-axis system, and even if there is an absorber positioned between the two apertures. The method is depicted in Fig. 2, where the etendue of the radiation that can be transferred between AB and CD is computed.

For a rotationally symmetric system, the etendue between two apertures (e.g., left side of Fig. 2) has been shown by Ref. 23 to be $(\pi^2/4)(AD-AC)^2$. Reference 24 has provided a generalized treatment of the 3D case.

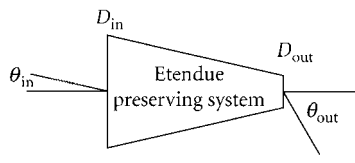


FIGURE 1 Graphical representation of etendue preservation.

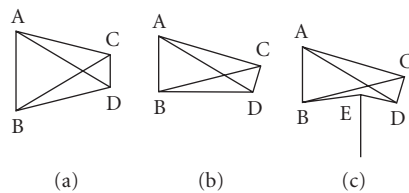


FIGURE 2 Hotell's crossed string relationship for computing etendue. This figure shows a symmetric case (a) and an asymmetric case (b) where the 2D etendue between aperture AB and CD is $(AD-AC + BC-BD)$. In (c), radiation is blocked below point E and the etendue is $(AD-AC + BC-BED)$.

Solid Angle and Projected Solid Angle

Solid angle is the ratio of a portion of the area of a sphere to the square of the sphere radius (see Chap. 3, and Refs. 24 and 25) and is especially important when working with sources that radiate into more than a hemisphere. Solid angle is directly applicable to point sources; however, when the size of the sphere is large enough, the solid angle is still used to characterize extended sources. For solid angle to be applied, the general rule of thumb is that the sphere radius should be greater than 10 times the largest dimension of the source, although factors of 30 or more are required in precision photometry (e.g., Ref. 26, pp. 4.29–4.31).

The etendue is sometimes called the *area-solid angle product*; however, this term can often cause confusion because etendue is actually the projected-area-solid-angle product [$\cos \theta dA d\Omega$] or area-projected-solid-angle product [$dA \cos \theta d\Omega$]. Reference 27 discusses some distinctions between projected solid angle (PSA) and solid angle. An important implication of the cosine factor is that the solid angle for a cone with half angle θ is

$$\text{Solid angle}_{\text{cone}} = 2\pi[1 - \cos\theta] = 4\pi \sin^2(\theta/2) \quad (6)$$

but the projected solid angle for the cone is

$$\text{Projected solid angle}_{\text{cone}} = \pi \sin^2 \theta \quad (7)$$

For a hemisphere, the solid angle is 2π and the PSA is π . This factor of 2 difference is often a source of confusion. PSA and solid angle are pictured in Fig. 3.

If the luminance at the receiver is constant, then PSA times luminance provides illuminance. Nonuniform luminance can be handled using weighted averages.²⁸

In the 2D case, the projected solid angle analog is simply projected angle, and is $2 \sin \theta$ for angles between $\pm\theta$ or $|\sin\theta_1 - \sin\theta_2|$ for angles between θ_1 and θ_2 .

Concentration

Concentrators can be characterized by the ratio of the output area to the input area.²⁹ For an *ideal* system, the etendue at the input aperture and the output aperture are the same, which leads to the ideal concentration relationships

$$\text{Concentration}_{2D} = n_{\text{out}} \sin\theta_{\text{out}} / (n_{\text{in}} \sin\theta_{\text{in}}) \quad (8)$$

$$\text{Concentration}_{3D} = n_{\text{out}}^2 \sin^2\theta_{\text{out}} / (n_{\text{in}}^2 \sin^2\theta_{\text{in}}) \quad (9)$$

where the input and output distributions are clipped Lambertians, θ_{in} is the maximum input angle, and θ_{out} is the maximum output angle. Maximum concentration occurs when $\sin\theta_{\text{out}} = 1$, so that the

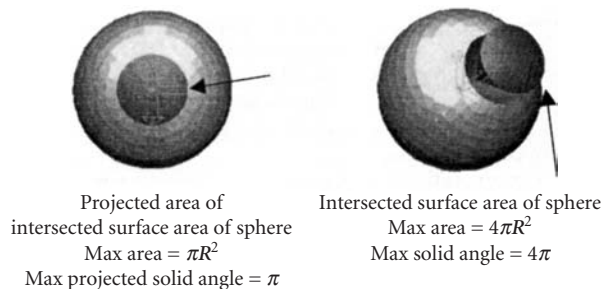


FIGURE 3 PSA (*left*) versus solid angle (*right*).

maximum concentration ratios for 2D and 3D symmetric systems in air are $1/\sin \theta_{in}$ and $1/\sin^2 \theta_{in}$, respectively. Concentrating the flux within a dielectric medium further increases the concentration ratio by a factor of n_{out}^2 in the 3D case and n_{out} in the 2D case.

In general, maximum concentration occurs when the phase space of the receiver is completely filled (i.e., all portions of the receiver are illuminated from all angles). Maximum concentration can be obtained if the etendue of the input is greater than the etendue of the output. In this case, however, the efficiency of the system is less than 1. One example of this is provided by a two-stage system where the second stage is an ideal maximum concentrator and the etendue of the first stage is greater than the etendue of the second stage.^{30,31} Such a system is not conventionally considered ideal because not all of the flux from the first stage is coupled into the entrance aperture of the second stage.

Dilution

Dilution is a term that is used to describe the situation in which the phase space of the receiver is unfilled. Consider the situation where a Lambertian disk of radius 1 is coupled directly to a receiver of radius 2. In this case, the area of the receiver is underfilled and represents *spatial dilution*. An example of *angular dilution* occurs when flux is transferred to a receiver using multiple discrete reflector elements. If there are gaps between the reflector elements, then the angular distribution of the flux incident on the receiver possesses gaps. The phrase *etendue loss* is also used to describe situations in which dilution is present.

39.3 SOFTWARE MODELING OF NONIMAGING SYSTEMS

The computer simulation of nonimaging systems is quite different from the computer simulation of imaging systems. Modeling of nonimaging systems typically requires three major items that imaging designers either do not need or tend to use infrequently. These are nonsequential ray tracing, spline surfaces, and modeling extended sources. Extended sources are often modeled using Monte Carlo techniques. Other differences include the need to model more complex surface properties including scattering surfaces; new methods to present simulation results including luminance, illuminance, and intensity distributions; and improved visualization because of the more complex geometries typically involved in nonimaging systems.

Nonsequential Ray Tracing

Nonsequential ray tracing means that the order of the surfaces with which the rays interact is not predetermined. In fact, rays can hit a single surface multiple times, which is especially common when skew rays are investigated. Nonsequential surfaces are especially important in the software modeling of lightpipes^{32,33} and prisms.

Spline Surfaces

Spline surfaces have been successfully applied to the design of nonimaging optical systems. Design work in this area includes genetic algorithm-based approaches,^{34,35} neural networks,³⁶ and variational principles.³⁷ Automotive headlamp designs also use spline surfaces routinely. Spline surfaces are also extremely important in computer-aided drafting (CAD) data exchange formats (e.g., IGES, STEP, and SAT) where Non-Uniform Rational B Splines (NURBS) have found tremendous utility.

Monte Carlo Techniques

Monte Carlo simulations are used to determine the intensity and/or illuminance distribution for optical systems. Typically, Monte Carlo techniques are used to model the source and/or surface properties of nonspecular surfaces. Reference 38 discusses some of the issues regarding Monte Carlo simulations. Some nonimaging systems that have been evaluated using Monte Carlo ray traces include an end-to-end simulation of a liquid crystal display (LCD) projector system,³⁹ light-emitting diodes,⁴⁰ integrating spheres,^{41,42} scintillating fibers,^{43,44} stray light,^{45,46} and automotive headlamps.⁴⁷

Source Modeling

Monte Carlo modeling of sources typically falls into three categories: point sources, geometric building blocks, and multiple camera images. Point sources are the simplest approach and are often used in the early stages of a design because of simplicity. Measured data, sometimes called *apodization data*, can be applied to a point source so that the intensity distribution of the source matches reality. The projected area of the source as seen from a given view angle can often be used to estimate the apodization data^{8,48} when measurements are unavailable. Since the luminance of a point source is infinite, point source models are of limited utility when etendue limitations must be considered. Geometric models create the source through the superposition of simple emitters (e.g., disks, spheres, cylinders, cubes, ellipsoids, toroids). The simple emitters are typically either Lambertian surface emitters or isotropic volume emitters, and the emitters can have surface properties to model effects that occur when flux propagates back through the emitter. The simple emitters are then combined with models of the nonemissive geometry (e.g., bulb walls, electrodes, scattering surfaces) to create an accurate source model. Apodization of the spatial and/or angular distributions of sources can help to improve the accuracy of the source model. In the extreme, independent apodization files can be used for a large number of view angles. This often results in the use of multiple source images,^{49,50} which has seen renewed interest now that charge-coupled device (CCD) cameras have matured.^{51–55}

The most typical Monte Carlo ray trace is one in which the rays traced from the source are independent of the optical system under investigation. When information about the portions of the source's spatial and angular luminance distribution that are contributing to a region of the intensity/illuminance distribution is available, importance sampling can be used to improve the efficiency of the ray trace. An example of such an approach occurs when an $f/1$ condenser is used to collect the flux from a source. Use of importance sampling means that rays outside of the $f/1$ range of useful angles are not traced.

Backward Trace

When a simulation traces from the source to the receiver, rays will often be traced that land outside of the region of the receiver that is under investigation. A backward ray trace can eliminate those rays by only tracing from the observation point of interest. The efficiency of the backward ray trace becomes dependent upon knowing which ray directions will hit the receiver. Depending upon the details of the system and prior knowledge of the system, backward tracing can often provide far more efficient use of traced rays when only a few observation points are to be investigated. Use of the backward ray trace and examples for a number of cases have been presented.²⁸ Backward ray tracing has also been called the *aperture flash mode*.⁵⁶

Field Patch Trace

Another trace approach is the field patch, where rays are traced backward from a point within the optical system and those rays are used to determine which rays to trace forward. Reference 56 describes the prediction of a distribution by summing results from multiple field patches. This

type of approach is especially useful if the effect of small changes in the system must be quantified because shifting the distributions and resuming can approximate small changes.

Software Products

There are many optical software vendors that have features that are primarily aimed toward the illumination market. A detailed comparison of the available products is difficult because of the speed at which software technology progresses. A survey of optical software products has been performed by Illuminating Engineering Society of North America (IESNA) and published in *Lighting Design and Application*.^{57,58}

39.4 BASIC BUILDING BLOCKS

This section highlights a number of the building blocks for designing systems to collect flux when the angles involved are large. Many of these building blocks are also important for imaging optics, which reflects the fact that image formation is not a fundamental requirement for nonimaging optics but that image formation is often useful.

Spherical Lenses

Spherical lenses are used in nonimaging systems although the aberrations can limit their use for systems with high collection angles. The aplanatic points of a spherical surface⁵⁹ are sometimes used to convert wide angles to lower angles and remain free of all orders of spherical aberration and low orders of coma and astigmatism. There are three cases⁶⁰ in which a spherical surface can be aplanatic:

1. The image is formed at the surface (e.g., a field lens).
2. Both object and image are located at the center of curvature.
3. The object and image are located on radii that are rn/n' and rn'/n away from the center of curvature of the spherical surface.

Case 3 is used to construct hyperhemispheric lenses that are often used in immersed high-power microscope objectives. Such lenses can suffer from curvature of field and chromatism (see Ref. 60, pp. 258–262, and Ref. 61). Some example aplanats are shown in Fig. 4.

Hyperhemispheric lenses can also be used in LED packages⁶² and have been used in photographic-type objectives and immersed IR detectors. Aplanatic surfaces for use as concentrators have been investigated⁶³ and can be nearly ideal if the exit surface is nonplanar (Ref. 4, pp. 37–38).

Aspheric Lenses

Spherical lenses have been the workhorse of the imaging industry because of the simplicity and accuracy with which they can be manufactured. Design requirements, especially the speed of the system, often necessitate the use of aspheric surfaces in nonimaging optics. Fortunately, the accuracy with which the surface figure must be maintained is often less severe than that of an imaging system. Aspheric lenses are often used as condensers in projection systems including LCD projectors, automotive headlamps, and overhead projectors.

The classic piano aspheric examples are the conic lenses where the eccentricity is equal to $1/n$ with the convex surface facing the collimated space and eccentricity = n with the piano surface facing the collimated space (see, for example, Ref. 64, pp. 112–113, and Ref. 65, pp. 100–103).

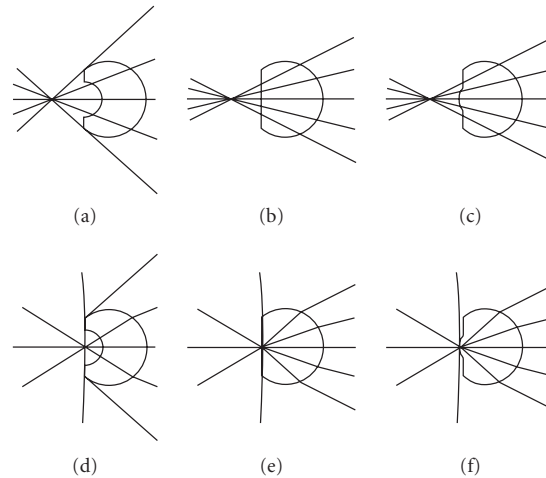


FIGURE 4 Several aplanats used to concentrate a virtual source. (a–c) The rays superimposed on top of the aplanats. (d–f) the ray trace. (d) A meniscus lens with hyperhemispheric outer surface and a spherical surface centered about the high concentration point. (e) A hyperhemisphere. (f) A hyperhemisphere with curved output surface.

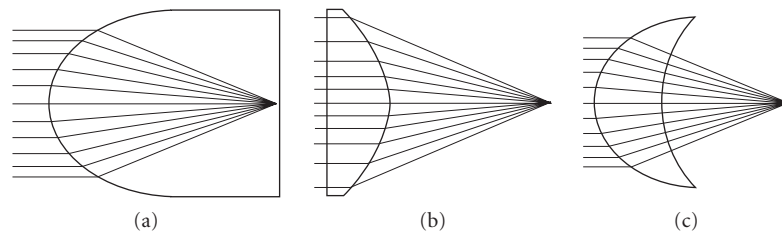


FIGURE 5 Conic collimators. (a) $e = 1/n$ and $f = Rn/(n - 1)$. (b) $e = n$ and $f = R/(n - 1)$. (c) $e = 1/n$ with other surface being spherical. $e =$ eccentricity, $n =$ index of refraction, $k = -e^2$.

Examples are shown in Fig. 5. The paraxial focal lengths are $Rn/(n - 1)$ and $R/(n - 1)$, respectively, where R is the radius of curvature. These examples provide correction for spherical aberration but still suffer from coma.

Refractive aspheric surfaces are often used in the classic Schmidt system.^{60,66} Reference 67 has described a procedure for determining the aspheric profile for systems with two aspheric surfaces that has been adapted to concentrators by Minano.⁶⁸

Fresnel Lenses

The creation of a lens using a symmetric arrangement of curved prisms is typically attributed to A. J. Fresnel and commonly called a *Fresnel lens*. General descriptions of Fresnel lenses are available.^{69,70}

Fresnel lenses provide a means to collect wide angular distributions with a device that can be easily molded. Standard aspheric lenses can be molded, but the center thickness compared to the edge thickness adds bulk and fabrication difficulties.

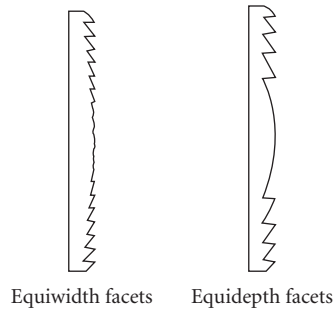


FIGURE 6 Two main types of Fresnel lenses: equiwidth and equidepth.

In imaging applications, Fresnel lenses are used, but the system designer must be careful to consider the impact of the facets on the image quality. The facet structure is also important in nonimaging applications where the use of facets introduces losses either by blocking flux or by underfilling the aperture, which typically results in an etendue loss.

Fresnel lenses have two main variations: constant-depth facet structures and constant-width facet structures (see Fig. 6). As the width of the facets becomes smaller, the effects of diffraction should be considered.⁷¹ Although Fresnel lenses are typically created on planar substrates, they can be created on curved surfaces.^{72,73}

The design of a Fresnel lens requires consideration of the riser angle (sometimes called the *draft angle*) between facets.^{74,75} A typical facet is depicted in Fig. 7. The riser angle should be designed to minimize its impact on the rays that are controlled by the facets that the riser connects.

Total Internal Reflection Fresnel Lenses As the Fresnel lens collection angle increases, losses at the outer facets increase.⁷⁶ One solution is the use of total internal reflection (TIR) facets. In this case, the flux enters the substrate material, hits the TIR facet, and then hits the output surface. The entering, exiting, and TIR facets can also have power in more sophisticated designs. The basic TIR Fresnel lens idea has been used in beacon lights since at least the 1960s.⁷⁷ TIR Fresnel lenses have received more recent design attention for applications including fiber illumination, LEDs, condensers, and concentrators.^{76,78–82} Vanderwerf⁸³ investigates achromatic issues with catadioptric Fresnel lenses.

In some applications, the TIR Fresnel lens degenerates into one refractive facet and two or more TIR facets. Combined TIR and refractive designs have been proposed for small sources with significant attention to LEDs.^{84–86}

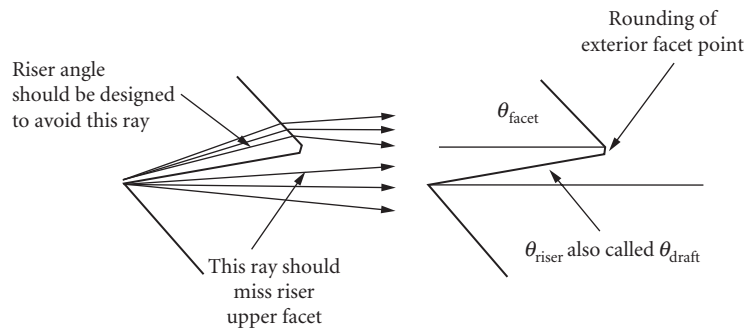


FIGURE 7 Facet terminology and design issues.

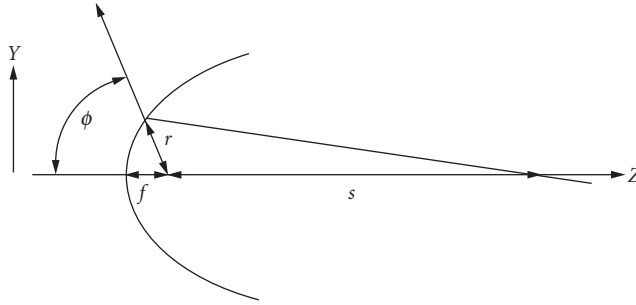


FIGURE 8 Conic reflector showing r , f , s , and ϕ .

Conic Reflectors

Reflectors made from conic sections (parabolas, ellipses, and hyperbolas) are commonly used to transfer radiation from a source to a receiver. They provide spherical-aberration-free transfer of radiation from one focal point to the other. The problem is that the distance from a focal point to a point on the surface of the reflector is not constant except for the case of a sphere. This introduces a nonuniform magnification (coma) that can result in a loss of concentration.

Many textbooks describe conic surfaces using equations that have an origin at the vertex of the surface or at center of symmetry. Reference 17 (page 1.34) shows a number of different forms for a conic surface. In nonimaging systems, it is often convenient to define conic reflectors using an origin shifted away from the vertex. Reference 87 describes these surfaces in Cartesian coordinates as *apo-vertex surfaces*. They can also be described in polar coordinates where the focal point is typically the origin. The polar definition of a conic surface⁸ is

$$r = \frac{f(1+e)}{1+e \cos \phi}$$

where r is the distance from the foci to the surface, f is the distance from one focus to its nearest vertex, ϕ is the angle from the foci to the vertex, e is $s/(s+2f)$, and s is the distance between the two foci. In Cartesian coordinates, $z = r \sin \phi$ and $y = r \cos \phi$. An example is shown in Fig. 8.

Macrofocal Reflectors

One can consider the standard conic reflector to map the center of a sphere to a single point. What can sometimes be more valuable for nonimaging optics is to map the edge of the sphere to a single point. Reflectors to provide this mapping have been called macrofocal reflectors.⁸⁸ They have also been called extinction reflectors⁸ because of the sharp edge in the illuminance distribution that they can produce. Reference 89 describes an illumination system that provides a sharp cutoff using tilted and offset parabolic curves.

Involute

A reflector that is used in many nonimaging systems is an *involute*. An involute reflector sends tangential rays from the source back onto themselves. One way to produce an involute for a circle is to unwind a string that has been wrapped about the circle. The locus of points formed by the string equals the involute. In Cartesian coordinates, the equation for the involute of a circle is⁹⁰

$$x = r(\sin \theta - \theta \cos \theta) \quad (10)$$

$$y = -r(\cos \theta + \theta \sin \theta) \quad (11)$$

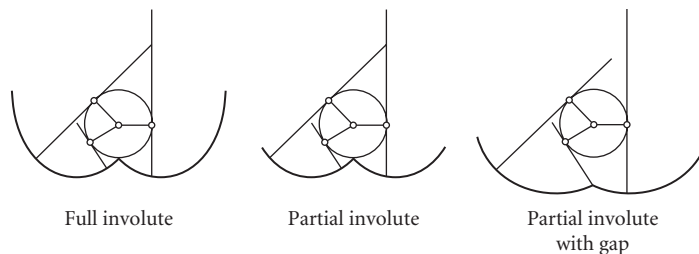


FIGURE 9 Full involute, partial involute, and involute with gap. Rays tangent to the source reflect back toward the same location on the source.

where r is the radius of the circle, $\theta = 0$ at the cusp point, and the origin is at the center of the circle. An example is shown in Fig. 9, where rays are drawn to highlight the fact that rays tangent to the circle reflect off of the involute and retrace their path.

The term *full involute* has been used to describe an involute that is continued until the output aperture is tangent to the surface of the source. A *partial involute* is an involute that is less than a full involute. An involute with gap is also shown in Fig. 9, and is sometimes called a *modified involute*.⁹¹ Involute segments can also be applied to noncircular convex shapes and to disjoint shapes.⁹²

39.5 CONCENTRATION

The transfer of radiation in an efficient manner is limited by the fact that the luminance of the radiation from a source in a lossless system cannot increase as the radiation propagates through the system. A consequence is that the etendue cannot be reduced and maximizing the luminance means avoiding dilution. A less understood constraint that can impact the transfer of radiation is that rotationally symmetric optical systems conserve the skewness of the radiation about the optical axis (see for example Ref. 4, Chap. 2.8, and more recently Ref. 93). This section describes some of the basic elements of systems used to transfer radiation with an emphasis on those that preserve etendue.

Many papers on nonimaging optics can be found in the Society of Photooptical Instrumentation Engineers (SPIE) proceedings. SPIE conferences specifically focused on nonimaging optics include Refs. 94–99. Reference 100 also provides a large number of selected papers. Textbooks on nonimaging optics include Refs. 4 and 101. Topical articles include Refs. 102 and 103.

Discussions of biological nonimaging optic systems have also been written.^{104–108}

Tapered Lightpipes

Lightpipes can be used to transport flux from one location to another. This can be an efficient method to transport flux, especially if total internal reflection (TIR) at the lightpipe surface is utilized to provide lossless reflections. Lightpipes can provide annular averaging of the flux. In addition, if the lightpipe changes in size from input to output over a given length (e.g., the lightpipe is tapered), then the change in area produces a corresponding change in angles. If the taper occurs over a long enough distance, then in most cases the etendue will be preserved (see, for example, Ref. 109). If the taper occurs over too short a length, the etendue will not be preserved, which can result in an increase in angles or possibly rays lost, such as when they are reflected back toward the input.

An interesting geometric approach to determining if meridional rays in a conical lightpipe can propagate from input to output end was shown by Williamson,¹¹⁰ and is informally called the Williamson construction or a tunnel diagram. An example is shown in Fig. 10, where the duplicate copies of the primary lightpipe are all centered about a single point. Where the rays cross the copies

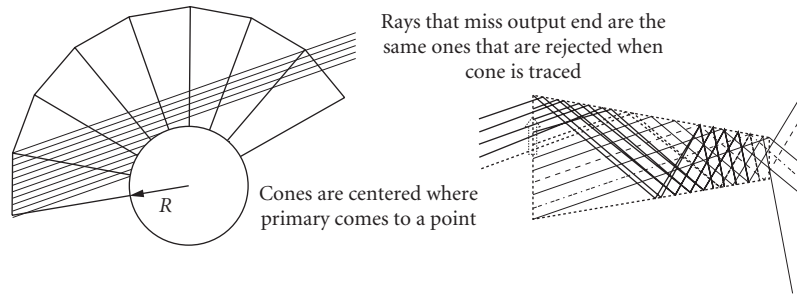


FIGURE 10 Williamson construction. Multiple copies of the primary cone are drawn. Where rays cross, these copies are the same as where they hit the actual cone. Rays that do not cross the output end are reflected back to the input end.

shows the position and at what angle the rays will hit the primary lightpipe. As seen in Fig. 10, the rays that do not cross the output end of any of the copies of the cone are also the rays that are reflected back out the input end of the cone. Other analyses of tapered geometries have been performed by Welford⁴ (pp. 74–76), and Meyer¹¹¹ provides numerous references.

Witte¹¹² has performed skew ray analysis using a variation of the Williamson construction. Witte also shows how a reference sphere centered at the point where the cone tapers to a point can be used to define a pupil when cones are used in conjunction with imaging systems analysis. Burton¹¹³ simplified the expressions provided by Witte.

Vector flux investigations have shown that a cone is an ideal concentrator for a sphere (Ref. 114, p. 539). An important implication of this result, which is similar to the information presented by Witte,¹¹² is that a cone concentrator can be analyzed by tracing from the cone input aperture to the surface of a sphere. If the rays hit the sphere surface, then they will also hit the sphere after propagating through the lightpipe. This is true for both meridional and skew rays.

CPC

A compound parabolic collector (CPC) can be used to concentrate the radiation from a clipped Lambertian source in a nearly etendue-preserving manner. Welford⁴ provides a thorough description of CPCs. An example CPC is shown in Fig. 11, where the upper and lower reflective surfaces are tilted parabolic surfaces. The optical axis of the parabolic curves is not the same for the upper and lower curves, hence the use of the term *compound*.

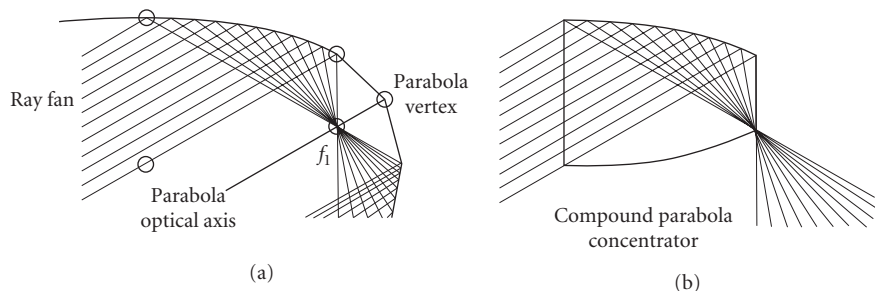


FIGURE 11 Compound parabolic concentrator (CPC). (a) Upper parabolic curve of CPC. Circles are drawn to identify the parabola vertex and the edges of the CPC input and output apertures. (b) Compound parabola concentrator.

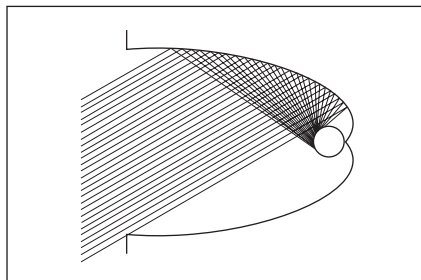


FIGURE 12 A CPC-type reflector for a tubular source.

In 2D, the CPC has been shown to be an ideal concentrator.⁴ In 3D, a small fraction of skew rays are rejected and bound the performance. For a CPC, this skew ray loss increases as the concentration ratio increases. Optimization procedures have produced designs that provide slightly better performance than the CPC in 3D.³⁷ Molledo¹¹⁵ has investigated a crossed cylindrical CPC configuration.

The length of a CPC concentrator is

$$\text{Length} = (R_{\text{in}} + R_{\text{out}}) / \tan(\theta_{\text{in}}) \quad (12)$$

A particularly important aspect of the CPC is that the length is minimized for the case where the input port can see the exit port. If the distance between the input and output ports is made smaller than the CPC length, then rays higher than the maximum design angle could propagate directly from input to output aperture. The reduction in concentration is often small, so many practical systems use truncated CPC-like designs.^{90,116–118} A cone that is the same length as an untruncated CPC will provide lower concentration, but as the cone is lengthened, the performance can exceed the performance of the CPC.

In the literature, CPC is sometimes used to denote any ideal concentrator. For nonplanar receivers, the reflector shape required to create an ideal concentrator is not parabolic. The convention is to use CPC-type.

An example of a CPC-type concentrator for a tubular receiver is shown in Fig. 12. Part of the reflector is an involute and the remainder reflects rays from the extreme input angle to the tangent of the receiver. The design of CPC-type concentrators for nonplanar sources has been described.^{4,116,119,120} Designs with gaps between reflector and receiver^{121–123} and with prisms attached to the receiver^{124,125} have also been described.

CPC-type geometries are also used in laser pump cavities.¹²⁶

Arrays of CPCs or CPC-like structures have also been applied to the liquid crystal displays,^{127,128} illumination,⁸⁹ and solar collection.¹²⁹

CEC

When the source of radiation is located at a finite distance away from the input port of the concentrator, a construction similar to the CPC can be used; however, the reflector surface is now elliptical.¹³⁰ In a similar manner, if the source of radiation is virtual, then the reflector curvature becomes hyperbolic. These two constructions are called compound elliptical concentrators (CEC) and compound hyperbolic concentrators (CHC). A CEC is shown in Fig. 13, where the edge of the finite size source is reimaged onto the edge of the CEC output aperture. Hottel strings can be used to compute the etendue of the collected radiation (see Sec. 39.2). The CPC is similar to the CEC, except the edge of the source is located at infinity for the CPC.

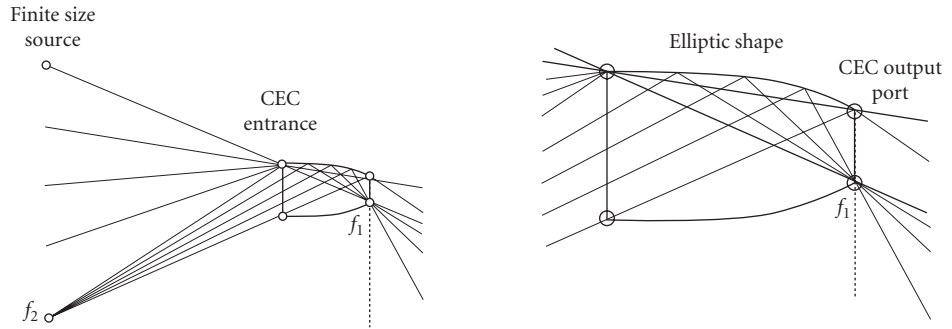


FIGURE 13 Compound elliptical concentrator (CEC). Rays originating at finite size source and collected at the CEC entrance are concentrated at the CEC output port. Edge of source is imaged onto edge of output aperture.

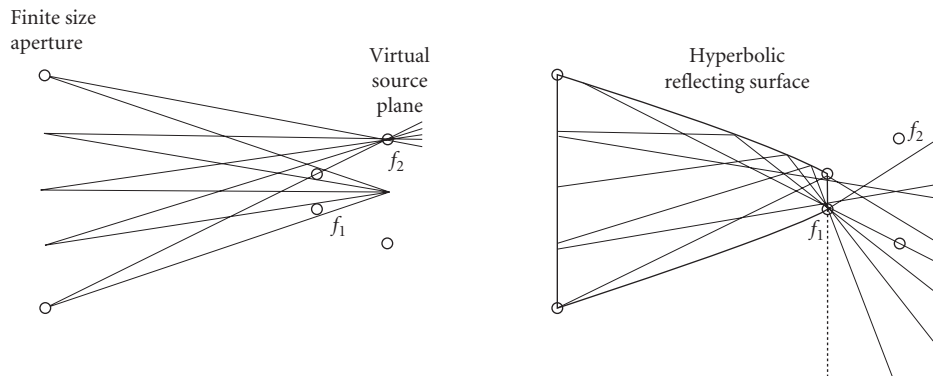


FIGURE 14 Compound hyperbolic concentrator for use with virtual source. Edge of virtual source is imaged onto the edge of the output port.

CHC

A CHC is shown in Fig. 14 where the rays that pass through the finite-size aperture create a virtual source. The CHC concentrates the flux. Arrangements where the CHC is just a cone are possible.¹³¹ Gush¹³² also describes a hyperbolic cone-channel device.

If the virtual source is located at the output aperture of the concentrator, then an ideal concentrator in both 2D and 3D can be created (see Fig. 15). Generically, this is a hyperboloid of revolution and is commonly called a *trumpet*.¹³³ The trumpet maps the edge of the virtual source to $\pm 90^\circ$. The trumpet can require an infinite number of bounces for rays that exit the trumpet at nearly 90° , which can limit the performance when finite-reflectivity mirrors are used.

DCPC

The CPC construction can also be applied when the exit port is immersed in a material with a high index of refraction. Such a device is called a *dielectric compound parabolic concentrator* (DCPC).¹³⁴ This allows concentrations higher than those found when the exit port is immersed in air. In this immersed exit port case, the standard CPC construction can still be used, but the exit port area is now n times smaller in 2D and n^2 smaller in 3D. Another motivation for a DCPC is that the reflectivity

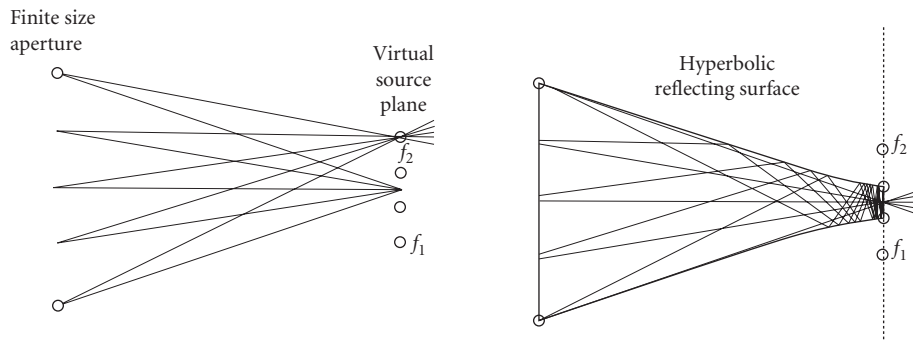


FIGURE 15 Trumpet with exit aperture located at virtual source plane. Rays “focused” at edge of virtual source exit the trumpet at $\pm 90^\circ$.

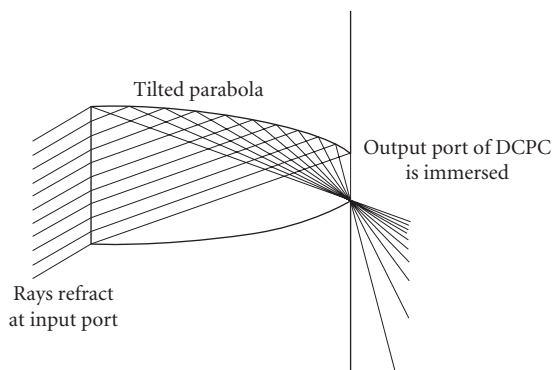


FIGURE 16 Dielectric CPC: same as CPC, except output port is immersed.

losses for a hollow CPC can be minimized because TIR can provide lossless reflections. A DCPC is shown in Fig. 16. A DCPC that uses frustrated TIR at the DCPC to receiver interface has been investigated.¹³⁵

A hollow CPC can also be attached to a DCPC to minimize the size of the DCPC.¹³⁶ Extraction of flux from a high-index medium into a lower-index medium has been investigated^{137,138} using faceted structures.

Multiple Surface Concentrators

CPC-type concentrators can produce a long system when θ_{in} is small. To avoid excessively long systems, other optical surfaces are typically added. Often, a CPC-type concentrator is combined with a primary element such as a condensing lens or a parabolic reflector. There are also designs where the optical surface of the primary and the optical surface of the concentrator are designed together, which often blurs the distinction between primary and concentrator.

Lens + Concentrator There are numerous examples of lenses combined with a nonimaging device in the literature. One approach is to place a lens with the finite-size aperture shown in Figs. 13, 14, or 15 and illuminate the lens with a relatively small angular distribution. The lens produces higher angles that are then concentrated further using the conventional CPC, CEC, CHC, or trumpet. Such a combination typically produces a shorter package length than a CPC-type concentrator.

A lens/CHC can be attractive because the CHC is placed between the lens and lens focal plane, whereas the CEC is placed after the lens focal plane. Thus, the lens/CHC provides a shorter package than a lens/CEC. Both the lens curvature and the CHC reflector can be optimized to minimize the effects of aberrations, or the combination can be adjusted so that the CHC turns into a cone. Investigations of lens/cone combinations include Williamson,¹⁰ Collares-Pereira,⁷⁵ Hildebrand,¹³⁹ Welford,⁴ and Keene.¹⁴⁰ Use of a cone simplifies fabrication complexity.

After the theoretical implications of the trumpet were recognized, Winston¹⁴¹ and O’Gallagher¹³³ investigated the case of a lens/trumpet.

When TIR is used to implement the mirrors, the mirror reflection losses can be eliminated, with the drawback that the package size grows slightly. A design procedure for a lens mirror combination with maximum concentration was presented by Ning,¹⁴² who coined the term *dielectric total internal reflecting concentrator* (DTIRC) and showed significant package improvements over the DCPC.

Eichhorn¹⁴³ describes the use of CPC, CEC, and CHC devices in conventional optical systems, and uses explicit forms for the six coefficients of a generalized quadric surface.¹⁴⁴ A related “imaging” lens and reflector configuration is the Schmidt corrector system.¹⁴⁵

Arrays of lenses with concentrators have also been investigated for use in illumination.¹⁴⁶

Mirror + Concentrator There have been numerous investigations of systems that have a parabolic primary and a nonimaging secondary.^{137,147–150} Other two-stage systems have been investigated.^{29–31,151–153} Kritchman¹⁵⁴ has analyzed losses due to aberrations in a two-stage system.

Asymmetric concentrators have also been investigated. Winston¹⁵² showed that an off-axis parabolic system can improve collection by tilting the input port of a secondary concentrator relative to the parabolic axis. Other asymmetric designs have been investigated.^{155,156}

Related two-mirror “imaging” systems (e.g., Ref. 60, Chap. 16) such as the Ritchey-Chretien can be used as the primary. Ries¹⁵⁷ also investigated a complementary Cassegrain used for concentration.

Simultaneous Multiple Surfaces Minano and coworkers have investigated a number of concentrator geometries where the edge of the source does not touch the edge of a reflector. The procedure has been called simultaneous multiple surfaces (SMS) and builds on the multisurface aspheric lens procedure described by Schultz.⁶⁷ Minano uses the nomenclature *R* for refractive, *X* for reflective (e.g., refleXive), and *I* when the main mode of reflection is TIR. The *I* surface is typically used for refraction the first time the ray intersects the surface and TIR for the second ray intersection. The *I* surface is sometimes mirrored over portions of the surface if reflection is desired but the angles do not satisfy the TIR condition. Illustrative citations include RR,¹⁵⁸ RX,^{159,160} and RXI.^{161,162} Example RX and RXI devices are shown in Fig. 17. SMS can be used to design all reflecting configurations.

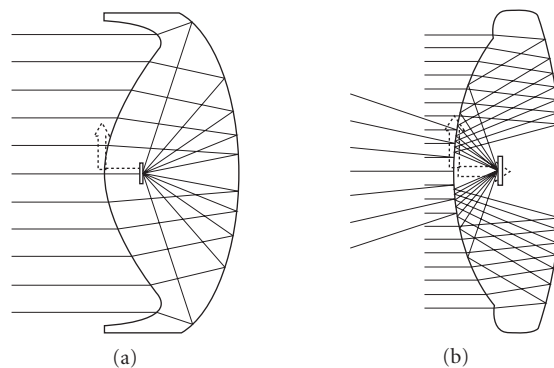


FIGURE 17 RX (a) and RXI (b) concentrators. The receiver is immersed in both cases. In the RXI case, the central portion of the refractive surface is mirrored.

Restricted Exit Angle Concentrators with Lenses

If the angular distribution of the flux is constant across an aperture and the centroid of the distribution is normal to the aperture, then the distribution is called *telecentric*. Many optical systems require the transfer of flux from one location to another. If the desired distribution of flux at the second location is also telecentric, then the system is called doubly telecentric. An afocal imaging system can be modified to provide a doubly telecentric imaging system when used with finite conjugates. An example is shown in Fig. 18.

In nonimaging systems, flux must often be transferred from one aperture to another and the angular distribution at the second aperture must be constant across the aperture. However, the point-by-point mapping is not required in nonimaging systems. A single lens can provide this transfer of radiation, as shown in Fig. 19. This type of lens is a collimator and is similar to a Fourier transform lens; it has been called a *beam transformer*.¹⁶³

The aberrations introduced by a lens can limit the etendue preservation of the aggregate flux collected by the lens. The increase in etendue tends to become more severe as range of angles collected by the lens increases. Somewhere around $f/1$, the change can become quite significant.

θ_1/θ_2

θ_1/θ_2 concentrator^{29,163,164} maps a limited-angle Lambertian distribution with maximum angle θ_1 into another limited-angle Lambertian with maximum angle θ_2 . One version of the θ_1/θ_2 is a compound parabolic construction with a cone replacing part of the small end of the CPC (Ref. 4, pp. 72–74, sec. 5.3). A picture is shown in Fig. 20. The length of a θ_1/θ_2 is given by the same equation as the CPC [e.g., Eq. (12)]. If the θ_1/θ_2 is hollow and $\theta_2 = 90^\circ$, then the cone disappears and the θ_1/θ_2 becomes a CPC.

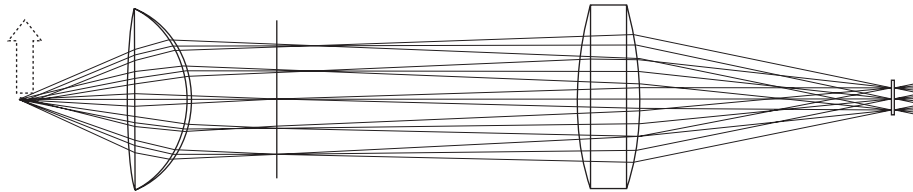


FIGURE 18 Doubly telecentric system showing a twofold increase in size and a twofold decrease in angles.

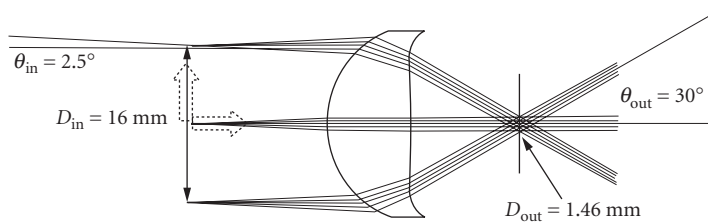


FIGURE 19 Condenser to provide an angle to area transformation and produce a distribution that is symmetric about the optical axis at both input and output planes. $D_{in} \sin \theta_{in} = D_{out} \sin \theta_{out}$.

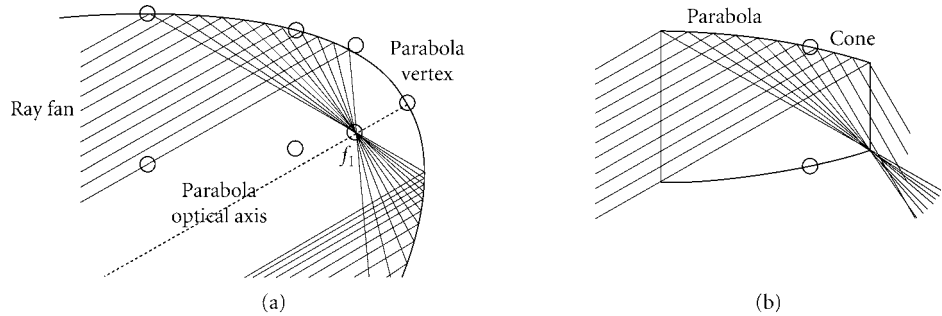


FIGURE 20 θ_1/θ_2 implemented using cone with compound parabolic. (a) Construction information. (b) θ_1/θ_2 , (b) compound parabola + cone.

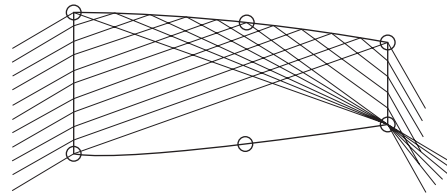


FIGURE 21 Dielectric θ_1/θ_2 implemented using cone with compound parabolic. The rays refract at the input port, total internally reflect (TIR) at the sidewall, and then refract at the output port.

To reduce losses introduced by nonunity mirror reflectivities, a θ_1/θ_2 converter can also be created using a solid piece of transparent material. The design is similar to the hollow θ_1/θ_2 and must take into account the refraction of rays at the input and output port. A dielectric θ_1/θ_2 converter is shown in Fig. 21.

If the θ_1/θ_2 is an all-dielectric device, nearly maximum concentration for an air exit port can be obtained without the mirror losses of a hollow CPC; however, Fresnel surface reflections at the input and output ports should be considered.

Similar to the lens-mirror combination for maximal concentration, the θ_1/θ_2 can be implemented with a lens at the input port. Going one step further, the output port of the θ_1/θ_2 can also incorporate a lens. Similar to a Galilean telescope, the θ_1/θ_2 with lenses at both the input and output ports can provide a short package length compared to embodiments with a flat surface at either port.¹⁶⁴

In many cases, a simple tapered cone can be used to implement a θ_1/θ_2 converter. One reason is that when the difference between θ_1 and θ_2 is small, the compound parabolic combined with a cone is nearly identical to a simple cone. By using an approximation to the ideal shape, some of the flux within θ_1 maps to angles outside of θ_2 . In many practical systems, the input distribution does not have a sharp transition at θ_1 , so a small mapping outside of θ_2 has minor impact. In general, the transition can be made sharper by increasing the cone length; however, there may be multiple local minima for a specific situation.¹⁶⁵

Uniformity at the output port of a θ_1/θ_2 compared to a CPC has been investigated.¹⁶⁶ Tabor¹¹⁸ describes hot spots that can occur with CPCs used in solar energy systems. Edmonds¹²⁵ uses a prism where hot spots that may occur are defocused after propagating through the prism. Emmons¹⁶⁷ discusses polarization uniformity with a dielectric 90/30 converter.

Two CPCs can be placed “throat-to-throat” to provide a θ_1/θ_2 converter that also removes angles higher than the desired input angle.¹⁶⁸ The removal of higher angles is a result of the CPCs ability to transfer rays with angles less than θ_1 and reject rays with angles greater than θ_1 . Two nonimaging concentrators have also been investigated for use in fiber-optic coupling.¹⁶⁹

Tight lightpipe bends are also made possible by inserting a bent lightpipe between the throats of the two concentrators.³³ In the special case of $\theta_1 = \theta_2 = 90^\circ$, a bent lightpipe has also been used as an ideal means to transport the flux for improved packaging.¹⁷⁰

2D versus 3D Geometries

Skew Ray Limits One of the standard design methods for concentrators has been to design a system using a meridional slice and then simply rotate the design about the optical axis. In some cases, such as the CPC used with a disk receiver, this may only introduce a small loss; however, there are cases where the losses are large.

The standard CPC-type concentrator for use with a circular receiver^{90,171} is an example where this loss is large. To assess the loss, compare the 2D and 3D cases. In 2D a concentrator can be designed where the input port is $2R_{in}$ with $a \pm \theta_{in}$ angular distribution; the size of the tube-shaped receiver is $2\pi R_{tube} = 2R_{in} \sin \theta_{in}$. The 2D case provides maximum concentration. In the 3D case where the 2D design is simply spun about the optical axis, the ratio of the input etendue to the output etendue is $\pi R_{in}^2 \pi \sin^2 \theta_{in} / (4\pi R_{tube}^2) = \pi^2 / 4 \sim 2.5$. The dilution is quite large. Feuermann¹⁷² provides a more detailed investigation of this type of 2D-to-3D etendue mismatch.

Ries⁹³ shows how skew preservation in a rotationally symmetric system can limit system performance even though the etendue of the source and receiver may be the same. The skewness distributions (etendue/dskewness as a function of skewness) for an example disk, sphere, and cylinder aligned parallel to the optical axis are shown in Fig. 22. The etendues of the disk, sphere, and cylinder are all the same. A second example where a cylinder is coupled to a disk is shown in Fig. 22b. This figure highlights losses and dilution. For skew values where the skewness distribution for the receiver is greater than the source, dilution occurs. Where the source is greater than the receiver, losses occur. If the size of the disk is reduced, then the dilution will also decrease, but the losses increase. If the disk is made bigger, then the losses decrease, but the dilution increases. Skew ray analysis of inhomogeneous sources and targets has also been considered.¹⁷³

Star Concentrator One way of avoiding the skew ray limit in a system with a rotationally symmetric source and a rotationally symmetric receiver is to avoid the use of rotationally symmetric optics to couple the flux from source to target. A device called a *star concentrator* has been devised¹⁷⁴ and shown to improve the concentration. The star concentrator derives its name from the fact that the cross section of the concentrator looks like a star with numerous lobes. The star concentrator is designed using a global optimization procedure. Performance of the star concentrator, assuming a reflectivity of 1, provides performance that exceeds the performance possible using a rotationally

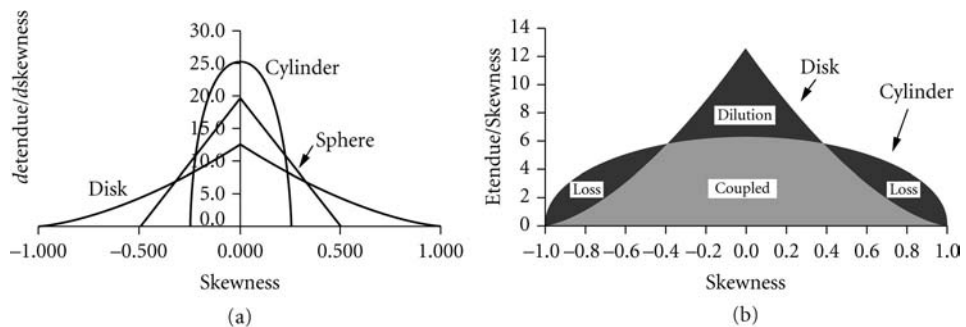


FIGURE 22 Skewness distributions (Ries⁹³). (a) Skewness distribution for a disk with unit radius, a cylinder with length = 2 and radius = 0.25, and a sphere of radius 0.5. The etendue for all three cases is constant, (b) Comparison of the coupling of a disk to a cylinder where the coupled flux is distinguished from the losses and the dilution.

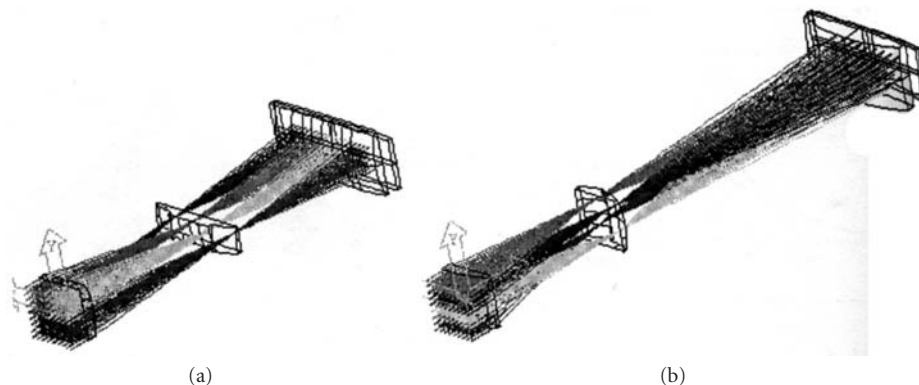


FIGURE 23 Square-to-rectangular mapping with an image dissector (a) and a lens array (b) where the output angular distribution is symmetric. The image dissector uses an array of discrete field lenses. The lens array uses a single field lens.

symmetric system. The star concentrator approach has also been investigated for the case of coupling from a cylindrical source to a rectangular aperture.^{7,175}

Image Dissectors The problem of coupling flux from a rotationally symmetric source to a spectrometer slit motivated some earlier efforts to overcome the skewness limit. Benesch¹⁷⁶ shows a system with two sets of facets that “dissect” an X by Y aperture into N regions of size X/N by Y and relay those images to create an X/N by NY region. The dissection can be done while preserving the NA of the input (e.g., the etendue is preserved but the skew is changed). Benesch calls this system an *optical image transformer* and it is similar in function to Bowen’s image slicer.¹⁷⁷ Refractive versions are also possible (see for example Ref. 178). One issue that must be addressed for some systems, which use these dissectors, is that the adjoining regions of the X/N by Y areas can produce a nonuniform illumination distribution. Laser arrays can also use dissectors to improve performance.^{179,180}

Tandem-lens-array approaches can be used that provide the image dissection and superimpose all the images rather than placing them side by side.^{181,182}

An example of a refractive image dissector compared to a lens array is shown in Fig. 23. The image dissector breaks an input square beam into four smaller square regions that are placed side by side at the output. The lens array breaks the input beam into four rectangles that are superimposed at the output. Both systems use field lenses at the output. The lens array does not have a gap between mapped regions at the output plane compared to the dissector.

Fiber bundles can also be used as image dissectors. For example, the fibers can be grouped in a round bundle at one end and a long skinny aperture at the other end. Feuermann¹⁸³ has provided a recent investigation of this idea. If the fibers in the bundle have a round cross section, tapering the ends of each of the fibers can help improve the uniformity at the output.¹⁸⁴ Special systems that use slowly twisted rectangular light guides have also been investigated for scintillators,^{185–187} document scanning,^{188,189} and microlithography.¹⁹⁰

Bassett¹⁹¹ also shows that arrays of ideal fibers can be used to create ideal concentrators, and Forbes¹⁹² uses the same approach to create a θ_1/θ_2 transformer.

Geometrical Vector Flux

The edge ray design method does not show how the jumble of rays from multiple reflections yields a Lambertian output. The geometric vector formalism was investigated in an attempt to understand the process. Summaries can be found in Bassett¹⁰¹ and Welford.⁴ The proof that a trumpet provides ideal performance has been performed using the flow line concept.¹¹⁴

The vector flux formalism provides the lines of flow and loci of constant geometric vector flux J . The vector J is the PSA for three orthogonal patches where the PSA on the front side of a patch can subtract from the PSA on the back side. The magnitude of J describes what the PSA would be for a surface aligned perpendicular to J . The flow line concept uses the idea that an ideal mirror can be placed along lines of flow without modifying the flow.

In the context of vector flux, the CPC, θ_1/θ_2 , and light cone are discussed in Winston.¹⁴¹ The light cone is demonstrated using an orange in Winston.⁹⁵ The 2D CEC and θ_1/θ_2 are discussed by Barnett.^{193,194} Gutierrez¹⁹⁵ uses Lorentz geometry to show that other rotationally symmetric ideal concentrators can be obtained using surfaces that are paraboloid, hyperboloid, and ellipsoid. Greenman¹⁹⁶ also describes vector flux for a number of ideal concentrators.

Edge Rays

The design of nonimaging optics has been facilitated by the use of the edge ray method^{4,68,101,197} and is sometimes called the string method.^{4,8} In the design of an ideal concentrator, the main principle is that extreme rays at the input aperture must also be extreme rays at the output aperture. Minano¹⁹⁸ describes the edge rays as those that bound the others in phase space. In many cases, the result is that the edge of the source is imaged to the edge of the receiver, and the reflector surface to perform the mapping is a compound conic. When the location of the source edge or the receiver edge is not constant (e.g., a circular source), the required reflectors are not simple conies (e.g., see the preceding section). Gordon¹⁹⁹ describes a complementary construction that highlights the fact that what is conventionally the “outside” of a concentrator can also be used in nonimaging optics.

Davies²⁰⁰ further explores the continuity conditions required by the edge-ray principle. Ries²⁰¹ refines the definition of the edge ray principle using topology arguments. Rabl²⁰² shows examples where the analysis of specular reflectors and Lambertian sources can be performed using only the edge rays.

Tailoring a reflector using edge rays is also described in the next section.

Inhomogeneous Media

A medium with a spatially varying index of refraction can provide additional degrees of freedom to the design of nonimaging optics. In imaging optics, examples of inhomogeneous media include Maxwell’s fisheye and the Luneburg lens.^{4,18} In concentrator design, 2D and 3D geometries have been investigated.^{198,203}

39.6 UNIFORMITY AND ILLUMINATION

In many optical systems, an object is illuminated and a lens is used to project or relay an image of the illuminated object to another location. Common examples include microscopes, slide projectors, overhead projectors, lithography, and machine vision systems. In a projection system the uniformity of the projected image depends upon the object being illuminated, the distribution of flux that illuminates the object, and how the projection optics transport the object modulated flux from object to image. The portion of the system that illuminates the object is often called the *illumination subsystem*.

Beam forming is another broad class of illumination systems. In beam forming, there is often no imaging system to relay an image of the illuminated object. Since the human eye is often used to view the illuminated object, the eye can be considered an imaging subsystem. Application examples include commercial display cases, museum lighting, room lighting, and automotive headlamps.

Approaches to provide uniform object illumination are discussed in this section. The discussed approaches include Kohler/Abbe illumination, integrating cavities, mixing rods, lens array, tailored

optics, and faceted structures. Many of the approaches discussed in this section can be designed to preserve etendue, but the design emphasis is on uniformity so they have been placed in this section rather than the preceding one.

Classic Projection System Uniformity

In classical projection systems, there is typically a source, a transparency (or some other item to be illuminated), and a projection lens. There are two classic approaches to obtaining uniformity in those systems: Abbe and Kohler illumination (see Fig. 24). In both cases, a source illuminates a transparency that is then often imaged onto a screen. *Transparency* is a general term that includes items such as spatial light modulators, film, slides, gobos, and microscope slides.

Some general references regarding Abbe/Kohler illumination include Jones,²⁰⁴ Kingslake,²⁰⁵ O'Shea,²⁰⁶ Wallin,²⁰⁷ Weiss,²⁰⁸ Ray (pp. 455–461),⁶⁶ Bradbury,²⁰⁹ Inoue,²¹⁰ and Born.¹⁸ A description of Kohler with partial coherence is included in Lacombat.²¹¹

Abbe Abbe (also called Nelsonian and Critical illumination) images the source onto the transparency. When the source is spatially uniform, the Abbe approach works fine. A frosted incandescent bulb is a good example. Sometimes a diffuser is inserted in the system and the diffuser is imaged onto the transparency. Fluorescent tubes are another case where a light source provides excellent spatial uniformity. The luminance of sources that are spatially uniform is typically lower than that of sources that are not uniform.

Kohler Even though the source's illuminance distribution may be nonuniform, there are often uniform regions of the source's intensity distribution. This phenomenon is used in the design of Kohler-type systems where the source is imaged into the projection lens aperture stop and a mapping of the source intensity distribution illuminates the transparency. One reason this approach works well is that far from the source, the illuminance distribution over a small area is independent of where the flux originated. This tends to provide uniformity over portions of the source's intensity distribution.

A nearly spherical reflector is sometimes added behind the source to improve the collection efficiency of a Kohler illumination system. Typically, the spherical reflector produces an image of the source that is aligned beside the primary source. The direct radiation from the source and the flux

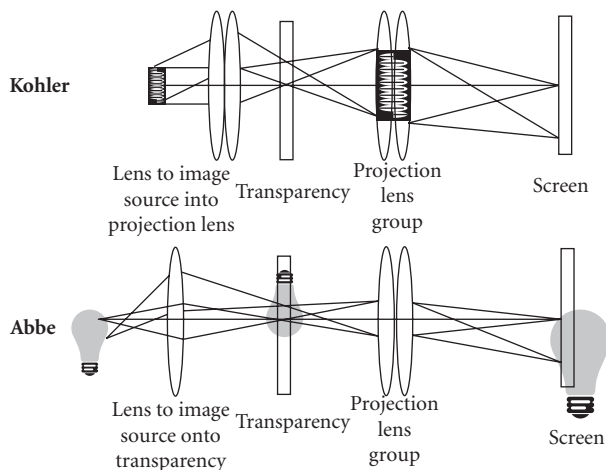


FIGURE 24 Abbe versus Kohler layout comparison.

from the image of the source then become the effective source for the Kohler illuminator. If the source is relatively transparent to its own radiation, then the image of the source can be aligned with the source itself. Blockage by electrodes and leads within the bulb geometry often limit the flux that can be added using the spherical reflector. Scattering and Fresnel surface reflections also limit the possible improvement. Gordon²¹² provides a recent description of a back reflector used with filament sources.

In Kohler illuminators, tailoring the illuminance distribution at the transparency to provide a desired center-to-edge relative illumination can be performed using the edge ray approach described by Zochling.²¹³ Medvedev⁸² also applied the edge ray principle to the design of a projection system. The other tailoring approaches described later in this section may also be relevant.

Integrating Cavities

Integrating cavities are an important example of nonimaging optics that are used in precision photometry,¹⁶ backlighting systems, collection systems (Basset,¹⁰¹ Sec. 7, and Steinfeld²¹⁴), scanners,⁴¹ and reflectometers.²¹⁵ They provide a robust means for creating a distribution at the output port that is independent of the angular and spatial luminance distributions at the input port, thereby removing “hot spots” that can be obtained using other methods. Integrating cavities are limited by a trade-off between maximizing the transmission efficiency and maximizing the luminance at the output port. The angular distribution exiting the sphere is nominally an unclipped Lambertian, but there are a number of techniques to modify the Lambertian distribution.²¹⁶

Nonuniformities with Integrating Spheres Integrating cavities, especially spherically shaped cavities, are often used in precision photometry where the flux entering the sphere is measured using the flux exiting the output port of the sphere. For precision photometry, the angular distribution of the exiting flux must be independent of the entering distribution. In general, the flux is scattered multiple times before exiting the sphere, which provides the required independence between input and output distributions. The two main situations where this independence is not obtained are when the output port views the input port and when the output port views the portion of the sphere wall in which the direct illumination from the input port illuminates the sphere wall.

Figure 25 shows an integrating sphere with five ports. A source is inserted into the upper port and the other four ports are labeled. Figure 26 shows the intensity distribution that exits output port 1 of the sphere shown in Fig. 25. The structure in the exiting intensity distribution is roughly Lambertian, but there are nonuniformities in the distribution that occur at angles that can “see” the source or the other output ports. To highlight this effect, rotated views of the sphere are included in Fig. 26. The bright peak in the intensity distribution occurs at the view angle where the output port

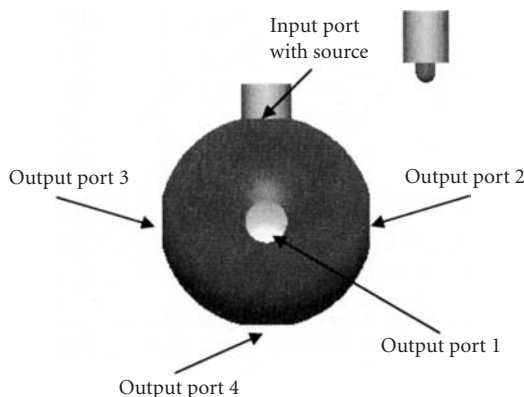


FIGURE 25 Spherical integrating cavity with one input port and four output ports.

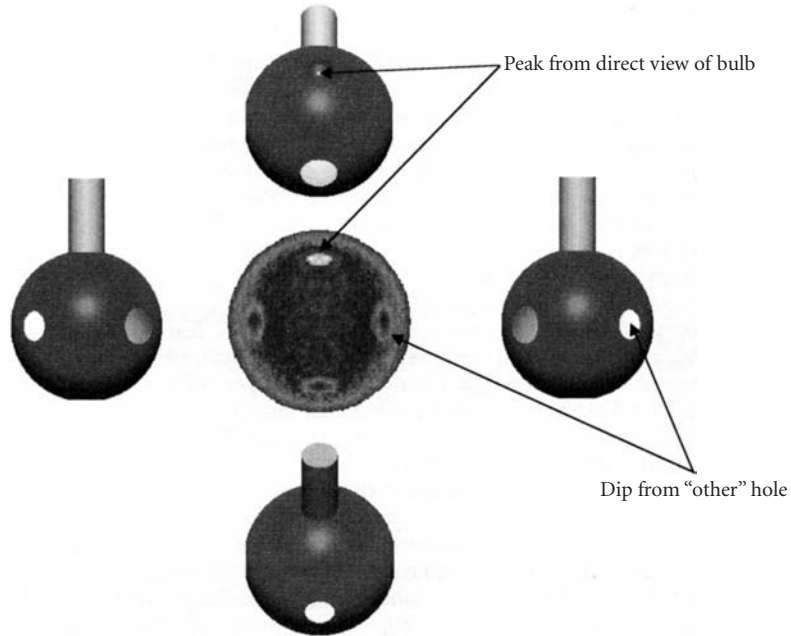


FIGURE 26 Relationship between the intensity distribution from the flux that exits an integrating sphere to what you see when looking back into the sphere.

can “see” direct radiation from the source. The dips in the intensity distribution occur at view angles where the output port can “see” the unused ports in the sphere.

If the source is replaced by a relatively collimated source that illuminates only a small region of the cavity interior (e.g., a collimated laser), then the output port intensity also shows a peak as shown in Fig. 27.

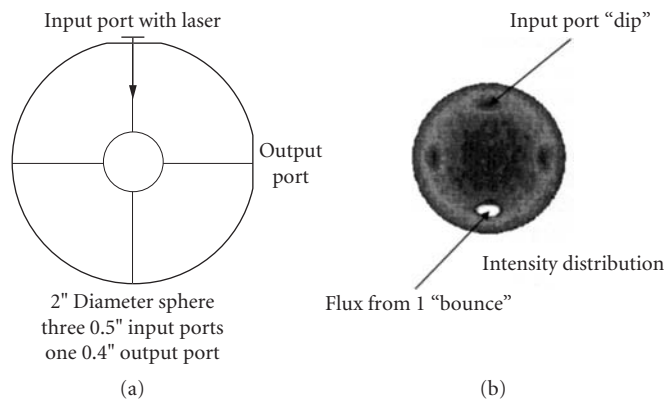


FIGURE 27 Intensity distribution (a) for an integrating sphere (b) that is illuminated with a highly collimated source. The flux that exits the sphere with only one bounce produces a higher intensity than the flux that rattles around and exits in a more uniform manner. The intensity distribution is shown normal to the orientation of the output port.

Three methods that are often used in precision photometry to create an output intensity distribution that is independent of the input involve the use of baffles, satellite spheres, and diffusers. Baffles block radiation from exiting the sphere directly, thereby forcing the exiting flux to scatter multiple times before exiting the sphere. A *satellite sphere* is a secondary sphere that is added to the primary sphere. Diffusers superimpose flux from multiple angles to provide an averaging effect that removes observable structure.

Finite-thickness walls at the exit port will also alter the view angle of external ports and can improve the on-axis intensity because light that reflects off of the walls can go back into sphere and have a chance to add to the output.

Efficiency versus Luminance An important aspect of an integrating cavity is that structure in the angular and spatial distribution of the source's light flux are removed by multiple reflections from the diffuse reflector. Under the assumption that the flux inside of an integrating sphere is scattered in a Lambertian manner, Goebel²¹⁷ has provided a description of the relationship between the fraction of flux that enters the sphere and the fraction of flux that exits the sphere (see also the Technical Information chapter in the LabSphere catalog). Goebel's results can be used to show that the ratio of the flux that exits the sphere via one port to the flux that exits out another port is

$$\text{Efficiency} = \eta = \frac{f_{\text{output}} R_{\text{sphere}}}{1 - f_{\text{input}} R_{\text{input}} - f_{\text{output}} R_{\text{output}} - f_{\text{sphere}} R_{\text{sphere}}} \quad (13)$$

where $f_{\text{sphere}} = (\text{total sphere surface area} - \text{port areas})/\text{total sphere surface area}$

$f_{\text{output}} = \text{output port area}/\text{total sphere surface area}$

$f_{\text{input}} = \text{input port area}/\text{total sphere surface area}$

$R_x = \text{reflectivity of } x$

When the input and output have zero reflectivity, the result is

$$\text{Efficiency} = \eta = \frac{f_{\text{output}} R_{\text{sphere}}}{1 - f_{\text{sphere}} R_{\text{sphere}}} \quad (14)$$

Averaging the input flux over a hemisphere gives a ratio of the luminance at the output port to the luminance at the input port of

$$\frac{\text{Output luminance}}{\text{Input luminance}} = \frac{\eta / \text{area}_{\text{output}}}{1 / \text{area}_{\text{input}}} = \frac{\eta f_{\text{input}}}{f_{\text{output}}} = \frac{R_{\text{sphere}} f_{\text{input}}}{1 - f_{\text{input}} R_{\text{input}} - f_{\text{output}} R_{\text{output}} - f_{\text{sphere}} R_{\text{sphere}}} \quad (15)$$

When the input and output ports have zero reflectivity, then the result is

$$\frac{\text{Output luminance}}{\text{Input luminance}} = \frac{R_{\text{sphere}} f_{\text{input}}}{1 - f_{\text{sphere}} R_{\text{sphere}}} \quad (16)$$

The efficiency and the ratio of input to output luminances are shown in Fig. 28 for the case of $R_{\text{sphere}} = 98.5$ percent.

When the fractional input port size equals the fractional output port size, the sphere transmission is less than 50 percent and the output luminance is less than 50 percent of the input luminance. Increasing the output port size provides higher transmission, but the luminance of the output port is now less than 50 percent of that of the input port. Reducing the output port size can increase the output luminance to values greater than the input luminance; however, the transmission is now less than 50 percent. Thus, an integrating cavity offers the possibility of improving uniformity at the expense of either luminance or efficiency.

Integrating Cavity with Nonimaging Concentrator/Collectors The flux exiting an integrating cavity tends to cover a hemisphere. Adding a nonimaging collector, such as a CPC, to the output port can

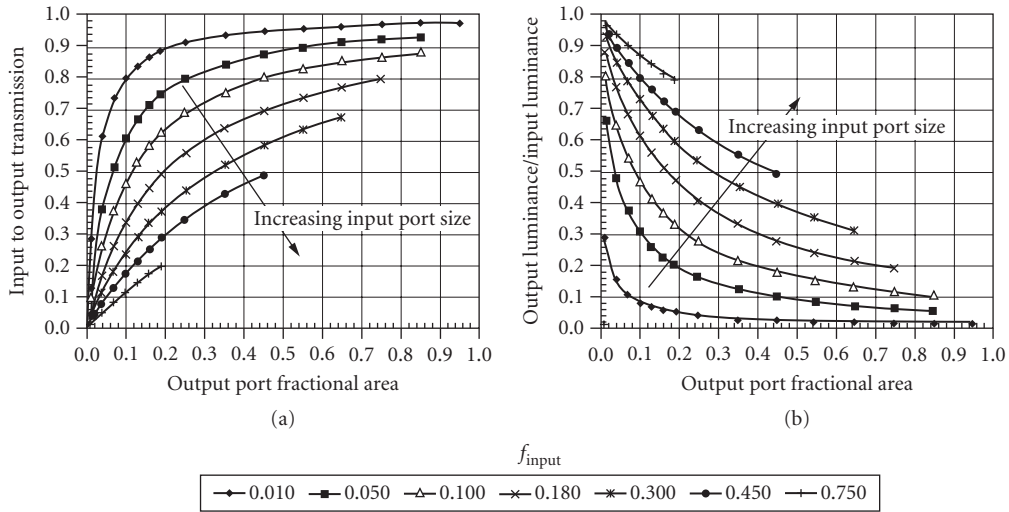


FIGURE 28 Input to output transmission (a) and relative output luminance (b) for the case of sphere reflectivity = 0.985. The different curves are for various input port fractional areas.

reduce the angular spread. The integrating cavity can also include a concentrator at the input port so as to minimize the size of the input port if the angular extent of the radiation at the input port does not fill a hemisphere.

A CEC can be added to the output port of an integrating cavity so that the portion of the sphere surface that contributes to the output distribution is restricted to a region that is known to have uniform luminance.²¹⁵ Flux outside of the CEC field of view (FOV) is rejected by the CEC. CPCs and lenses with CHCs have also been investigated.^{215,218,219}

Modifying the Cavity Output Distribution Prismatic films can also be placed over the exit port of an integrating cavity, which is an approach that has been used by the liquid crystal display community (e.g., the BEFII films commercialized by 3M). The prismatic films increase the intensity at the angles that exit the sphere by reflecting the undesired angles back into the sphere so that they can be rescattered and have another chance to exit the sphere. For systems that require polarized light, nonabsorbing polarizing films can also be used to return the undesired polarization state back into the sphere.

Another approach to selectively controlling the exit angles is the addition of a hemispherical or hemiellipsoid reflector with a hole in the reflector. Webb²²⁰ described this approach for use with lasers.

In a similar manner, reflecting flux back through the source has been used with discharge sources, and is especially effective if the source is transparent to its own radiation.

Mixing Rods (Lightpipes)

When flux enters one end of a long lightpipe, the illuminance at the other end can be quite uniform. The uniformity depends upon the spatial distribution at the lightpipe input, the intensity distribution of the input flux, and the shape of the lightpipe.³⁰⁶ Quite interestingly, straight round lightpipes often do not provide good illuminance uniformity, whereas some shapes like squares and hexagons do provide good illuminance uniformity. This is shown in Fig. 29.

Shapes (Round versus Polygon) Square lightpipes have been used very successfully to provide extremely uniform illuminance distributions. To understand how a lightpipe works, consider the

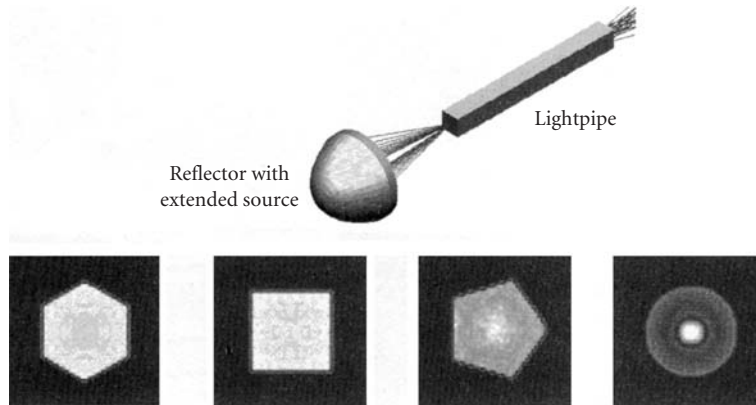


FIGURE 29 When flux from a source is coupled into a mixing rod, the illuminance at the output end of the mixing rod is significantly affected by the shape of the lightpipe. Square and hexagonal lightpipes work much better than round or many other regular polygons.

illuminance distribution that would result at the output face of the lightpipe if the lightpipe's sidewalls were removed. Including the sidewalls means that various regions of the illuminance distribution with no lightpipe are superimposed. In the case of a square, the multiple regions add in a controlled manner such that there is improved uniformity by virtue of the superposition of multiple distributions. If the source distribution is symmetric about the optical axis, then distributions with a positive slope are superimposed onto distributions with a negative slope. This means that the uniformity from N distributions can be better than the addition of N uncorrelated distributions that would provide uniformity with a $1/\sqrt{N}$ standard deviation.

One way to view the operation of a rectangular mixing rod is to consider the illuminance distribution that would occur at the output face of the mixing rod if no sidewall reflections were present. The sidewalls force subregions of this unreflected illuminance distribution to superimpose in a well-controlled manner.²²¹ This is shown in Fig. 30, where a Williamson-type construction is shown on the left. The illuminance distribution at the lightpipe output with no sidewalls is also shown, as is the superposition of the subregions when sidewall reflection is included. Every other subregion is inverted because of the reflection at the sidewall. This inversion provides a difference between a lightpipe used for homogenization and a lens array.

Figure 31 shows the results of a Monte Carlo simulation where a source with a Gaussian illuminance distribution and a Gaussian intensity distribution are propagated through a square lightpipe. The illuminance at the lightpipe input and output faces is shown in the top of the figure. The illuminance distribution in the lower right shows the illuminance distribution if no sidewall reflections occurred. Lines are drawn through the no-sidewall reflection illuminance distribution to highlight the subregions that are superimposed to create the illuminance distribution when the lightpipe is present. The lower left illuminance distribution results from propagating the flux from the output end of the lightpipe back to the input end. Multiple images of the source are present, similar to the multiple images observed in a kaleidoscope. This kaleidoscope effect is one reason why lightpipe mixing is sometimes considered to be the superposition of multiple virtual sources.

In general, if the cross-sectional area is constant along the lightpipe length, and the area is covered completely by multiple reflections with respect to straight sidewalls of the lightpipe, then excellent uniformity can be obtained. The shapes that provide this *mirrored tiling* include squares, rectangles, hexagons, and equilateral triangles.²²² A square sliced along its diagonal and an equilateral triangle sliced from the center of its base to the apex will also work. Pictures of these shapes

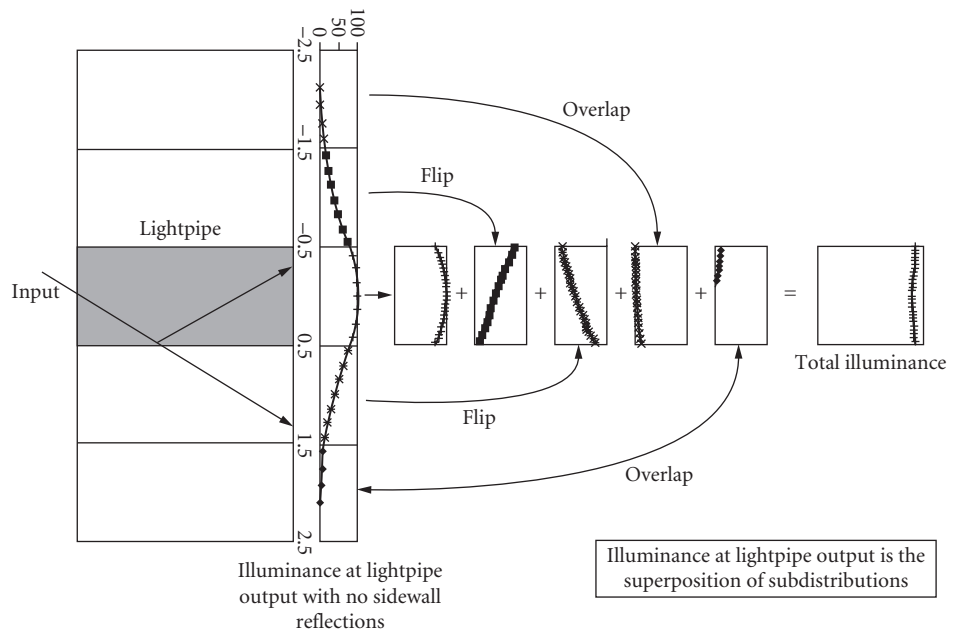


FIGURE 30 Flip-and-fold approach²²¹ for use with rectangular mixing rods. Subregions of the illuminance at the lightpipe output with no sidewall reflections are superimposed to create a uniform distribution. The lightpipe is shown shaded and a Williamson construction is included. The path of one ray is shown with and without sidewall reflection.

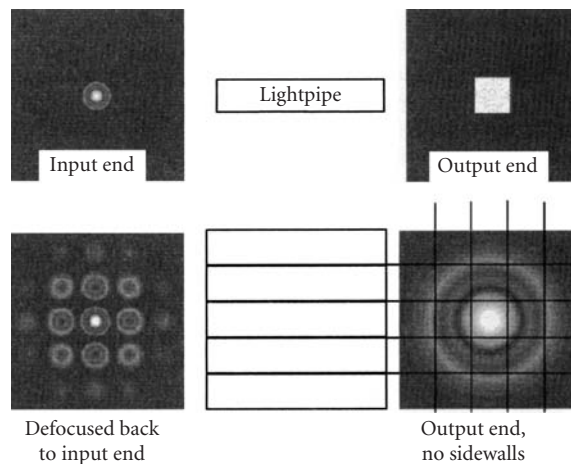


FIGURE 31 Example with a source that has a Gaussian intensity distribution. If the source is propagated to the end of the lightpipe without sidewall reflections, the illuminance distribution at the output can be sliced into distinct regions. These regions are superimposed to create the illuminance distribution that actually occurs with the sidewall reflections. If the flux at the lightpipe output is propagated virtually back to the input end, the multiple images can be observed.

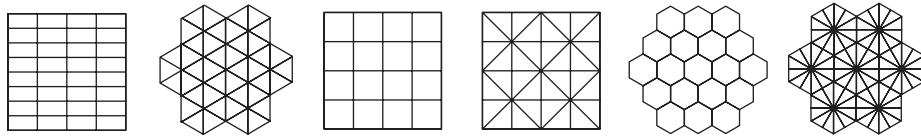


FIGURE 32 Mirror-tiled shapes that can ensure uniformity with adequate length.

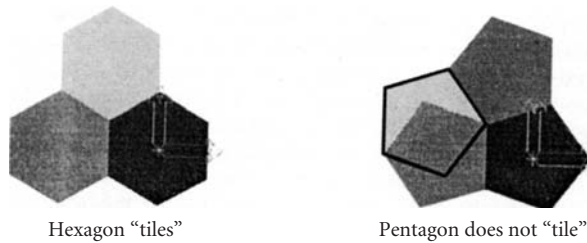


FIGURE 33 Diagram showing why pentagon does not tile like the hexagon. Simulation results are shown in Fig. 29.

arranged in a mirror-tiled format are shown in Fig. 32. The lightpipe length must be selected so as to provide an adequate number of mirrored regions. In particular, the regions near the edge that do not fill a complete mirror-tiled subregion must be controlled so that they do not add a bias to the output distribution.

Shapes that do not mirror-tilde may not provide ideal uniformity as the lightpipe length is increased, at least for a small clipped Lambertian source placed at the input to the lightpipe. One explanation is shown in Fig. 33.

Round lightpipes can provide annular averaging of the angular distribution of the flux that enters the input end. Because of this, the illuminance distribution along the length of the lightpipe tends to be rotationally symmetric; however, the illuminance distribution tends to provide a peak in the illuminance distribution at the center of the lightpipe. The magnitude of this peak relative to the average varies along the length of the lightpipe. If the distribution entering the lightpipe is uniform across the whole input face, and the angular distribution is a clipped Lambertian, then the round lightpipe maintains the uniformity that was present at the input end. However, if either the input spatial distribution or the input angular distributions are not uniform, then nonuniform illuminance distributions can occur along the lightpipe length. This presents the interesting effect that a source can be coupled into a square lightpipe of adequate length to provide a uniform illuminance distribution at the lightpipe output face; however, if the angular distribution is nonuniform and coupled into a round lightpipe, then the illuminance at the output face of the round lightpipe may not be uniform.

Periodic Distributions If the unmirrored illuminance distribution has a sinusoidal distribution with certain periods that match the width of the lightpipe, then the peaks of the sinusoid can overlap. A Fourier analysis of the lightpipe illuminance distribution can help to uncover potential issues.²²¹

Length The length of the lightpipe that is required to obtain good uniformity will depend upon the details of the source intensity distribution. In general, overlapping a 9×9 array of mirrored regions provides adequate uniformity for most rotationally symmetric distributions. For an $f/1$ distribution with a hollow lightpipe, this means that the lightpipe aspect ratio (length/width) should be greater than 6:1. For a lightpipe made from acrylic, the aspect ratio needs to be more like 10:1.

Solid versus Hollow In addition to a solid lightpipe requiring a longer length to obtain a given uniformity, the design of the solid case should consider the following: Fresnel surface losses at the input

and output ends, material absorption, cleanliness of the sidewalls (heat shrink Teflon can help here), and chipping of the corners. For the hollow case, some considerations include dust on the inside, coatings with reflectivities of less than 100 and angular/color dependence, and chips in mirrors where the sidewalls join together.²²¹

With high-power lasers and high-power xenon lamps, a solid coupler may be able to handle the average power density, but the peak that occurs at the input end may introduce damage or simply cause the solid lightpipe to shatter.

Hollow lightpipes can also be created using TIR structures.^{223,224}

Angular Uniformity Although the illuminance at the output end of a lightpipe can be extremely uniform, the angular distribution may not be, especially the “fine” structure in the distribution. One method to smooth this fine structure issue is to add a low-angle diffuser at the lightpipe output end.²²⁵ A diffuser angle that is about the angular difference between two virtual images is often sufficient. Since multiple images of the source are usually required to obtain illuminance uniformity at the lightpipe output, the diffuser does not increase the etendue much. Making the lightpipe longer reduces the angular difference between neighboring virtual images and the required diffuser angle is reduced accordingly. Combinations of lightpipes and lens arrays have also been used.²²⁶

Tapered Lightpipes If the lightpipe is *tapered*, meaning that the input and output ends have different sizes, excellent uniformity can still be obtained. Since the lightpipe tends to need a reasonable length to ensure uniformity, the tapered lightpipe can provide reasonable preservation of etendue. This etendue preservation agrees with Garwin’s³ adiabatic theorem. In general, for a given light pipe length, better mixing is obtained by starting with high angles at the light pipe input and tapering to lower angles than if the input distribution starts with the lower angles.²²⁷ A Williamson construction can be used to explain why the tapered pipe has more virtual images than a straight tapered case. Tapered lightpipes have been used with multimode lasers²²⁸ and solar furnaces.²²²

An example with a small Lambertian patch placed at the entrance to a tapered lightpipe and a straight lightpipe is shown in Fig. 34. The input NA for the untapered case is 0.5/3 for the 100-mm-long

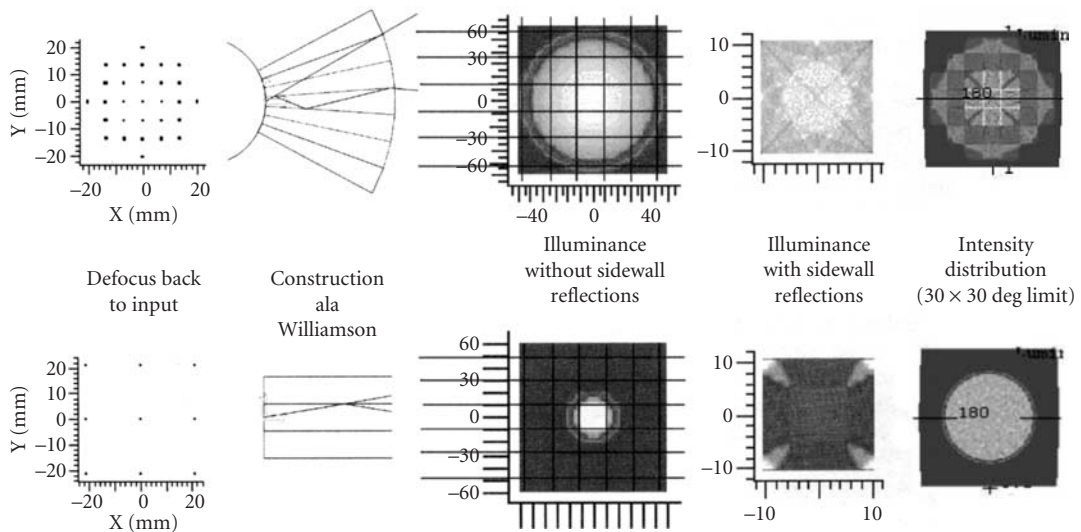


FIGURE 34 Tapered lightpipe example. The upper and lower rows are for the tapered and untapered cases, respectively. Moving left to right, the drawings show the source when defocused from the lightpipe output back to the lightpipe input, a Williamson construction, the no-sidewall illuminance at the lightpipe output, the illuminance at the lightpipe output, and the angular distribution of the flux exiting the lightpipe.

straight lightpipe and 0.5 for the 100-mm-long tapered lightpipe. The input size of the straight lightpipe is a square with a 30-mm diagonal. The input size of the tapered lightpipe is a 10-mm diagonal and the output size is a 30-mm diagonal. The Williamson construction shows that the extreme rays for the 0.5-NA meridian have almost three bounces for the tapered case, but only slightly more than one bounce for the extreme ray in the 0.5/3-NA straight lightpipe case. The illuminance distribution that would occur at the lightpipe output if the lightpipe sidewalls were removed is shown in Fig. 34. Gridlines are superimposed on the picture to highlight that the tapered case has numerous regions that will be superimposed, but the superposition will be small for the straight lightpipe case. The tapered case provides superior uniformity, although the output NA in the tapered case is higher than 0.5/3 because the angular distribution has “tails.” These results are obtained using a source that is a small Lambertian patch. The comparison is typically far less drastic when the input face of the lightpipe is better filled.

Applications of Mixing Rods Mixing rods have found use in a number of applications. Examples include projection systems,²²⁹ providing uniformity with a bundle of fibers using discharge lamps,²²⁷ fiber-to-fiber couplers,²³⁰ and solar simulators.²²²

Mixing rods have also been used with coherent sources. Some applications include semiconductor processing;²³¹ hyperpigmented skin lesions;^{232,233} lithography, laser etching, and laser scanning;²²⁸ and high-power lasers.²³⁴ If the coherence length of the laser is not small compared to the path length for overlapping distributions, then speckle effects must be considered.^{231,234} Speckle effects can be mitigated if the lightpipe width is made large compared to the size of the focused laser beam used to illuminate the input end of the lightpipe (e.g., Refs. 231 and 236). A negative lens for a laser illuminated lightpipe is described by Grojean;²³⁵ however, this means that etendue of the source is less than the etendue of the lightpipe output.

Speckle can be controlled in lightpipes by adding a time-varying component to the input distribution. Some methods include placing a rotating diffuser at the lightpipe input end and moving the centroid of the distribution at the input end. In principle, the diffuser angle can be small so that the change in etendue produced by the diffuser can also be small.

Coherent interference can also be used advantageously, such as in the use of mixing rods of specific lengths that provide sub-Talbot plane reimaging. This was successfully applied to an external cavity laser.²³⁷

In addition to uniformity and the classic kaleidoscope, mixing rods can also be used in optical computing systems by reimaging back to the lightpipe input.^{238,239}

Lens Arrays

Lens arrays are used in applications including Hartmann sensors,²⁴⁰ diode laser array collimators,²⁴¹ beam deflection,²⁴² motion detectors,²⁴³ fiber-optic couplers,²⁴⁴ commercial lighting, and automotive lighting. In commercial and automotive lighting, lens arrays are often used to improve uniformity by breaking the flux into numerous subregions and then superimposing those subregions together. This superposition effect is illustrated in Fig. 35, where fans of rays illuminate three lenslets and are superimposed at the focal plane of a condensing lens. The condensing lens is not used if the area to be illuminated is located far from the lens array.

If the incident flux is collimated, then the superposition of wavefronts provides improved uniformity by averaging the subregions. However, when the wavefront is not collimated, the off-axis beams do not completely overlap the on-axis beams (see Fig. 35*b*). Tandem lens arrays add additional degrees of design freedom and can be used to eliminate the nonoverlap problem (see Fig. 35*c*).

Literature Summary Some journal citations that describe the use of lens arrays with incoherent sources include Zhidkova's²⁴⁵ description of some issues regarding the use of lens arrays in microscope illumination. Ohuchi²⁴⁶ describes a liquid crystal projector system with an incoherent source and describes the use of two tandem lens arrays with nonuniform sizes and shapes for the second lens array. This adjustment to the second lens array provides some adjustment to the coma introduced by the conic reflector.

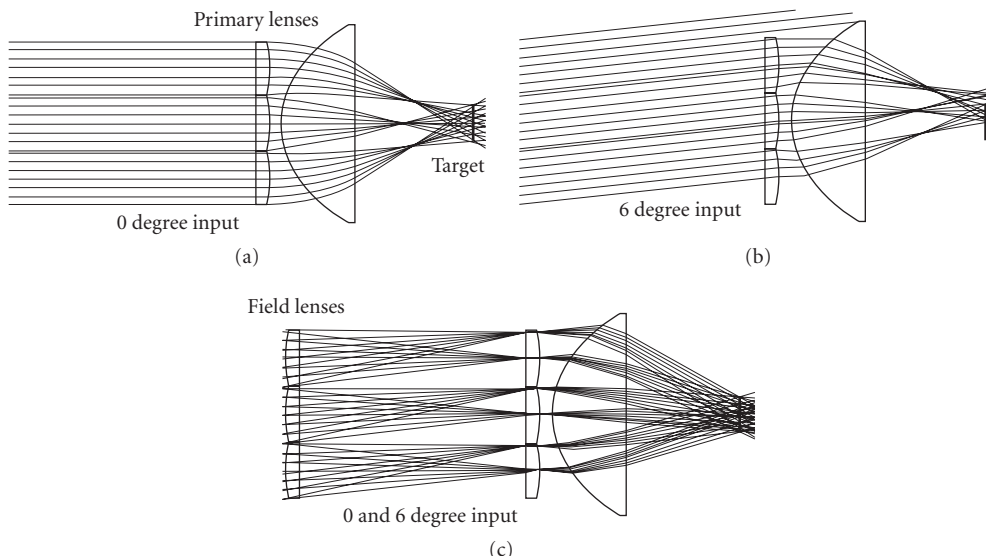


FIGURE 35 Lens arrays used for illumination. (a) Superposition of three subregions of the input wavefront. (b) Off-axis angles are displaced relative to the on-axis angles. (c) Tandem lens array case, where the addition of an array of field lenses allows the superposition to work with a wide range of angles.

Rantsch²⁴⁷ shows one of the earliest uses of tandem lens arrays. Miles²⁴⁸ also uses a tandem lens array system and is concerned about the distribution in the pupil. Ohtu²⁴⁹ uses two sets of tandem lens arrays in a lithography system where one array controls the spatial distribution of the image and the other controls the angular (pupil) distribution of the image. Kudo²²⁶ combines a lens array and lightpipe to provide the same type of control over the illuminance and intensity distributions. Matsumoto²⁵⁰ adds a second set of tandem lens arrays with a sparse fill to correct the fact that an odd spatial frequency term in the illuminance distribution can result in a final nonuniformity, independent of the number of lenslets. Van den Brandt²⁵¹ shows numerous configurations using tandem lens arrays, including the use of nonuniform sizes and shapes in the second lens array. Watanabe²⁵² provides design equations for tandem lens arrays in a converging wavefront.

Journal publications that include lens arrays often involve the use of lens arrays with coherent sources rather than incoherent sources. Deng²⁵³ shows a single lens array with a coherent source and shows that the distribution at the target is a result of diffraction from the finite size of each lenslet and interference between light that propagates through the different lenslets. Nishi²⁵⁴ modifies the Deng system to include structure near the lenslet edges so as to reduce the diffraction effects. Ozaki²⁵⁵ describes the use of a linear array of lenses with an excimer laser, including the use of lenslets where the power of the lenslets varies linearly with distance from the center of the lens array. Glockner²⁵⁶ investigates the performance of lens arrays under coherent illumination as a function of statistical variation in the lenslets. To reduce interference effects with coherent sources, Bett²⁵⁷ describes the use of a hexagonal array of Fresnel zone plates where each zone incorporates a random phase plate. Kato²⁵⁸ describes random phase plates.

Single-Lens Arrays Single-lens-array designs for uniformity tend to have two extremes. One occurs when adding the lens array significantly changes the beam width. This is the beam-forming category where the shape of the beam distribution is essentially determined by the curvature and shape of the lenslets. The other extreme occurs when adding the lens array does not significantly change the beam width. This is the beam-smearing category and is useful for removing substructure in the beam, such as structure from filament coils or minor imperfections in the optical elements. There are many designs that lie between these two extremes.

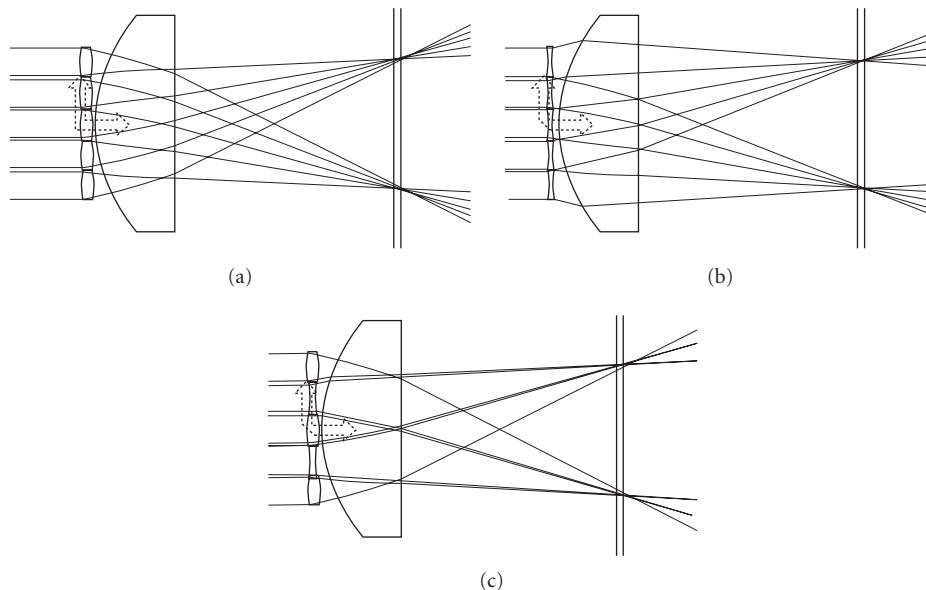


FIGURE 36 Superposition using an array of positive lenses (a), negative lenses (b), and also a hybrid system (c) that uses negative and positive lenses. A collimated wavefront passing through each lenslet is mapped over the target. The hybrid system offers the possibility of a continuous surface with no slope discontinuities.

Figure 36 depicts the beam-forming case where the source divergence is assumed to be small. Positive and/or negative lenses can be used to create the distribution. Hybrid systems where mixtures of positive and negative lenses are used provide the added benefit that the slope discontinuity between neighboring lenses can be eliminated. Slope discontinuities can degrade performance in many systems.

Tandem-Lens Arrays The use of a single-lens array has two issues that sometimes limit its use: First, the lens array increases the etendue because the angles increase but the area stays the same. Second, the finite size of the source introduces a smearing near the edges of the “uniform” distribution. Adding a second lens array in tandem with the first can minimize these limitations.

Figure 37a shows a 1:1 imager in an Abbe (critical) illumination system. Figure 37b shows the same imager with two-lens arrays added. Adding the lens arrays has changed the system from Abbe to Kohler (see the preceding text). In addition, as shown in Fig. 38, the lenslets provide channels that can be designed to be independent of one another. There are three main reasons why two tandem-lens arrays provide uniformity:

1. Each channel is a Kohler illuminator, which is known to provide good uniformity assuming a slowly varying source-intensity distribution.
2. Adding all the channels provides improved uniformity as long as the nonuniformities of one channel are independent of the nonuniformities of the other channels (e.g., a factor of $1/\sqrt{N}$ improvement by adding independent random variables).
3. For symmetric systems, the illuminance from channels on one side of the optical axis tends to have a slope that is equal but opposite to the slope on the other side of the optical axis. The positive and negative slopes add to provide a uniform distribution.

Reflector/Lens-Array Combinations Using a reflector to collect the flux from the source can provide a significant increase in collection efficiency compared to the collection efficiency of a single lens. The

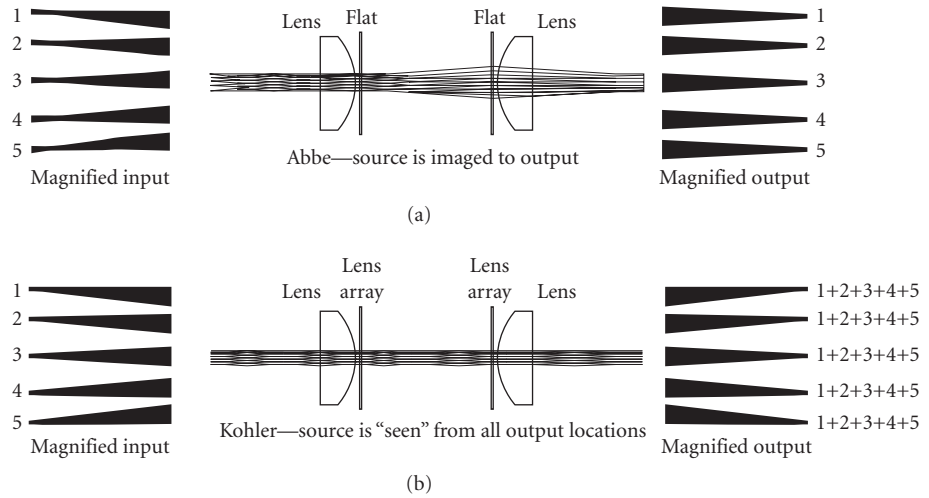


FIGURE 37 Critical illumination (Abbe, *a*) compared to Kohler (*b*). In Abbe, the source is imaged to the output. In Kohler, most of the source is “seen” from all output points.

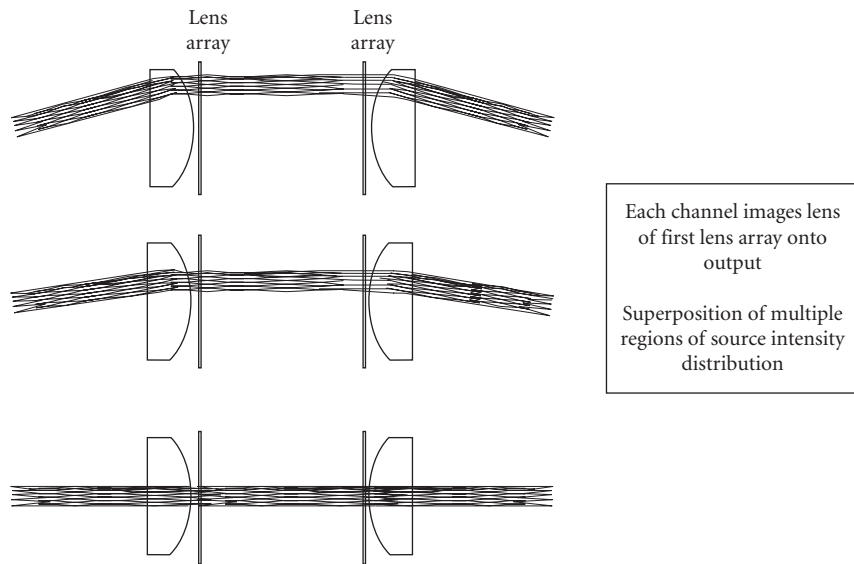


FIGURE 38 Multiple channels in a tandem-lens array. Each channel images a subregion of the flux at the first lens array. The flux from the different channels is superimposed.

drawback is that a single conic reflector does not generally provide uniform illuminance. The use of a lens array in this situation allows the illuminance distribution to be broken into multiple regions, which are then superimposed on top of each other at the output plane.

Examples showing the illuminance distribution at the first lens array, the second lens array, and the output plane are shown in Fig. 39 for the case of a lens collector, Fig. 40 for the case of a parabola

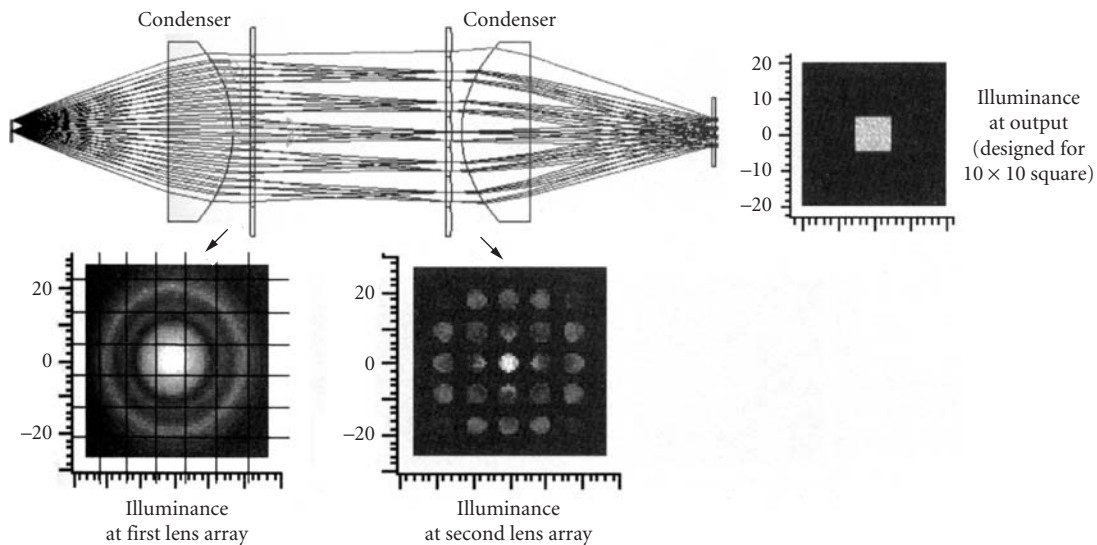


FIGURE 39 Source/condenser/tandem lens arrays/condenser configuration. Illuminance at first lens array is broken into subregions that are superimposed at the target by the second lens array/second condenser. All units are millimeters.

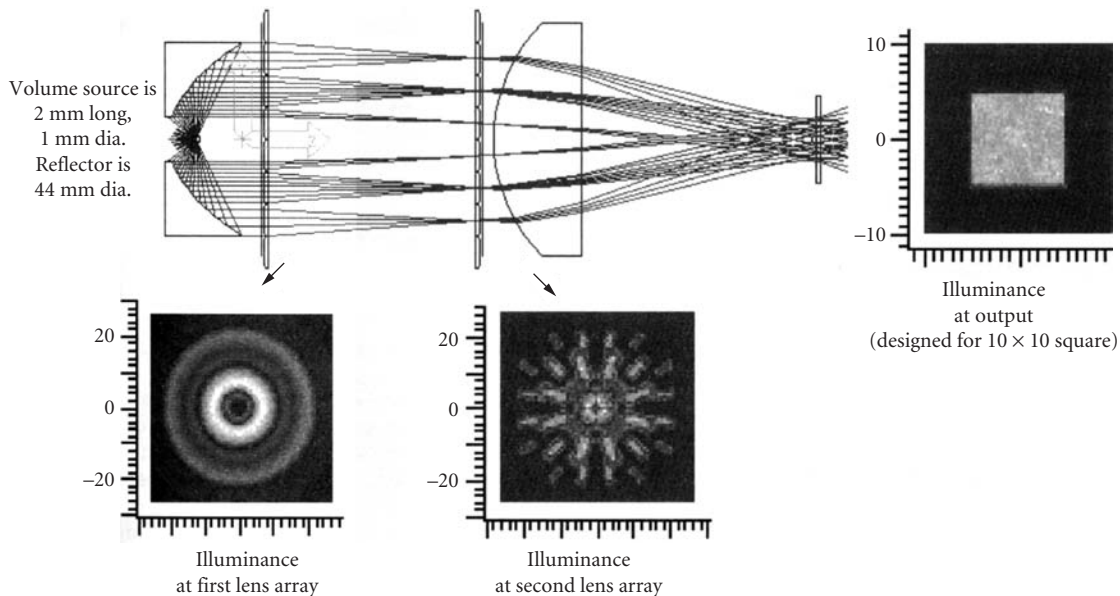


FIGURE 40 Parabola/tandem lens arrays/condenser configuration. Nonuniform magnification from the reflector (coma) is seen in the array of images at the second lens array. The illuminance at the target is the superposition of subregions at the first lens array. All units are millimeters.

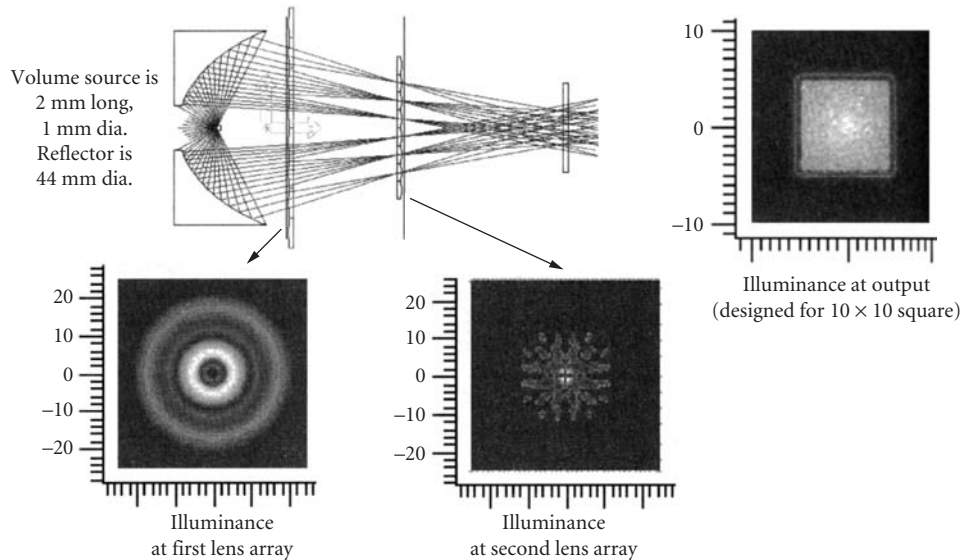


FIGURE 41 Ellipse/tandem lens arrays/condenser configuration. Nonuniform magnification from the reflector (coma) is seen in the array of images at the second lens array. The illuminance at the target is the superposition of subregions at the first lens array. The spacing between lenses is different for the two lens arrays. All units are millimeters.

collector, and Fig. 41 for the case of an elliptical collector. These examples show that nonuniform illuminance distribution at the first lens array is broken into regions. Each region is focused onto the second set of lenses, which creates an array of images of the source. The lenslets in the second array image the illuminance distribution at the first lens array onto the output plane. Enhanced uniformity is obtained because the distributions from each of the channels are superimposed at the target. The discrete images of the source that are seen at the second lens array indicate that the pupil of the output distribution is discontinuous.

Improvements to this basic idea include the use of lens arrays where the lenslets in the second lens array do not all have the same size/shape.²⁵¹ The curvature of each lenslet in the first array is decentered to “aim” the flux toward the proper lenslet in the second array. When nonuniform lenslet size/shapes are used, the etendue loss that results from the use of a conic reflector can be minimized. Such a system can be used to create asymmetric output distributions without a correspondingly misshapen pupil distribution (see earlier section on image dissectors).

Tailored Optics

The tailoring of a reflector to provide a desired distribution with a point or line source has been explored.^{8,259–265} Some simple cases are available in closed form.^{88,266,267} Refractive equivalents have also been investigated.^{268,269}

In a manner similar to the lens configurations shown in Fig. 36, tailored reflectors can produce converging or diverging wavefronts. Examples are shown in Fig. 42. The converging reflector creates a smeared image of the source between the reflector and the target. The diverging reflector creates a smeared image of the source outside of the region between the reflector and target. This terminology was developed assuming that the target is far from the reflector. The terms *compound hyperbolic* and *compound elliptic* are also used to describe these two types of reflector configurations.²⁷⁰

If the flux from the two sides of the reflector illuminates distinct sides of the target, then there are four main reflector classifications (e.g., Ref. 26, pp. 6.20–6.21; Refs. 91 and 260). Examples of the

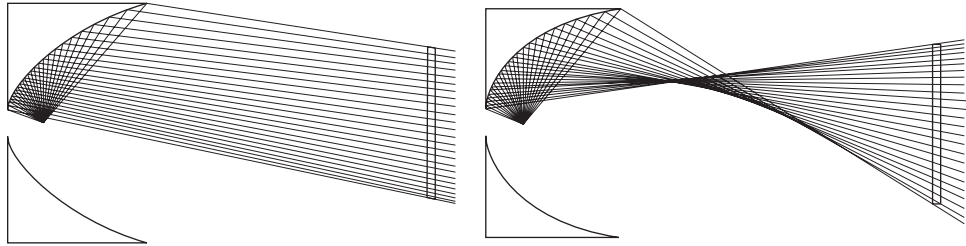


FIGURE 42 Two main classes of reflector types—diverging (*left*) and converging (*right*). The target is shown on the right of each of the two figures.

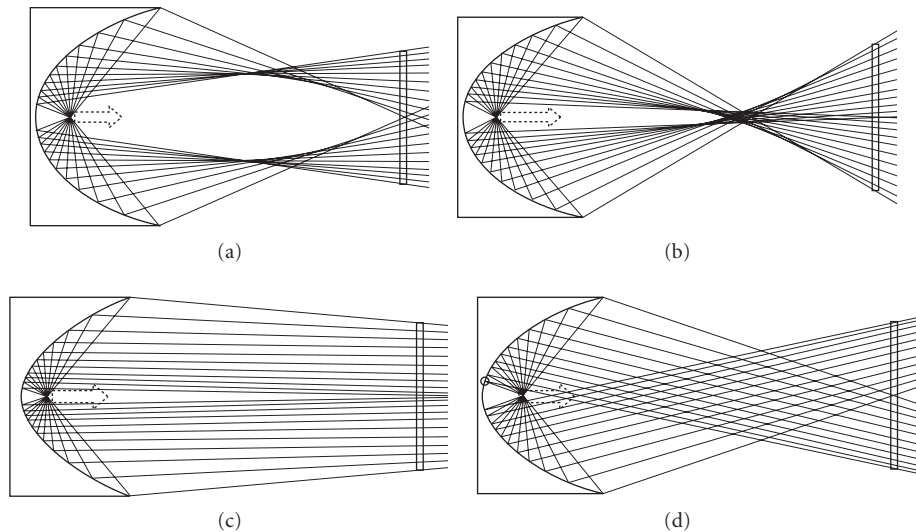


FIGURE 43 Four types of one-to-one mapping reflectors where the flux from each side of the reflector does not overlap at the target. (a) Uncrossed converging. (b) Crossed converging. (c) Uncrossed diverging. (d) Crossed diverging.

four cases are pictured in Fig. 43. They are the crossed and uncrossed versions of the converging and diverging reflectors, where *crossed* is relative to the center of the optical axis. Elmer²⁷¹ uses the terms *single* and *divided* for uncrossed and crossed, respectively. The uncrossed converging case tends to minimize issues with rays passing back through the source. The four cases shown provide one-to-one mappings because each point in the target receives flux from only one portion of the reflector.

Tailored Reflectors with Extended Sources Elmer⁸ describes some aspects of the analysis/design of illumination systems with extended sources using what are now called *edge rays* (see also Ref. 271). Winston²⁶⁶ sparked renewed interest in the subject. The following is a brief summary of some of the edge-ray papers that have been published in recent years. Some concentration related citations are included. Edge rays were also discussed earlier.

Gordon²⁷² adds a gap between the source and reflector to offset the cosine cubed effect and obtain a sharp cutoff. Gordon²⁷³ explores the case where the extreme direction is a linear function of polar angle. Friedman¹⁴⁷ uses a variable angle to design a secondary for a parabolic primary. Gordon²⁷⁴ describes a tailored edge ray secondary for a Fresnel primary, in particular for a heliostat field. Winston²⁷⁰ describes the tailoring of edge rays for a desired functional rather than maximal

concentration. Ries²⁷⁵ describes the tailoring based on a family of edge rays. Rabl²⁷⁶ tailors the reflector by establishing a one-to-one correspondence between target points and edge rays for planar sources. Ong^{277,278} explores tailored designs using full and partial involutes. Jenkins²⁷⁹ also describes a partial involute solution and generalizes the procedure as an integral design method.²⁸⁰ Ong⁹¹ explores the case of partial involutes with gaps. Gordon²⁸¹ relates the string construction to tailored edge ray designs for concentrators.

Faceted Structures

Nonuniform distributions can be made uniform through the use of reflectors where the reflector aims multiple portions of the flux from the source toward common locations at the target. In many cases, the resulting reflector has a faceted appearance. If the source etendue is significantly smaller than the etendue of the illumination distribution, then smearing near the edges of the uniform region can be small. However, if the source etendue is not negligible compared to the target, then faceted reflectors tend to experience smearing similar to the case of one lens array.

The term *faceted reflector* identifies reflectors composed of numerous distinct reflector regions. Elmer²⁷⁰ uses the term *multiphase* and Ong⁹¹ uses the term *hybrid* to describe these types of reflectors. If there are more than two regions, then the regions are often called *facets*. As with lens arrays, there are two extremes for faceted reflector designs. One is the beam-smearing category, where the facets remove substructure in the beam. The other extreme is the beam-forming category, where each facet creates the same distribution. If all facets create the same distribution, Elmer²⁷¹ calls the reflector *homogeneous*. If the distributions are different, then the term *inhomogeneous* is used.

For a given meridional slice, each facet can create an image of the source either in front of or behind the reflector. Using the convergent/divergent terminology to distinguish where the blurred image of the source occurs (see Fig. 42), a meridional slice of a faceted reflector can be composed of an array of convergent or divergent facets. These two cases are shown in Fig. 44 *a* and *b*. Flat facets are

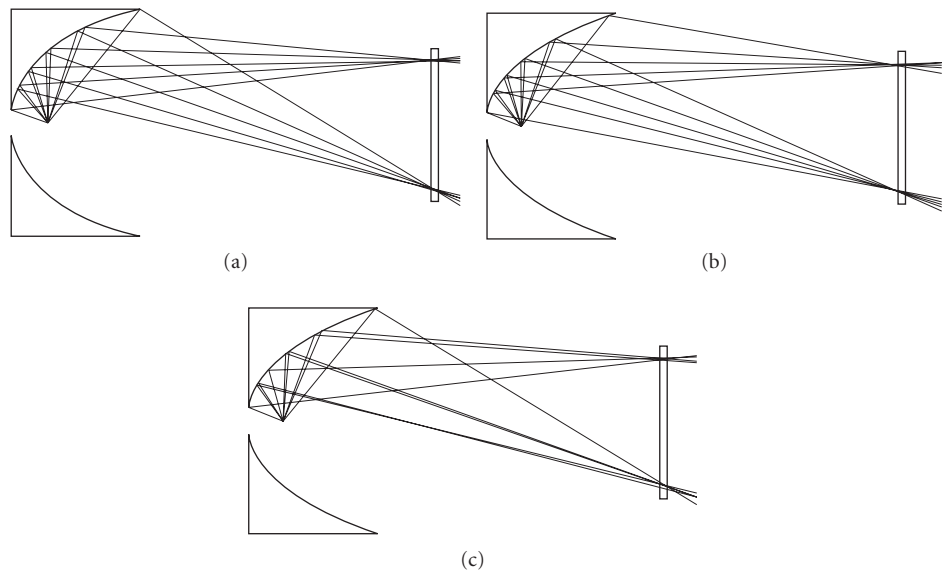


FIGURE 44 Some faceted reflector profiles. Convergent (*a*), divergent (*b*), and mixed convergent/divergent (*c*) profiles are all shown with five facets for the upper half of the reflector. Rays from the source that hit near the edges of the facets are shown. The rays for the upper facets are highlighted for all three cases. The rays for the second facet are highlighted in the mixed case.

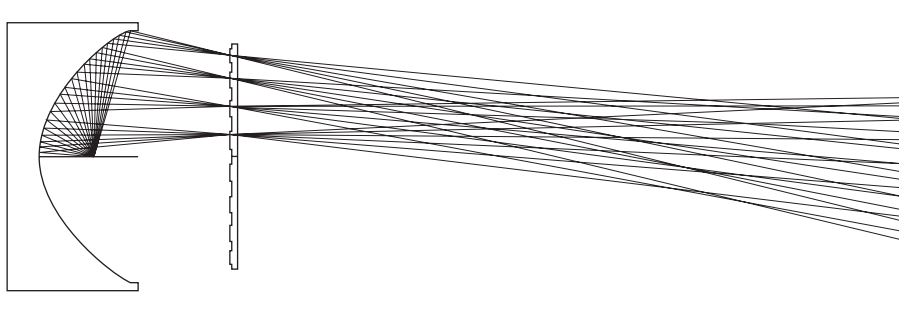


FIGURE 45 Reflector with array of concave facets and an array of lenslets. The lenslets image the facets.

often desired because of the simplicity of fabrication and fall under the divergent category because the image of the source is behind the reflector.

The intersection between facets of the same type introduces a discontinuity in the slope of the reflector. Mixed convergent/divergent versions are possible, as shown in Fig. 44c. Mixed versions offer the possibility of minimizing discontinuities in the slope of the reflector curve.

If all the distributions are superimposed, then improved uniformity can arise through superposition if sufficient averaging exists. Consider a manufacturing defect where a subregion of the reflector that is about the size of a facet is distorted. If a superposition approach is used, then the effect of the defect is small because each facet provides only a small contribution to target distribution. However, if a one-to-one mapping approach is used, then the same defect distorts an entire subregion of the target distribution because there is no built-in redundancy.

The effects of dilution should be considered when faceted optics are used. A single-lens array provides insufficient degrees of freedom to provide homogeneous superposition without introducing dilution, which is why a second lens array is often added (e.g., the tandem-lens-array configuration). Schering²⁸² shows lenslets placed on a reflector combined with a second array of lenslets. A straightforward extension of lenslets on a reflector is the use of concave facets, as pictured in Fig. 45.

Some Literature on Faceted Structures If the source is highly collimated (e.g., a laser) and has a Gaussian distribution, then only a couple of regions of the wavefront need to be superimposed to provide excellent uniformity. Descriptions of optical implementations that emphasize nearly Gaussian wavefronts include a four-surface prism by Kawamura,²⁸³ a wedge prism by Lacomat,²¹¹ and a multipass wedge prism by Sameshima.²⁸⁴ Austin²⁸⁵ has described a prism system for an excimer laser. Brunsting²⁸⁶ has described a system with nonplanar wedges. Bruno²⁸⁷ and Latta²²⁸ have discussed examples of mirrored systems with very few superimposed wavefronts.

For the more general case, uniformity can be obtained by superimposing a number of different regions of the beam. Dourte²⁸⁸ describes a 2D array of facets for laser homogenization. Ream²⁸⁹ discusses a system with a convex array of flat facets for a laser. Doherty²⁹⁰ describes a segmented conical mirror that can be used to irradiate circular, annular, and radial surfaces with a laser source. Dagenais²⁹¹ describes the details of a system with nonplanar facets for use with a laser in paint stripping. Dickey²⁹² shows a laser-based system composed of two orthogonal linear arrays of facets and includes curvature in the facet structures to minimize the effects of diffraction. Experimental data for two orthogonal linear arrays of facets [called a *strip mirror integrator* (SMI)] are provided by Geary.²⁹³ Lu²⁹⁴ has described a refractive structure using a multiwedge prism structure. Henning²⁹⁵ and Unnebrink²⁹⁶ show curved facets for use with UV laser systems.

Laser-based systems provide the degree of freedom that the wavefront divergence can be small, although the effects of diffraction and coherence must be considered in those designs. The divergence of nonlaser sources is often an important issue.

Jolley²⁵⁹ describes a faceted mirror design with disc and spherical sources. David²⁹⁷ describes faceted reflectors like the ones used in commercial lighting products and highlights the fact that facets on a

conic reflector base can improve uniformity with very little increase in “spot size” for long linear sources. Coaton²⁹⁸ provides a brief description of a faceted reflector design with a point source. Donohue²⁹⁹ describes faceted structures in headlamps.

There are many patents that describe the use of faceted structures for incoherent sources. Examples include Wiley,³⁰⁰ who uses the facets to remove filament image structure; Laudenschlager,³⁰¹ who describes design equations for a faceted reflector; Heimer,³⁰² who uses a faceted reflector in a lithography system; and Hamada,³⁰³ who uses a faceted structure with a linear lamp.

Ronnelid³⁰⁴ uses a corrugated reflector to improve the angular acceptance of a 2D CPC in solar applications. Receiver uniformity can also be improved in CPCs using irregular surfaces (Ref. 4, pp. 153–161).

Faceted structures have often been used in complex lightpipes such as those used in the illumination of backlit liquid-crystal displays and instrument panels. For instrument-panel application, more recent efforts have been directed toward using aspheric surfaces.³⁰⁵

39.7 ACKNOWLEDGMENTS

Special acknowledgement is given to Doug Goodman, whose list of uniformity references and OSA course notes provided an important starting point for portions of this chapter. The support of Optical Research Associates (ORA) during the manuscript preparation is also appreciated, as are the editorial and terminology discussions with many of the technical staff at ORA. Simulation and ray trace results were all generated using either LightTools® or CodeV®.

39.8 REFERENCES

1. H. Hinterberger, L. Lavoie, B. Nelson, R. L. Sumner, J. M. Watson, R. Winston, and D. M. Wolfe, “The Design and Performance of a Gas Cerenkov Counter with Large Phase-Space Acceptance,” *Rev. Sci. Instrum.* **41**(3):413–418 (1970).
2. D. A. Harper, R. H. Hildebrand, R. Stiening, and R. Winston, “Heat Trap: An Optimized Far Infrared Field Optics System,” *Appl. Opt.* **15**:53–60 (1976) (Note: References are on page 144 of vol. 15).
3. R. L. Garwin, “The Design of Liquid Scintillation Cells,” *Rev. Sci. Instrum.* **23**:755–757 (1952).
4. R. Winston, J. C. Miñano, and P. Benítez, *Nonimaging Optics*, Elsevier Academic Press, San Diego, CA, 2005.
5. W. J. Cassarly and J. M. Davenport, “Fiber Optic Lighting: The Transition from Specialty Applications to Mainstream Lighting,” SAE 1999-01-0304 (1999).
6. A. E. Rosenbluth and R. N. Singh, “Projection Optics for Reflective Light Valves,” *Proc. SPIE* **3634**:87–111 (1999).
7. J. Bortz, N. Shatz, and R. Winston, “Advanced Nonrotationally Symmetric Reflector for Uniform Illumination of Rectangular Apertures,” *Proc. SPIE* **3781**:110–119 (1999).
8. W. Elmer, *The Optical Design of Reflectors*, 3rd ed., Academic Press, New York, 1989.
9. J. Palmer, Frequently Asked Questions, www.optics.arizona.edu/Palmer/rpfaq/rpfaq.htm, 1999.
10. J. M. Palmer, “Radiometry and Photometry: Units and Conversions,” in *OSA Handbook of Optics*, 3rd ed., vol. III, chap. 3, McGraw-Hill, New York, 2009.
11. A. Stimson, *Photometry and Radiometry for Engineers*, John Wiley & Sons, 1974.
12. F. Grum and R. J. Becherer, *Optical Radiation Measurements*, vol. 1, *Radiometry*, Academic Press, 1979.
13. W. Budde, *Optical Radiation Measurements*, vol. 4, *Physical Detectors of Optical Radiation*, Academic Press, 1983.
14. A. D. Ryer, *Light Measurement Handbook*, International Light, Inc., 1997.
15. W. L. Wolfe, *Introduction to Radiometry*, *Proc. SPIE Press*, 1998.
16. W. R. McCluney, *Introduction to Radiometry and Photometry*, Artech, 1994.
17. D. Goodman, “Geometrical Optics,” in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 1, McGraw-Hill, New York, 2009.

18. M. Born and E. Wolf, *Principles of Optics*, Cambridge Press, pp. 522–525, 1980.
19. M. S. Brennessoltz, “Light Collection Efficiency for Light Valve Projection Systems,” *Proc. SPIE* **2650**:71–79 (1996).
20. B. A. Jacobson, R. Winston, and P. L. Gleckman, “Flat-Panel Fluorescent Backlights with Reduced Illumination Angle: Backlighting Optics at the Thermodynamic Limit,” *SID 92 Digest*, 423–426 (1992).
21. H. Ries, “Thermodynamic Limitations of the Concentration of Electromagnetic Radiation,” *J. Opt. Soc. Am.* **72**(3):380–385 (1982).
22. H. Hottel, “Radiant Heat Transmission,” in *Heat Transmission*, W. H. McAdams, ed., 3rd ed., McGraw Hill, New York, 1954.
23. R. Winston, “Cone-Collectors for Finite Sources,” *Appl. Opt.* **17**(5):688–689 (1978).
24. G. H. Derrick, “A Three-Dimensional Analogue of the Hottel String Construction for Radiation Transfer,” *Opt. Acta* **32**:39–60 (1985).
25. C. Weiner, *Lehrbuch der darstellenden Geometrie*, vol. 1, Leipzig, 1884 (see F. E. Nicodemus, *Self Study Manual on Optical Radiation Measurements*, part 1, *Concepts*, NBS, Washington, March 1976).
26. M. S. Rea (ed.), *Lighting Handbook*, Illuminating Engineering Society of North America, 1993.
27. F. O. Bartell, “Projected Solid Angle and Black Body Simulators,” *Appl. Opt.* **28**(6):1055–1057 (March 1989).
28. D. G. Koch, “Simplified Irradiance/Illuminance Calculations in Optical Systems,” *Proc. SPIE* **1780**:226–242, 1992.
29. A. Rabl and R. Winston, “Ideal Concentrators for Finite Sources and Restricted Exit Angles,” *Appl. Optics* **15**:2880–2883 (1976).
30. E. Harting, D. R. Mills, and J. E. Giutronich, “Practical Concentrators Attaining Maximal Concentration,” *Opt. Lett.* **5**(1):32–34 (1980).
31. A. Luque, “Quasi-Optimum Pseudo-Lambertian Reflecting Concentrators: An Analysis,” *Appl. Opt.* **19**(14): 2398–2402 (1980).
32. R. M. Saraiji, R. G. Mistrick, and M. F. Modest, “Modeling Light Transfer through Optical Fibers for Illumination Applications,” *IES* 128–139 (Summer 1996).
33. T. L. Davenport, W. J. Cassarly, R. L. Hansler, T. E. Stenger, G. R. Allen, and R. F. Buelow, “Changes in Angular and Spatial Distribution Introduced into Fiber Optic Headlamp Systems by the Fiber Optic Cables,” *SAE*, Paper No. 981197, 1998.
34. S. Doyle and D. Corcoran, “Automated Mirror Design Using an Evolution Strategy,” *Opt. Eng.* **38**(2):323–333 (1999).
35. I. Ashdown, “Non-imaging Optics Design Using Genetic Algorithms,” *J. Illum. Eng. Soc.* **3**(1):12–21 (1994).
36. Y. Nakata, “Multi B-Spline Surface Reflector Optimized with Neural Network,” *SAE* **940638**:81–92 (1994).
37. N. E. Shatz and J. C. Bortz, “Inverse Engineering Perspective on Nonimaging Optical Design,” *Proc. SPIE* **2538**:136–156, 1995.
38. C. Gilray and I. Lewin, “Monte Carlo Techniques for the Design of Illumination Optics,” *IESNA Annual Conference Technical Papers*, Paper #85, pp. 65–80, July 1996.
39. J. C. Schweyen, K. Garcia, and P. L. Gleckman, “Geometrical Optical Modeling Considerations for LCD Projector Display Systems,” *Proc. SPIE* **3013**:126–140 (1997).
40. D. Z. Ting and T. C. McGill, “Monte Carlo Simulation of Light-Emitting Diode Light-Extraction Characteristics,” *Opt. Eng.* **34**(12):3545–3553 (1995).
41. M. Kaplan, “Monte Carlo Calculation of Light Distribution in an Integrating Cavity Illuminator,” *Proc. SPIE* **1448**:206–217 (1991).
42. B. G. Crowther, “Computer Modeling of Integrating Spheres,” *Appl. Opt.* **35**(30):5880–5886 (1996).
43. R. C. Chaney, “Monte Carlo Simulation of Gamma Ray Detectors using Scintillation Fibers,” *EUV, X-ray, and Gamma-Ray Instrumentation for Astronomy and Atomic Physics*; Proceedings of the Meeting, San Diego, California, Aug. 7–11, 1989 (A90-50251 23–35), SPIE 1989.
44. J. A. Bamberg, “Scintillation Detector Optimization Using GUERAP-3 Radiation Scattering in Optical Systems, Proceedings of the Seminar, Huntsville, Alabama, September 30–October 1, 1980 (A81-36878 16-74) p. 86–93, SPIE 1981.
45. B. K. Likeness, “Stray Light Simulation with Advanced Monte Carlo Techniques,” *Stray-Light Problems in Optical Systems*; Proceedings of the Seminar, Reston, VA, April 18–21, 1977 (A78-40270 17-35) pp. 80–88, SPIE 1977.

46. E. R. Freniere, "Simulation of Stray Light in Optical Systems with the GUERAP III," *Radiation Scattering in Optical Systems*; Proceedings of the Seminar, Huntsville, AL, September 30–October 1, 1980 (A81-36878 16-74) pp. 78–85, SPIE 1981.
47. N. Shatz, J. Bortz, and M. Dassanayake, "Design Optimization of a Smooth Headlamp Reflector to SAE/DOT Beam-Shape Requirements," SAE 1999.
48. T. Hough, J. F. Van Derlofske, and L. W. Hillman, "Measuring and Modeling Intensity Distributions of Light Sources in Waveguide Illumination Systems," *Opt. Eng.* **34**(3):819–823 (1995).
49. R. E. Levin, "Photometric Characteristics of Light-Controlling Apparatus," *Illum. Eng.* **66**(4):202–215 (1971).
50. R. J. Donohue and B. W. Joseph, "Computer Synthesized Filament Images from Reflectors and through Lens Elements for Lamp Design and Evaluation," *Appl. Opt.* **14**(10):2384–2390 (1975).
51. I. Ashdown, "Near Field Photometry: A New Approach," *J. IES* 163–180 (Winter 1993).
52. R. D. Stock and M. W. Siegel, "Orientation Invariant Light Source Parameters," *Opt. Eng.* **35**(9):2651–2660 (1996).
53. M. W. Siegel and R. D. Stock, "Generalized Near-Zone Light Source Model and Its Application to Computer Automated Reflector Design," *Opt. Eng.* **35**(9):2661–2679 (1996).
54. P. V. Shmelev and B. M. Khana, "Near-Field Modeling versus Ray Modeling of Extended Halogen Light Source in Computer Design of a Reflector," *Proc. SPIE* (July 27–28, 1997).
55. R. Rykowski and C. B. Wooley, "Source Modeling for Illumination Design," 3130B-27, *Proc. SPIE* (July 27–28, 1997).
56. T. P. Vogl, L. C. Lintner, R. J. Pegis, W. M. Waldbauer, and H. A. Unvala, "Semiautomatic Design of Illuminating Systems," *Appl. Opt.* **11**(5):1087–1090 (1972).
57. "1998 IESNA Software Survey," *Lighting Design and Application* **28**(10):53–62 (1998).
58. "1999 IESNA Software Survey," *Lighting Design and Application* **29**(12):39–48 (1999).
59. W. T. Welford, *Aberrations of Optical Systems*, Adam Hilger, Bristol, 1986, pp. 158–161.
60. W. J. Smith, *Modern Lens Design*, McGraw Hill, New York, 1992.
61. Y. Shimizu and H. Takenaka, "Microscope Objective Design," in *Advances in Optical and Electron Microscopy*, vol. 14, Academic Press, San Diego, CA, 1994.
62. G. P. Smestad, "Nonimaging Optics of Light-Emitting Diodes: Theory and Practice," *Proc. SPIE* **1727**: 264–268 (1992).
63. W. T. Welford and R. Winston, "Two-Dimensional Nonimaging Concentrators with Refracting Optics," *J. Opt. Soc. Am.* **69**(6):917–919 (1979).
64. R. Kingslake, *Lens Design Fundamentals*, Academic Press, San Diego, CA, 1978.
65. E. Hecht and A. Zajac, *Optics*, Addison-Wesley, 1979.
66. S. F. Ray, *Applied Photographic Optics*, Focal Press, 1988.
67. G. Schultz, "Achromatic and Sharp Real Image of a Point by a Single Aspheric Lens," *Appl. Opt.* **22**(20): 3242–3248 (1983).
68. J. C. Minano and J. C. Gonzalez, "New Method of Design of Nonimaging Concentrators," *Appl. Opt.* **31**(16): 3051–3060, 1992.
69. E. A. Boettner and N. E. Barnett, "Design and Construction of Fresnel Optics for Photoelectric Receivers" *J. Opt. Soc. Am.* **41**(11):849–857 (1951).
70. Fresnel Technologies, Fort Worth, TX. See <http://www.fresneltech.com/pdf/FresnelLenses.pdf>
71. S. Sinzinger and M. Testorf, "Transition Between Diffractive and Refractive Microoptical Components," *Appl. Opt.* **34**(26):5670–5676 (1995).
72. F. Erismann, "Design of Plastic Aspheric Fresnel Lens with a Spherical Shape," *Opt. Eng.* **36**(4):988–991 (1997).
73. D. J. Lamb and L. W. Hillman, "Computer Modeling and Analysis of Veiling Glare and Stray Light in Fresnel Lens Optical Systems," *Proc. SPIE* **3779**:344–352 (1999).
74. J. F. Goldenberg and T. S. McKechnie, "Optimum Riser Angle for Fresnel Lenses in Projection Screens," U.S. Patent 4,824,227, 1989.
75. M. Collares-Pereira, "High Temperature Solar Collector with Optimal Concentration: Nonfocusing Fresnel Lens with Secondary Concentrator," *Solar Energy* **23**:409–419 (1979).

76. W. A. Parkyn and D. G. Pelka, "Compact Nonimaging Lens with Totally Internally Reflecting Facets," *Proc. SPIE* **1528**:70–81 (1991).
77. R. I. Nagel, "Signal Lantern Lens," U.S. Patent 3,253,276, 1966.
78. W. A. Parkyn, P. L. Gleckman, and D. G. Pelka, "Converging TIR Lens for Nonimaging Concentration of Light from Compact Incoherent Sources," *Proc. SPIE* **2016**:78–86 (1993).
79. W. A. Parkyn and D. G. Pelka, "TIR Lenses for Fluorescent Lamps," *Proc. SPIE* **2538**:93–103 (1995).
80. W. A. Parkyn and D. G. Pelka, "New TIR Lens Applications for Light-Emitting Diodes," *Proc. SPIE* **3139**:135–140 (1997).
81. V. Medvedev, W. A. Parkyn, and D. G. Pelka, "Uniform High-Efficiency Condenser for Projection Systems," *Proc. SPIE* **3139**:122–134 (1997).
82. V. Medvedev, D. G. Pelka, and W. A. Parkyn, "Uniform LED Illuminator for Miniature Displays," *Proc. SPIE* **3428**:142–153 (1998).
83. D. F. Vanderwerf, "Achromatic Catadioptric Fresnel Lens," *Proc. SPIE* **2000**:174–183 (1993).
84. J. Spigulis, "Compact Illuminators, Collimators and Focusers with Half Spherical Input Aperture," *Proc. SPIE* **2065**:54–59 (1994).
85. D. Silvergate, "Collimating Compound Catadioptric Immersion Lens," U.S. Patent 4,770,514, 1988.
86. McDermott, "Angled Elliptical Axial Lighting Device," U.S. Patent 5,894,196, 1999.
87. D. Korsch, *Reflective Optics*, Academic Press, San Diego, CA, 1991.
88. D. E. Spencer, L. L. Montgomery, and J. F. Fitzgerald, "Macrofocal Conies as Reflector Contours," *J. Opt. Soc. Am.* **55**(1):5–11 (1965).
89. D. R. Philips, "Low Brightness Louver," U.S. Patent 2,971,083, 1961.
90. H. P. Baum and J. M. Gordon, "Geometric Characteristics of Ideal Nonimaging (CPC) Solar Collectors with Cylindrical Absorber," *Solar Energy* **33**(5):455–158 (1984).
91. P. T. Ong, J. M. Gordon, and A. Rabl, "Tailored Edge-Ray Designs for Illumination with Tubular Sources," *Appl. Opt.* **35**(22):4361–4371 (1996).
92. I. M. Bassett and G. H. Derrick, "The Collection of Diffuse Light onto an Extended Absorber," *Optical Quantum Electronics* **10**:61–82 (1978).
93. H. Ries, N. Shatz, J. Bortz, and W. Spirkl, "Performance Limitations of Rotationally Symmetric Nonimaging Devices," *J. Opt. Soc. Am. A* **14**(10):2855–2862 (1997).
94. M. Ruda (ed.), "International Conference on Nonimaging Concentrators," *Proc. SPIE* **441** (1983).
95. R. Winston and R. L. Holman (eds.), "Nonimaging Optics: Maximum Efficiency Light Transfer," *Proc. SPIE* **1528** (1991).
96. R. Winston and R. L. Holman (eds.), "Nonimaging Optics: Maximum Efficiency Light Transfer II," *Proc. SPIE* **2016** (1993).
97. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer III," *Proc. SPIE* **2538** (1995).
98. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer IV," *Proc. SPIE* **3139** (1997).
99. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer V," *Proc. SPIE* **3781** (1999).
100. R. Winston (ed.), "Selected Papers on Nonimaging Optics," *Proc. SPIE Milestone Series* **106** (1995).
101. I. M. Bassett, W. T. Welford, and R. Winston, "Nonimaging Optics for Flux Concentration," in *Progress in Optics*, E. Wolf (ed.), 1989, pp. 161–226.
102. R. Winston, "Nonimaging Optics," *Sci. Am.* 76–81 (March 1991).
103. P. Gleckman, J. O'Gallagher, and R. Winston, "Approaching the Irradiance of the Sun Through Nonimaging Optics," *Optics News*: 33–36 (May 1989).
104. M. F. Land, "The Optical Mechanism of the Eye of Limulus," *Nature* **280**:396–397 (1979).
105. M. F. Land, "Compound Eyes: Old and New Optical Mechanisms," *Nature* **287**:681–685 (1980).
106. R. Levi-Seti, D. A. Park, and R. Winston, "The Corneal Cones of Limulus as Optimized Light Concentrators," *Nature* **253**:115–116 (1975).
107. D. A. Baylor and R. Fettiplace, "Light Path and Photon Capture in Turtle Photoreceptors," *J. Physiol.* **248**(2):433–464 (1975).
108. R. Winston and J. Enoch, "Retinal Cone Receptor as Ideal Light Collector," *J. Opt. Soc. Am.* **61**(8):1120–1121 (1971).

109. R. L. Garwin, "The Design of Liquid Scintillation Cells," *Rev. Sci. Instr.* **23**:755–757 (1952).
110. D. E. Williamson, "Cone Channel Condensor," *J. Opt. Soc. Am.* **42**(10):712–715 (1952).
111. J. H. Myer, "Collimated Radiation in Conical Light Guides," *Appl. Opt.* **19**(18):3121–3123 (1980).
112. W. Witte, "Cone Channel Optics," *Infrared Phys.* **5**:179–185 (1965).
113. C. H. Burton, "Cone Channel Optics," *Infrared Phys.* **15**:157–159 (1975).
114. R. Winston and W. T. Welford, "Ideal Flux Concentrators as Shapes That Do Not Disturb the Geometrical Vector Flux Field: A New Derivation of the Compound Parabolic Concentrator," *J. Opt. Soc. Am.* **69**(4):536–539 (1979).
115. A. G. Molledo and A. Luque, "Analysis of Static and Quasi-Static Cross Compound Parabolic Concentrators," *Appl. Opt.* **23**:2007–2020 (1984).
116. R. Winston and H. Hinterberger, "Principles of Cylindrical Concentrators for Solar Energy," *Solar Energy* **17**:255–258 (1975).
117. W. R. McIntire, "Truncation of Nonimaging Cusp Concentrators," *Solar Energy* **23**:351–355 (1979).
118. H. Tabor, "Comment—The CPC Concept—Theory and Practice," *Solar Energy* **33**(6):629–630 (1984).
119. A. Rabl and R. Winston, "Ideal Concentrators for Finite Sources and Restricted Exit Angles," *Appl. Opt.* **15**:2880–2883 (1976).
120. A. Rabl, N. B. Goodman, and R. Winston, "Practical Design Considerations for CPC Solar Collectors," *Solar Energy* **22**:373–381 (1979).
121. W. R. McIntire, "New Reflector Design Which Avoids Losses through Gaps between Tubular Absorber and Reflectors," *Solar Energy* **25**:215–220 (1980).
122. R. Winston, "Ideal Flux Concentrators with Reflector Gaps," *Appl. Opt.* **17**(11):1668–1669 (1978).
123. R. Winston, "Cavity Enhancement by Controlled Directional Scattering," *Appl. Opt.* **19**:195–197 (1980).
124. F. Bloisi, P. Cavaliere, S. De Nicola, S. Martellucci, J. Quartieri, and L. Vicari, "Ideal Nonfocusing Concentrator with Fin Absorbers in Dielectric Rhombuses," *Opt. Lett.* **12**(7):453–155 (1987).
125. I. R. Edmonds, "Prism-Coupled Compound Parabola: A New Look and Optimal Solar Concentrator," *Opt. Lett.* **11**(8):490–492 (1986).
126. J. D. Kuppenheimer, "Design of Multilamp Nonimaging Laser Pump Cavities," *Opt. Eng.* **27**(12):1067–1071 (1988).
127. D. Lowe, T. Chin, T. L. Credelle, O. Tezucar, N. Hariston, J. Wilson, and K. Bingaman, "SpectraVue: A New System to Enhance Viewing Angle of LCDs," *SID 96 Applications Digest* 39–42 (1996).
128. J. L. Henkes, "Light Source for Liquid Crystal Display Panels Utilizing Internally Reflecting Light Pipes and Integrating Sphere," U.S. Patent 4,735,495, 1988.
129. R. Winston, "Principles of Solar Concentrators of a Novel Design," *Solar Energy* **16**:89–94 (1974).
130. R. Winston, "Cone-Collectors for Finite Sources," *Appl. Opt.* **17**(5):688–689 (1978).
131. M. Collares-Pereira, A. Rabl, and R. Winston, "Lens-Mirror Combinations with Maximal Concentration," *Appl. Opt.* **16**(10):2677–2683 (October 1977).
132. H. P. Gush, "Hyperbolic Cone-Channel Condensor," *Opt. Lett.* **2**:22–24 (1978).
133. J. O'Gallagher, R. Winston, and W. T. Welford, "Axially Symmetric Nonimaging Flux Concentrators with the Maximum Theoretical Concentration Ratio," *J. Opt. Soc. Am. A* **4**(1):66–68 (1987).
134. R. Winston, "Dielectric Compound Parabolic Concentrators," *Appl. Opt.* **15**(2):291–292 (1976).
135. J. R. Hull, "Dielectric Compound Parabolic Concentrating Solar Collector with a Frustrated Total Internal Reflection Absorber," *Appl. Opt.* **28**(1):157–162 (1989).
136. R. Winston, "Light Collection within the Framework of Geometrical Optics," *J. Opt. Soc. Am.* **60**(2):245–247 (1970).
137. D. Jenkins, R. Winston, R. Bliss, J. O'Gallagher, A. Lewandowski, and C. Bingham, "Solar Concentration of 50,000 Achieved with Output Power Approaching 1kW," *J. Sol. Eng.* **118**:141–144 (1996).
138. H. Ries, A. Segal, and J. Karni, "Extracting Concentrated Guided Light," *Appl. Opt.* **36**(13):2869–2874 (1997).
139. R. H. Hildebrand, "Focal Plane Optics in Far-Infrared and Submillimeter Astronomy," *Opt. Eng.* **25**(2):323–330 (1986).
140. J. Keene, R. H. Hildebrand, S. E. Whitcomb, and R. Winston, "Compact Infrared Heat Trap Field Optics," *Appl. Opt.* **17**(7):1107–1109 (1978).

141. R. Winston and W. T. Welford, "Geometrical Vector Flux and Some New Nonimaging Concentrators," *J. Opt. Soc. Am.* **69**(4):532–536 (1979).
142. X. Ning, R. Winston, and J. O'Gallagher, "Dielectric Totally Internally Reflecting Concentrators," *Appl. Opt.* **26**(2):300–305 (1987).
143. W. L. Eichhorn, "Designing Generalized Conic Concentrators for Conventional Optical Systems," *Appl. Opt.* **24**(8):1204–1205 (1985).
144. W. L. Eichhorn, "Generalized Conic Concentrators," *Appl. Opt.* **21**(21):3887–3890 (1982).
145. G. H. Smith, *Practical Computer-Aided Lens Design*, Willmann-Bell, VA, 1998, pp. 379–383.
146. K. W. Beeson, I. B. Steiner, and S. M. Zimmerman, "Illumination System Employing an Array of Microprisms," U.S. Patent 5,521,725, 1996.
147. R. P. Friedman, J. M. Gordon, and H. Ries, "New High-Flux Two-Stage Optical Designs for Parabolic Solar Concentrators," *Solar Energy* **51**:317–325 (1993).
148. P. Gleckman, J. O'Gallagher, and R. Winston, "Concentration of Sunlight to Solar-Surface Levels Using Non-imaging Optics," *Nature* **339**:198–200 (May 1989).
149. H. Ries and W. Sprickl, "Nonimaging Secondary Concentrators for Large Rim Angle Parabolic Trough with Tubular Absorbers," *Appl. Opt.* **35**:2242–2245 (1996).
150. J. O'Gallagher and R. Winston, "Test of a 'Trumpet' Secondary Concentrator with a Paraboloidal Dish Primary," *Solar Energy* **36**(1):37–14 (1986).
151. A. Rabl, "Comparison of Solar Concentrators," *Solar Energy* **18**:93–111 (1976).
152. R. Winston and W. T. Welford, "Design of Nonimaging Concentrators as Second Stages in Tandem with Image-Forming First-Stage Concentrators," *Appl. Opt.* **19**(3):347–351 (1980).
153. X. Ning, R. Winston, and J. O'Gallagher, "Optics of Two-Stage Photovoltaic Concentrators with Dielectric Second Stages," *Appl. Opt.* **26**(7):1207–1212 (1987).
154. E. M. Kritchman, "Second-Stage Concentrators—A New Formalism," *J. Opt. Soc. Am. Lett.* **73**(4):508–511 (1983).
155. D. R. Mills and J. E. Giutronich, "New Ideal Concentrators for Distant Radiation Sources," *Solar Energy* **23**:85–87 (1979).
156. D. R. Mills and J. E. Giutronich, "Asymmetrical Non-imaging Cylindrical Solar Concentrators," *Solar Energy* **20**:45–55 (1978).
157. H. Ries and J. M. Gordon, "Double-Tailored Imaging Concentrators," *Proc. SPIE* **3781**:129–134 (1999).
158. J. C. Minano and J. C. Gonzalez, "New Method of Design of Nonimaging Concentrators," *Appl. Opt.* **31**(16):3051–3060 (1992).
159. J. C. Minano, P. Benitez, and J. C. Gonzalez, "RX: A Nonimaging Concentrator," *Appl. Opt.* **34**(13):2226–2235 (1995).
160. P. Benitez and J. C. Minano, "Ultrahigh-Numerical-Aperture Imaging Concentrator," *J. Opt. Soc. Am. A* **14**(8):1988–1997 (1997).
161. J. C. Minano, J. C. Gonzalez, and P. Benitez, "A High-Gain, Compact, Nonimaging Concentrator: RXI," *Appl. Opt.* **34**(34):7850–7856 (1995).
162. J. L. Alvarez, M. Hernandez, P. Benitez, and J. C. Minano, "RXI Concentrator for 1000× Photovoltaic Conversion," *Proc. SPIE* **3781**:30–37 (1999).
163. J. F. Forkner, "Aberration Effects in Illumination Beams Focused by Lens Systems," *Proc. SPIE* **3428**:73–89 (1998).
164. X. Ning, "Three-Dimensional Ideal Θ_1/Θ_2 Angular Transformer and Its Uses in Fiber Optics," *Appl. Opt.* **27**(19):4126–4130 (1988).
165. A. Timinger, A. Kribus, P. Doron, and H. Ries, "Optimized CPC-type Concentrators Built of Plane Facets," *Proc. SPIE* **3781**:60–67 (1999).
166. J. P. Rice, Y. Zong, and D. J. Dummer, "Spatial Uniformity of Two Nonimaging Concentrators," *Opt. Eng.* **36**(11):2943–2947 (1997).
167. R. M. Emmons, B. A. Jacobson, R. D. Gengelbach, and R. Winston, "Nonimaging Optics in Direct View Applications," *Proc. SPIE* **2538**:42–50 (1995).
168. D. B. Leviton and J. W. Leitch, "Experimental and Raytrace Results for Throat-to-Throat Compound Parabolic Concentrators," *Appl. Opt.* **25**(16):2821–2825 (1986).

169. B. Moslehi, J. Ng, I. Kasimoff, and T. Jansson, "Fiber-Optic Coupling Based on Nonimaging Expanded-Beam Optics," *Opt. Lett.* **14**(23):1327–1329 (1989).
170. M. Collares-Pereira, J. F. Mendes, A. Rabl, and H. Ries, "Redirecting Concentrated Radiation," *Proc. SPIE* **2538**: 131–135 (1995).
171. A. Rabl, "Solar Concentrators with Maximal Concentration for Cylindrical Absorbers," *Appl. Opt.* **15**(7): 1871–1873 (1976). See also an erratum, *Appl. Opt.* **16**(1):15 (1977).
172. D. Feuermann and J. M. Gordon, "Optical Performance of Axisymmetric Concentrators and Illuminators," *Appl. Opt.* **37**(10):1905–1912 (1998).
173. J. Bortz, N. Shatz, and H. Ries, "Consequences of Etendue and Skewness Conservation for Nonimaging Devices with Inhomogeneous Targets," *Proc. SPIE* **3139**:28 (1997).
174. N. E. Shatz, J. C. Bortz, H. Ries, and R. Winston, "Nonrotationally Symmetric Nonimaging Systems that Overcome the Flux-Transfer Performance Limit Imposed by Skewness Conservation," *Proc. SPIE* **3139**:76–85 (1997).
175. N. E. Shatz, J. C. Bortz, and R. Winston, "Nonrotationally Symmetric Reflectors for Efficient and Uniform Illumination of Rectangular Apertures," *Proc. SPIE* **3428**:176–183 (1998).
176. W. Benesch and J. Strong, "The Optical Image Transformer," *J. Opt. Soc. Am.* **41**(4):252–254 (1951).
177. I. S. Bowen, "The Image-Slicer, a Device for Reducing Loss of Light at Slit of Stellar Spectrograph," *Astrophysical J.* **88**(2):113–124 (1938).
178. W. D. Westwood, "Multiple Lens System for an Optical Imaging Device," U.S. Patent 4,114,037, 1978.
179. J. Bernges, L. Unnebrink, T. Henning, E. W. Kreutz, and R. Poprawe, "Novel Design Concepts for UV Laser Beam Shaping," *Proc. SPIE* **3779**:118–125 (1999).
180. J. Endriz, "Brightness Conserving Optical System for Modifying Beam Symmetry," U.S. Patent 5,168,401, 1992.
181. T. Mori and H. Komatsuda, "Optical Integrator and Projection Exposure Apparatus Using the Same," U.S. Patent 5,594,526, 1997.
182. J. A. Shimizu and P. J. Janssen, "Integrating Lens Array and Image Forming Method for Improved Optical Efficiency," U.S. Patent 5,662,401, 1997.
183. D. Feuermann, J. M. Gordon, and H. Ries, "Nonimaging Optical Designs for Maximum-Power-Density Remote Irradiation," *Appl. Opt.* **37**(10):1835–1844 (1998).
184. F. C. Genovese, "Fiber Optic Line Illuminator with Deformed End Fibers and Method of Making Same," U.S. Patent 4,952,022, 1990.
185. P. Gorenstein and D. Luckey, "Light Pipe for Large Area Scintillator," *Rev. Sci. Instrum.* **34**(2):196–197 (1963).
186. W. Gibson, "Curled Light Pipes for Thin Organic Scintillators," *Rev. Sci. Instrum.* **35**(8):1021–1023 (1964).
187. H. Hinterberger and R. Winston, "Efficient Design of Lucite Light Pipes Coupled to Photomultipliers," *Rev. Sci. Instrum.* **39**(3):419–420 (1968).
188. H. D. Wolpert, "A Light Pipe for Flux Collection: An Efficient Device for Document Scanning," *Lasers Appl.* 73–74 (April 1983).
189. H. Karasawa, "Light Guide Apparatus Formed from Strip Light Guides," U.S. Patent 4,824,194, 1989.
190. D. A. Markle, "Optical Transformer Using Curved Strip Waveguides to Achieve a Nearly Unchanged F/Number," U.S. Patent 4,530,565, 1985.
191. I. M. Bassett and G. W. Forbes, "A New Class of Ideal Non-imaging Transformers," *Opt. Acta* **29**(9):1271–1282 (1982).
192. G. W. Forbes and I. M. Bassett, "An Axially Symmetric Variable-Angle Nonimaging Transformer," *Opt. Acta* **29**(9):1283–1297 (1982).
193. M. E. Barnett, "The Geometric Vector Flux Field within a Compound Elliptical Concentrator," *Optik* **54**(5):429–432 (1979).
194. M. E. Barnett, "Optical Flow in an Ideal Light Collector: The Θ_1/Θ_2 Concentrator," *Optik* **57**(3):391–400 (1980).
195. Gutierrez, J. C. Minano, C. Vega, and P. Benitez, "Application of Lorentz Geometry to Nonimaging Optics: New 3D Ideal Concentrators," *J. Opt. Soc. Am. A* **13**(3):532–540 (1996).
196. P. Greenman, "Geometrical Vector Flux Sinks and Ideal Flux Concentrators," *J. Opt. Soc. Am. Lett.* **71**(6):777–779 (1981).

197. W. T. Welford and R. Winston, "On the Problem of Ideal Flux Concentrators," *J. Opt. Soc. Am.* **68**(4):531–534 (1978). See also addendum, *J. Opt. Soc. Am.* **69**(2):367 (1979).
198. J. C. Minano, "Design of Three-Dimensional Nonimaging Concentrators with Inhomogeneous Media," *J. Opt. Soc. Am. A* **3**(9):1345–1353 (1986).
199. J. M. Gordon, "Complementary Construction of Ideal Nonimaging Concentrators and Its Applications," *Appl. Opt.* **35**(28):5677–5682 (1996).
200. P. A. Davies, "Edge-Ray Principle of Nonimaging Optics," *J. Opt. Soc. Am. A* **11**:1256–1259 (1994).
201. H. Ries and A. Rabl, "Edge-Ray Principle of Nonimaging Optics," *J. Opt. Soc. Am.* **10**(10):2627–2632 (1994).
202. A. Rabl, "Edge-Ray Method for Analysis of Radiation Transfer among Specular Reflectors," *Appl. Opt.* **33**(7):1248–1259 (1994).
203. J. C. Minano, "Two-Dimensional Nonimaging Concentrators with Inhomogeneous Media: A New Look," *J. Opt. Soc. Am. A* **2**(11):1826–1831 (1985).
204. M. T. Jones, "Motion Picture Screen Light as a Function of Carbon-Arc-Crater Brightness Distribution," *J. Soc. Mot. Pic. Eng.* **49**:218–240 (1947).
205. R. Kingslake, *Applied Optics and Optical Engineering*, vol. II, Academic Press, New York, 1965, pp. 225–226.
206. D. O'Shea, *Elements of Modern Optical Design*, Wiley, New York, 1985, pp. 111–114 and 384–390.
207. W. Wallin, "Design of Special Projector Illuminating Systems," *J-SMPTE* **71**:769–771 (1962).
208. H. Weiss, "Wide-Angle Slide Projection," *Inf. Disp.* 8–15 (September/October 1964).
209. S. Bradbury, *An Introduction to the Optical Microscope*, Oxford University Press, 1989, pp. 23–27.
210. R. Oldenbourg and M. Shribak, "Microscopes," in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 28, McGraw-Hill, New York, 2009.
211. M. Lacomat, G. M. Dubroeuq, J. Massin, and M. Brevignon, "Laser Projection Printing," *Solid State Technology* **23**(115): (1980).
212. J. M. Gordon and P. T. Ong, "Compact High-Efficiency Nonimaging Back Reflector for Filament Light Sources," *Opt. Eng.* **35**(6):1775–1778 (1996).
213. G. Zochling, "Design and Analysis of Illumination Systems," *Proc. SPIE* **1354**:617–626 (1990).
214. A. Steinfeld, "Apparent Absorptance for Diffusely and Specularly Reflecting Spherical Cavities," *Int. J. Heat Mass Trans.* **34**(7):1895–1897 (1991).
215. K. A. Snail and L. M. Hanssen, "Integrating Sphere Designs with Isotropic Throughput," *Appl. Opt.* **28**(10):1793–1799 (1989).
216. D. P. Ramer and J. C. Rains, "Lambertian Surface to Tailor the Distribution of Light in a Nonimaging Optical System," *Proc SPIE* **3781** (1999).
217. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**(1):125–128 (1967).
218. L. Hanssen and K. Snail, "Nonimaging Optics and the Measurement of Diffuse Reflectance," *Proc. SPIE* **1528**:142–150 (1991).
219. D. B. Chenault, K. A. Snail, and L. M. Hanssen, "Improved Integrating-Sphere Throughput with a Lens and Nonimaging Concentrator," *Appl. Opt.* **34**(34):7959–7964 (1995).
220. R. H. Webb, "Concentrator for Laser Light," *Appl. Opt.* **31**(28):5917–5918 (1992).
221. D. S. Goodman, "Producing Uniform Illumination," OSA Annual Meeting, Toronto, October 5, 1993.
222. M. M. Chen, J. B. Berkowitz-Mattuck, and P. E. Glaser, "The Use of a Kaleidoscope to Obtain Uniform Flux over a Large Area in a Solar or Arc Imaging Furnace," *Appl. Opt.* **2**:265–271 (1963).
223. L. A. Whitehead, R. A. Nodwell, and F. L. Curzon, "New Efficiency Light Guide for Interior Illumination," *Appl. Opt.* **21**(15):2755–2757 (1982).
224. S. G. Saxe, L. A. Whitehead, and S. Cobb, "Progress in the Development of Prism Light Guides," *Proc. SPIE* **692**:235–240 (1986).
225. K. Jain, "Illumination System to Produce Self-Luminous Light Beam of Selected Cross-Section, Uniform Intensity and Selected Numerical Aperture," U.S. Patent 5,059,013, 1991.
226. Y. Kudo and K. Matsumoto, "Illuminating Optical Device," U.S. Patent 4,918,583, 1990.
227. W. J. Cassarly, J. M. Davenport, and R. L. Hansler, "Uniform Lighting Systems: Uniform Light Delivery," *SAE 950904*, **SP-1081**:1–5 (1995).

228. M. R. Latta and K. Jain, "Beam Intensity Uniformization by Mirror Folding," *Opt. Comm.* **49**:27 (1984).
229. C. M. Chang, K. W. Lin, K. V. Chen, S. M. Chen, and H. D. Shieh, "A Uniform Rectangular Illuminating Optical System for Liquid Crystal Light Valve Projectors," *Euro Display '96*:257–260 (1996).
230. L. J. Coyne, "Distributive Fiber Optic Couplers Using Rectangular Lightguides as Mixing Elements," *Proc. FOC '79*:160–164 (1979).
231. M. Wagner, H. D. Geiler, and D. Wolff, "High-Performance Laser Beam Shaping and Homogenization System for Semiconductor Processing," *Meas. Sci. Technol. UK* **1**:1193–1201 (1990).
232. K. Iwasaki, Y. Ohyama, and Y. Nanaumi, "Flattening Laserbeam Intensity Distribution," *Lasers Appl.* **76** (April 1983).
233. K. Iwasaki, T. Hayashi, T. Goto, and S. Shimizu, "Square and Uniform Laser Output Device: Theory and Applications," *Appl. Opt.* **29**:1736–1744 (1990).
234. J. M. Geary, "Channel Integrator for Laser Beam Uniformity on Target," *Opt. Eng.* **27**:972–977 (1988).
235. R. E. Grojean, D. Feldman, and J. F. Roach, "Production of Flat Top Beam Profiles for High Energy Lasers," *Rev. Sci. Instrum.* **51**:375–376 (1980).
236. B. Fan, R. E. Tibbetts, J. S. Wilczynski, and D. F. Witman, "Laser Beam Homogenizer," U.S. Patent 4,744,615, 1988.
237. R. G. Waarts, D. W. Nam, D. E. Welch, D. R. Scifres, J. C. Ehlert, W. Cassarly, J. M. Finlan, and K. Flood, "Phased 2-D Semiconductor Laser Array for High Output Power," *Proc. SPIE* **1850**:270–280 (1993).
238. J. R. Jenness Jr., "Computing Uses of the Optical Tunnel," *Appl. Opt.* **29**(20):2989–2991 (1990).
239. L. J. Krolak and D. J. Parker, "The Optical Tunnel—A Versatile Electrooptic Tool," *J. Soc. Motion Pict. Telev. Eng.* **72**:177–180 (1963).
240. T. D. Milster and T. S. Tkaczyk, "Miniature and Micro-Optics," in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 22, McGraw-Hill, New York, 2009.
241. R. A. Sprague and D. R. Scifres, "Multi-Beam Optical System Using Lens Array," U.S. Patent 4,428,647, 1984.
242. K. Flood, B. Cassarly, C. Sigg, and J. Finlan, "Continuous Wide Angle Beam Steering Using Translation of Binary Microlens Arrays and a Liquid Crystal Phased Array," *Proc. SPIE* **1211**:296–304 (1990).
243. W. Kuster and H. J. Keller, "Domed Segmented Lens System," U.S. Patent 4,930,864, 1990.
244. W. J. Cassarly, J. M. Davenport, and R. L. Hansler, "Uniform Light Delivery Systems," SAE, Paper No. 960490 (1996).
245. N. A. Zhidkova, O. D. Kalinina, A. A. Kuchin, S. N. Natarovskii, O. N. Nemkova, and N. B. Skobeleva, "Use of Lens Arrays in Illuminators for Reflected-Light Microscopes," *Opt. Mekh. Promst.* **55**:23–24 (September 1988); *Sov. J. Opt. Technol.* **55**:539–541 (1988).
246. S. Ohuchi, T. Kakuda, M. Yatsu, N. Ozawa, M. Deguchi, and T. Maruyama, "Compact LC Projector with High-Brightness Optical System," *IEEE Trans. Consumer Electronics* **43**(3):801–806 (1997).
247. K. Rantsch, L. Bertele, H. Sauer, and A. Merz, "Illuminating System," U.S. Patent 2,186,123, 1940.
248. J. R. Miles, "Lenticulated Collimating Condensing System," U.S. Patent 3,296,923, 1967.
249. M. Ohtu, "Illuminating Apparatus," U.S. Patent 4,497,013, 1985.
250. K. Matsumoto, M. Uehara, and T. Kikuchi, "Illumination Optical System," U.S. Patent 4,769,750, 1988.
251. A. H. J. Van den Brandt and W. Timmers, "Optical Illumination System and Projection Apparatus Comprising Such a System," U.S. Patent 5,098,184, 1992.
252. F. Watanabe, "Illumination Optical System," U.S. Patent 5,786,939, 1998.
253. X. Deng, X. Liang, Z. Chen, W. Yu, and R. Ma, "Uniform Illumination of Large Targets Using a Lens Array," *Appl. Opt.* **25**:377–381 (1986).
254. N. Nishi, T. Jitsuno, K. Tsubakimoto, M. Murakami, M. Nakatsuma, K. Nishihara, and S. Nakai, "Aspherical Multi Lens Array for Uniform Target Illumination," *Proc. SPIE* **1870**:105–111 (1993).
255. Y. Ozaki and K. Takamoto, "Cylindrical Fly's Eye Lens for Intensity Redistribution of an Excimer Laser Beam," *Appl. Opt.* **28**:106–110 (1989).
256. S. Glocker and R. Goring, "Investigation of Statistical Variations between Lenslets of Microlens Arrays," *Appl. Opt.* **36**(19):4438–4445 (1997).
257. T. H. Bett, C. N. Danson, P. Jinks, D. A. Pepler, I. N. Ross, and R. M. Stevenson, "Binary Phase Zone-Plate Arrays for Laser-Beam Spatial-Intensity Distribution Conversion," *Appl. Opt.* **34**(20):4025–4036 (1995).

258. Y. Kato, K. Mima, N. Miyanaga, S. Arinaga, Y. Kitagawa, M. Nakatsuka, and C. Yamanaka, "Random Phasing of High-Power Lasers for Uniform Target Acceleration and Plasma-Instability Suppression," *Phys. Rev. Lett.* **53**(11):1057–1060 (1984).
259. L. B. W. Jolley, J. M. Waldram, and G. H. Wilson, *The Theory and Design of Illuminating Engineering Equipment*, John Wiley & Sons, New York, 1931.
260. D. G. Burkhard and D. L. Shealy, "Design of Reflectors Which Will Distribute Sunlight in a Specified Manner," *Solar Energy* **17**:221–227 (1975).
261. D. G. Burkhard and D. L. Shealy, "Specular Aspheric Surface to Obtain a Specified Irradiance from Discrete or Continuous Line Source Radiation: Design," *Appl. Opt.* **14**(6):1279–1284 (1975).
262. O. K. Kusch, *Computer-Aided Optical Design of Illuminating and Irradiating Devices*, ASLAN Publishing House, Moscow, 1993.
263. T. E. Horton and J. H. McDermit, "Design of Specular Aspheric Surface to Uniformly Radiate a Flat Surface Using a Nonuniform Collimated Radiation Source," *J. Heat Transfer*, 453–458 (1972).
264. J. S. Schruben, "Analysis of Rotationally Symmetric Reflectors for Illuminating Systems," *J. Opt. Soc. Am.* **64**(1):55–58 (1974).
265. J. Murdoch, *Illumination Engineering, from Edison's Lamp to the Laser*, MacMillan Publishing, 1985.
266. R. Winston, "Nonimaging Optics: Optical Design at the Thermodynamic Limit," *Proc. SPIE* **1528**:2–6 (1991).
267. D. Thackeray, "Reflectors for Light Sources," *J. Photog. Sci.* **22**:303–304 (1974).
268. P. W. Rhodes and D. L. Shealy, "Refractive Optical Systems for Irradiance Redistribution of Collimated Radiation: Their Design and Analysis," *Appl. Opt.* **19**:3545–3553 (1980).
269. W. A. Parkyn, "Design of Illumination Lenses via Extrinsic Differential Geometry," *Proc. SPIE* **3428**: 154–162 (1998).
270. R. Winston and H. Ries, "Nonimaging Reflectors as Functionals of the Desired Irradiance," *J. Opt. Soc. Am. A* **10**(9):1902–1908 (1993).
271. W. B. Elmer, "The Optics of Reflectors for Illumination," *IEEE Trans. Industry Appl.* **IA-19**(5):776–788 (September/October 1983).
272. J. M. Gordon, P. Kashin, and A. Rabl, "Nonimaging Reflectors for Efficient Uniform Illumination," *Appl. Opt.* **31**(28):6027–6035 (1992).
273. J. M. Gordon and A. Rabl, "Nonimaging Compound Parabolic Concentrator-Type Reflectors with Variable Extreme Direction," *Appl. Opt.* **31**(34):7332–7338 (1992).
274. J. Gordon and H. Ries, "Tailored Edge-Ray Concentrators as Ideal Second Stages for Fresnel Reflectors," *Appl. Opt.* **32**(13):2243–2251 (1993).
275. H. R. Ries and R. Winston, "Tailored Edge-Ray Reflectors for Illumination," *J. Opt. Soc. Am. A* **11**(4): 1260–1264 (1994).
276. A. Rabl and J. M. Gordon, "Reflector Design for Illumination with Extended Sources: The Basic Solutions," *Appl. Opt.* **33**(25):6012–6021 (1994).
277. P. T. Ong, J. M. Gordon, A. Rabl, and W. Cai, "Tailored Edge-Ray Designs for Uniform Illumination of Distant Targets," *Opt. Eng.* **34**(6):1726–1737 (1995).
278. P. T. Ong, J. M. Gordon, and A. Rabl, "Tailoring Lighting Reflectors to Prescribed Illuminance Distributions: Compact Partial-Involute Designs," *Appl. Opt.* **34**(34):7877–7887 (1995).
279. D. Jenkins and R. Winston, "Tailored Reflectors for Illumination," *Appl. Opt.* **35**(10):1669–1672 (1996).
280. D. Jenkins and R. Winston, "Integral Design Method for Nonimaging Concentrators," *J. Opt. Soc. Am. A* **13**(10):2106–2116 (1996).
281. J. M. Gordon, "Simple String Construction Method for Tailored Edge-Ray Concentrators in Maximum-Flux Solar Energy Collectors," *Solar Energy* **56**:279–284 (1996).
282. H. Schering and A. Merz, "Illuminating Device for Projectors," U.S. Patent 2,183,249, 1939.
283. Y. Kawamura, Y. Itagaki, K. Toyoda, and S. Namba, "A Simple Optical Device for Generating Square Flat-Top Intensity Irradiation from a Gaussian Laser Beam," *Opt. Comm.* **48**:44–46 (1983).
284. T. Sameshima and S. Usui, "Laser Beam Shaping System for Semiconductor Processing," *Opt. Comm.* **88**:59–62 (1992).
285. L. Austin, M. Scaggs, U. Sowada, and H.-J. Kahlert, "A UV Beam-Delivery System Designed for Excimers," *Photonics Spectra* 89–98 (May 1989).

286. A. Brunsting, "Redirecting Surface for Desired Intensity Profile," U.S. Patent 4327972, 1982.
287. R. J. Bruno and K. C. Liu, "Laserbeam Shaping for Maximum Uniformity and Minimum Loss," *Laser Appl.* 91–94 (April 1987).
288. D. D. Dourte, C. Mesa, R. L. Pierce, and W. J. Spawr, "Optical Integration with Screw Supports," U.S. Patent 4,195,913, 1980.
289. S. L. Ream, "A Convex Beam Integrator," *Laser Focus*: 68–71 (November 79).
290. V. J. Doherty, "Design of Mirrors with Segmented Conical Surfaces Tangent to a Discontinuous Aspheric Base," *Proc. SPIE* 399:263–271 (1983).
291. D. M. Dagenais, J. A. Woodroffe, and I. Itzkan, "Optical Beam Shaping of a High Power Laser for Uniform Target Illumination," *Appl. Opt.* 24:671–675 (1985).
292. F. M. Dickey and B. D. O'Neil, "Multifaceted Laser Beam Integrators: General Formulation and Design Concepts," *Opt. Eng.* 27:999–1007 (1988).
293. J. Geary, "Strip Mirror Integrator for Laser Beam Uniformity on a Target," *Opt. Eng.* 28(8):859–864 (August 1989).
294. B. Lu, J. Zheng, B. Cai, and B. Zhang, "Two-Dimensional Focusing of Laser Beams to Provide Uniform Irradiation," *Opt. Comm.* 149:19–26 (1998).
295. T. Henning, L. Unnebrink, and M. Scholl, "UV Laser Beam Shaping by Multifaceted Beam Integrators: Fundamentals and Principles and Advanced Design Concepts," *Proc. SPIE* 2703:62–73 (1996).
296. L. Unnebrink, T. Henning, E. W. Kreutz, and R. Poprawe, "Optical System Design for Excimer Laser Materials Processing," *Proc. SPIE* 3779:413–422 (1999).
297. S. R. David, C. T. Walker, and W. J. Cassarly, "Faceted Reflector Design for Uniform Illumination," *Proc. SPIE* 3482:437–446 (June 12, 1998).
298. J. R. Coaton and A. M. Marsden, *Lamps and Lighting*, 4th ed., John Wiley & Sons, New York, 1997.
299. R. J. Donohue and B. W. Joseph, "Computer Design of Automotive Lamps with Faceted Reflectors," *IES* 36–42 (October 1972).
300. E. H. Wiley, "Projector Lamp Reflector," U.S. Patent 4,021,659, 1976.
301. W. P. Laudenschlager, R. K. Jobe, and R. P. Jobe, "Light Assembly," U.S. Patent 4,153,929, 1979.
302. R. J. Heimer, "Multiple Apparent Source Optical Imaging System," U.S. Patent 4,241,389, 1980.
303. H. Hamada, K. Nakazawa, H. Take, N. Kimura, and F. Funada, "Lighting Apparatus," U.S. Patent 4,706,173, 1987.
304. M. Ronnelid, B. Perers, and B. Karlsson, "Optical Properties of Nonimaging Concentrators with Corrugated Reflectors," *Proc. SPIE* 2255:595–602 (1994).
305. D. J. Lamb, J. F. Van Derlofske, and L. W. Hillman, "The Use of Aspheric Surfaces in Waveguide Illumination Systems for Automotive Displays," SAE Technical Paper 980874, 1998.
306. W. J. Cassarly and T. L. Davenport, "Non-Rotationally Symmetric Mixing Rods," International Optical Design Conference, *Proc. SPIE* 6342: (June 2006).

This page intentionally left blank.

Anurag Gupta

*Optical Research Associates
Tucson, Arizona*

R. John Koshel

*Photon Engineering LLC and
College of Optical Sciences
University of Arizona
Tucson, Arizona*

40.1 GLOSSARY

Illuminance. Luminous flux incident on a surface per unit projected area in the direction of emission relative to the surface normal. $1 \text{ lux} = 1 \text{ lumen/m}^2$.

Intensity. Luminous flux emitted by a source per unit solid angle in a given direction. $1 \text{ candela} = 1 \text{ lumen/steradian}$.

Luminance. Luminous flux emitted in a given direction per unit solid angle per unit projected area in the direction of emission relative to the surface normal. $1 \text{ nit} = 1 \text{ lumen}/(\text{m}^2 \times \text{steradian})$.

Luminaire. Lamp or lighting fixture that includes optical components, baffles, housing, and electronics.

Display. It refers to many things: a computer monitor, a projected image, and a piece of art or a decorative item.

Backlighting. It refers to illumination of an object from behind. The object can be opaque, translucent, or transparent. The light source is usually large in size and diffuse.

40.2 INTRODUCTION

Lighting is an area of science that includes the interaction between light and people in their daily lives. The primary goals of lighting are to provide the illumination to perform tasks, direct people to desired locations and provide a sense of security. Additionally, lighting has a profound effect on mood and sleep-wake cycles of all living beings. As mankind has advanced technologically, the needs fulfilled by lighting have also increased to additionally provide relaxation, alter moods, attract people, provide entertainment, create virtual environments and improve human productivity.

Unlike most of other fields in optics, lighting is a more subjective field than objective—it is based on the emission aspects and how it is perceived within its surroundings. The subjective nature is based upon our interactions with lighting shaped by vision biology and brain perception. The capabilities of the human eye largely determine the detectable range and variations (in both time and space) in colors and brightness. Perception depends upon the brain's interpretation of the input from the eye and is influenced by past experiences. For example, although brightness and the color

of lit objects are mathematically represented by luminance metrics as a function of distance, direction and wavelength, the perception of these quantities is context dependent based upon the observer's experiences and environment. Thus, the same illumination levels in two different environments can be perceived as two drastically different lighting outcomes.

The importance of understanding lighting goes beyond the basic need to illuminate objects or surroundings. The importance of representing lighting accurately has been exemplified in art over the centuries.¹ Understanding lighting in the context of human perception and the ability to represent and simulate lighting on computers has attracted attention in fields such as virtual reality for training especially in the fields of medicine, aviation and computer-aided design, and entertainment such as video games and movies. As our understanding and computational capabilities have increased over time, so has the sophistication and quality of virtual environments and the entertainment media.

An effective lighting design is contextual, cultural and is well integrated into its surroundings. Light interacts intimately with everything it impinges upon. Every object, including humans, plants and architectural elements, has distinct scatter, reflective, transmissive, and absorptive properties that are dependent upon wavelength and direction of light. As a result, the lit objects contribute to the appearance of the scene. Lighting therefore must complement in concert with the architecture and its surroundings in form, composition, and style to meet our expectations for the lit environment. For example, our expectations for the lighting environment of a casino are quite different to that of a sports stadium, a retail store, an office or a hospital. In each case, there is a distinct purpose that directs the layout of the light sources with respect to the objects. Therefore, the lighting must change in each case as dramatically as the differences between the contexts of lighting. If it does not, then responses to the lit environment are often not positive, typically making, for example, a poor work environment for an office, a lackluster sales venue for retail, or an uncomfortable room in someone's home. This subjectivity is the most difficult aspect of lighting design. After all, the success of any lighting scheme depends upon being able to meet the expectations of its users. Simply said, we know a good lighting design as soon as we see it, but it has limited quantifiable metrics.

The history of modern lighting traces its origin back to the advent of artificial light sources: the incandescent light bulb, more than a century ago. Currently there is a wide variety of sources to choose from: incandescent (includes halogen), discharge (fluorescent, high-intensity discharge, sodium vapor), lasers, electroluminescent (includes LEDs and OLEDs), and daylight. The choice of light source depends on operating characteristics, cost, efficiency, and safety. The luminaire optics play a critical role in shaping the light distribution from the source. Thus the choice of the source and luminaire optics is an integral part of the lighting design process.

Any lighting design must comply with the relevant government regulations for the specific purpose of lighting. For example, safety is critical in transportation lighting, so there are stringent regulations on the relative intensity distribution from street lights and automobile headlamps to minimize glare and achieve the desired visibility. There are increasing numbers of government mandates on using efficient sources in many countries, such as the banning of inefficient incandescent sources in favor of fluorescent and LED light sources. In addition to regulations, there are guidelines for best lighting practices in various situations. These guidelines are published by national and international committees such as CIE (Commission Internationale de l'Éclairage or The International Commission on Illumination), IESNA (Illumination Engineering Society of North America), SAE International (Society of Automotive Engineers), CIBSE (Chartered Institution of Building Services Engineers), and many others. These groups publish guidelines on many aspects of lighting in the interiors and exteriors of homes, offices, educational institutions, hospitals, entertainment facilities, malls, industrial complexes, sports stadiums, theatres, museum, streets, parks, transportation, and even underwater. The goals of these guidelines are to help, at a minimum, to create designs that are functional, efficient and provide safe and comfortable lighting. It is in the hands of the lighting designer to add to this mandated objective illumination to achieve the desired subjective lighting—in other words, the aesthetics.

In this chapter, we provide an introduction to many facets of lighting by touching upon the perceptual and biological factors that guide lighting design (Sec. 40.3), design elements and methods

to create functional and aesthetically pleasing designs (Sec. 40.4), technology of sources and design of relevant optics (Sec. 40.5), and measurement of lighting conditions (Sec. 40.6). We end the chapter with application examples on interior and exterior lighting in Sec. 40.7. These sections provide comprehensive data on source selection and guidelines for best practices in general and in specific application areas such as lighting for offices, homes, healthcare, retail, and transportation. Although very important aspects of lighting engineering, we do not discuss electronic control mechanisms, maintenance, and commissioning of a lighting design due to limited space.

We make extensive use of terms used in Radiometry and Photometry. The reader is encouraged to refer to the chapters on various aspects of radiometry and photometry in Chaps. 34 to 39 in this volume.

40.3 VISION BIOLOGY AND PERCEPTION

The lighting design is guided by our understanding of perception and vision biology. In this section, we touch upon various aspects of biology of vision and perception.

Vision Biology

Biological aspects of human visual system response² that are parameterized and used in lighting design are visual acuity (resolution, vernier, recognition, and stereoscopic), color sensitivity, accommodation, field of view, and adaptability to color and brightness changes.

Once adapted, the human visual system response does not change appreciably with time. The human visual system responds to over eleven orders of magnitude of luminance. However, at any given instance, only 2 to 3 orders of magnitude are adapted to by the eye. It takes up to 60 minutes to fully adapt to lighting conditions. Light adaptation takes place via change in the eye pupil size that controls the amount of light entering the eye (2 mm to 8 mm), photochemical processes in the retinal cells and neural processes that respond to change in the luminance below 600 nits when cone cells have not fully bleached. Neural adaptation occurs quickly, within the first 200 ms, and allows us to adapt to 2 to 3 orders of magnitude of luminance fluctuation, such as in lit spaces of building interiors. Pupil adaptation via change in the pupil size takes up to few minutes and allows us to adapt up to 1 to 2 orders of magnitude of luminance fluctuation. For tasks that require good color discrimination, several minutes to an hour are needed for this color adaptation. The higher the luminance, the shorter is the adaptation time. Based upon luminance, three vision conditions exist:

1. Photopic vision occurs when bright illumination conditions in the visible (luminance >3 nits) exist. Both the rod and the cone cells in the retina are excited, but the rods are saturated and thus effectively ignored except for peripheral vision. Full color vision with highest resolution is possible. Most indoor lighting conditions ensure photopic enabling lighting conditions.
2. Scotopic vision occurs when low illumination level in the visible (luminance <0.001 nits) exists. Only the rod cells in the retina are excited. No color vision occurs and only low-resolution peripheral vision is possible. Lighting design is usually not done for the scotopic domain but it must take into account our peripheral vision capability wherever possible to help guide the individual toward relevant objects.
3. Mesopic vision is described by the transition between photopic and scotopic. The rod cells are excited and the cone cells are partially excited. The eye has limited color discrimination and resolution capabilities. Most outdoor artificial lighting conditions operate in this region.

The luminance and color distributions should be such that it includes individuals with impaired vision that cannot be corrected such as reduced illumination at the retina (50 percent for a 50-year old as compared to a 20-year old) due to smaller pupil size and reduced transmission efficiency of the eye, reduced contrast, increased glare sensitivity, loss of accommodation and reduced field of view caused by ageing, macular degeneration, cataract, or glaucoma. For most such cases

increased illuminance, slow transients from the light to dark regions and an environment with reduced glare help considerably.³

Perception

It is the perception of lit environment that eventually determines the acceptability and adequateness of lighting conditions. Lighting perception is ultimately determined by the human visual system response and brain processes that are individual dependent. It is the latter that brings considerable subjectivity to lighting perception. Lighting design criteria for specific environments are often guided by studies that determine average perceptual responses. These guidelines are continually evolving, in response to geopolitical factors, research on vision biology and technology.

Depending upon an object's optical and physical characteristics and how it is viewed in relation to its surroundings, lighting can alter its perceived visual attributes. Visual attributes of physical objects are defined in terms of brightness, lightness, hue, saturation, transparency, and glossiness. We explain each of these attributes below and also show that not every object has each of these attributes.

Brightness is the perceptual correlate of luminance. Although, in the absence of a background, the perception of brightness of an object is proportional to its luminance in a logarithmic manner (log of brightness is proportional to log of luminance), the perception of brightness is dependent upon adaptation to the surroundings and the relative hue and saturation of the object. Figure 1 shows an example of how the brightness perception is altered by surroundings.⁵⁰ Different portions of a uniform luminance bar appear to have different brightness depending upon the local background. Similarly, car headlights appear much brighter in the dark than in daylight. Both vision biology and perception play a role here. Due to low ambient lighting, the pupil dilates and admits much more light (up to 16 times in night than in the day). Therefore, any bright source leads to sudden pupil contraction and ensuing discomfort. The sources appear even brighter due to a dark background, an effect similar to Fig. 1. Color perception is similarly affected. Color-saturated objects tend to appear brighter and vice versa.

Lightness is the perceptual correlate of diffuse reflectance. Our perception of an object being dark or light depends upon our estimate of its reflectance. For example, a piece of white paper appears white irrespective of the illumination falling upon it as we know from prior experience that it has a high reflectance. We tend to perceive it whiter than other objects of lesser reflectance even when the flux reflected by white paper in relatively low illumination is lower than a grey object that is placed under higher illumination. This confusion would not occur if we view a small region of the object through an aperture without knowing anything about the object. The aperture masks the contextual information and forces us to make objective judgments.

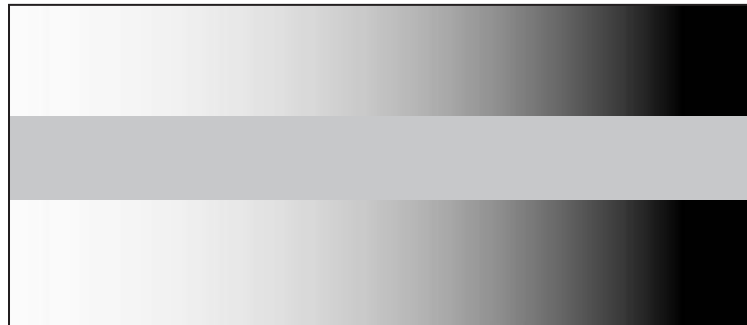


FIGURE 1 Perception of brightness. The rectangular bar in the center has constant luminance yet it seems brighter against a darker background at the right hand side.

Hue is the perception of dominant wavelength of the color spectrum transmitted or reflected by the object. It helps us judge if an object appears close to a known color such as red, blue, yellow, purple, green, colorless, or a combination of multiple colors such as bluish-green. Saturation is the perception of the extent of color purity. Highly saturated colors have a narrow band of wavelengths, centered on the hue. Transparency is the perception of the degree of light penetration into the object. Glossiness is the perception of smoothness of a surface relative to a matte finish.

Depending upon the properties of an object, the lighting conditions and how it is viewed, one or more visual attributes can be determined. For example, if an object is viewed through an aperture that hides the information of its surroundings or is illuminated such as nothing other than the object is visible then only attributes like brightness, hue, and saturation can be observed. A self-illuminated object or an object that gives the appearance of being so (such as a lit computer monitor) does not display lightness or glossiness, but if it stops emitting light, these attributes can be observed under external illumination.

Our brain attempts to maintain a perceptual constancy of shape, color, size, and lightness of lit objects based on past experiences and the context in which an object exists. Perceptual constancy allows us to maintain a certain perception of an object or a scene under changing viewing conditions. For example, we can recognize an apple as such when viewed at different angles, surroundings, or lit conditions. We can perceive a skyscraper as a tall building even if seen from far away and having a small image on the retina. Similarly, a tilted book appears rectangular although the image in the retina is trapezoidal. We can perceive the colors of common objects around us fairly correctly while wearing mildly colored glasses. Lighting can help maintain or alter perceptual constancies. To maintain perceptual constancies the lighting conditions must be such as there is adequate ambient lighting with high-color-rendering sources and without any disability glare. We discuss the terms such as color rendering, glare, and the impact of ambient lighting in the next section. The position of the light sources must be obvious to the observer even if not directly visible to establish the directionality of illumination. Direction of illumination is important as we tend to expect light to come from above and our perception of lit objects takes that into account at a subconscious level.

The relative distribution of light can impact the perception of the space itself and as such can impact human behavior.⁴ For example, there is a general tendency among humans to be attracted to brighter regions of space. It can be demonstrated by viewing people that there is a clear trend toward walking on brighter areas of pathways or facing brightly illuminated areas of a restaurant. That is why brightly lit shopping malls are quite effective in attracting traffic as compared to a poorly lit shopping complex. The knowledge of human behavior toward lit space is also used in making more effective lit object displays and navigational aids. There are four distinct categories of light distribution⁵ that affect the perception of space: privacy, relaxation, visual clarity, and spaciousness. An effect of privacy can be created by utilizing nonuniform high-brightness illumination across the vertical surfaces in the room with dark spaces in the occupant domains and low ambient luminance. An effect of relaxation can be created by using nonuniform warm (correlated color temperature (CCT) <3500 K) ambient light across the room. Visual clarity in the environment is emphasized with cool light (CCT >4000 K), high and uniform brightness near the center of the room and at all task planes. In addition, higher emphasis is given to ceiling and horizontal surfaces. A sense of spaciousness can be implemented with uniform room illumination and relatively higher levels of brightness on the walls and ceilings.

Summary

Understanding perception and biology of vision helps us develop models to describe desired lighting conditions. An example is the development of the relative visual performance (RVP) model. This model has been extensively developed by Mark Rea and his colleagues over the last few decades⁶ to obtain the relative visual performance of a given task under different lighting conditions and establish lighting guidelines. RVP is provided as a probability of performing a visual task successfully under given lighting conditions. Task performance depends upon both visual and nonvisual aspects

of the task. To obtain the true impact of lighting conditions on task performance, it is necessary to isolate those tasks for evaluation that are dominated by the visual component. The impact of the nonvisual components is minimized by quantifying their effect to the fullest possible extent and subtracting it from the overall task performance. A key finding of the RVP model is that the visual performance improves rapidly as the luminance contrast between the task and the background increases up to 40 percent. Beyond this value, the improvement in visual performance is negligible with increase in contrast. Luminance contrast in this context is defined as

$$\text{Luminance contrast} = 100(L_{\text{Task}} - L_{\text{background}}) / L_{\text{background}} \quad (1)$$

Visual performance curves (performance metric versus luminance contrast) can be evaluated for various task sizes and background luminance. A major limitation of the existing RVP model is its limited validity to only those tasks that are quantifiable by task size, luminance contrast, and background luminance and only under the conditions of foveal vision. We need more sophisticated task performance models to cover a larger range of tasks.

In the sections to follow, design guidelines make use of the understanding of lighting perception and vision biology to justify their development.

40.4 THE SCIENCE OF LIGHTING DESIGN

The lighting design process begins with identifying the needs to be addressed. An understanding of the functions of lighting and knowledge of basic building blocks helps us in making a preliminary design. The quality of lighting is determined by its ability to fulfill human needs in an economical and environmental friendly manner while at the same time complementing the architecture in form, composition, style, codes, and standards. We discuss the lighting design process in the following subsections:

Design considerations. We discuss the factors involved in creating a lighting environment for an application. The categories discussed are: goals, context, illuminance, color, visual discomfort, trespass and light induced damage of objects.

Functions of lighting. Here we discuss the four primary functions of lighting: ambient, task, decorative and accent.

Lighting geometries to achieve specific functions of lighting. Here we discuss the building blocks for lighting design.

Properties of objects and their impact on lit scene are discussed.

Modeling. Here we discuss the techniques used to simulate a lit environment in order to achieve the best design.

Design Considerations

Goals Lighting for any application must take into account the needs (both human and nonhuman such as plants or animals), lighting economics, environment impact, and architectural aspects of the application. Human needs include the desired degree of visibility, comfort, ability to perform the needed tasks, social communication, ambience, and aesthetics.

Context Lighting helps create a perceptual environment or ambience to suit a specific application such as office, home, lobby, restaurant, casino, or a sports stadium. Appropriate selection of lighting schemes and luminaires that complement the architecture and interior design helps in achieving the desired ambience.

TABLE 1 IESNA Guidelines on Illumination Categories and Average Illuminance Levels

Category	Average Illuminance (lx)
Public spaces	30
Simple orientation or short visits in a new environment	50
Working spaces where simple visual tasks are performed	100
Performance of visual tasks of high contrast and large size	300
Performance of visual tasks of high contrast and small size, or visual tasks of low contrast and large size	500
Performance of visual tasks of low contrast and small size	1,000
Performance of visual tasks near vision threshold	3,000–10,000

Large size: Object's projected solid angle subtense at the eye $>4.0 \times 10^{-6}$ sr.

Small size: Object's projected solid angle subtense at the eye $\leq 4.0 \times 10^{-6}$ sr but not near the visual acuity limit.

Low contrast: ≤ 0.3 but greater than visual threshold.

High contrast: >0.3 .

Illuminance—Horizontal and Vertical Horizontal and vertical illuminances refer to illuminance distribution on horizontal and vertical planes respectively. Table 1 describes IESNA guidelines on illumination categories and average illuminance levels needed for each.⁷ These guidelines do not apply to special situations that involve setting up a particular ambience or focusing on an object for emphasis.

The uniformity of luminance/illuminance is generally defined by the ratio of maximum to minimum luminance/illuminance. The need for uniformity across the field of view and across the entire space depends upon the application. The human eye is a brightness detector and is thus responsive to changes in luminance. To calculate luminance, illuminance, and the reflectivity of surfaces must be taken into account. For Lambertian surfaces (surfaces of constant luminance, independent of the viewing direction),

$$\text{Luminance} = (\text{illuminance} \times \text{surface reflectivity}) / \pi \quad (2)$$

Generally, a luminance uniformity of 0.7 within the field of view across the task is considered adequate as the eye is unable to detect these variations. Not all tasks require high luminance uniformity. For example, in tasks involving inspection of 3D objects, nonuniform illumination is able to highlight the geometrical features, especially surface textures much better due to being able to provide better depth perception. In lit environments, a nonuniform light distribution is used to provide perceptions of privacy or exclusivity, for example in retail lighting or in restaurants. A luminance ratio (maximum luminance: minimum luminance) of greater than 15:1 is generally considered undesirable. In the applications section, we discuss the recommended luminance ratios for several scenarios in a variety of applications.

Color Lighting influences the color appearance of lit objects. For most lighting applications, white light is the standard form of illumination: exteriors (landscape, roadways, buildings, city, and stadiums) and interiors (homes, offices, restaurants, museums, industrial complexes, and shopping malls). Saturated colors in lighting are used only in special applications like indicators or signals, displays, color-specific industrial applications such as those involving color discrimination, special effects in casinos, hotels, malls, or discotheques, and so forth. The chapter on colorimetry in this *Handbook* (Vol. III, Chap. 10) provides an excellent introduction to the subject of color.

Depending upon needs such as aesthetics, task performance, and color-dependent reflective properties of objects, an appropriate light spectrum must be selected. The desired light spectrum can be achieved by using light sources that emit in that spectrum, by using static filters on the sources to tailor the spectrum or by spatial and time-averaged color mixing. Spatial color mixing involves using multiple light sources that emit in different portions of the desired spectrum but are laid out in such a manner that the lit environment or object appear to be illuminated by light

having the combined spectrum of individual light sources. Time-averaged color mixing involves high-frequency (>60 Hz) mixing of different portions of the light spectrum in different proportions. If the frequency of color mixing is high enough, the brain perceives a specific color based on the time average of the varying spectra used in their respective strengths. For example, in modern digital projectors, a color wheel is used that has various segments of different spectral transmission. When it rotates through those segments, one can create the appearance of any color within the color gamut of the light source. In the section on LEDs, we discuss color mixing to create white light or any desired color. Spatial and temporal color mixing can be used together to create a variety of effects. Although color mixing can achieve the visual perception of any desired color in emission, its ability to color render the object it illuminates depends upon the product of incident light spectrum and wavelength-dependent reflectance of the object. In this section we discuss how white light is specified and how its ability to render objects is estimated.

It is tedious to choose light sources if we have to analyze the spectrum and its impact on common objects. For white light, the color rendering index (CRI) and the correlated color temperature (CCT) help provide a quick estimate of the appearance of light and its color rendering of lit objects. For example, at a CCT of 2700 to 6500 K, a CRI ≥ 70 is adequate for common situations such as in offices and homes, a CRI ≥ 50 is sufficient for most industrial tasks and a CRI ≥ 90 is needed for stringent color discrimination tasks such as hospital surgery, paint mixing, or color matching. In the future, it is likely that different metrics based on color-appearance models would be in use. This is especially true with the increasing use and availability of a variety of light sources, especially LEDs, which have significantly different emission spectra from incandescent sources.

The CCT is the temperature of the planckian (perfect blackbody) radiator in Kelvin whose perceived color most closely resembles that of a given stimulus at the same brightness and under specified viewing conditions.⁸ To find the CCT, the nearest point on the planckian locus is considered but only in a perceptually uniform color space. The isotherms across the planckian locus in a uniform color space are represented as normal to locus curve but in a nonuniform color space such as XYZ, these are no longer perpendicular to the locus.

Color rendering is defined as the effect of an illuminant on the color appearance of objects by conscious or subconscious comparison with their color appearance under a reference illuminant.⁸ Two functional reference illuminants are currently used for calculating CRI: (1) for a source CCT up to 5000 K, a blackbody at the same color temperature is used and (2) for a source CCT above 5000 K, one of the phases of daylight is used. The phase of daylight selected is such that its chromaticity is within $1/(15E6 \text{ K})$ to that of the test source. It is defined by a mathematical formula based on the CCT of the test source.⁹ For all cases, a CRI value of 100 is considered to be a perfect match between the test source and the illuminant. The precise steps and calculations needed to calculate the CRI are recommended by the CIE.¹⁰ Here we summarize the results:

$$\text{CRI}(R_a) = 100 - 4.6 \overline{E}_{UVW} \quad (3)$$

where R_a refers to the general CRI and \overline{E}_{UVW} is the average of the Euclidean distances between the color coordinates of the reflected reference and test light sources from the first 8 out of the 14 CIE-prescribed test samples (see Ref. 10). The color coordinates of the test source are chromatically adapted to the reference source and are expressed in the CIE 1964 color space. The CRI calculation is valid only when the color difference between the test source and the reference source is not large.

Although CRI is a useful metric and is widely used, it has several shortcomings. The existing method of calculating CRI is not perceptually well correlated. The high CRI predicted for sources with extreme CCT do not have good color-rendering properties. CRI is not valid for sources such as discrete spectral LEDs where CCT cannot be defined. CRI cannot be used to evaluate white-light LEDs with nonuniform spectra; the CRI predicted is quite low although the quality of white light appears to be better. The practice of using only eight test samples, none of which are saturated creates situations where high CRI sources do not color render color-saturated objects correctly. The CRI formulation can be modified to include the impact of different reference illuminants for different source types, color spaces, test sample set, different chromatic adaptation formulas or even a

reduced focus on absolute color fidelity. CIE reviews various propositions in this regard and updates its recommendations.

Visual Discomfort Visual discomfort can be caused by a variety of reasons.¹¹ Many of these reasons are context dependent where the expectation of the nature of lit environment determines the suitability of lighting. It also depends on cultural differences between various groups of people. For example, lighting flicker in a dance club may be desirable as opposed to almost all other situations. Similarly, the preference of color among different cultural groups varies considerably. Visual discomfort occurs when lighting creates perceptual confusion. For example, if the pattern of illumination is such that surfaces of higher reflectance reflect less light than the surfaces of lower reflectance, perceptual confusion may result. The causes of visual discomfort that are more specific to lighting are summarized as follows:

Insufficient lighting Insufficient lighting, especially for task performance, results in eye strain besides reduced task performance. For different tasks, there are different levels of horizontal and vertical illuminance necessary to be considered adequate. Table 1 describes suggested horizontal illuminance levels needed for certain situations. Lighting communities across the world have established guidelines for illuminance levels for a wide variety of specific tasks and environments.

Uniformity Visual discomfort occurs when the uniformity across the field of view is not as expected. A high uniformity can be as undesirable as a high degree of nonuniformity especially when considered across the entire visual field of view. Both can cause severe eye strain. In the section on applications, we provide examples of preferred luminance ratios between task and its vicinity.

Glare Glare results when there are regions of unexpected very high levels of luminance in the field of view. Glare can be direct or indirect. Direct glare occurs when a light source or a portion of it is visible such as in overhead lamps, automobile headlights, or direct sunlight. Indirect glare occurs when the light is reflected or scattered directly into the eye, such as the sky reflecting off a lake surface and obscuring the view beneath the water surface. Glare comes in many forms:¹²

1. Flash blindness: This is caused by a sudden onset of bright light leading to temporary bleaching of retinal pigment.
2. Paralyzing glare: This is caused by sudden illumination with bright light that can temporarily “freeze” the movements of the observer.
3. Distracting glare: This is caused by flashing bright sources of light in the peripheral vision field.
4. Retinal damage: When the light is bright enough to cause retinal damage.
5. Saturation or dazzle: When a large portion of the vision field is dominated by bright source(s), which can be alleviated by wearing low transmittance eye glasses.
6. Adaptation: When one enters from a low ambient illumination region to a bright ambient illumination region without a transition region to help in vision adaptation.
7. Disability: This is caused by intraocular light scattering which reduces the luminance contrast [See Eq. (1)] of the task image at the retina. The impact is loss in task performance due to reduced visibility.
8. Discomfort: When the glare causes discomfort or distraction but does not affect the task visibility to the extent of limiting its performance

Note that several forms of glare can exist concurrently, especially discomfort with the other types. Discomfort and disability glare are the most commonly experienced glare forms. There are several mechanisms available to estimate the impact of these forms of glare. We discuss briefly the CIE-recommended models for disability and discomfort glare.

Disability glare Disability glare occurs when the luminance contrast [Eq. (1)] of the task image at the retina falls due to the superposition of intraocular scattered light on the retinal image.

This background noise from scattered light can be thought of as viewing through a veil. Therefore, it makes sense to define an equivalent glare (or veiling) luminance (EVL) that mimics the impact of disability glare. The luminance contrast C is described as

$$C = \left| \frac{(L_b + L_{\text{EVL}}) - (L_{\text{object}} + L_{\text{EVL}})}{(L_b + L_{\text{EVL}})} \right| = \left| \frac{L_b - L_{\text{object}}}{L_b + L_{\text{EVL}}} \right| \quad (4)$$

where L is object luminance and b is background.

As the equivalent veiling luminance L_{EVL} increases, contrast at the retina C reduces. The contrast reduction is especially severe during difficult viewing conditions such as fog or nighttime when the object luminance is low. That is why car headlights are a much stronger glare source in the night than in the day.

Equivalent veiling luminance L_{EVL} is defined as:¹³

$$L_{\text{EVL}} = \sum_i \left[\frac{10}{\theta_i^3} + \left\{ \frac{5}{\theta_i^2} + \frac{0.1p}{\theta_i} \right\} \cdot \left\{ 1 + \left(\frac{A}{62.5} \right)^4 \right\} + 0.0025p \right] \cdot E_i \quad (5)$$

where E_i = illuminance at the eye due to i th glare source

θ_i = angle of the glare source (in degrees) from the line of sight, $0.1^\circ < \theta_i < 100^\circ$

p = eye pigmentation factor (0 for black eyes, 0.5 for brown eyes, 1.0 for light eyes, and 1.2 for very light-blue eyes)

A = age of the viewer in years

For young adults (<35 years of age) and for glare source angle, $1^\circ < \theta_i < 30^\circ$, Eq. (5) approximates to

$$L_{\text{EVL}} = \sum_i 10(E_i/\theta_i^2) \quad (6)$$

The EVL as described above is strictly due to intraocular scatter. In practice, it is necessary to add to this the luminance from external scatterers, such as fog or dust in the atmosphere, to yield the correct retinal contrast.

Equation (5) is currently the most sophisticated recommended treatment for disability glare. It can be applied toward a wide variety of circumstances, including indoor lighting, street lighting, and bright sky at a tunnel's exit.

For road lighting, the CIE recommendation on disability glare¹⁴ is given by a percentage threshold increment (TI) described by Eq. (7). TI is limited between 10 and 15 percent.

$$TI = 65(L_{\text{EVL}}/L^{0.8}) \quad (7)$$

where L_{EVL} = equivalent veiling luminance as described by Eq. (5)
 L = average road surface luminance

Discomfort glare There are many formulations that describe discomfort glare. Each model is prescribed for well defined geometries and sources. Discomfort glare has been described by the visual comfort probability (VCP) model¹⁵ in North America, British Glare Index system¹⁶ (CIBSE) and the European glare limiting system.¹⁷⁻¹⁹ Each of these systems has validity under specific constraints. Many luminaire manufacturers in North America and Europe provide VCP or glare index tables for worst case scenarios. CIE has proposed a Unified Glare Rating (UGR) model²⁰ to replace these systems. We describe here the formulation recommended by CIE.

The UGR formula is described by Eq. (8). UGR values range from 5 to 30, the higher values signifying a greater level of discomfort. For home and offices, UGR is specified at <20 and for industrial application it is >20. This formula is valid for source areas between 0.005 and 1.5 m². For smaller

TABLE 2 Glare Specification for Large Sources Such as an Illuminated Ceiling

Maximum Average Illuminance (lx)	UGR
300	13
600	16
1000	19
1600	22

sources, UGR overestimates the glare and for larger sources, UGR underestimates the glare. Therefore, for ranges outside the validity of UGR, CIE has provided detailed prescriptions to tackle small, large, and complex luminaires.²¹ For source areas smaller than 0.005 m², Eq. (9a) is recommended. In practice, any bare incandescent lamp, frosted or clear qualifies as small. For source areas larger than 1.5 m², but not as large as an illuminated ceiling or uniform indirect lighting (see Section, "Lighting Geometries," for a definition of indirect lighting), UGR is modified into a large room glare rating (GGR) and is described by Eq. (9b). The same UGR and GGR values represent an identical level of discomfort. For very large sources, only the maximum average illuminance values correspond to a specific UGR rating, as shown in Table 2. For nonuniform indirect source, CIE has provided guidelines.²¹ Each of the glare formulations discussed in this section are independent of the light spectrum.

Equations (9a) and (9b) are expressed for a single small and a single large source, respectively. Equation (9) is valid for viewing angles greater than 5° from the line of sight. For a combination of sources of different sizes, Eq. (8) must be modified to include glare from sources specified by equations (9a) and (9b).

$$\text{UGR} = 8 \log_{10} \left(\frac{0.25}{L_b} \right) \sum_i \frac{L_i^2 \omega_i}{P_i^2} \quad (8)$$

$$\text{UGR}_{\text{SingleSmallSource}} = 8 \log_{10} \left(\frac{0.25}{L_b} \right) \frac{200(I^2/R^2)}{P_i^2} \quad (9a)$$

$$\text{GGR}_{\text{SingleLargeSource}} = \text{UGR} + \{1.18 - (0.18/\text{CC})\} 8 \log [2.55\{1 + (E_d/220)\} / \{1 + (E_d/E_i)\}] \quad (9b)$$

where L_b = average luminance of the field of view without the luminaire or glare source

L_i = luminance of the i th luminaire in the observer's direction

ω_i = solid angle of the i th luminaire subtended at the observer's eye

P_i = Guth Position index²² of the i th luminaire. [It is a function of angular deviations (vertical and horizontal) from the line of sight and valid up to 53° of deviation from the line of sight. See Eq. (10).]

I = luminous intensity of the small source expressed in lumens per steradians. The source must be >5° away from the line of sight

E_d = direct illuminance at the eye due to the source

E_i = indirect illuminance at the eye = πL_b

CC = ceiling coverage = (area projected by the source at nadir)/(area lit by the source)

$$P_i = \exp\{(35.2 - 0.31889\alpha - 1.22e^{-2\alpha/9})10^{-3}\beta + (21 + 0.26667\alpha - 0.002963\alpha^2)10^{-5}\beta^2\} \quad (10)$$

where α = angle of the plane containing the observer's line of sight and a line from the observer to the source from the vertical direction. [Vertical direction is the height above (orthogonal to) the floor on which the observer is positioned.]

β = angle between the observer's line of sight and the line from the observer to the source

For roadway lighting, the glare rating is affected by driver fatigue, vehicular speed, and whether the person in the car is driving or not. All these effects must be dealt with comprehensively to formulate a single glare rating model that is largely driver independent. There is ongoing research in this field to develop better models for evaluating glare for road and vehicular lighting.^{23,24} Current CIE recommendation for limiting the discomfort glare for road lighting is identical to disability glare as described by Eq. (7).

Veiling reflections Veiling reflections are reflections from the task surface that result in the luminance contrast [Eq. (1)] reduction of the task itself. Veiling reflections can be evaluated from the source-task-eye geometry. It has been found that 90 percent of test subjects find a luminance contrast reduction of 25 percent as acceptable.²⁵ Veiling reflections are sometimes used as highlights to reveal the specular nature of a display (a lit object).

Miscellaneous Design Issues Other issues include designing against unacceptable light pollution or trespass, light-induced degradation or damage of objects and flicker from light sources. UV and IR filters are used with luminaires when there is a potential of damage to artwork or object displays. Flicker is more noticeable with high levels of the percentage modulation, area of visual field that is impacted by it, or the adaptation luminance. The impact of flicker can be reduced by using high frequency electronic ballasts or multiphase power supplies for different sources.

Functions of Lighting

There are four functions of lighting: ambient, task, decorative, and accent.²⁶ For each function, there are several implementation geometries. For most applications, more than one of these functions is necessary. We first discuss each of these lighting functions and then the lighting geometries used to accomplish these functions.

Ambient lighting fills up the space and is integral to almost any lighting scheme and yet is commonly ignored. It reduces the difference between the magnitudes of vertical and horizontal illuminance. As a result, it reduces glare, softens shadows, and provides a well-lit appearance. A common mistake is to consider any light that is illuminating the space as ambient light. For example, ceilings with recessed down lights with a narrow angular spread cause harsh shadows of objects on the ground. Even facial features show unflattering shadows. This is a result of insufficient level of vertical illuminance. So although there seems to be enough light to illuminate the space, it does not achieve a proper balance between vertical and horizontal illuminances resulting in cast shadows. To achieve proper ambient illumination, it is necessary to use those lighting geometries that spread light into large angles and from many directions. Ambient lighting is provided by large overhead luminaires with diffusers, torchieres, wall sconces, cove lighting, cornice lighting, valence and wall slot lighting, illuminated ceilings and wall washings. Figure 2 illustrates some of these schemes. Figure 4 illustrates ambient lighting with wall sconces.

Task lighting is used to provide sufficient illumination at the task plane such as a desk or work plane. Task lighting should be free of glare and shadows caused by the illuminated objects such as shadows from the hand and body or shadows from machine parts. Several lamp types and lighting geometries are available to reduce the impact of shadows. Lamps such as Banker or Bouillotte (Fig. 24) or large overhead luminaires with diffusers are good choices. The light emanating from these lamps is spread over a wide range of angles from an effectively large source. In a Banker lamp, a significant portion of light from the source undergoes multiple reflections from the inside of the luminaire before exiting. Large fluorescent light sources in a vertical configuration in the Bouillotte lamp provide excellent vertical and horizontal illuminance. Lamps with batwing lenses achieve cross illumination or lighting from two different directions overlapping in the task region. A batwing lens has a linear prism array (Fig. 23c) at the front of the source that leads to spreading of the light in predominantly two directions from each source point. Side lights installed in the vicinity of the task region increase vertical illuminance, and when used with overhead lighting provide excellent task lighting. Sometimes backlighting is necessary for certain tasks that involve transparent objects. Task lights must have

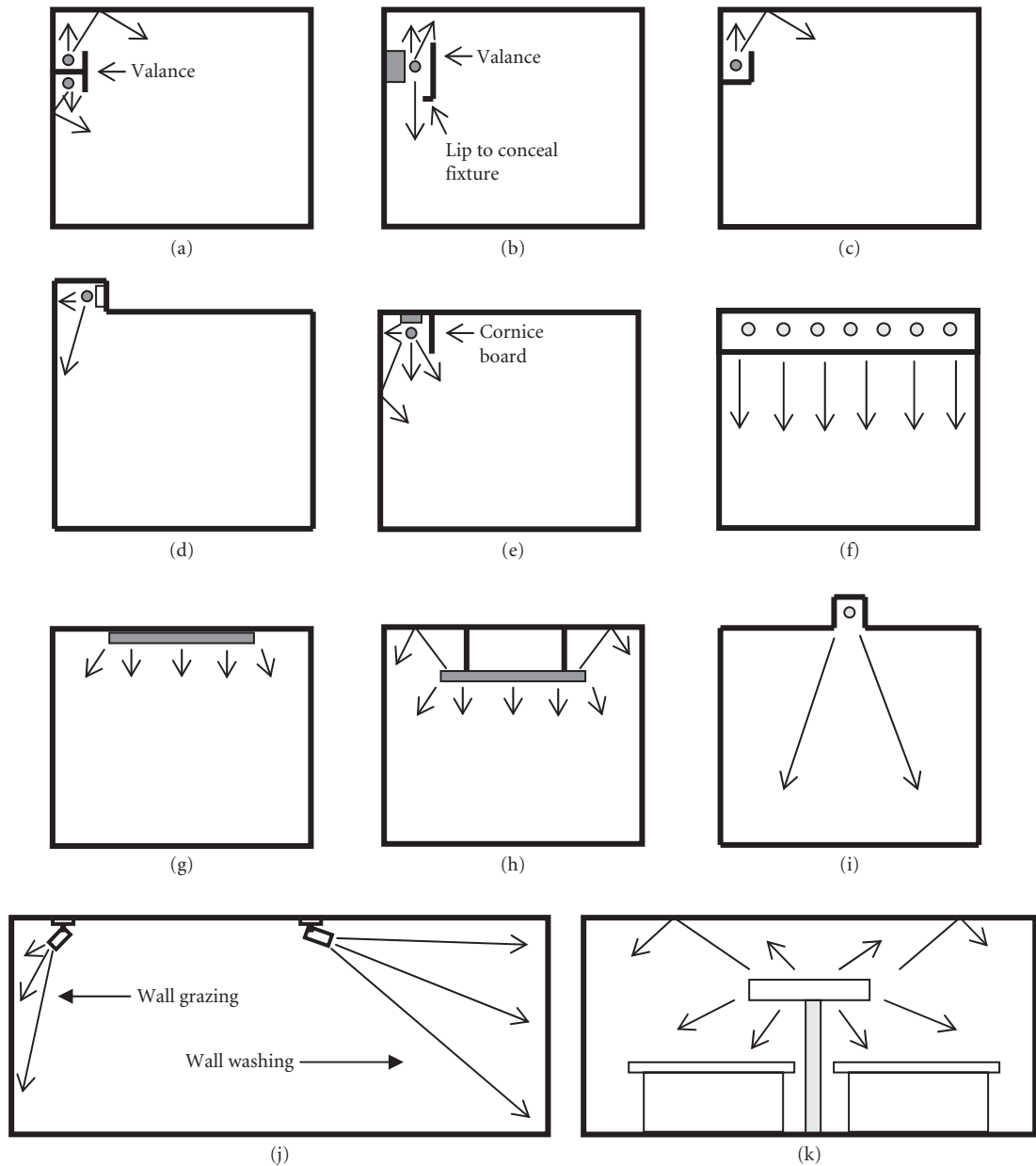


FIGURE 2 Various lighting geometry components. (a) and (b) Valance lighting. (c) Cove lighting. (d) and (e) Wall slot and Cornice lighting for wall illumination. It can be used to create an effect of a floating ceiling. (f) Illuminated ceiling. There is a diffuser below the row of line sources. This scheme can also be used to create illuminated wall panels. Backlighting is another technique. (g) Overhead luminaire. (h) Suspended luminaire. (i) Recessed downlight. (j) Wall-grazing and wall-washing illumination emphasizes and flattens the wall textures respectively. (k) Furniture-integrated lighting system. In cases (a) to (i), the sources are line sources such as fluorescent tubes along the length of the wall. In each case, the neighboring surfaces whether wall or shielding surfaces are coated with high reflectance diffuse paints. The distance of light source from the ceiling/wall determines the uniformity of illumination across the respective regions.



FIGURE 3 Accent lighting. (See also color insert.) (Courtesy of Pegasus Associates Lighting, www.pegasusassociates.com.)

the light sources well shielded by the luminaire and baffles to avoid direct glare. For large overhead sources, louvers (Fig. 24e) are used to limit the direct source view. Veiling reflections are avoided by ensuring that the source-task-eye geometry does not direct the reflections off the task into the eye.

Decorative lighting, as the term implies, is used to add sparkle to the lit environment. Decorative lighting is most effective when it appears to provide most of the lighting in a scene. Using decorative lighting to provide other functions of lighting by increasing the lamp brightness is not a suitable solution, although it increases the illumination in the region. It draws too much attention toward the lamp itself and creates undesirable levels of source brightness against the background leading to glare and unsightly shadows; therefore, decorative lighting must be immersed in an environment with good ambient light. Decorative lighting is provided by low-wattage table or floor lamps, gas lights, chandeliers, sconces, bare light sources, backlighting, light art, and torchieres. See the use of chandelier in Fig. 5.

Accent lighting primarily provides highlighting of objects such as artwork (Fig. 3), plants and decorative objects within a lit environment. Track lights, adjustable recessed lights, uplights, and backlights are commonly used to provide accent lighting.

Lighting Geometries

Lighting geometries can essentially be broken down into four broad classifications. These classifications are not mutually exclusive.

1. **Direct Lighting:** Most of the light (>90 percent) from a luminaire is targeted toward a certain region such as in downlighting in an office by overhead light fixtures. Applications of direct lighting include all the lighting functions such as ambient, task, decorative, and accent.
2. **Indirect Lighting:** An object or a region is illuminated by light that is not directly coming from the source. For example, the light from a lamp is directed toward the ceiling, wall, or even a



FIGURE 4 Wall sconces for providing ambient lighting and the much needed vertical illumination in various situations. (See also color insert.) (Courtesy of Lightcrafters Inc., www.lightcrafters.com.)

diffuse reflecting region of the luminaire. The reflections illuminate the space around the luminaire. Indirect lighting tends to give a more spacious appearance and eliminates shadows. Its primary application is to provide ambient lighting. See Figs. 4, 5, and 24c.

3. **Diffuse Lighting:** When lighting does not appear to come from any specific direction. Examples of diffuse lighting include indirect lighting and certain direct lighting geometries such as overcast skies, large area lighting fixtures with diffusing or prismatic optics or large spatial extent fluorescent lamps. It is mostly used for ambient illumination but can also be used to create local areas of high and uniform brightness for task lighting or accent lighting.
4. **Direct-Indirect/Semi-Direct/Semi-Indirect Lighting:** These terms typically apply to lamps that distribute some portion of their light toward the target and the remaining portion toward a surface that reflects (specular and/or diffuse) light toward the target. Semi-direct typically refers to the case where 60 to 90 percent of the light is directed toward the target while semi-indirect refers to the case where 60 to 90 percent of the light is directed away from the target. Applications include all the lighting functions: ambient, task, accent, and decorative. See Fig. 24d for an example of direct-indirect lighting fixture.

Figure 2 illustrates several implementations of these lighting geometries⁵ within the confines of indoor lighting. A practical lighting layout is likely to include one or more of these concepts.



FIGURE 5 Indirect lighting with cove lighting in a restaurant using light strips. The chandelier provides the decorative lighting without significantly contributing to any other lighting function. (See also color insert.) (Courtesy of Pegasus Associates Lighting, www.pegasusassociates.com.)

Properties of Objects and Their Impact on Lit-Scene

Objects in a lit-scene have geometrical (shape, location, and orientation) and optical properties (wavelength dependent absorption, reflection, transmission, and diffraction). Reflection and transmission include both diffuse and specular components. Lit-scene characteristics are then in-part determined by geometrical and optical properties of the objects and also in part determined by the interaction of these characteristics with the light sources in the environment. The object-light source interaction gives rise to such phenomena as glare, nonuniform illumination and color change. Major indoor features that impact a lighting environment are walls, floors (including large area carpeting) room partitions and ceilings. Minor features consist of temporary objects such as wall coverings, furniture, partitions, art, displays and plants. Major outdoor features that impact lighting environment are ground, buildings, vegetation, distant and close landscape features. Minor outdoor features are temporary objects.

System Layout and Simulation

Based upon the lighting design criteria, appropriate lamps with desired lighting schemes are selected. However, to obtain the desired illuminance distribution over time, appropriate calculations and simulations are needed. In this section, we discuss the tools for system simulation.

For any given system layout, a certain light level is needed. Equation (11) provides an approximation of the number and type of luminaires needed to obtain a certain horizontal illuminance over a work plane. This equation is useful when expressed as a summation of individual luminaires with their specific constants.

$$E_{\text{maintained}} = FnLLFCU/A \quad (11)$$

where $E_{\text{maintained}}$ = average illuminance maintained
 F = total rated luminaire lumen
 n = number of luminaires

CU = coefficient of utilization (It defines the percentage of light from the lamp reaching the work plane. CU depends upon the relative placement of lamp and work plane and illuminance distribution at the work plane corresponding to the geometry.)

A = work plane area

LLF = total light loss factor

LLF²⁷ accounts for the lamp output reduction over time. LLF has both recoverable and nonrecoverable components. Nonrecoverable factors are due to permanent degradation of the luminaire surface, reduction of output due to deviation from ideal operating temperature and environment (convection and ambient temperature), inefficiency of the electronic drive components and deviation in the operating position (tilt) from the ideal position. Recoverable factors are those whose effects can be mitigated by regular cleaning of the luminaire, operation in a clean environment, and regular replacement of bulbs or sources upon their natural degradation with time. LLF is a product of all these factors. Many of these factors also affect the intensity profile of the lamp output and can therefore affect the CU. The *Lighting Handbook* by IESNA²⁸ lists a detailed procedure for calculating the room surface dirt depreciation factor and luminaire dirt depreciation factor. Other factors can be obtained either by lamp manufacturer specification data or by measurement in an as-used configuration of each lamp before use.

System simulation can be done in two ways: manually or using specialized software. To do it manually, the lumen method [see Eq. (11)] and the zonal cavity method are used. The zonal cavity method²⁹ involves modifying the CU by calculating the effects of room geometry, wall reflectance, luminaire intensity, luminaire suspension distance and workplane height. The impact of various parameters is available in the form of look-up tables that can be consulted to provide estimates. These methods are quite powerful but outside the allowable space in this chapter, please consult the mentioned references.

System simulation with specialized software allows unparalleled accuracy and flexibility. Modern software tools allow easy modeling of 3D geometries, material and source properties. Sophisticated analysis tools allow for calculation of luminance, intensity, illuminance and chromaticity at any location, and they also provide photorealistic rendering of lit models that account for specular as well as diffuse reflections. Accurate system modeling involves the following steps:

1. Model the 3D geometry of the lit region. Windows, skylights, and light shelves must be modeled with their coverings: blinds or glazings with their optical properties.
2. Model the surfaces and paints. The reflectance is a combination of specular and diffuse components. For those surfaces and paints that cannot be simply described by specular or Lambertian properties, the bidirectional reflectance distribution function (BRDF) is used.³⁰ The dependency of BRDF is composed of the incident direction and the direction of observation. This concept has been explained in much more detail in Chap. 7, "Control of Stray Light," in this volume. BRDF measurements are mostly obtained experimentally. These measurements are available sometimes by paint vendors or makers of specific surfaces. It is also possible to simulate BRDF approximately by assigning texture and reflective properties to the surface and performing a ray trace. Once the reflectance properties of each surface are available, they are assigned to the surfaces in the optical model and allow accurate photorealistic rendering of the model for all cases of source and viewer locations.
3. Model the lamps: luminaire geometry and source model. Source models are increasingly being made available by the lamp vendors. If the model is not available, source measurements may be needed. Source models consist of a collection of rays that represent the output from the source. Each ray is described by its position and direction cosines in 3D space. Daylight can be simulated by creating two sources outside the model: sun and sky and assigning diffuse reflective properties (10 to 20 percent) to the ground outside the model. There are a variety of ways to model the sun. One way is to represent it as a Lambertian disk of its angular extent (0.52°) and a luminance value that provides the insolation on the earth's surface at the geographic location ($\sim 1000 \text{ W/m}^2$). The final step is to locate the Sun's position relative to the model. The sky is modeled as a large Lambertian disk of a prescribed luminance. For example, a clear sky is represented by 8000 nits and an overcast sky is represented by 2000 nits.
4. Decide upon the location of the measurement surfaces. These could be real or virtual. A photorealistic scene rendering helps in realizing the most promising layouts.

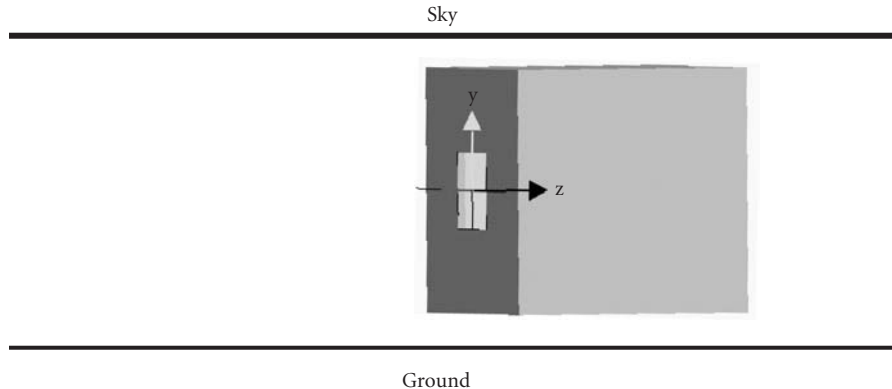


FIGURE 6 System layout. (Source: Clear sky 8000 cd/m^2 . Room Size: $3 \text{ m} \times 3 \text{ m}$. Window size: $1 \text{ m} \times 1 \text{ m}$. Window location: center of the wall. Wall reflectance: 35 percent. Roof Reflectance: 60 percent. Floor reflectance: 40 percent. Ground reflectance exterior to room: 20 percent.)

Figures 6 and 7 show an example that illustrates the above steps.³¹ We simulate a room with one window. All the light received in the room is daylight from a clear sky. In this example, we calculate the light distribution across the floor for two cases: with and without assuming that the room surfaces are diffuse reflective. Figure 7 shows the impact of including reflectances from various objects which results in a wider spread of illumination across the floor.

It is easy to make the model more complex by adding internal light sources, objects with a variety of surface properties, skylights, influence of sun and so forth. Next, we discuss the software tools available to perform such simulations.

Software Tools There are extensive software tools to assist the lighting designer in simulating illumination systems. The software includes computer-aided design (CAD), source modeling, optical analysis and design, and computer graphics. Each of these plays a crucial role in the design process, so they are described in the following subsections. Some of the software areas have applicability in a number of aspects of the design process. Finally, we do not mention the explicit names of the software packages since they are continuously evolving, and we make use of certain ones in our daily lives, so we do not want to bias the discussion presented here.

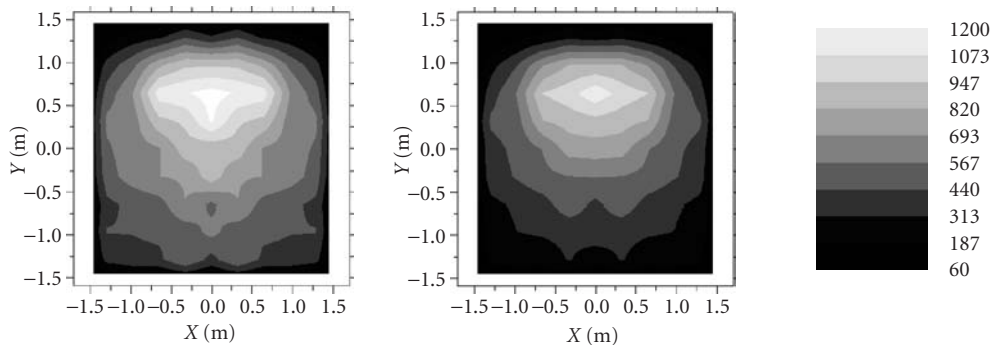


FIGURE 7 Horizontal illuminance (lux) on the floor (a) with and (b) without taking into account the room reflectance.

CAD software Computer-aided design software is used to build the geometry of any system and is then interfaced with machine tools to fabricate the components. Lighting design makes great use of CAD software to ease the process of integration of complex optical components and the mechanics that hold them and the electronics. Not only are there self-standing, mechanical CAD software packages, but native CAD geometry generation capabilities in the optical analysis software packages discussed later. The latter provide a wealth of tools for the generation of complex illumination systems, but they do not have the range of tools that mechanical CAD software provides. The software can be broken down into two subsets: surface based and solid based. The former implies that each surface is defined separately (e.g., a cube is made up of six separate surfaces), while the latter implies that each object is defined (e.g., a cube is made from one function call). Of course tools are provided in each code such that the design process is simplified (e.g., in a surface-based code a macro could generate all of the six surfaces of the cube with one function call). Solid-based codes tend to provide a more efficient process to enter the geometry of the system, but surface-based codes have a longer market history. None of the mechanical CAD packages provide optical analysis and design tools, but they do provide hooks to integrate into them. In fact, some of the optical analysis and design software companies have developed plug-ins that allow the user to specify optical characteristics such as materials and surfaces. These tools assist with the design process by requiring only one iteration of assigning such properties and simplify the transfer of the geometry to the optical analysis software. The transfer of the geometry is typically accomplished by two formats: International Graphics Exchange Specification (IGES) and Standard for the Exchange of Product model data (STEP). Other protocols including proprietary ones, DXF, DWG, STL, and SAT are understood by certain optical design and analysis tools.

IGES IGES is a standard first published in 1980 by the National Bureau of Standards (now the National Institute of Standards and Technology, NIST) as NBSIR 80-1978, and then approved by an ANSI committee as Version 1.0.³² It was first known as Digital Representation for Communication of Product Definition Data. It is essentially a surface/curve-based method to represent the geometry that comprises a component or system. The IGES standard was updated through the years, with Version 5.3 in 1996 being the last published version.³³ STEP (see the next section) was to replace IGES, but IGES remains the prevalent neutral-based method to transfer geometry data. Most optical analysis codes can read and interpret various versions of IGES, but the surface-based optics codes tend to have better performance (i.e., fewer import errors).

STEP STEP is a standard first published in 1994 by the International Standards Organization (ISO) as ISO 10303 with the goal of replacing IGES. It is essentially a solid-based method to represent a system, but it also has 2D and database aspects.³⁴ It has been updated through the years and is now comprised of many part and application protocols. STEP is developed and maintained by the ISO technical committee TC 184, Technical Industrial automation systems and integration, subcommittee SC4 Industrial data.³⁵ Most optical analysis codes can read and interpret various versions of STEP, but the solid-based optics codes tend to have better performance (i.e., fewer import errors).

Source modeling software In order to perform accurate simulation of a lighting system it is imperative that an accurate source model is used. There are essentially three methods to accomplish this:

- Generation of ray data based on manufacturer data sheets
- Experimental measurements of the emitted radiation to create ray sets
- Modeling of the geometry of the source and the physics of emission to create ray sets

The first method creates either ray sets or Illuminating Engineering Society (IES) intensity distributions that can be used to model the performance of a lighting system that incorporates the prescribed source. The second, experimental measurement uses a goniometer or similar device to measure the output of the source both as a function of position and angle; therefore, it measures the luminance distribution of the source. The last is based on modeling of the source components, such as filament, base, and glass envelope in an incandescent bulb; electrodes, glass envelope, and

base for an HID lamp; and die, epoxy dome, and reflector cup for a LED. Based on the physics of the source, rays are assigned to the emission areas. A number of optical analysis software companies are providing geometrical and ray source model libraries to their customers. All three methods are based on the data, measurements, or model of a single source, called the nominal source. Thus, there is the opportunity for error between the nominal source and what is used in your fabricated system.

Source manufacturers and architectural lighting software tend to use the IES source files to specify their sources. These files are not as precise as the other two methods, but they provide a good enough and fast method to implement the source emission characteristics into the design process. Companies that make these accurate experimental measurements keep libraries of the data, so that their customers can include them in their lighting system models. The CAD software allows a designer to make a complex model of the source, and then the optical analysis software allows the rays to be generated. This method also provides accurate source models, but with the only time expenditure to develop such models. It has been found that methods that employ both the geometry and experimental emission measurements provide the highest level of accuracy while also giving an avenue to model the tolerances of the emission.^{36,37}

Optical modeling software As previously stated, optical analysis and design software not only allows for the modeling of optical systems, but it also provides tools for inclusion of geometry, source modeling, and rendering (discussed later). There are both solid-based and surface-based codes. Most optical phenomena occur at the surface interfaces, such as reflection, refraction, scattering, and diffraction, but there are volume effects, such as scattering, absorption, and emission. Therefore, though a single software packages is based around solids or surfaces, it must be able to effectively model the other type of phenomena. A number of the codes are generic in nature, such that they can handle virtually any type of system or application, from backlights to luminaires to biomedical applications. There are application specific codes in a number of areas, especially external automotive lighting and architectural illumination. The software is increasingly adding tools such as source modeling macros, optimization, tolerancing, and rendering.

There are essentially three types of software packages that can be used to model a lighting system: optically based ray tracing, lighting-based radiosity, and computer graphics rendering. Ray tracing is simply the tracing of a multitude of rays from sources through the optical system. It is quite accurate, only limited both by the characterization of the geometry and the assigned optical properties within the model and the number of rays that are traced. Radiosity algorithms are essentially scatter-based methods that propagate approximate wavefronts from one object to another. Initially, Lambertian scatter properties of all objects were assumed, but more recently radiosity implementations propagation based on the bidirectional surface distribution function (BSDF), generally, or the respective reflective (BRDF) or transmissive (BTDF) forms, specifically, have been developed.³⁸ Ray tracing can handle both specular and diffuse reflections, but radiosity is limited to diffuse reflections. Ray tracing has a number of benefits including accuracy and utility from the near field to far field. Radiosity is for the most part limited to far field calculations, where the approximations of the propagation model are minimized. As the distance between the source and target is reduced, the limitations of the scatter-based propagation inhibit accuracy. The primary limitation of ray tracing is the calculation time, which is several orders of magnitude more than radiosity. Thus, hybrid methods that employ both ray tracing and radiosity are in use. The goal of hybrid methods is to obtain the benefits of both ray tracing and radiosity in a single algorithm.

In the next three subsections the three types of software packages are discussed in more detail. Notably, the lines between these three types of software packages is disappearing, especially between the lighting design and computer graphics sectors. In each of these sections the applicability to lighting design and modeling is provided. These software packages are seeing rapid growth, so consultation of the literature on active research on future development is suggested.

Optical design and analysis software There are two types of optical design and analysis software: imaging system software and general analysis software. The first is typically called lens design software, and has limited utility to the design of lighting systems. The second uses nonsequential ray tracing from the source to the target. This process allows the illumination distribution at the

target to be accurately determined. This type of software often includes, at the discretion of the user, such features as spectrum, coherence, polarization, and so forth. Thus, the accuracy is only limited by the user input. The design of the actual luminaire is best done with this type of software. It allows the designer to design efficient systems that effectively couple and broadcast the emission from the light source. These codes also provide tools for optimization and tolerancing of the luminaire. However, due to lengthy computation time, optical design and analysis codes have limited (but increasing, due to the advances in computer speeds) utility in determination of a lit scene (i.e., rendering).

Lighting design software Unlike optical design and analysis software, lighting design software typically makes use of radiosity algorithms.³⁹ Radiosity codes are quite fast and have acceptable accuracy in the far field. The use of Lambertian or BSDF scatter properties allows the diffuse reflection from objects to be quickly ascertained and propagated further into the system. The diffuse emission is both collected at the observation location and is cascaded to other objects in the scene. In order to obtain a higher convergence speed objects are typically parameterized with polygons, while optical phenomena such as refraction or specular reflection are approximated or even ignored. This type of software thus provides at worst a first-order approximation of the lit appearance such that architects and lighting designers can view the results of their design work. Lit-scene rendering or illumination in the far field of a luminaire is the best use of lighting design software. More advanced software packages can then be employed, such as those that employ some semblance of ray tracing—see the previous and next sections.

Computer graphics software In the past two decades the computer graphics community has grown rapidly. The tools they develop and employ are useful in the illumination community of optics. Foremost they have tools to model the simulated look of an unlit or lit lighting system, called unlit-and lit-appearance modeling respectively. These tools are important in illumination since acceptance of a system is often based on subjective criteria such as appearance. Thus, these tools provide such before potential costly and time-consuming manufacture. Additionally, the computer graphics community uses both ray-based and scatter-based radiosity methods, while also employing hybrid methods. These methods are especially geared to the rendering of scenes in video games, movies, and other types of visual media. Thus, they tend to have the least amount of accuracy since the completion of numerous images in a timely manner is demanded.

Computer graphics software makes direct use of forward ray tracing (from the source to observer) and reverse ray tracing (from the observer to the source). The latter is especially useful for the rendering of scenes where there is a discrete viewpoint, like that of a virtual observer. These codes and algorithms are increasingly being used to model the lit appearance of illumination systems such as luminaires and lightpipes. This process involves some form of ray tracing and/or radiosity calculations and then employs vision biology (Sec. 40.3). As an example, consider a star-shaped taillight as shown in Fig. 8.^{40,41} The taillight is oriented at several angles such that the effects of changing ones aspect to the lit taillight is taken into account. The combination of these different angular views of the taillight provides the luminance distribution, or essentially what an observer would see by walking around the lamp. A ray tracing method employing a pupil collection of 15° is used for each of the plots within Fig. 8. Note that saturation of the retinal cones is included, which is evidenced by the whitish appearance of the filament at the center of the lit patterns.

Furthermore, the resulting intensity pattern for the lit-appearance model can be projected into a scene to provide a rendering of what the illumination from the lamp will look like. For example consider Fig. 9, which shows three view aspects of an automobile headlight designed to meet standards (see section on vehicular lighting): (a) the driver's perspective, (b) 20 m above and behind the drive, and (c) the bird's eye view.⁴² These renderings are quite accurate since the goal is to completely mimic the lit-scene appearance prior to fabrication. These types of renderings can be extended to any scene that involves sources and objects with accurate optical characteristics applied to them. Figure 10a shows the rendering of a south-facing office room⁴³ in Tucson, Arizona. The illumination was modeled from average direct and diffuse insolation data for November 15 at noon for this location. The CAD model was generated from architectural blueprints of the facility. Surface reflectances were determined from first principles. The scene outside the window was created from a

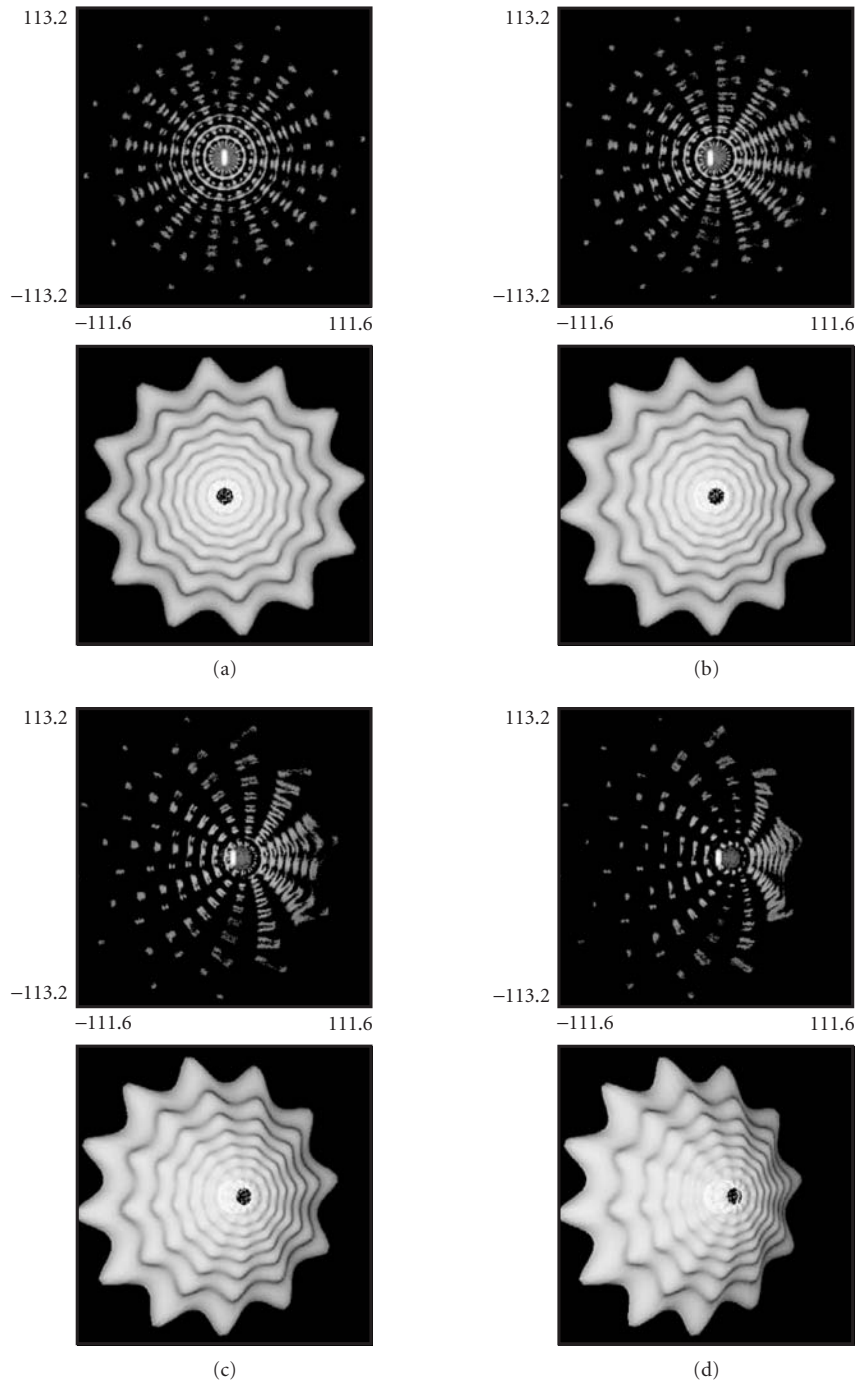


FIGURE 8 Views of the lit appearance (upper) of a star-shaped taillight (lower) at four horizontal angles of (a) 0°; (b) 10°; (c) 20°; and (d) 30°. (See also color insert.) (Used with permission from SPIE;⁴⁰ Developed with Advanced Systems Analysis Program from Breault Research Organization.)

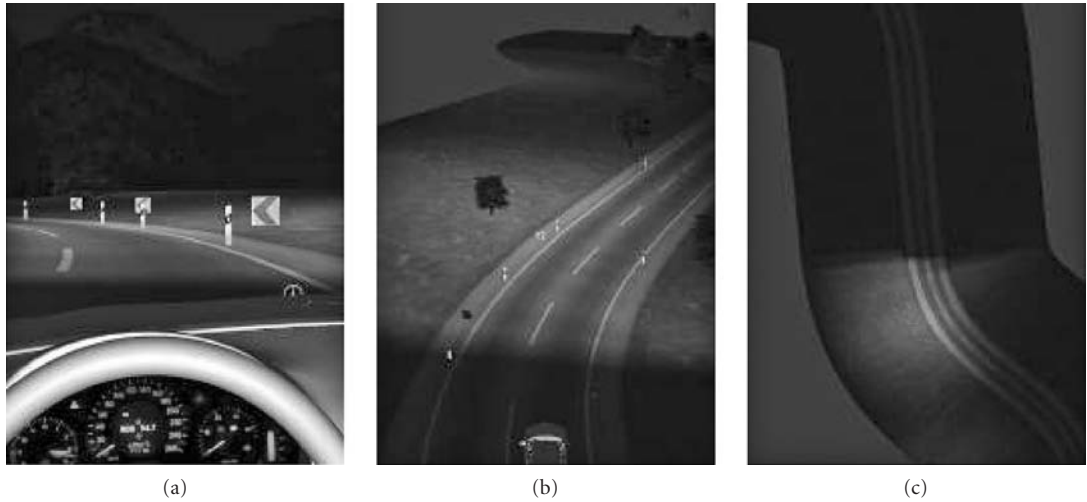


FIGURE 9 Three perspectives of lit-scene renderings from a low-beam headlamp: (a) driver's view; (b) 20 m above and behind automobile; and (c) bird's eye view. (See also color insert.) (Developed with *LucidShape* and *LucidDrive* from Brandenburg, GMBH.)

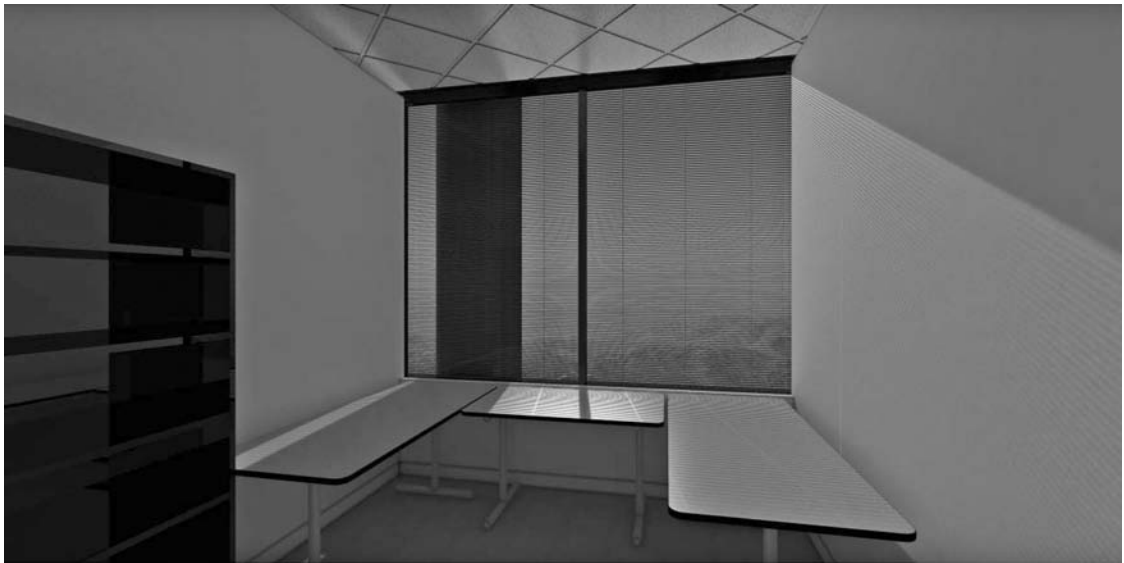


FIGURE 10a Rendering of a lit office room. (See also color insert.) (Developed with *LightTools* from Optical Research Associates.)

digital photograph taken in the mountains outside of Tucson. 100,000,000 rays were traced from the source throughout the model. Similarly, Fig. 10b shows the rendering of a desk surface lit by interior, incandescent lighting.⁴⁴ All surfaces, except the three objects on the desk (shown in wireframe to ease view through the objects), are diffuse Lambertian reflectors. The three objects: wine glass, ice cube, and crystal ball, display the effects of specular refraction and total internal reflection.

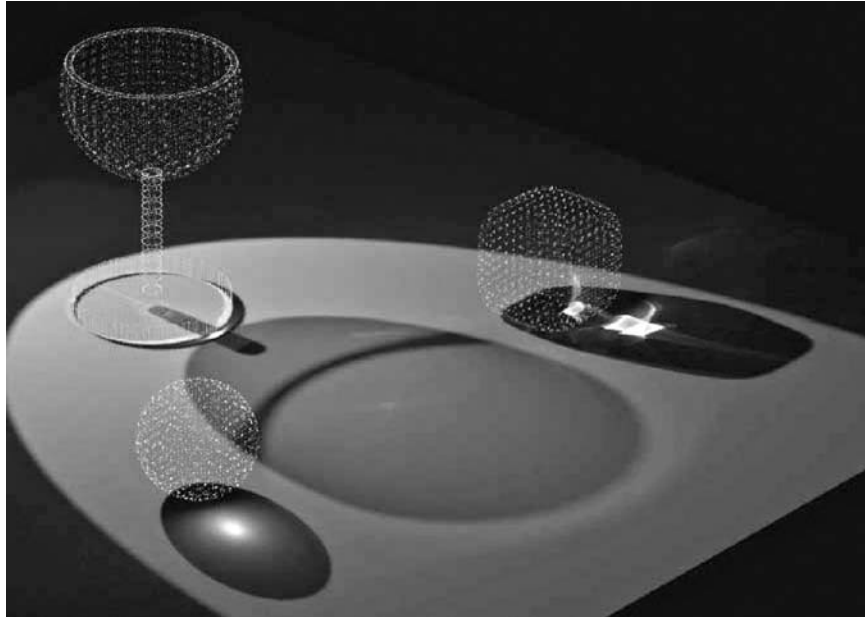


FIGURE 10b Rendering of a lit desk with three objects located on it (wine glass, ice cube, and crystal ball) to show both diffuse and specular effects. (See also color insert.) (*Developed with FRED from Photon Engineering.*)

40.5 LUMINAIRES

A luminaire is a packaged light source consisting of emitter, optics to baffle or redirect light, fixture and electrical components. Luminaires form an integral part of the lighting design by the virtue of the illumination they provide and also by their appearances. In this section, we discuss the optical components of luminaires: types of light sources, both artificial and natural (daylight), and the design of luminaires with optical components such as reflectors, lenses, lightguides, fibers, windows, skylights and baffles to achieve the desired light distribution.

Types of Light Sources

Light sources are optically characterized by their luminous spectrum (lumens as a function of wavelength), intensity (in candela), and efficiency (lumens per watt, lpw). For white light, CCT and CRI are derived from the spectrum of the source and are typically reported on the data sheets to give an estimate of its appearance and its ability to color render lit objects. Other source characteristics are cost, safety, government regulations, package size/type, lamp sockets, electrical driver requirements and constraints related to environment or operating conditions. Cost is defined in terms of luminous flux per dollar per hour of use, replacement, and initial installation costs.

We discuss daylighting and the following artificial light sources: incandescent, fluorescent, high-intensity discharge (HID), light-emitting diode (LED), electrodeless HID, electroluminescent, nuclear, laser, bare discharge, low-pressure sodium (LPS), and short arc sources. In the next few subsections we discuss various light sources in terms of operating principles, construction, and packaging. Refer to Table 3, which provides the performance comparison of various lamps; and Table 4, which

TABLE 3 General Lamp Characteristics for Most Lighting Applications

	Watts	Efficacy lpw	CCT K	CRI	Lifetime K hours	Notes
Standard Incandescent	40–100 ^a	10–17 ^T	2700	>95	0.75	Undesired IR radiation, fire hazard. Naturally low flicker, instant-on and dimming to 0.
Tungsten-Halogen	300 ^a	20	2850–3200	>95	<6	Same as standard incandescent.
Fluorescent/CFL	5–55	15–100 60–70 ^T	3000–6500 4100 ^T	50–98 70–85 ^T	5–20 ^b	Complex ballasts for good dimming range, short start-up, low hum, low flicker. Hg disposal issue, EMI.
HID-Hg	50–1000	30–65	3900–5700 ^c	15–20 higher with phosphor coatings	16–24	Complex ballasts for start-up and flicker control, poor dimming, high flicker, explosion, fire and UV hazard, lamp disposal issue, high start-up and re-strike intervals, poor color stability with time, operating position affects performance.
HID-MH	30–18 K 50–1000 ^T	75–125	2500–6000	60–70	7.5–20	Same as HID-Hg. Lifetime <1 K hours for >10 kW lamps.
HID-HPS	175–1000	45–150	1900–2700	22–85	7.5–24	Same as HID-Hg, CRI is inversely proportional to efficacy. 110 lpw corresponds to CRI 20, CCT 1900–2100 K.
HID-CMH	20–400	70–90	3000–4200	80–96	7.5–20	Same as HID-Hg except for good color stability with time and performance independent of operating position.
Electrodeless	4	60+	2700–6500	50–98	15–30 ^d , <100 ^e	Complex ballast, EMI.
LPS LED^f	20–100	80–150	1800	0	14–18	LEDs (large chip or arrays) are likely to replace most of the existing lamp sources. Efficacy can reach up to and beyond 200 ^g lpw with the theoretical limit being the CIE standard observer luminous efficacy curve, lifetimes up to 100 K hours. Key advantages are: fast turn-on times, color tunability, high dimming range, low voltage operation, no Hg or Lead, no UV or IR, no catastrophic failures and scalable packages. Phosphor coated LEDs can provide white light at desired CRI and CCT.

HID—high-intensity discharge, Hg—mercury, MH—metal halide, HPS—high-pressure sodium, CMH—ceramic mercury halide, Xe—xenon, LED—light-emitting diode, CFL—compact fluorescent lamps, ^T—typical.

^aAlso available in kW. Lifetimes shrinks to a few hundred hours, ^bRegular Fluorescent have typical lifetimes >10 K hours, Compact Fluorescent lamps have lifetimes >5 K hours, ^cCCT 5700 K at CRI 15, CCT 3900 K at CRI 50, ^ddeveloping technology, ^ewith integrated ballasts, ^fseparate ballast, ^gonly light generation efficiency inside the LED die.

lists various lamp types, currently available packages and their applications. Figure 11 shows various lamp packages. The alphabetic designation is explained in Fig. 11. The numeric designation with the letters is the diameter specified in 1/8th of an inch. For example, MR16 refers to a multifaceted reflector, 2 in. in diameter.

Incandescent Sources These are thermal sources that emit electromagnetic radiation from a heated filament, which is a phenomenon known as *incandescence*. Incandescent light sources are being steadily replaced with fluorescent lamps, LEDs, and potentially with electrodeless lamps. Other sources are able to provide similar or improved performance at higher efficiency and lifetime leading to reduced operational costs.

Modern day incandescent sources use a tungsten filament.⁴⁵ Tungsten has a high melting point (3382°C), high ductility, high conductivity, and low thermal expansion that make it a preferred material for use as a lamp filament. Tungsten is alloyed with tiny amounts of potassium (60 ppm), aluminum oxide (10 ppm), and silicon (1 ppm) to give it high strength near its melting point. This allows lamp operation close to the melting point, thus improving its efficacy. Sometimes tungsten

TABLE 4 Common Lamp Packages and Applications

	Available Packaging	Current Applications
Tungsten	A-line, elliptical, decorative (B, C, CA, F, G, M), PAR, reflector (R, BR, ER), appliance and indicators (S), tubular (T)	General purpose, 3-way, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight and spot light), outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, bug light, yard stake), appliances, colored lamp, display, exit sign, freshwater and saltwater aquarium, heat lamp, marine, nightlight, party, recreation vehicle, plant, rough service, sewing machine, shatter resistant, terrarium, vacuum cleaner, airport, emergency, city lighting, projection, photoflood, filmstrip, retail display, restaurants.
Tungsten-Halogen	A-line, decorative (B, G, F, T10), PAR, AR, MR, single ended, double ended	General purpose, 3-way, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight), outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, yard stake), camera light, microfilm, curio cabinet, display, enlarger & printer, equipment, fiber optics, landscape lighting, projection, stage & studio, torchiere, airport, emergency, city lighting, special service, monuments, museums, heat lamp.
Fluorescent (Linear)	Straight linear (T5, T6, T8, T10, T12), circular (T9), U-shaped (T8, T12), grooved (PG17)	Kitchen, bath, shop, work light, Appliances, blacklight, blacklight blue, cold temperature, colored lamp, freshwater and saltwater aquarium, plant, shatter resistant, terrarium, stage & studio, projection, diazo reprographic, germicidal, gold UV blocking, superstores, warehouses.
Compact Fluorescent	Plug-in (2-pin, 4-pin, proprietary), self-ballasted (decorative, reflectors, proprietary)	General purpose, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight and spot light), outdoor (post & lantern, pathway, garden & deck, bug light), appliances, blacklight blue, facilities, hospitality, office, plant, restaurant, retail display, saltwater aquarium, terrarium, torchiere, warehouse.
HID - Hg	A-line, elliptical, reflector	Street lighting.
HID - MH	Elliptical, PAR, single ended, double ended, tubular	Street lighting, sports lighting, decorative lighting of architectural wonders.
HID - HPS	Elliptical, double ended, tubular	Street lighting, horticulture.
HID - CMH	Elliptical, par, single ended, double ended, tubular	Track and Recessed (Indoor floodlight and spot light), retail display.
Miniature	B, G, R, RP, S, T, TL, discharge	Outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, yard stake), automotive (headlamp, fog, daytime running, parking, directional front & rear, tail, stop, high mount stop, side-marker front & rear, backup/cornering, instrument, license plate, glove compartment, map, dome, step/convenience, truck/cargo and under hood), flashlight, landscape lighting, low voltage, marine, telephone, traffic signal, emergency.
Sealed Beam	PAR, rectangular	Outdoor (motion-sensing & security, yard stake), automotive headlamp, railway, shatter resistant, stage & studio, directional lighting, aircrafts, tractors, airport, emergency, city lighting.
LPS Electrodeless	Tubular T, P	Street lighting, parking lots. Any application where long lifetimes and high efficacy are needed due to high replacement costs and/or difficult access. Road signs, warehouses.
LED	MR16, miniature (2, 3, 5 mm)	LEDs can potentially replace most of other light sources being used for various applications. Currently being used in building interiors (homes, offices, commercial places such as health clubs, hospitals rooms) stairways and pathways, flashlights, traffic lights, signs, digital projection, instrument indicator panels, backlighting applications, decorative, display.

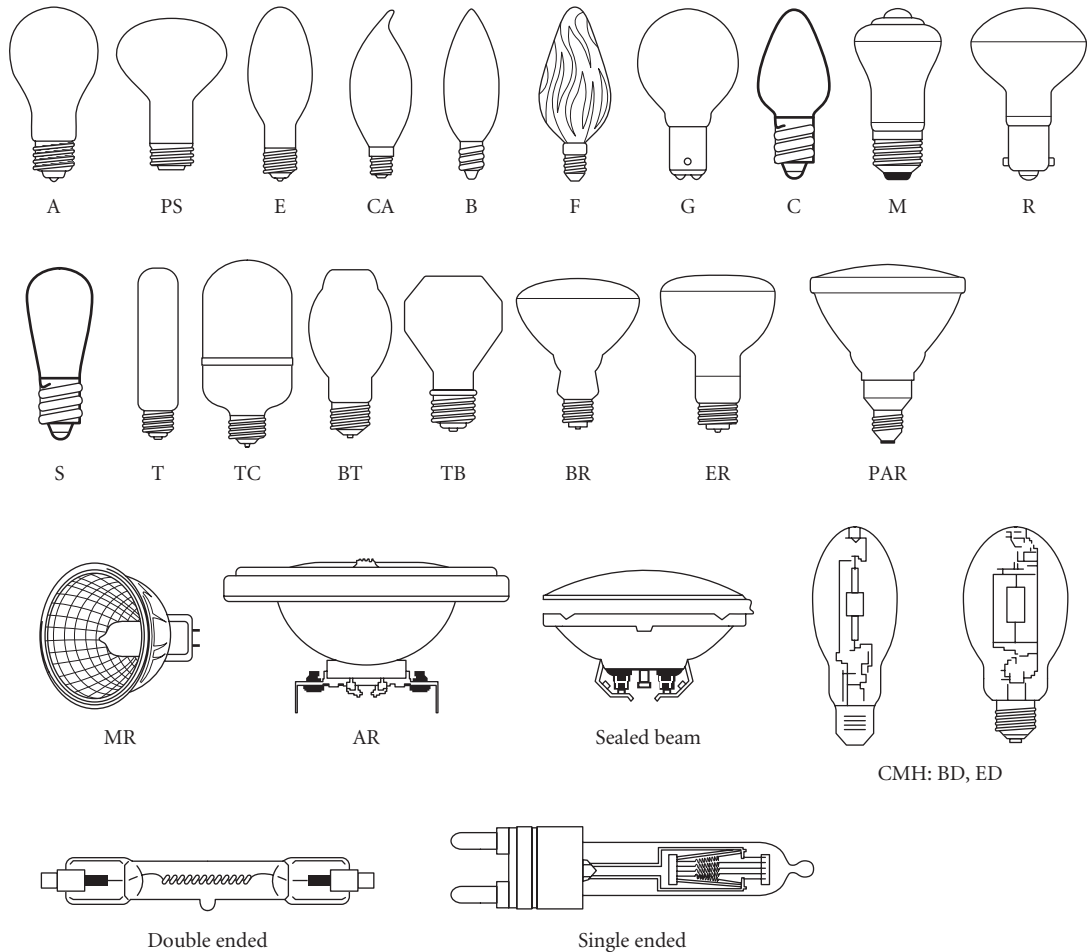
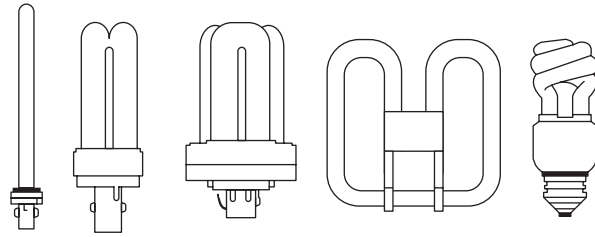


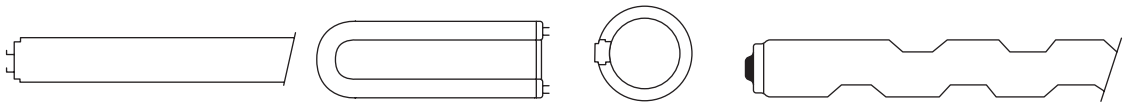
FIGURE 11a Various lamp package representations. Shape and sizes are not to scale. Various sizes are available for each package type. For each bulb shape, a variety of bases are available. A—arbitrary spherical shape tapered to narrow neck, B—bulged or bullet shape, BT—bulged tubular, C—conical, CA—conical with blunt tip, E—elliptical, blunt tip, ED—elliptical with dimple in the crown, F—flame shaped, decorative, G—globe, M—mushroom-shaped with rounded transitions, MR—multifaceted reflector, PS—pear shaped with straight neck, PAR—parabolic aluminized reflector, BD—bulged with dimple in crown, S—straight, T—tubular, TB—Teflon bulb, TL—tubular with lens in crown. (Illustration courtesy of General Electric Company.)

is alloyed with rhenium (3 to 25 percent) to make it more ductile at low temperatures and achieve higher recrystallization temperatures thereby giving the lamp a longer life. Alloying tungsten with thorium provides increased strength, better machinability and high recrystallization temperatures. Such filaments are used for very high voltage applications.

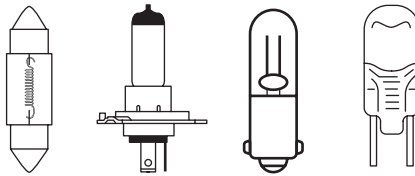
Figure 12 shows the radiating characteristic of tungsten at 3000 K and its comparison with a blackbody emitter. It differs from the blackbody due to wavelength-dependent low emissivity. A perfect blackbody has an emissivity of one for all wavelengths. The hotter the filament, the higher are the luminous flux radiated per watt, the percentage of luminous flux of the total radiation, and the CCT. However, the lifetime is inversely proportional to the filament temperature as filament evaporation is



CFLs: Biax, double biax, triple biax, 2D, spiral



Fluorescent T, U-line, circline and grooved PG



Miniature: neon, festoon, automotive, TL

FIGURE 11b Various lamp package representations. Shape and sizes are not to scale. (Illustration courtesy of General Electric Company.)

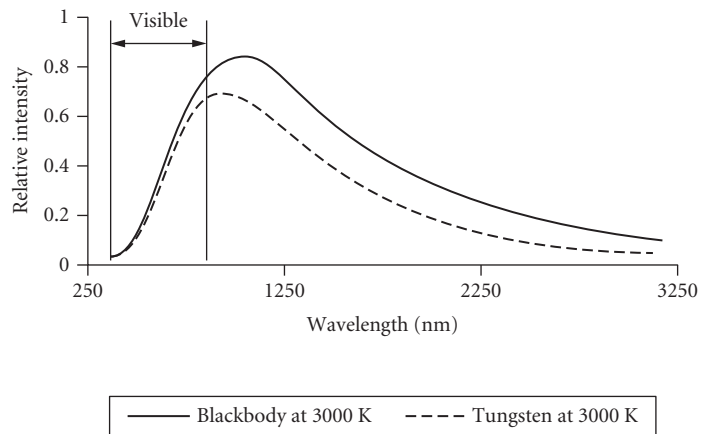


FIGURE 12 Blackbody at 3000 K versus tungsten filament at 3000 K.

the primary mode of lamp failure. The best engineered solution consists of the hottest possible lamp filament at an acceptable lifetime, voltage rating, and packaging. Near its melting point, an uncoiled tungsten wire has a luminous efficacy of 53 lm/W. To achieve an acceptable lifetime, the tungsten filament is operated at much lower temperatures. The luminous efficacy of typical incandescent lamps with tungsten filaments ranges from about 5 to 20 lm/W for lamp wattages 5 to 300, respectively.

A typical incandescent lamp consists of an evacuated glass envelope; filament, with or without fill gases; filament leads; and base. Figure 13 shows the key characteristics.

The bulb material is application dependent. For most applications soda lime glass is used. For applications requiring heat resistant glass, borosilicate, quartz, or aluminosilicate is used. When the bulb envelope is frosted from the inside or coated with powdered silica, it provides diffuse illumination from the bulb surface and masks the bright filaments from direct view. The bulb envelope material can be used to filter the radiation to alter the CCT. Daylight application bulbs filter out the longer wavelengths to provide a higher CCT.

Lead-in wires are made of borax coated dumet (alloys of nickel, copper, and iron). Dumet is able to form a glass-metal seal. It is important to match the thermal expansion of the lead-in wires with the envelope material. When high bulb-envelope temperatures are involved, molybdenum strips bonded to lead-in wires are used for glass-metal seals. The bulb base is cemented to the bulb envelope and is designed to withstand the operating temperatures. The filament itself comes in various configurations depending upon the application. Various filament configurations are a straight wire (designated as S),

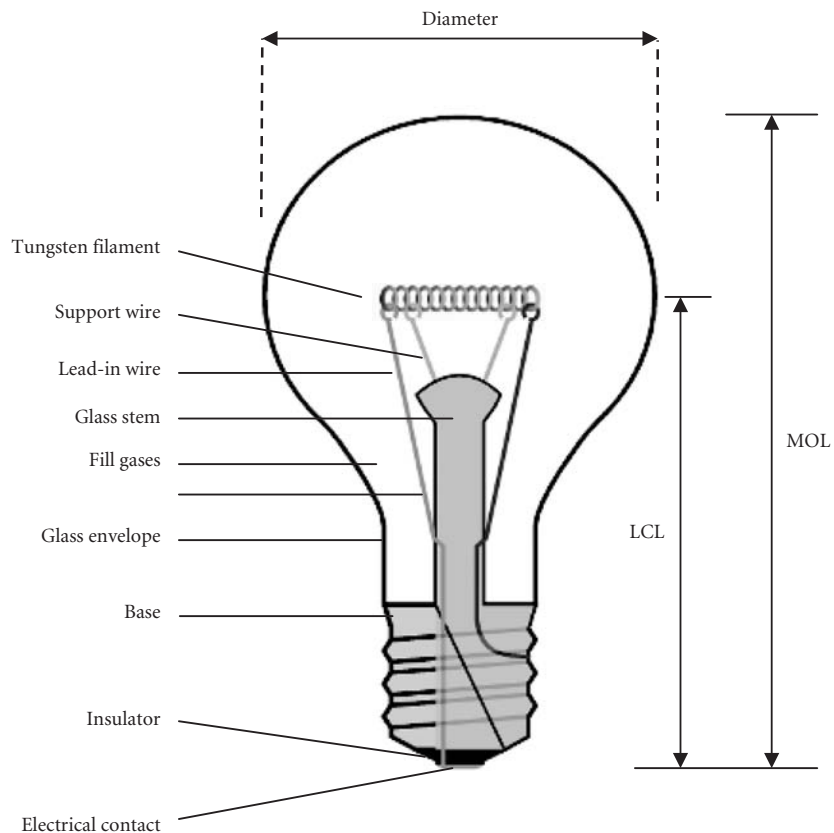


FIGURE 13 An incandescent bulb. (MOL—maximum overall length, LCL—light center length.) (Adapted from Wikimedia Commons.)

a coiled wire (designated as C) or a coiled coil (the coiled wire is further coiled onto itself, designated as CC). See Fig. 16 in Chap. 15, "Artificial Sources," in this volume. Coiling increases the surface area per unit length of the filament as well as volume packing density. This allows higher operation temperature at higher efficiency due to relative reduction in heat lost to convection. Multiple supports are used to reduce filament vibration. The number and type of supports depend upon the bulb operating characteristics. A higher number of supports is needed for rough/vibration service lamps. These lamps have low wattages and efficiency and operate in environments that involve shock and vibration for the lamp. Heavier filaments withstand vibrations much better and need fewer filament supports.

Fill gases at low pressure are used to prolong the lamp life at high operating temperatures (high efficacy and power) by reducing the evaporation rate of tungsten. For most applications, mixtures of argon and nitrogen are used as fill gases. More expensive gases such as krypton and xenon are added to the mixture when the higher efficacy justifies the cost increase. Bromine is added to create a new class of lamps called tungsten-halogen lamps which we discuss next.

Tungsten-halogen lamps operate at substantially higher temperatures leading to higher CCT and efficacy. In addition, they have longer lifetimes. At high temperatures, evaporated tungsten combines with the halogen gas and forms a gaseous compound. This gaseous compound circulates throughout the bulb via convection. When this gaseous compound comes in contact with the hottest parts of the lamp such as electrodes the halogen compound breaks down. Tungsten is re-deposited on the electrodes and the halogen is freed up to take part in the halogen regenerative cycle. For the halogen regenerative cycle to work, the bulb envelope must reach a high temperature ($>250^{\circ}\text{C}$). At lower temperatures, the tungsten will deposit on the bulb envelope instead and lead to lamp blackening and filament thinning, a common failure mode of incandescent lamps. Operating a tungsten-halogen lamp at less than the rated voltage inhibits the halogen regenerative cycle and takes away the longevity advantage. Tungsten halogen lamps have relatively compact bulb envelopes to allow for high bulb-envelope temperatures. The bulb envelope may also have an IR reflective coating to enhance the heat density inside the bulb. Compact bulb sizes make them good candidates for use with reflectors (such as a PAR) to provide directional properties for the lamp emission. The bulb envelopes for such lamps are made of heat resistant material to withstand high temperatures. The high filament temperature in halogen lamps generates UV. The UV rays must be blocked by a UV absorbing lamp cover or UV absorbing but heat resistant bulb envelope such as high-silica or aluminosilicate.

Flicker in incandescent lamps is naturally very low due to slow response of filament temperature to voltage fluctuations caused by power supply frequency or noise.

Incandescent lamps fail when bulb blackening or filament notching (thinning of the filament by evaporation) reduces the output substantially or the filament breaks either by vibration or by complete evaporation of some filament portion. With the exception of halogen lamps, operating the incandescent lamps at less than rated voltage dramatically extends the lifetime (by reducing the filament evaporation rate) at the cost of reduced efficacy, CCT, and luminous flux. In situations, where long lifetimes are needed such as when the lamps are located in a hard to replace areas and/or under tough environmental conditions, heavy filament lamps that are operated at less than the rated voltage are used. But using lower operating voltage to extend the lifetime is not always economical once the increased cost of electricity due to reduced efficacy and compensation of the reduced flux by using more bulbs is taken into account.

Fluorescent Lamps These are luminous sources based on light emission by excited states of phosphors, a phenomenon known as *fluorescence*. These phosphors are typically excited by UV emission due to spectral line transitions across gases such as mercury vapor and/or rare gases such as Xe and Ar. Most commonly available fluorescent lamps are mercury-vapor-based although mercury-free fluorescent sources are available. Phosphors are now available that are excitable by visible light. Such phosphors are also being used for LED-based light sources to produce white light. Fluorescent lamps are actively replacing incandescent light sources and in many cases are themselves being replaced by LEDs.

As shown in Fig. 14, a fluorescent lamp consists of a closed tubular fluorescent material coated glass envelope filled with low-pressure mercury vapor, electrodes at each end of the tube and a ballast to provide high-strike voltage across the electrodes and limit the current during operation. A high-strike voltage across the electrodes initiates a gas discharge in the mercury vapor with light

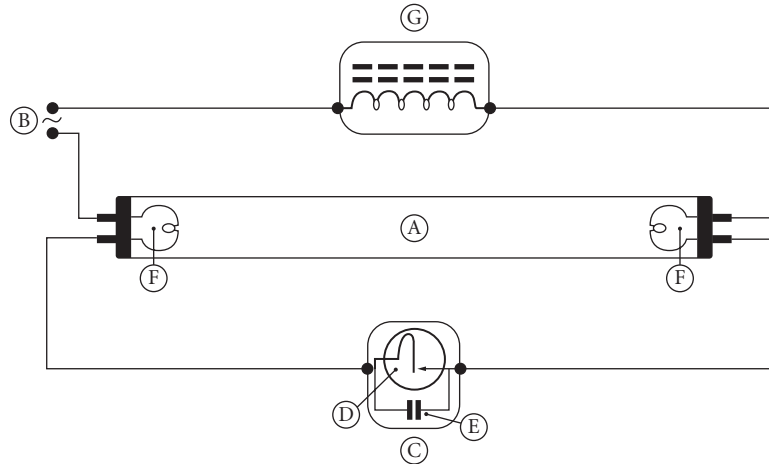


FIGURE 14 A preheat fluorescent lamp circuit using an automatic starting switch. A—fluorescent tube with fill gases, B—power, C—starter, D—switch (bimetallic thermostat), E—capacitor, F—filaments, and G—ballast. (Courtesy of Wikipedia Commons.)

emissions across various wavelengths. The UV emission lines of mercury, mostly at 254 nm, excite the phosphors coated on the inside of the tube, cause them to fluoresce and lead to emission in the visible. The fill gases consist of mercury vapor at low pressure (10^{-5} atm) for UV emission and mixtures of inert gasses such as argon, krypton, neon, or xenon at a relatively higher pressure (10^{-3} atm). Inert gases help in lowering the strike voltage across the electrodes as discussed later.

Mercury-free fluorescent lamps mostly use excimers (excited dimers of rare gases and/or their halides with Xe being popular) to produce UV to excite the phosphors. High-wattage operation and high lifetime is achievable but the efficacy is low as compared to the mercury-based lamps. Xenon-filled fluorescent lamps are also available but far less efficient than excimer based. Both excimer and Xe-based fluorescent lamps have additional advantages: instant-on, instant restrike, and color stability.

Fluorescent lamp envelope material is made of soft soda lime glass. This glass material blocks all UV. The length and diameter of the fluorescent tubes affect the efficacy and operating characteristics (voltage, current, and temperature). Longer tubes need higher voltage and power but provide higher efficacy. The efficacy is optimized for a certain tube diameter. To reduce the angular extent of emission, fluorescent tubes are coated across a prescribed region for high reflection in the visible. In such cases, only the uncoated regions of the tube emit light.

The fluorescent coating consists of different mixtures of phosphor salts.^{46,47} Phosphor salts consist of oxides, sulfides, selenides, halides, or silicates of zinc, cadmium, manganese, aluminum, silicon, or rare earth materials. Inclusions to the base phosphor material help tailor the emission characteristics. There is a large variety of phosphors available. Each phosphor emits light in one or more narrow wavelength bands. A fluorescent coating consists of one or more phosphor materials to produce a desired CRI and CCT for white light emission or specific spectral characteristics. Halophosphates (wide band) and triphosphors (blends of three narrow band red, green, and blue rare-earth phosphors) are commonly used for general lighting applications. Figure 15 shows an example of the spectrum from a halophosphate phosphor. Figure 17 shows the spectrum of a standard fluorescent lamp.

Due to the variety of fluorescent lamp spectra available, it is possible to achieve the desired CCT, CRI, and color requirements of an application. For example, an application may require a specific blend of white light such as “warm white” (CCT 3000 K, CRI 53), “cool white” (CCT 4100 K, CRI 62), “daylight” (CCT 6500 K, CRI 79), or special requirements such as aquarium lighting for providing optimal plant and coral growth and color enhancing of the display. CFLs use triphosphor coating to produce a high CRI (>90) in order to effectively replace the incandescent lamps. Special application

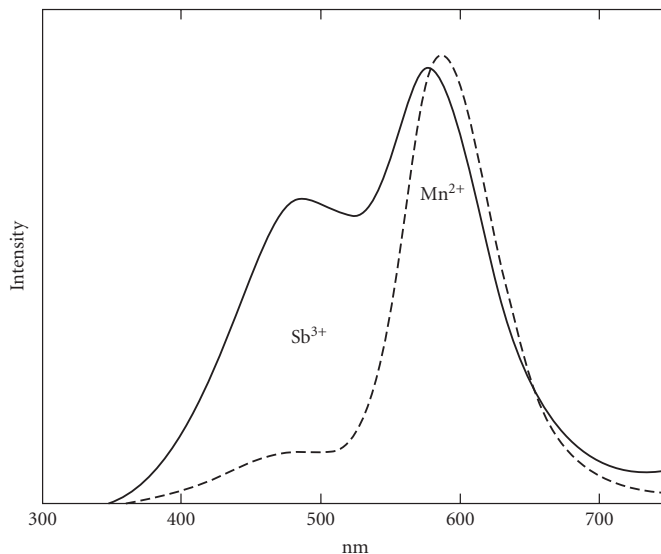


FIGURE 15 Emission spectrum of halophosphate phosphor:⁴⁷ Sb^{3+} , Mn^{2+} activated $\text{Ca}^5(\text{PO}_4)^3(\text{Cl}, \text{F})$. The ratio of Sb^{3+} and Mn^{2+} species can be adjusted to adjust the spectral distribution.

fluorescent lamps include backlight and tanning lamps. These lamps have phosphors that convert short-wave UV into long-wave UV for applications in tanning (UVA and UVB), detecting materials that fluoresce at long UV (urine, paints, or dyes), or attracting insects.

The lamp electrodes are coated with materials such as oxides of barium, calcium, and strontium to provide low-temperature thermionic emission (electron emission from a heated electrode). Under a potential, these electrons accelerate and ionize the inert gas atoms by impact ionization. Each ionization event generates more electrons that are available to accelerate and further ionize leading to an avalanche that rapidly lowers the gas conductivity. Eventually mercury atoms are ionized. The operating voltage drops and a steady current is established. UV is emitted by transitions across the excited states of mercury atoms. The electrodes can be operated in either of the two modes: cold cathode (arc mode) or hot cathode (glow mode). In the hot cathode mode, the electrodes are preheated to improve the thermionic emission and lower the strike voltage. To enable preheating, electrodes at each end are in an incandescent bulb-like configuration with a tungsten filament (straight wire or coiled). The ballast circuitry allows for preheating (up to 1100°C) before the voltage strike. Hot cathode operation allows operation at relatively higher power and over larger tube sizes. In contrast, cold cathode can be a single cylindrical pin electrode at each end. The strike voltage across the electrodes is relatively higher (~10×). The high voltage strips the electrons from the cathode at ambient temperature and initiates the breakdown process of the gas. The coatings on the electrode surface amplify production of secondary electrons, which are produced when high-energy ions and electrons collide against the cathode. These electrons further increase the gas conductivity. Cold cathode fluorescent lamps (CCFL) are compact (~3-mm diameter) and are used in applications such as thin monitors (i.e., LCDs), backlights, timers, photocells, dimmers, closets, and bathrooms. CCFLs are less efficient (<50 lm/W) but provide instant-on and have long lifetimes (50,000 hours). CCFLs operate at high surface temperatures and require complex power supplies which are fairly compact.

The primary functions of the lamp ballast are to provide a high-strike voltage to start the lamp, give regulated current supply during lamp operation and sometimes provide for cathode preheating for rapid restart applications. Often starter circuitry is deployed to preheat the electrodes. It is critical to use the correct lamp-ballast combination for proper operation. Lamp ballast can be a current limiting resistor, magnetic ballast or electronic high-frequency ballast (most modern ballasts). Due to high-frequency operation of electronic ballasts, flicker in fluorescent lamps is reduced to

almost unnoticeable. Flicker is caused by fast response of the fluorescent lamps to the voltage fluctuations in the lamp power supply caused by noise or operating power supply frequency.

Lamp failure mode consists of thermionic emission electrode coating degradation (a function of strike voltage and the number of strikes), phosphor degradation, mercury loss (diffusion or absorption by lamp materials) and ballast malfunction. Fluorescent lamps fail to operate far outside the ambient temperature range for which they are designed. Most fluorescent lamps are designed to operate in an ambient temperature of $\sim 20^{\circ}\text{C}$. For operation at low temperatures, special cold start circuitry and mercury amalgams are needed.

High-Intensity Discharge (HID) and Low-Pressure Sodium (LPS) Lamps These sources emit light across the spectral line transition of enclosed gases by electrical discharge. The source spectrum also includes background thermal radiation due to the heated electrodes and plasma. These sources are similar to fluorescent sources in basic physics involving discharge and emission of light in the UV and visible due to transitions between the excited states of gas atoms. Unlike fluorescent lamps, the electrodes are separated by less than 1 mm up to a few inches. The enclosed gases are at a pressure that is three orders of magnitude higher than in fluorescent sources. As a result, HID sources are far brighter than fluorescent sources with much higher lumen output. The following types of HID lamps are commonly available: mercury vapor (Hg), metal halides (MH), high-pressure sodium (HPS or nicknamed as White SON), and ceramic metal halides (CMH). CMH combines the advantages of MH and HPS technologies. Each HID lamp technology has very different performance and operating characteristics. An LPS lamp is similar to HID lamps in construction and operation with some differences that are identified as we describe the HID lamps below. Figure 16 describes the construction of various HID and LPS lamps. Figure 17 shows the relative spectrum of various HID lamps and comparison with the fluorescent lamp spectrum.

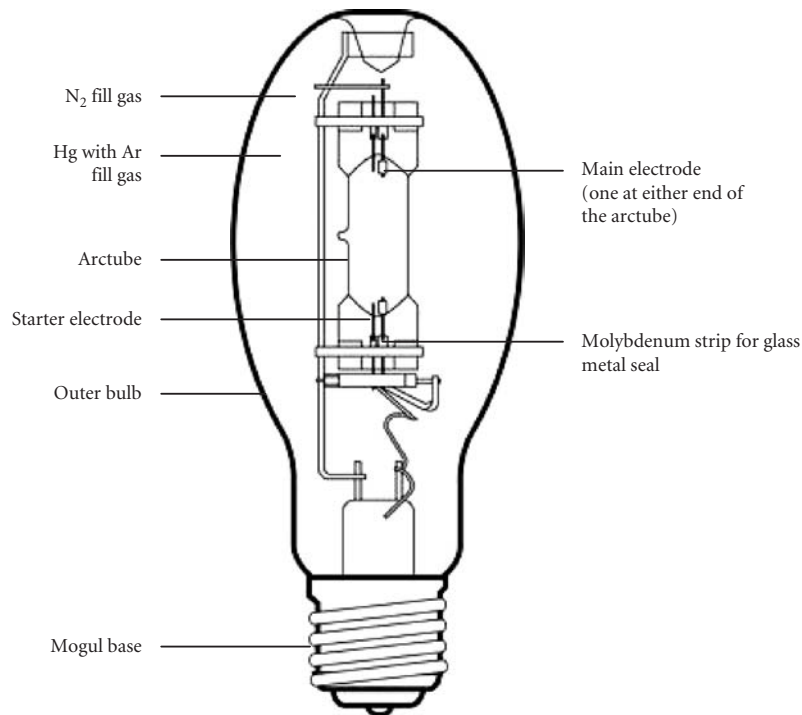


FIGURE 16a Mercury lamp construction. (Illustration courtesy of General Electric Company.)

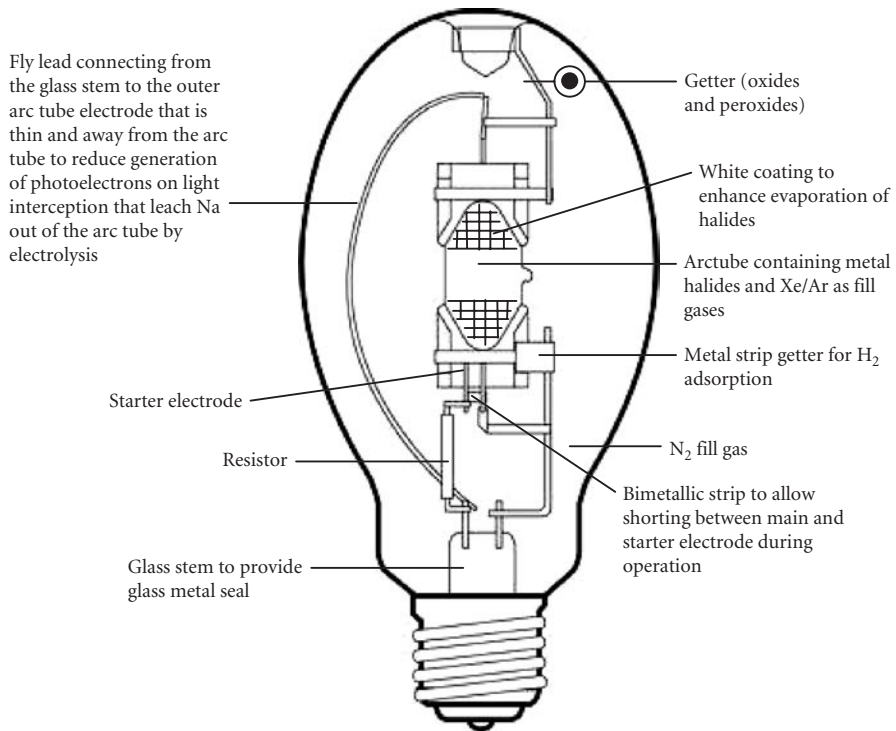


FIGURE 16b Mercury halide lamp construction. (Illustration courtesy of General Electric Company.)

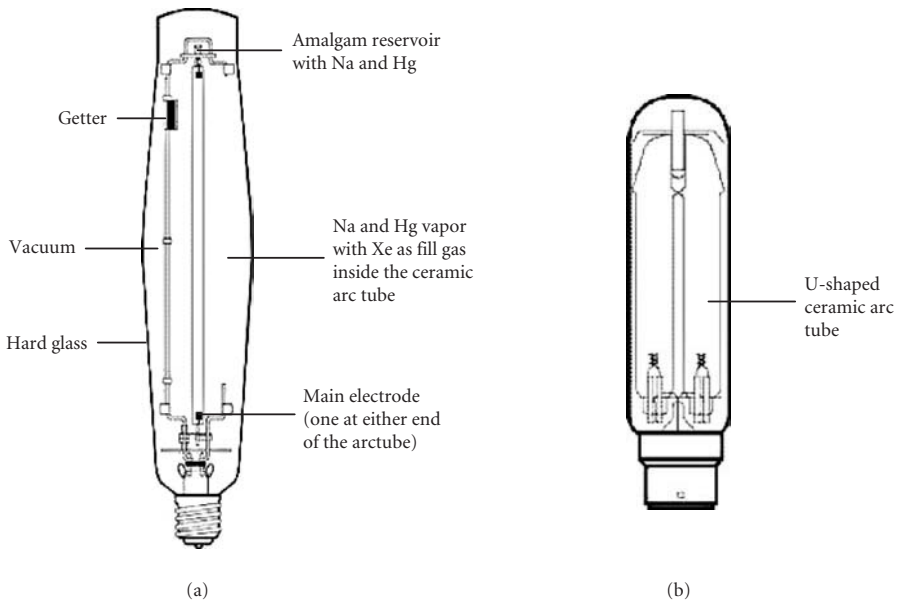


FIGURE 16c (a) High-pressure sodium lamp construction and (b) low-pressure sodium lamp construction. (Illustration courtesy of General Electric Company.)

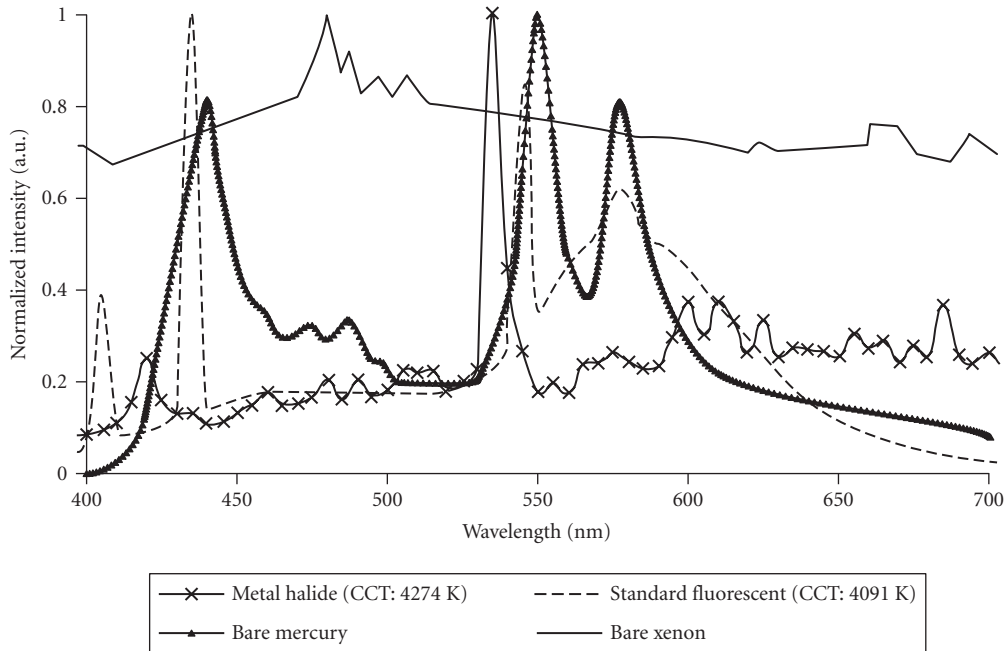


FIGURE 17 Representative lamp spectra of standard fluorescent, bare mercury, bare xenon, and metal halide lamps.

An HID lamp consists of two glass envelopes: inner and outer. The inner glass envelope or the arc tube is made of quartz for MH and Hg and alumina ceramic for HPS, LPS, and CMH lamps. The arc tube houses the discharge gases at high pressure (several atmospheres) and the tungsten electrodes. The outer glass envelope is made of borosilicate and sometimes with soft glass in Hg lamps. It absorbs UV and insulates the arc tube from outer convection currents and from large ambient temperature ranges. It houses the lead-in wires, circuitry to help initiate the high-voltage discharge, getters in case of MH to absorb impurities and has a vacuum (HPS) or low-pressure nitrogen (MH and Hg) to prevent shorting of the lead-in wires. The outer-envelope Hg lamp is sometimes coated with phosphors to provide white light at high CRI and CCT like in regular fluorescent lamps.

Tungsten electrodes are coated with various oxides in a tungsten matrix (except in MH lamps where gases can react with such electrodes) to slow down evaporation and assist in thermionic emission when heated. Starter electrodes, when used in MH and Hg lamps, assist in arc initiation via an electric field between the starter electrode and the adjacent main electrode. During operation the starter electrode is removed from the active circuitry with a bimetallic strip, otherwise premature lamp failure results.

An LPS lamp is similar in construction and operation to HID lamps. Key differences are lower arc tube pressure (0.7 atm), long arc length with a U-shaped arc tube. The arc tube gases include sodium vapor and small amounts of Ne, Ar, or Xe as startup gases.

Different HID lamp types have different gas mixtures. An easily ionizable gas such as Argon (Ar), Neon (Ne), or Xenon (Xe) is used in the arc tube to help in arc initiation. Similarly, mercury is used in most HID lamps to achieve high pressure and improved color rendering. An HPS lamp used sodium (Na)-Hg amalgam. Mercury-free HPS lamps are also available. An LPS lamp uses sodium vapor. The MH lamps use halides of metals in addition to mercury, argon, and xenon.

A low-pressure sodium lamp (LPS) has a CRI of zero due to spectral emission only at 589.0 nm and 589.6 nm (sodium-D line). An HPS lamp has broader spectrum and thus better CRI due to

pressure broadening of the sodium-D line. At very high pressures (27 atm), the sodium-D line is self absorbed by the cooler outer layers of the arc leading to a narrow spectral hole around 589 nm. Mercury atoms help in further broadening of the red end of the pressure broadened sodium-D line due to Van Der Waals forces. A CRI from 22 to 85 is achievable depending upon the pressure which can be greater than 90 atm. MH lamps achieve a rich spectrum due to the line spectra of metals like sodium, tin, dysprosium, holmium, thulium, scandium, iron, or cesium. An MH lamp may have mixtures of halides of one or more metals to achieve the desired efficacy, CRI, and CCT.

An MH lamp using a single metal halide compound can also be used to generate discrete spectral output: orange (sodium), green (thallium), blue (indium), and UV (iron). The metal-halide compound is stable at low temperatures and does not react with the arc tube material unlike some metals. At high temperatures, near the arc, the metal-halide compound breaks down and provides the spectral line emission from the metals. The operating temperature is lower than would be the case where the metals are evaporated to see the spectral emission lines. CRI is usually traded for lifetime and efficacy. CMH lamps combine the advantages of MH and HPS by using a poly crystalline alumina as the bulb envelope material. This material does not allow diffusion of metals, especially sodium or reaction of metals with the bulb material. The bulb is operated at a much higher temperature and pressure than the MH lamps. These advantages lead to high color stability with high CRI, uniformity, and efficacy over the lifetime in spite of the bulb material allowing only ~90 percent transmission.

HID lamps require several minutes of start-up time (time to reach stable output) and restrike. A long start-up is due to the time taken to reach a stable operating temperature and pressure within the arc tube. The restrike time interval results from the need to have low pressure inside the arc tube for arc initiation. Complex ballasts are needed to provide startup, restrike, and stable operation with constant current. To improve startup, restrike, and operations at low voltage, a high voltage-low current pulsed start is used. Sometimes, multiple arc tubes within the outer bulb envelope are used to provide faster restrike. Only one arc tube operates at a time in such lamps. Xe-based HID lamps are capable of instant-on and restrike. Automotive HID lamps use Xe with metal halides to improve the start and restrike times dramatically.

The physical orientation of HID lamps such as Hg and MH during operation is far more important than with the other lamps. Due to convection, within and outside the lamp, different portions of any lamp, not just HID, are heated to different temperatures. The lamp engineering must take this into account and ensure that the hottest regions do not constitute a failure mode either by design or by providing instructions to the user for best operating configuration. In MH and Hg HID, the high convection roll within the long arc tube has an overwhelming effect on the arc shape and position under gravity. It can make the arc shape curved and lead to nonuniform degradation of the electrode tips and impact the light output, lifetime, and light distribution patterns. HPS lamps, however, can be operated in any position primarily due to a compact arc tube at high gas pressure (5 to 27 atmospheres).

Lamp degradation and failure occur due to electrode degradation by evaporation, arc tube blackening due to electrode material deposition, loss of gas pressure, and selective diffusion of gases leading to change in the lamp color. Arc tube blackening also leads to the rise in the arc tube temperature leading to an analogous rise in pressure and operating voltage. The effect is especially pronounced in HPS where lamp cycling can occur: as the lamp cools down, it is able to restrike but after some time temperature rises to the point that it shuts down.

Electrodeless Lamps As the name suggests, electrodeless lamps do not have any electrodes internal to the bulb envelope. As a result lifetime is not limited by electrode degradation. The concept behind these sources is over a century old.⁴⁸ These sources are being sought to replace conventional light sources where high flux is needed at low operational costs (long lifetimes and high efficiency). There are two kinds of electrodeless lamps: induction lamps (IL) also known as *electrodeless fluorescent lamps* and microwave powered lamps also known as *electrodeless sulfur lamps* (ESL). Extraordinary high bulb life times (>25,000 hours) are possible due to lack of electrodes that degrade under operation. The causes of lamp failure are due to electrical components rather than the bulb itself, which implies a greater lifetime.

In each case, the goal is to excite a discharge with an EM field without the need of electrodes inside the bulb. Alternating magnetic fields in IL or microwaves in ESL initiate the discharge by

accelerating the free electrons of a gas with low ionization potential, such as argon or krypton. Free electrons are created in the gas by a spark from a high-voltage pulse across two electrodes in the vicinity of the bulb. These free electrons ionize the gas atoms by impact ionization. Ionization yields more free electrons and ions and the process resumes, eventually resulting in plasma formation. Excited states of gas atoms produce light via spectral transitions across various wavelengths. In case of IL, mercury is present in addition to argon/krypton to produce UV from excited mercury atoms. The UV excites the phosphor coating on the inner surface of the bulb envelope and emits white light just like a regular fluorescent lamp.

In ESL, microwaves are used to produce an intense plasma inside a rotating quartz ball containing argon/krypton and sulfur. The rotation of the quartz ball helps in stabilizing the fill for uniform emission as well as convective cooling with a fan to prevent its meltdown. Initially, the microwaves create a high-pressure (several atmospheres) noble gas plasma. This heats sulfur to a high temperature resulting in brightly emitting plasma. Light emission by sulfur plasma is due to the molecular emission spectra of the sulfur dimer molecules (S_2). The spectrum is continuous across the visible and has >70 percent of its emission in the visible. It peaks at 510 nm, giving a greenish hue. The resultant CRI is 79 at a CCT 6000 K. The lamp spectrum can be modified with additives such as calcium bromide, lithium iodine, or sodium iodide or by using an external color filter.

Electroluminescent Sources Electroluminescent sources are materials that emit light in response to an electric field. Examples of such materials include powdered ZnS doped with copper or silver, thin film ZnS doped with manganese, natural blue diamond (pure diamond with boron as a dopant), III-V semiconductors or inorganic LED materials such as AlGaAs, phosphors coated on a capacitor plane and powered by pulsating current, organic LED (OLED) also known as light-emitting polymer or organic electroluminescent. The sources can operate at low electrical power with simple circuitry.

Electroluminescent sources are commonly used for providing illumination across small regions such as indicator panels. LEDs are already a major lighting source that is rapidly replacing incandescent and fluorescent light sources. The remainder of this section discusses LEDs and OLEDs.

LEDs emit light by electron-hole pair recombination across the P-N junction of a diode. The wavelength of the emitted light corresponds to the band gap (energy gap between valence and conduction bands) across which the electron hole pair is created. The degeneracy in the valence and conduction bands leads to a closely spaced band of wavelengths that constitute light from the LED. Narrow spectral bandwidth enables applications that require saturated colors. Although LEDs are not available for every desired color, it is possible to combine LEDs of different colors to create any color within the color gamut defined by these LEDs. As such, color mixing has become an important field. LEDs emitting in specific bandwidths can be combined with sources with continuous spectra to either enhance a certain spectral region or to provide an easy dynamic color control. One easy method of combining multiple LEDs is using lightguides. Lightguides with rippled surface texture along the cross section are particularly efficient in combining multiple colors with excellent uniformity in a short path length.^{49,50}

LED lamps may have an array of small LED chips or a single large chip to achieve the desired power levels. The directionality is controlled by appropriately mounted optics. DC operation of LEDs makes them flicker free sources. Figure 18 shows the structure of a simple packaged LED. Chapter 17 in this volume is dedicated to the subject of LEDs.

LED packages come in various sizes and shapes. Surface mount LEDs (SMD) have minimum packaging and are almost a bare die. LED packages are also offered in multicolored die formats.

Over the years, a variety of materials for LEDs have been used with the goal of obtaining higher efficiencies and different colors across the visible spectrum. LED technology is fast evolving with ever increasing brightness, lifetime, colors, and materials and decreasing costs. The available materials, at the time this chapter was written, are listed in Table 5.

Most lighting applications require white light. Some of the processes for making white-light LEDs are listed below:

- Arrays of small red, green, and blue dies placed in close proximity in a single LED package. Good color mixing takes place in angular space.
- Color-mixing of red, green, and blue colors using lightguide or other optical means.^{50,52}

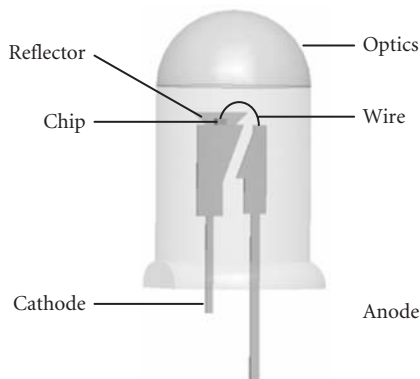


FIGURE 18 Structure of a simple packaged LED.

TABLE 5 LED Materials and Emitted Colors

Material	Color
AlGaAs	red, infra
AlGaP	green
AlGaInP	higher brightness orange, orange-red, yellow, and green
GaAsP	orange, orange-red, orange, and yellow
GaN	green and blue
InGaN	blue (450–470 nm), near UV, bluish-green, and blue
SiC (as substrate)	blue
Si (as substrate)	blue
Sapphire (as substrate)	blue
ZnSe	blue
Diamond	UV
AlN, AlGaN, AlGaInN	UV(<210 nm) ⁵¹
Organic light-emitting diodes (OLED)	Red, green, and blue

- Phosphor excitation by blue or UVLEDs.
- Novel techniques like quantum dot blue LEDs or homoepitaxially grown ZnSe blue LEDs on a ZnSe substrate. The active region emits blue light while the substrate emits yellow light.^{53,54}

There is a continuous push to improve the LED efficiency and brightness while keeping the lifetimes high. High brightness LEDs became possible due to large area chips, efficient heat extraction and better light extraction from the chip. Internal quantum efficiency of LEDs can be increased by placing emitters inside a cavity⁵⁵ to increase the radiative recombination rate. Due to the high internal Fresnel reflections and lateral waveguiding, a lot of light fails to exit the chip. Techniques such as texturing the surface with photonic crystals assist in increasing the light extraction from large dies.^{56–58} Figure 19 shows the internal structure of a photonic crystal LED.

Organic LEDs (OLEDs)^{59,60} in contrast to inorganic LEDs are size-scalable light sources with richer color spectra. OLEDs can be used to create flexible transparent lighting solutions as they can be printed on a malleable substrate with transparent electrodes. Currently OLEDs are being used for displays and are competing with LCD flat panels. Conceptually, OLEDs are no different from inorganic

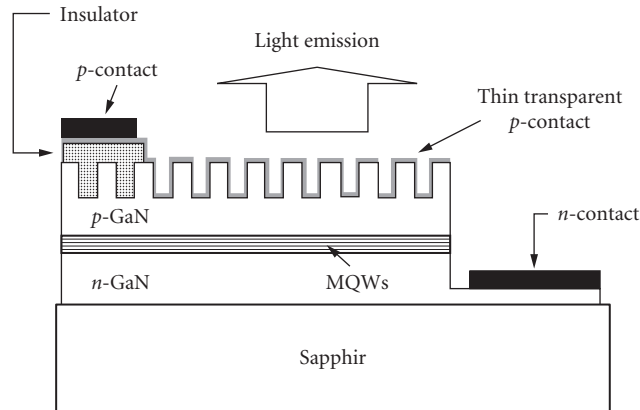


FIGURE 19 Internal Structure of a Photonic Crystal LED. MQW refers to “multiple quantum wells.” (Courtesy: Seoul National University, Korea, http://optics.org/cws/article/research/23635/1/sem_image.)

semiconductor based LEDs. An OLED deploys layers of organic materials on polymer substrates to form conductive and emissive layers connected to a cathode and anode respectively. Much of the OLED research is aimed at making them brighter and longer lasting.

LED failure causes include damage due to degradation of the active layers with time (spontaneously or in operation); plastic package degradation due to ambient UV; electrostatic discharge; current crowding or inhomogeneous current distribution across the junction leading to hot spots; and thermal stresses causing rupture of the LED package, diffusing of the metal contact material into the die material at high currents, high output leading to facet melting and phosphor degradation in white LEDs, and degradation of organic layers in OLEDs.

Miscellaneous Artificial Light Sources *Neon signs* are essentially cold cathode-like operation of a fluorescent tube without phosphors. A low-pressure mixture of noble gases such as neon, argon, helium, xenon and a small amount of mercury is used in the discharge tube. Neon emits a reddish-orange color; argon emits blue; and krypton, helium, and xenon emit over a wider spectrum. Colored filter glass can be used for making different colors.

Short arc sources function almost identical to HID lamps with the only difference being that the electrodes are much closer (less than 1 mm to 12 mm). The gases are mercury (with argon), mercury-xenon, pure xenon, or metal halides (with mercury and argon). These lamps are primarily used in illuminating high loss systems where the source étendue needs to be as small as possible. Projectors, medical optical instruments, metrology instruments, and daylight or solar simulators use such lamps. These lamps have lifetimes from a few hundred hours up to 10,000 hours. The light sources are typically used with a reflector (parabolic or ellipsoidal). Sometimes the reflector is an integral part of the source/lamp package.

Pure *Xe arc lamps* have an instant-on capability and provide high CRI (>80) at high CCT (>6000 K). Digital cinema projectors and flash tubes (for warning signs, entertainment applications, camera flash lights, and warning or emergency signs and indicators) often deploy pure Xe-arc sources.

Lasers are used for visual displays for entertainment. The subject of Lasers is discussed in Chap. 16 in this volume.

Nuclear sources are self luminous light sources that function by phosphor excitation caused by beta radiation from radioactive materials such as tritium. These light sources are used to illuminate tiny spaces such as watches or displays of instrument panels in very low ambient light.

Glow lamps are low wattage arc sources with gases such as argon emitting in the UV to excite UV-excitable materials or neon to emit orange light to be used as indicator lights.

Carbon arc sources are now obsolete but still find applications in illuminating small areas with bright light under demanding environmental conditions (such as outer space). An arc is struck across a pair of carbon rods and the incandescence from the heated carbon rods provides the light.

Gas lights that operate by the burning of gases like methane, ethylene or hydrogen are used with appropriate lanterns primarily for decorative applications.

Natural Sources: Daylight Daylight can be utilized in the lighting design of buildings to provide a pleasing environment that enhances physiological well-being and productivity and also energy savings during the day by reducing the need for artificial lighting and solar influx contribution to building heating. Daylight is primarily used for ambient lighting. It can be used for task lighting when integrated with electrical lighting. Daylight constitutes direct sunlight, scattered sunlight from the atmosphere, reflected sunlight from the clouds, and reflected light from the surroundings such as ground (especially snow) and objects such as buildings. The solar spectrum changes with atmospheric conditions and so does scattered light from the sky or reflected light from the ground. The CCT of Sun is 1000 to 5500 K, clear blue sky is 10000 to 100000 K, overcast sky is 4500 to 7000 K and clear sky with sunlight is 5000 to 7000 K. Figure 20 shows an example of the solar spectrum at the ground, at noon, at Golden, Colorado. Note the IR content in the spectrum and the shift in spectrum from noon to evening.

Designing for daylight requires close attention to many factors:

- Goals of providing daylight: physiological well being of occupants and/or energy savings.
- Intended distribution of light inside the building during the day and in the night. Penetration of daylight into the interiors and the impact of reflectivity from various surfaces.
- Impact of daylight on materials such as wall paints, artwork, plastic materials, furniture, and plants. The UV component of daylight is generally harmful to most materials via solarization of plastics or fading of paints and stains. If the impact of UV is not known and acceptable, then the UV should be rejected by the optical system components through coatings or materials.
- Integration of daylighting-based building design with other controls such as for electrical lighting, cooling and any automated systems.

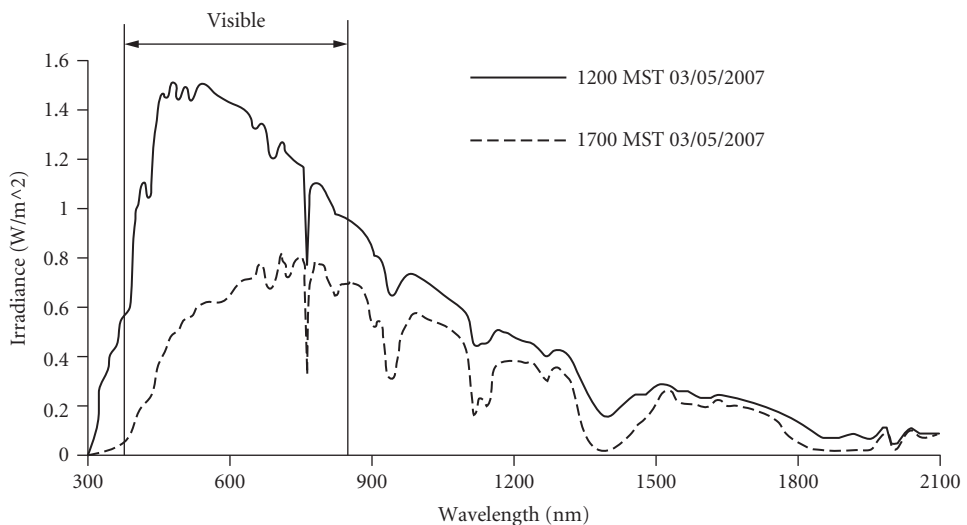


FIGURE 20 Solar spectrum recorded at ground at Golden, Colorado at two different times during the day. (Courtesy of National Renewable Energy Laboratory, USA.)

- Location (latitude and longitude) and orientation of the building relative to the surroundings. Simulating the daylight entering the building during various times of the day and the year.
- Design for reduction or elimination of glare from daylight sources. Internal layout of the building during intended use will play a key role.
- Outside view requirements via windows. A minimum window size^{61–63} is needed for a given geographic location, and a certain window location is dictated by aesthetic reasons including the desired view from the window. The minimum window size can also depend upon the prevalent building regulations based on wall size or floor size. Under such circumstances, the daylight entering must be controlled or balanced with interior lighting and cooling for uniformity, glare, and heat management. Daylight can be controlled by providing appropriate window shades, glazings, or even additional windows or skylights at appropriate places.

It is not always possible to obtain energy savings with daylight due to increased heat load, especially during summers and the need to supplement daylight with electrical lighting during the night. Infrared rejection is achieved by using special glazings. To limit glare, glazing such as high reflectance, low transmittance, electrically controlled transmission, tilted glazings with angularly dependent transmittance or IR reflecting films are used. Adequate consideration must be given to the window view impact due to any of the daylight control mechanisms. A very low-transmittance glazing can make the outside view appear gloomy on a bright day.

Luminaire Design

The source is the starting point of the luminaire design. The optics performs at a minimum two functions: first to capture the light from source and second to transfer light efficiently to the desired distribution at the target (i.e., illuminance, intensity, and/or luminance). The choice of the optics is also a challenging process. First, the designer must select if refractive, reflective, or both types of optics (i.e., hybrid) are to be used in the design. Reflective optics have been a standard for most sources, with flat refractive optics used at the output aperture in order to protect the other optics. Until recently reflectors were standard conics such as ellipses, hyperbolas, parabolas, or spheres. Standard refractors in use include cover plates (still called a lens), pillow lens arrays, Fresnel lenses, and other faceted designs. However, with advances in LEDs, refractive optics are increasingly being used for advanced function. Solid-state lighting optics are either plastic or glass that surround individual or a limited number of LEDs. They are hybrid optics that use at least total internal reflection (TIR) and refraction at the input and output facets, but they can also use reflection, scattering, and even diffraction. With the advent of better technology, especially manufacturing capabilities and software for modeling, faceted and continuously smooth surfaces parameterized with nonuniform rational B-splines (NURBS) are being developed. For high-performance, injection molding of plastic or glass is being increasingly used. For refractive optics, the surfaces are left bare, but for reflective optics there is vacuum metallization of the surfaces. Reflectors are also made with stamping methods, but these optics tend to have lower performance.

Baffles (louvers) are often considered by many to be optics of the system; however, many, rather than shaping a distribution through redirection of the light, block or even absorb incident radiation. Thus, while some baffles are reflective and provide some shaping of the illumination distribution, they most often achieve shaping through subtraction rather than addition. They are primarily added to alleviate glare, trespass, and pollution, but they are also included for aesthetic reasons, correcting errors in the design process, and to hide structure within the luminaire.

The source coupling aspect is the focus of the next section, following that is a discussion of the design of the optics and subsections on baffling in the form of luminaire cutoff classification and types of luminaires. Following this process from the source to the optics to the baffle to the target, while always considering the perception aspects of the illumination (see Sec. 40.3), means that aesthetically pleasing while technically sound lighting can be developed.

Étendue and Source Coupling For reflective optics, the source-coupling components also act as the transfer optics, often in conjunction with a front, protective lens. With the advances in source

technology that use hybrid, dielectric components, coupling of the sources, is increasingly important in order to improve upon system efficiency. Typically, individual LEDs are placed in recesses in the dielectric optic, and by obeying the conditions for TIR, all of the emitted light can be captured by the coupler and then transferred to following optics that shape the emitted distribution.

In all cases, the term étendue describes the flux transfer characteristics of the optics, starting with the coupling optics, of an optical system, such as a luminaire. Étendue is a geometrical quantity that is the integrated product of an area and solid angle. In paraxial form it is the Lagrange Invariant, but in nonparaxial form it is given by

$$\mathcal{E} = n^2 \iint_{\text{pupil}} \cos \theta dA d\Omega \quad (12)$$

where \mathcal{E} is the étendue, dA is the differential area, $d\Omega$ is the differential solid angle, θ is the angle with respect to the surface of interest (i.e., normal), and n is the index of refraction of the space. The limits of integration in area are over some aperture (e.g., a lens clear aperture or a reflector exit aperture), while the solid angle integration is over the limits that are passed by the aperture. For example, consider a source of area A_s that emits into a half angle of θ_0 from every point on the surface. The étendue for this source is

$$\mathcal{E} = n^2 A_s \int_0^{2\pi} \int_0^{\theta_0} \cos \theta \sin \theta d\theta d\phi = \pi n^2 A_s \sin^2 \theta_0 \quad (13)$$

In lossless optical systems, étendue is conserved. Thus, in order to design the most efficient luminaire, one must continue to match the étendue as one progresses through the optical components of the system. For example, if the source of Eq. (13) is used for a luminaire, one must keep this étendue quantity consistent. If one desires to reduce the angular spread of the output from an optic ($\theta_0 > \theta_{\text{optic}}$), then the area of the optic must be increased ($A_s < A_{\text{optic}}$). The counter also holds true: to reduce the real extent ($A_s > A_{\text{optic}}$), one must increase the angle ($\theta_0 < \theta_{\text{optic}}$). An expression for conservation of étendue in a generalized form is

$$dx dy dp dq = dx' dy' dp' dq' \quad (14)$$

where the dx and dy terms are the differential position terms and the dp and dq terms are the differential optical direction cosine terms, which are equivalent to ndL and ndM respectively. More information about étendue can be found in Chap. 39, “Nonimaging Optics: Concentration and Illumination,” in this volume. Another factor related to étendue through a differential is skewness, which denotes the twist on individual rays of light in an optical system. Skewness is also invariant and implies that transfer from one source geometry (e.g., a square) cannot be transformed to another source geometry (e.g., a circle) without loss except if some rotational asymmetry is added to the optical system. Further information about étendue and associated terms like skewness can be found in the literature.⁶⁴

Luminaire Design Methods There are a multitude of design principles for the design of the optics of a luminaire. Fundamentally, most design methods are based on the basic conic shapes as listed in Table 6. Each of these shapes provides a basic intensity distribution at its output aperture. However, increasing demands of tailored light distributions and also increased efficiency require perturbations to these basic design forms. Furthermore, the topics of light trespass, light pollution, and glare are receiving a wealth of attention from ordinance and regulatory agencies. To reduce glare issues it is best to use diffuse optics with a well-defined cutoff. In the field of nonimaging optics (see Chap. 39), the edge-ray theorem provides a means to have a well delineated cutoff. The edge ray is defined by the maximum extent of the source, thus providing a maximum cone of light from the reflector designed around the source shape. However, the edge-ray principle is passive with respect to the luminance distribution of the source—it contends for the maximum extents but not the physical distribution of light in the radiation pattern.

Thus, tailoring methods have been developed. The tailoring methods specify the shape of the optics based upon the luminance distribution of the source and the desired illumination pattern

TABLE 6 Basic Conic Shapes, Their Conic Provide at Their Output Aperture Constant, and the Basic Intensity Distribution That They Provide at Their Output Aperture

Shape	Conic Constant (k)	Base Intensity Distribution
Hyperbolic	$k < -1$	Diverging; far-field applications
Parabolic	$k = -1$	Collimating; pseudocollimation applications
Elliptic	$-1 < k < 0$	Converging; near-field applications
Spheric	$k = 0$	Converging; self-imaging applications

(i.e., luminance, illuminance, and/or intensity) at the target. These methods are extensive and beyond the confines of this chapter. The reader is encouraged to consult the *Handbook* chapter on nonimaging optics (Chap. 39) for a brief introduction, the theoretical book on nonimaging optics by Winston, Miñano, and Benítez,⁶⁴ and the applied book on nonimaging optics by Chaves.⁶⁵

Luminaire Cutoff Classification A cutoff ensures that light from the luminaire is restricted above the horizon with respect to the lamp geometry. Cutoff is designed into the luminaire through the optics (i.e., edge-ray designs) and/or the integration of baffles. While most applications do not require cutoff classification, except for those used on the exterior, such as automotive, roadway, and landscape lighting, most designers include such to make effective lighting systems by alleviating potential glare, trespass, and pollution concerns. Often strict cutoff guidelines are mandated by governmental standards, such as for automotive, traffic signal, and roadway lighting. The goals are to provide the required lighting level to its users, while also alleviating light pollution and light trespass. Light pollution is light that is directed up into the atmosphere, causing sky glow, which is especially present in urban settings. The reduction of light pollution is a growing trend being addressed by the astronomy community. When light is incident on surfaces outside the intended illumination region, it is called *light trespass*. The impact of light trespass from roadway lighting is a major concern in residential areas.

For both trespass and pollution the luminaire cutoffs provide a protocol to reduce both. Automotive lighting does not use the criteria presented here, but rather use a set of governmental standards. Roadway and external lighting make the most use cutoff criteria as shown in Table 7. See Fig. 22 for a depiction of the angles listed in Table 7.⁶⁶

Luminaire Classification System In 2007 the IESNA published research results for refinement of the cutoff classification system of the previous section.⁶⁷ This study focused on light distribution in front of the luminaire (forward light), behind the luminaire (back light), and above the luminaire (uplight) as shown in Fig. 21. They found the photometric luminaire efficiency ($\eta_{\text{luminaire}}$) to be

$$\eta_{\text{luminaire}} = 100 \frac{\Phi_{\text{forward}} + \Phi_{\text{back}} + \Phi_{\text{uplight}}}{\Phi_{\text{source}}} \tag{15}$$

where Φ_{forward} , Φ_{back} , Φ_{uplight} , and Φ_{source} are the integrated fluxes in lumens over the solid angles shown in Fig. 21 for forward light, back light, uplight, and the bare source, respectively.

TABLE 7 Amount of Emitted Light Criteria for the Luminaire Cutoff Classification⁶⁶

Type	Horizon and above (90° and greater)	10° below Horizon (80° to 90°)	Remainder (0° to 80°)
Full Cutoff	0%	≤10%	≥90%
Cutoff	<2.5%	≤10%	≥87.5%
Semicutoff	<5%	≤20%	≥75%
Noncutoff	No restrictions over entire angular space		

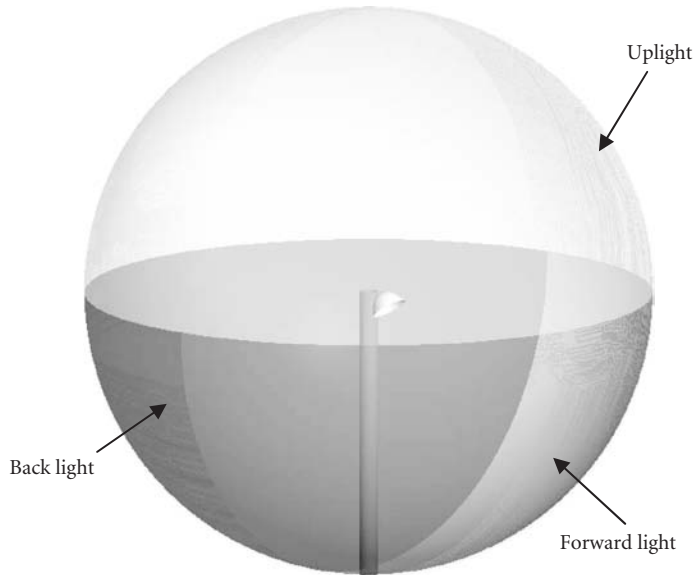


FIGURE 21 Lighting classification system zones for forward light, back light, and uplight based upon the exit aperture of the luminaire.

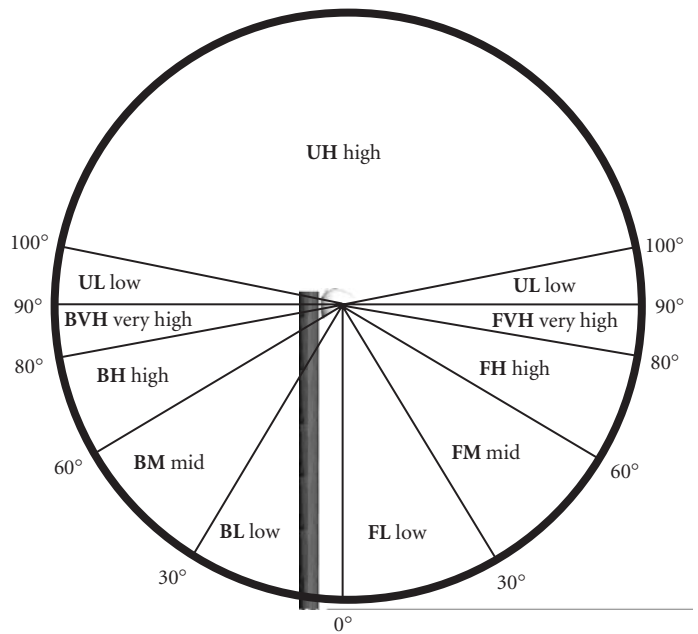


FIGURE 22 Layout of the light classification system subzones. (See also color insert.)

These three zones are then broken down into a total of 10 subzones as shown in Fig. 22. The lumens in each subzone is measured with respect to the bare source flux, as in Eq. (15), in each of these subzones and then reported as an evaluation of the luminaire. These subzones indicate the distribution of light over several regions rather than restricted to the 80° and greater as per the previous section. Standard goals include reduction of light above the horizon (i.e., all uplight subzones and the BVH and FVH subzones) and the desired uniformity over the other zones. This new classification system provides for better control of the illumination such that both vertical and horizontal surface illuminances can be addressed in the design process. Horizontal surface illuminance criteria are met by increasing the flux in the BM to FM subzones. Vertical surface illuminance criteria are met by increasing the flux in the BM, BH, FM, and FH subzones.

Luminaire Optics There are innumerable schemes to the design and availability of luminaire optics, both for artificial and natural sources. Methods employing reflectors, refractors, TIR optics, and combinations thereof have been developed. Additionally, the optics are both specular and diffuse or a combination of these two are used. In the next few subsections, we provide examples of commonly available luminaire optics. Finally, as per Table 6, the shapes of the optics are typically based on conic shapes.

Luminaire optics for artificial sources While the design of the optics of the luminaire is to provide a desired illumination distribution, we learned from the previous sections that light cutoff is important in the design of the luminaire. In fact, the ability to hide or shield the source from direct view is as equally important as obtaining the desired illumination distribution. A reflector as shown in Fig. 23a has a fairly wide direct view of the source, denoted as the shielding angle or similarly the cutoff angle. Room lamps are perfect examples of this dual requirement since it is typical to use a diffuse shade around the source. The source and shade (which also acts as a diffuse reflector) provide the desired general illumination, while the shade provides the requisite cutoff. Besides using the body of the reflector to hide the direct view of the source, louvers or baffles (see Fig. 23b), a prismatic or Fresnel lens (see Fig. 23c), or a bulb shield (see Fig. 23d) are

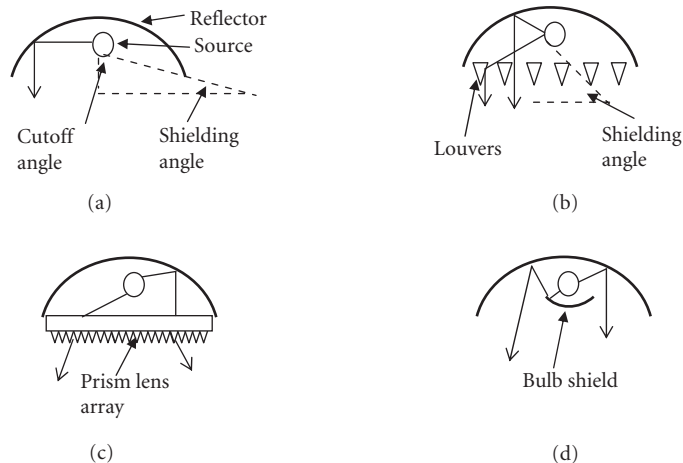


FIGURE 23 Glare issues of four luminaire geometries: (a) the reflector providing the cutoff; (b) louvers, flat or parabolic, increasing the shielding angle; (c) a lens array, prismatic or Fresnel, which directs more of the light emission downward thus reducing the direct intensity level of the source; and (d) the inclusion of a bulb shield to ensure that all emission strikes the primary reflector at least once. (Adapted from Ref. 5, p. 71.)

added to the luminaire. For the baffle case two standard options are typically used: vertical, strongly absorbing louvers, or specular reflecting parabolic louvers. The former inhibits direct view of the source in the luminaire while also absorbing the glare-inducing light. The latter also reduces the direct view angle of the source, but it uses the parabolic, specular louvers to direct the glare-producing light into directed radiation typically outside the direct view of an observer. The lens, prismatic or Fresnel, coupled to a reflector causes most of the emission to be directed downward, thus frustrating the direct imaging of the source. This obscuration of the source means glare effects are reduced. The bulb shield, which is typically spherical, ensures that all light from the source is incident on the primary reflector of the luminaire, thus completely hiding the source from direct view and in turn greatly reducing the glare potential.

The specifics of the lamp design depend on the application. Figure 24 shows some representative luminaire designs. Figure 24*a* shows a banker's desk lamp, which has multiple bounces within the glass envelope optic. It creates a wide illumination distribution, but the colored glass and the multiple bounces softens the appearance of the bulb located within. The envelope optic can be positioned by the user to alleviate source glare while also providing the desired illumination over a desk surface. Figure 24*b* shows a Bouillotte table lamp, which uses two vertically oriented fluorescent lamps. The shade hides part of the bulbs from direct view, and the coating on the fluorescent tubes causes near-Lambertian emission; therefore, the illumination for this lamp is wide and glare issues are minimal. Figure 24*c* shows an indirect RLM lamp fixture that is suspended from a ceiling. The indirect nature

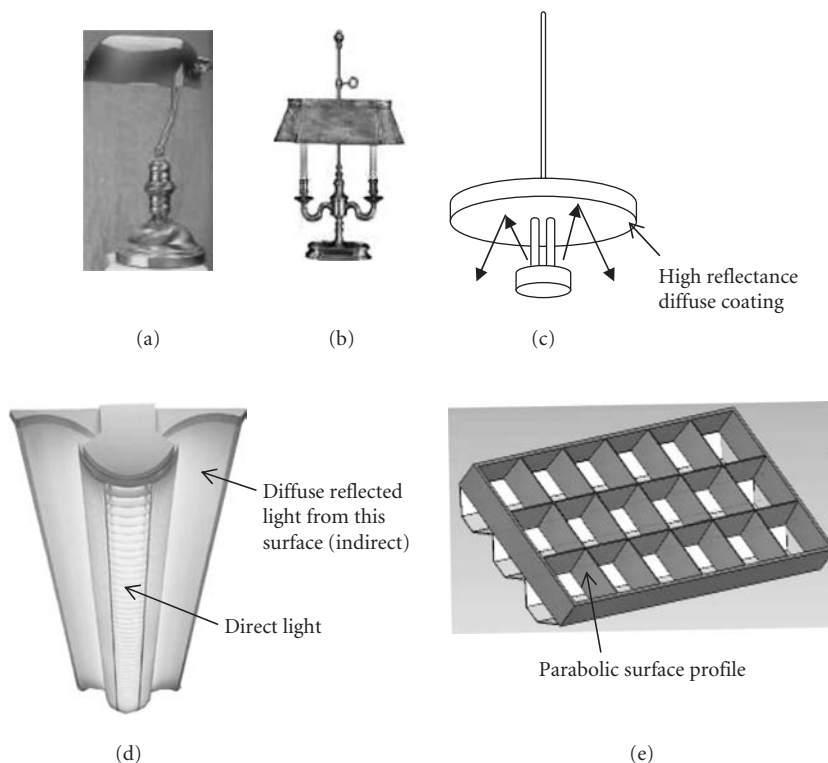


FIGURE 24 Depictions of luminaires: (a) Banker's lamp: multiple bounces inside the reflector create a wide angled uniform illumination; (b) Bouillotte lamp: vertical fluorescent tubes provide diffuse illumination; (c) indirect lighting with RLM fixture where the top surface reflects light into a wide angular range; (d) overhead direct-indirect lighting fixture using fluorescent tubular bulbs; and (e) parabolic louvered trough reflector for fluorescent tubes. (See also color insert.)

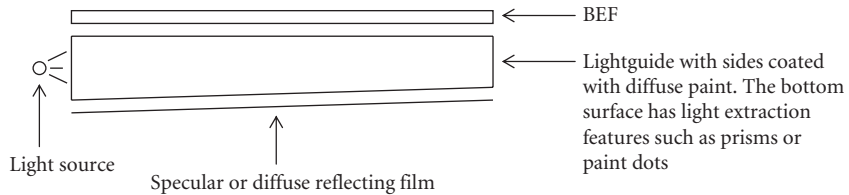


FIGURE 25 An edge lit backlight. (BEF: Brightness enhancing films.)

is due to the inclusion of a bulb shield and a highly reflective diffuse primary reflector thus alleviating glare. Figure 24*d* depicts a direct-indirect overhead, fluorescent lamp fixture. The white, diffuse wing structures provide indirect lighting, while the central section is louvered direct lighting. The direct lighting provides task illumination needs, while the indirect lighting provides a general light level for the room. Finally, Fig. 24*e* shows a standard overhead office luminaire: a series of parabolic troughs (called a troffer) in which long, tubular fluorescent lamps are located. The parabolic louvers reduce glare concerns by increasing the shielding angle, while also increasing the task illumination due to the specular reflectivity of the vanes. Instead of using the louvers, a pillow lens array or prismatic lens can be used; however, computer monitor glare issues can arise with these lenses. Other luminaire geometries not presented here include track lighting, recessed lights, chandeliers, and spot lamps. Please see Refs. 71 and 5 for more information.

Backlighting Backlighting is used extensively in photography to separate the background from the subject and create 3D effects. It can be used for the same purpose in interior lighting to illuminate displays from behind. Backlighting is used extensively in signage, to illuminate instrument panels and device display panels such as laptop and cellphone screens. A typical edge lit backlight operation is illustrated in Fig. 25. Light enters the planar lightguide⁶⁸ (for example, 50 mm × 50 mm × 5 mm) from the thin edges and bounces around. Carefully designed and positioned light extraction features such as prism or spheres deflect the light out of the lightguide. As a result the entire planar surface is lit and appears as a planar light source. Modern backlights use sources such as CCFLs and LEDs. With LEDs, dynamic multicolored backlit displays are possible as lightguides allow for efficient color mixing. Edge lit backlighting can be used in direct lighting of large spaces such as living rooms by creating illuminated ceilings, walls, or artificial windows. Figure 26 shows one such application where backlit ceiling tiles are used to create an illuminated ceiling or an artificial skylight. A picture of sky and vegetation is superimposed on the tiles to create an effect of natural sky with vegetation on the roof. Backlights could be replaced by OLEDs which provide not only illumination but also information content.

Luminaire optics for daylight sources Daylighting schemes^{69,70} involve careful layout of windows, skylights, skytubes, and controls such as shades, window overhangs, window depths, light shelves, and hybrids with electrical lighting. For effective daylight illumination, it is necessary to determine the access to sunlight by taking into account the sun path across the sky during the day and across different months and the impact of neighboring buildings and ground features. Daylight can be used for city lighting too with careful planning. Heliostats or large plane mirrors have been used atop buildings or even mountains to direct sunlight into the city interior.

Layout of windows and skylights can utilize schemes such as side lighting, top-level lighting, and clerestory lighting. Placement of nonview providing windows for providing daylight must be designed carefully for maximum daylight penetration and controlled glare. Side lighting allows daylight in from the walls usually at eye level, top lighting allows daylight in from skylights, and clerestory lighting allows daylight in from the side windows near the roof above the eye level. Clerestory windows provide more uniform ambient illumination over a larger region. However, proper attention must be given to glare from the sky or direct sun by using baffles. The depth of windows near the ceilings or window overhangs also helps in limiting such glare. Figure 27 shows several different



FIGURE 26 A conference room with artificial skylight made up of backlit ceiling image tiles. (See also color insert.) (Courtesy of *The Sky Factory, LLC.*)

layouts of windows and skylights that bring daylight into the interiors.^{71,72} The location and number of openings for daylight determine the penetration and uniformity of the illumination achieved. The height and the slope of the ceilings determine the penetration and the illumination gradient achieved inside the room. For example, for daylight from windows located near the ground, having a tall ceiling allows light to penetrate to the farthest ends of the room. Similarly a sloping ceiling with windows located near the ground allows a gentler illumination gradient from the window to interior.

The reflective properties of the walls and the ceiling must be taken into consideration when simulating the impact of daylight. High reflectance paints on the ceiling can be used to spread the daylight entering from the window portion near the ceiling into the building interiors. Window blinds, shades, and mechanical louvers are commonly used to control daylight. Sometimes partition walls around certain window sections are used to control glare or even limit the extent of illuminated region such as artwork in museums.

Light shelves are used interior or exterior to the windows to allow daylight from the sky in without glare. Light shelves consist of large horizontal sections hanging below the top edge of the window. A layout of mirrors or even just a plane aluminum sheet directs the daylight from the sky toward the ceiling or deep into the room without hitting the ground. A reflective ceiling will scatter daylight into the room. Combinations and variations of these schemes can be used to suit a particular situation. Figure 28 illustrates the concept of light shelves. Suncatchers are similar in concept as shown in Fig. 29.

Skytubes are lightguides that carry daylight into the building interiors. The inside walls are specular with high reflectance. These tubes may be straight or curved, as needed, to transport daylight to different regions. Curved lightguides have a symmetric cross section to maximize light transfer efficacy across the bends.

Hybrid daylighting systems consist of daylight integrated with electrical light as shown in Fig. 30.

Solar lighting systems include the use of tracked or untracked mirrors, lenses or apertures for collecting sunlight and channeling it into the building interior via skytubes, lightguides or fiber-optical cables. Heliostats or large plane mirrors are used at locations with access to the Sun and they direct that light into building interiors or even city interiors. Optimally, these mirrors track the Sun.

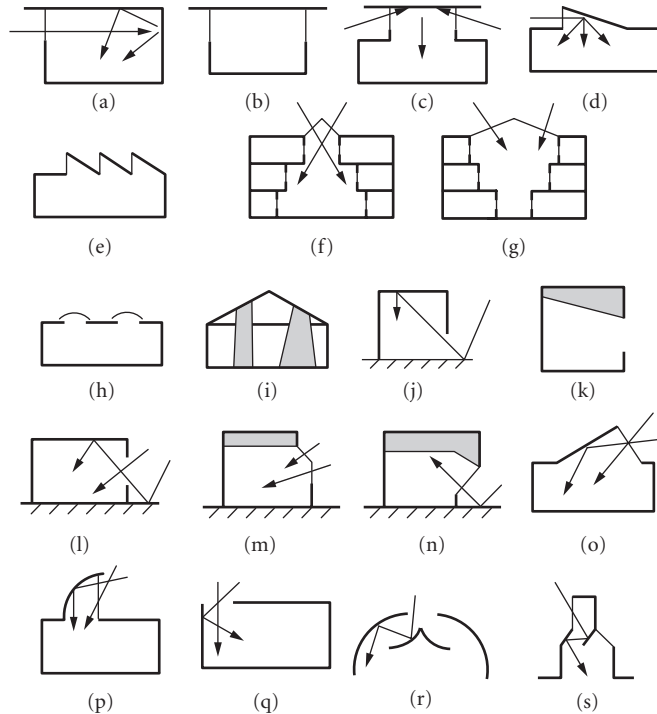


FIGURE 27 Windows for gathering daylight: (a) Unilateral, side lighting. (b) Bilateral, side lighting (c) Roof monitoring. (d) Clerestory. (e) Sawtooth skylighting. (f) Atrium. (g) Litrium. As opposed to atrium it provides best light distribution to adjacent spaces. (h) Top skylighting. (i) Skywells: straight and splayed sections. Shaded region shows the lit region. Splayed section distributes light farther, more evenly and with reduced peak brightness. (j) Window near the ground. Deep ceiling allow deeper light penetration. (k) Tilted ceiling reduced the illumination gradient from the window to interior. (l) Window in the mid-section allows direct skylight and ground reflected skylight. (m) “Greenhouse” opening for overcast sky. (n) “Overbite” opening for ground reflected skylight. (o) Tilted glazing with clerestory opening to allow sunlight from a wider range of elevations. (p) Capturing daylight via top lighting. (q) Using high reflectance wall adjacent to top lighting to provide glare free light. (r) and (s) Top lighting via lightguiding.

Figure 31 shows use of a solar light pipe (SLP) to bring sunlight deep into the building interior. The atrium is 140 ft deep, 60 ft long, and 9 ft across. Without the SLP, the view in the building interior was a dark concrete wall. A rooftop heliostat captures sunlight and directs in down into a prismatic glass cone. The glass refracts the incoming sunlight horizontally onto an outboard cylinder of open weave fabric; creating a glowing, translucent, 120-ft-long tube of diffuse sunlight. This display is visible from the 14 floors of atrium offices and also from the ground floor lobby, elevator lobbies, and the library. The SLP projects a 10-in diameter sunburst onto the lobby floor (Fig. 31*b*). Besides injecting daylight into dark spaces, the SLP’s unique design provides a compelling and a very dynamic visual focus for the atrium occupants. It constantly updates their understanding of the Sun, the sky, and the weather patterns and reconnects them to the otherwise solar rhythms of the day and seasons. At night, powerful searchlights use the “at rest” heliostat and SLP to inject a shifting palette of colored light into atrium.

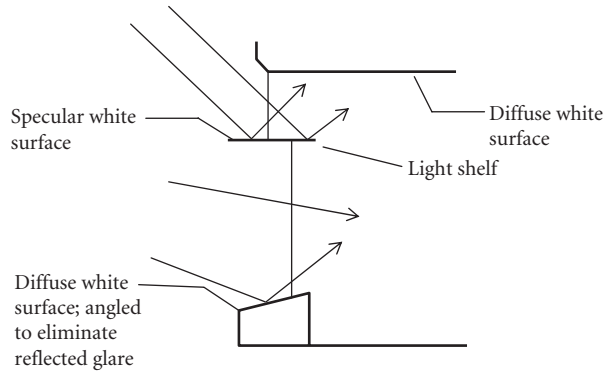


FIGURE 28 Light shelf to limit glare from direct skylight and redirect light to the ceiling.⁷² The light shelf can be curved in shape and moveable angularly. The position of the light shelf relative to the window can be interior, exterior, or both. Exterior light shelf provides shading while interior light shelf limits glare. Blinds or moveable shades can further help in glare control.

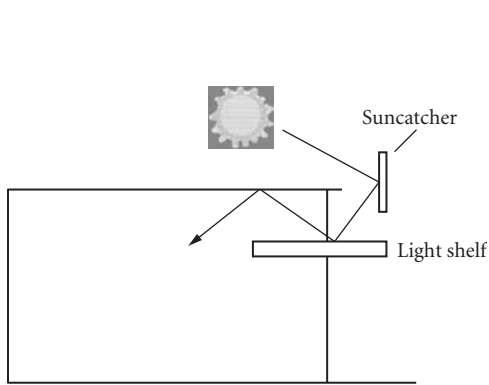


FIGURE 29 The concept of suncatcher. The surfaces of suncatcher and light shelf are of high reflectance specular or diffuse finish. Light shelf is optional. Suncatcher reduces the view, eliminates direct glare and increased illumination when the sun is at a particular location.

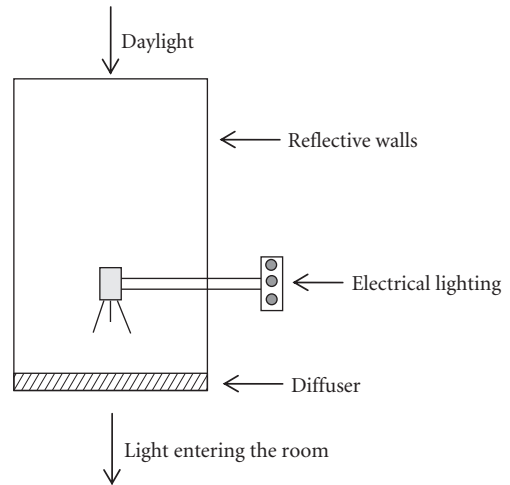


FIGURE 30 A skytube or skywell with integrated Daylight and electrical light.

The use of modern CAD software, especially radiosity based, makes it easy simulate the daylight (that includes sun-tracking, scattered and reflected light from the sky, ground, and neighboring buildings throughout the day and the year), its interaction with optical components such as optical fibers, lightguides, lenses, and diffusers; its penetration through the windows into the interior; impact of reflection from walls, ceilings or furniture, and interaction with interior lighting.

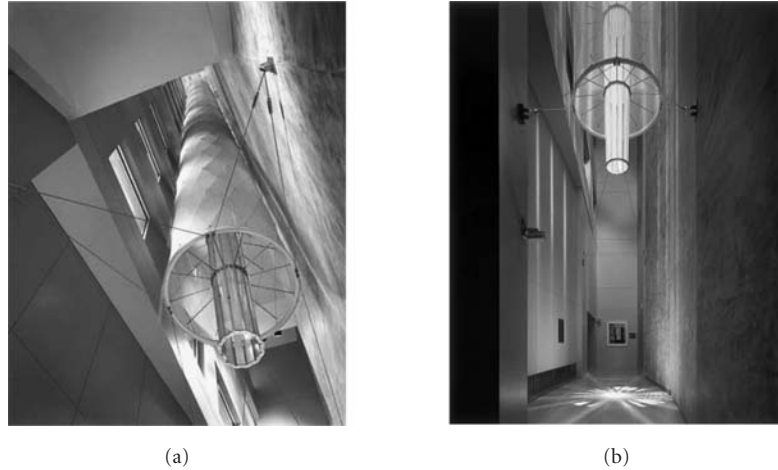


FIGURE 31 A Solar light pipe. (a) A 140-ft-tall light gathering and distributing device that presents daylight down into the core of a building that has no other access to daylight. (b) Light projected (10-in diameter) on the floor. (See also color insert.) (Photograph by Paul Warchol; Courtesy of Carpenter Norris Consulting.)

40.6 LIGHTING MEASUREMENTS

In this section, we discuss briefly the tools and techniques used for measuring light, 3D object profiles and optical properties of objects, in the context of lighting applications. The quantities most relevant for light measurement are horizontal/vertical illuminance and luminance of lamps and lit objects as a function of wavelength. To measure the optical properties of objects, luminance measurements in transmission and reflection as a function of wavelength, magnitude, angle of viewing, and direction and angle of incident light are needed. To model the objects in CAD software, measurement of 3D object profiles is needed. These measurements help in picking the lamps and their placement geometry to achieve desired lighting goals.

For measuring light from the source, we discuss the following instruments: illuminance meters and goniometers. For measuring the optical properties of objects, including luminaire optics, in reflection or in transmission, we discuss instruments such as reflectometers and luminance meters. For the measurement of luminaires conformance to the design, we present two methods: CMM and laser scanning.

Illuminance Meters

Illuminance meters or luxmeters measure illuminance in lumens/area. These are typically handheld devices, as shown in Fig. 32 that consist of a photodiode with a photopic correction filter and a cosine corrector on top of it. The photopic correction filter multiplies the incident light spectrum by the eye's photopic response to convert the incident energy in watts into lumens. The cosine corrector consists of a diffuser such as a plastic disk or flashed opal disk or a small integrating sphere with a knife edge entrance port. The goal of the diffuser is to provide a constant relative distribution of angles to the detector regardless of the incident angle of light on the diffuser. High Fresnel losses at high angles of incidence at the surface of the diffusers such as plastic disk or opal glass can still introduce errors. These errors can be minimized by allowing some light to leak through the edges of the diffuser and then using a screening ring to prevent additional error due to leaked light.⁷³



FIGURE 32 Handheld lighting measurement instruments. (a) Simple Luxmeter; (b) Luxmeter (from Labsphere) with an integrating sphere as a diffuser; and (c) luminance meter (from Konica Minolta).

Luminance Meters

Luminance meters measure luminance in lumen/area/solid angle. These are handheld instruments (Fig. 32) and are equivalent to using a lens and placing a luminance meter at the image location of the image formed by the lens of the target object. Thus it consists of a detector, a photopic correction filter or a color filter, a cosine corrector and an imaging lens. The lens is used to image the region of the object to be measured on the entrance to the detector (or the cosine corrector surface). Since the NA of the lens is known, the solid angle is known. This allows us to obtain the luminance by dividing the detected illuminance by the solid angle. Care must be taken to focus the lens only on the region whose average luminance has to be evaluated. Apertures with different field of view are provided to limit the region over which luminance is evaluated.

These instruments can be made more sophisticated by using a high-quality CCD detector array to image an entire scene and provide luminance distribution across it. The imaging lens has a flat field and is telecentric in image (detector) space.

The use of neutral density filter helps in expanding the dynamic range of the instrument. The measurement on color coordinates in various colors spaces is provided by using multiple detectors, each with a characteristic color filter. The incoming spectrum is multiplied by the transmission spectrum of the color filter for each detector. The relative signals on the detectors help in determining the color coordinates. Similarly the CCT can be evaluated.

Reflectometers

Reflectometers measure the reflectance from samples for cases such as reflectance as a function of wavelength, angle of incidence, angle of viewing and polarization. Depending upon the requirements, not all of these parameters are needed. The sample surfaces may be diffuse, specular, or a mix of two. When samples must be characterized by BRDF, scatterometers are used. Reflectometers are discussed in Chap. 35 in this volume and Scatterometers are discussed in Chap. 1 in Vol. V of this *Handbook*.

Goniometers

Goniometers are instruments that measure the irradiance or illuminance distribution of a source or luminaire. They accomplish this by measuring the illumination distribution at a number of points. By locating the detector in the far field with respect to the source or luminaire, one is actually measuring the intensity distribution. The far field is when the irradiance/illuminance distribution closely matches that of the respective intensity distribution, which means the inverse

TABLE 8 Three Types of Standard Goniophotometer with Their Respective Standard Coordinate Systems⁷⁴

	Type A	Type B	Type C
Polar axis	Vertical	Horizontal	Vertical
Vertical angle designation	Y	V	V
Horizontal angle designation	X	H	L
Range of vertical angles*	$Y \in [-90^\circ, 90^\circ]$	$V \in [-180^\circ, 180^\circ]$	$V \in [0^\circ, 180^\circ]$
Range of horizontal angles	$X \in [-180^\circ, 180^\circ]^\dagger$	$H \in [-90^\circ, 90^\circ]^*$	$L \in [0^\circ, 360^\circ]^\ddagger$
Straight ahead/down	Ahead: $Y = 0^\circ, X = 0^\circ$	Ahead: $V = 0^\circ, H = 0^\circ$	Down: $V = 0^\circ, L = 0^\circ$
Primary applications	Optical systems Automotive lighting	Floodlights	Indoor lighting Roadway lighting

*The lower angle is in the nadir direction while the upper angle is in the zenith direction.

†The lower angle is measured to the left from the perspective of the luminaire.

‡Measured from the primary axis of the luminaire.

square law applies. A good rule of thumb is ten times or greater the greatest extent of the emitter. For example, a luminaire with a 100-mm-diameter aperture indicates that measurement should be made at no closer than 1 m. Of course, better results are obtained as the distance between the source and detector is increased. By rotating either the source or detector with respect to the other, the full intensity distribution can be measured. The inclusion of a rotation device on either the luminaire or detector while the other remains fixed defines a goniometer. For the purposes of the lighting community, the luminous intensity is measured, so goniometers are better known as goniophotometers in the community.

There are three standard types of goniophotometers (Type A, Type B, and Type C), which are listed with their design criteria in Table 8. There are also three coordinate systems used for the purposes of the measurement reporting. These spherical coordinate systems are also labeled Types A, B, and C, and they typically conform to the type of goniophotometers.⁷⁴ However, individual goniophotometers may be of one type and report measurements in another coordinates system. Table 8 assumes that the goniophotometers types and coordinates systems agree.

For Types A and B goniophotometers the detector is held fixed while the luminaire is located in the rotating device. Typically guidelines or standards define the distance that the detector is located away from the luminaire. For example, for U.S. and European automotive headlamps this distance is 75 ft and 25 m, respectively. Type C goniophotometers have the detector rotate around the horizontal axis of the luminaire while the luminaire is rotated around its vertical axis. This setup is important for sources that have limitations to orientation (e.g., metal halide arc lamps). Figure 33 shows a Type C goniometer that is used to measure the luminance distribution from sources.

There are many other types of goniophotometers, especially those labeled as snapshot systems. Snapshot systems capture an “image” of the intensity distribution through one measurement. Examples include systems with several detectors; rapid scanning systems for smaller sources such as LEDs; camera-capture systems that incorporate an intermediate diffuse, reflective screen; and tapered-fiber bundles integrated to detector arrays.⁷⁴

Surface Measurement Systems

Not only are the illumination distributions, spectra, scatter distributions, and optical characteristics measured, but also the geometry of the optics. These measurements are done to characterize the geometry of the fabricated optic in regard to the design. This step is done, in conjunction with tests from the previous testing sections, when there is a disagreement between the laboratory and modeling results. There are essentially two methods in use: coordinate measurement method (CMM) and laser scanning.^{75,76} CMM uses a probe that is drawn across the surface of the optic, which provides the (x, y, z) data at a series of points on the surface. This method is analogous to using a spherometer

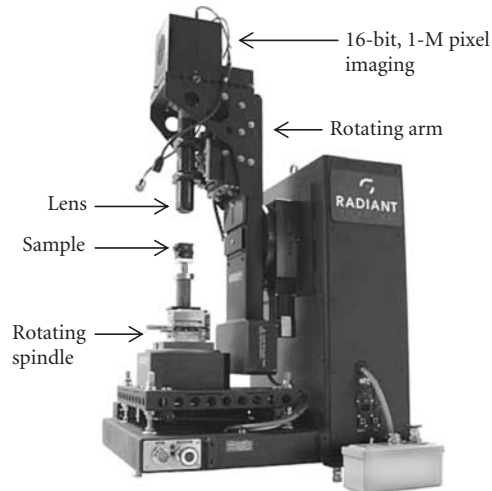


FIGURE 33 Photograph of a source measurement goniometer that is used to ascertain the luminance distribution of the source. The system wobble (electro-mechanical-software runout) is $15\ \mu\text{m}$ to allow for measuring small light sources like LED die. (See also color insert.) (Courtesy of Radiant Imaging, Inc., model SIG 400.)

to measure the sag and thus curvature of a lens surface. CMM is a contact method so the optic under test can be detrimentally affected (e.g., scratched) during measurement. Noncontact methods such as laser scanning essentially replace the mechanical probe with a laser beam. The position is measured through location of the reflected spot, triangulation, and/or time of flight as the laser is scanned across the part. Software can then determine the slopes of the test points, and thus rebuild the shape through a point cloud. If the component under test reflects well, then the surface must be coated or a fine, white powder can be placed across the surfaces. Laser scanning is a rapid method to ascertain the part shape. It provides a multitude of points in a short time; however, laser scanning accuracy is currently limited in comparison to CMM.

40.7 LIGHTING APPLICATION AREAS

In this section, we discuss several lighting application areas. We broadly divide the applications into two areas: interior lighting (office, residential, retail and healthcare) and exterior lighting (industrial and transportation). This set of applications is limited; we provide design concepts, metrics and data reference on major application areas. These ideas could be readily applied to other areas of lighting such as entertainment, sports, theatre, and so forth. Each field has its own region-specific and time-dependent guidelines and standards to implement.

The data provided here is primarily from IESNA suggested guidelines, thus exceptions are expected. For example, an office lobby is likely to have far dimmer ambient lighting than dictated by general office lighting guidelines. In essence high-end retail guidelines provide the design goals in order to provide a desired ambience. Although we have provided the guideline data to indicate the relevant specification parameters, by no means is this complete and the guidelines and standards are evolving with time, place, and technology advances. The reader must check the prevalent guidelines for the exact numbers to use. Previous sections of this chapter describe in detail the general aspects

of the lighting design process, principles, and techniques for the applications described here. Each of the following subsections provides insights into the specifics of the lighting design process for the given application.

Interior Lighting

Interior lighting includes the office, retail, residential, and health-care facility subfields. In each of these areas there are established guidelines to provide a pleasing environment to the users of the space. Of special importance are required illuminance levels, luminance ratios, and the reduction of glare. The following four subsections describe each of these application areas in more detail.

Office Lighting Modern office lighting has assumed greater importance than in the past as more people work in them. The goals of office lighting include efficient task performance, energy-efficient lighting, nonmonotonic ambient lighting that provides a balance between horizontal and vertical illuminance and minimum glare.

The task surface in offices has historically been the horizontal desk surface. However, computer monitors are ubiquitous these days and present a self-luminous vertical task surface. In addition it is specular and reflects light which can cause glare. Lighting goals are to provide adequate task illumination and adequate light in the task vicinity to eliminate eye strain due to varying brightness and disability glare. This is ensured by limiting the maximum luminance ratio of task to nontask regions. Veiling reflections (see section on glare) from computer monitor screens and paper surfaces must be avoided by taking into account source task and eye geometry. The glare from the computer monitor arises typical from ceiling source, and it is avoided by putting a limit upon the maximum ceiling luminance. Glare across the horizontal surface can be eliminated by using low luminance, wide-area lighting from overhead luminaires or special desk lamps (Fig. 24). Daylight from windows at desk level or near the ceiling can be a source of glare and steps such as daylight control via blinds or changing the task-eye geometry with respect to the glare source must be taken into account.

Balancing of daylight with electrical lighting, preferably with automatic controls, is not only energy efficient but also preferred by the workers. Therefore electrical lighting with high CRI (>70) is desirable as it mimics daylight illumination quality better.

Table 9 summarizes the common specifications and major guidelines on office lighting. The reflectivity of room parameters such as walls and ceiling is provided to aid in room modeling to aid in determining the layout of luminaires.

Retail Lighting The goal of lighting for the retail environment is to attract the customer, facilitate evaluation of the merchandise by the customer, and provide light for completing the transaction.⁷⁷ There are varying lighting methods dependent on the type of retail store, what goods are being illuminated, and the background. Table 10 provides various guidelines on lighting levels dependent on the type of store and what is being illuminated. Note that this table is not complete by any standard, so the reader is encouraged to consult the literature (see, for example, Ref. 77) for more specific information for a given retail lighting environment. Note that there are particular design issues that have varying importance levels for each type of retail outlet, such as, glare reduction in jewelry and china stores.⁷⁸

The circulation areas are those not typically used for the display of merchandise such as walkways, aisles, and foyers. The general areas are those for the generic display of merchandise. Perimeter areas are the walls where merchandise is placed for sale. Feature areas are where important displays are positioned. Horizontal illuminance values are listed for all of the columns in Table 10 except the perimeter areas, which provide vertical illuminance values. The feature areas have peak illuminance values of 5:1 to 10:1 compare to the respective general area peak illuminance. Of particular note in Table 10, the trends indicate

- Bulk stores, or those described as “big box,” tend to have much higher illuminance levels with higher uniformity across areas. This type of illumination provides the user with an abundance of light to fully inspect the merchandise, associated labeling, and compare to similar products.

TABLE 9 Summary of Specifications and Guidelines for Office Lighting

Specifications	Goals
Maximum luminance ratios:	
Task to neighboring areas	3:1 to 1:3
Task to distant regions	10:1 to 1:10
Max ambient illuminance (lx)	500
Max ceiling luminance (nits)	<1000 (in the presence of computer monitors) <425 nits does not cause glare
Reflectivity of room objects* (%)	
Ceiling	≥80
Walls	50–70
Partition	40–70
Furniture	25–45
Floors	20–40
Corridor Floors	>20 (of the illuminance of the adjacent areas)
CRI	>70
Lighting schemes	Direct, indirect and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, glare control, provide a spacious ambience
Ceiling uniformity (max: min)	<8:1 (for indirect lighting only)
Common light sources	Daylighting, LED, CMH, fluorescent, and CFL
Glare sources to be eliminated or reduced	Veiling reflections from direct sources or ceilings on computer monitor screen, reflections from any specular surface including walls and desks

*Valid for matte or diffuse finish.

TABLE 10 Suggested Illuminance Levels for Circulation, General, Perimeter, and Feature Areas for Various Types of Retail Stores⁷⁷

Type of Retail Store	Circulation Area Illuminance (lx)	General Area Illuminance (lx)	Perimeter Area Illuminance (lx)	Feature Area Illuminance (lx)
Warehouse	250–300	750–850	750–850	3750–8500
Supermarket	250–300	750–1000	750–1000	3750–5000
Discount	250–300	750–850	750–850	3750–8500
Mass merchant	200–250	500–600	750–850	2000–5000
Department	200–250	400–500	500–750	2000–3500
Upscale	150–200	300–400	400–800	1500–4000
Specialty	200–250	400–500	500–750	2000–3500
Upscale	150–200	300–400	400–800	1500–4000
Boutique	80–120	200–300	200–600	1000–3000
Jewelry	80–120	200–600	200–600	1000–6000
Upscale China	80–120	200–600	200–600	1000–6000
Drugstore	250–300	750–850	750–850	3750–8500
Home	200–250	400–500	500–750	2000–3500
Furniture	80–120	200–300	200–600	1000–3000

- More specialized or exclusive, the lighting levels tend to be reduced. This type of lighting provides a more intimate environment between the customer and the salesperson with illumination to highlight the item under evaluation.

Essentially, Table 10 can be broken down into three categories: mass merchandising, department, and exclusive stores. At the lower end, mass merchandising, one typically specifies ambient illuminance

TABLE 11 Suggested Illuminance Values for Areas within a Department Store⁷⁷

Department Store Area	Horizontal Illuminance (lx)	Vertical Illuminance (lx)	Very Important Design Issues
Alteration room	500	300	Color appearance and source geometry
Dressing area	300	50	Color appearance and object modeling
Fitting area	500	300	Color appearance and object modeling
Stock room	300	50	
Sales area	300		Direct glare
General merchandise area	500		Color appearance
Feature display	2,500		Appearance of area, color appearance, direct glare, and object modeling
Display window	2,000 (day) 500 (night)		Appearance of area, color appearance, daylighting, object modeling, and reflected glare
Feature	10,000 (day) 5,000 (night)		

between 750 and 1000 lx and a CCT of 3500 to 4100 K. Department stores are in the range of 400 to 600 lx with a color temperature around 3500 K. High-end stores have ambient light levels of 150 to 300 lx and color temperatures of 2700 to 3000 K.

The values provided in Table 10 are guidelines, and values are expected to vary dependent on the background, the items being illuminated, and any external lighting. The reasoning behind this is that observers see luminance rather than illuminance. Thus, the lighting designer must take into account the reflectance, both diffuse and specular, from the merchandise and background. Thus, Table 10 assumes that there is constant reflectance between the features and the background, such that illuminance ratios of 5:1 to 10:1 are specified, when in actuality luminance ratios in this range are prescribed.

Within any type of store there are different lighting levels dependent on the application. Consider, for example, a department store, which is made up of several different environments, from the entrance areas to fitting areas, to general displays areas. Table 11 provides illuminance level guidelines for typical areas in a department store. It also includes the very important design issues for each of the retail areas.

The methods of lighting a given area are dependent on the application of the space. Ambient lighting provides the baseline lighting level, while the addition of secondary lighting units provides the increased illuminance values as listed in Tables 10 and 11. Table 12 provides a synopsis of these other application space lighting demands, including the typical luminaire used and design issues.

As previously noted the lighting designer must remain cognizant of external lighting conditions when specifying the artificial retail lighting environment. A large facet of external lighting is better known as daylighting, and of particular concern are the varying light level that is provided through the day, direct glare through windows and doors, and the strong background to incorrectly situated merchandise. Other aspects of the retail lighting environment include

- The CRI should be 70 or greater for most environments.
- Transitions between spaces in stores should have luminance ratios of no greater than 3:1 for similar neighboring spaces, greater than 3:1 when there is a distinct transition between the neighboring spaces, and 5:1 to 10:1 for abrupt transitions.
- The perimeter area illuminance should be greater than the overhead area in order to draw the attention of the shopper to the merchandise.
- Calculating baseline lighting levels for retail spaces one should use diffuse reflection values of office spaces.

Residential Lighting The goals of residential lighting are to provide ambient illumination to create a pleasing ambience due to a well-lit environment; sufficient task lighting in workspaces such

TABLE 12 Specific Application Space Lighting for Illumination of Merchandise Locations

Application Space Lighting	Typical Luminaires Used	Design Issues
Ambient Light		
Mass merchandising department	Fluorescent and halogen fluorescent, recessed fluorescent and halogen	Uniformity
High end	Recessed fluorescent and halogen, track lighting	Flexibility for change
Perimeter	Fluorescent, incandescent, or HID wall wash, track or recessed spot lighting	Uniform vertical illuminance, hidden luminaires
Rack	Recessed or track dpot lighting	Direct glare reduction
Shelf	Ambient lighting with recessed surface lighting	Direct glare reduction, hidden luminaires
Counter	Focused downlighting	Direct glare reduction, 3 to 5 times ambient lighting
Mirror	Downlighting	Glare reduction, color appearance, object modeling, and consistent lighting with use of product
Showcase	Fluorescent and fiber-optic lighting	Hidden luminaires
Accent	Small (point-like) sources	Provide luminance ratios of 5:1 to 15:1
Decorative	Sconces, chandelier, table, and torchiere lamps	For high-end stores to provide a desired look and feel

as office, garage, kitchen, and workshop; and decorative lighting and accent lighting for lit object displays. A variety of lighting techniques are used to create layers of lighting that perform different lighting functions (see book by Whitehead on “Residential Lighting” in Ref. 26). Like retail lighting, residential lighting can be quite complex and creative. Although we summarize many common guidelines for residential lighting requirements in Table 13, other specialized guidelines are more relevant in some areas of the house, such as retail lighting guidelines for the kitchen, office lighting guidelines for the home office, industrial lighting guidelines for task lighting in garage or workshop and exterior lighting guidelines for landscaping and the house facade. Thus depending upon the purpose of a specific region of the house the lighting scheme must be tailored. However, the different lighting schemes used across the house must be gently blended into one another. For example, exterior lighting for residences includes the lighting of entry, walkways, and landscape. Although the primary goals are to provide direction, safety, identification and aesthetic appearance, achieving a balance between interior and exterior lighting makes the interior and exterior spaces extensions of one another. Residential lighting should be customized to maximize the comfort of the inhabitants especially when the inhabitants are disabled or elderly. For example, elderly people require much higher levels of illumination especially for task performance.

Other than those areas that require excellent task lighting such as offices or workshops, the limits on luminance ratios are relatively relaxed when compared to other lighting applications, such as the maximum residential luminance ratio of 5:1 between Task and neighboring areas is higher than that of office lighting of 3:1. The idea of limiting the maximum luminance ratio is to minimize visual discomfort caused by disability glare and adapting to variations in brightness. The integration of daylight with artificial lighting is highly desirable.

Health-Care Facility Lighting The lighting in health-care facilities, which includes hospitals, out-patient clinics, chronic and extended care centers, and other facilities, requires careful understanding of the lighting requirements for not only the patients but also the individuals working therein. The lighting must be pleasing and comforting to the patients in order to put them at ease and assist in their healing. The patients and visitors have a wide range of ages, with the majority being elderly;

TABLE 13 Summary of Specifications and Guidelines for Residential Lighting

Specifications	Goals
Maximum luminance ratios	
Task to neighboring areas	5:1 to 1:5 (in general) 3:1 for demanding tasks such as sewing
Task to distant regions	10:1 to 1:10
Luminaire to ceiling	20:1
Hallway or stair to adjacent area	1:5
Reflectivity of room objects* (%)	
Ceiling	60–90
Walls, curtains, draperies†	35–60
Floors†	40–70
Average luminance for luminaire (nits)	1700
Maximum luminance for luminaire (nits)	2700 (in utility areas)
CRI	>80 (kitchen and clothing closets)
Lighting schemes	Direct, indirect, and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, glare control, providing a spacious and cozy ambience as well cozy as desired
Common light sources	Daylighting, LED, fluorescent, and CFL.
Glare sources to be eliminated or reduced	Direct glare from light sources, veiling reflections, indirect glare from shiny objects

*Valid for matte or diffuse finish.

†Reflectance of walls and floor can be increased by 40% and 25%, respectively, to improve visual task lighting where needed.

therefore, there is a large variance in the response to lighting, especially increased glare sensitivity, loss of contrast sensitivity, the need for higher lighting levels, and slow adaptation to changes in brightness as one ages.⁷⁹ The lighting levels also need to provide for the medical staff such that they can effectively carry out their job—from meticulous and demanding surgery to patient interviews and diagnoses to manning the check-in desk. Finally, since in a number of these facilities patients may be there for extended periods, circadian system illumination levels that conform to the human biological clock are the norm. As can be seen, there is an extensive range of tasks, observers, and daily requirements for illumination in health-care facilities, thus, making the lighting design a challenging assignment.

Foremost is the need to specify the light requirements for a given location based on the tasks to be performed there. Table 14 provides a synopsis of the illuminance level guidelines and important design issues for a number of locations in health-care facilities; however, since these are only guidelines, controls for the area illumination are often available to the medical staff and/or patients. This flexibility in the control of lighting levels based on the specific function of the environment and even mood of the occupants is important in health-care facilities. Note that there are many other areas and functions than can be included in Table 14, so the reader is encouraged to see Ref. 79 for this additional information.

The operating room environment is especially challenging since the lighting is preferentially directed to the task area; however, this has the potential of creating large luminance ratio variation within the room. Practices suggest three regions within the operating room with the following luminance ratios: task area to the surgical table of 3:1 or less and task area to the background (i.e., walls) of 10:1 or less. Additionally, there are specific guidelines for the finishes for the surfaces within the surgery theatre as provided in Table 15. In most cases all surfaces are white or pastels with a matte finish to reduce reflected glare. The lighting in an operating room is provided by a multitude of sources including ambient and even daylighting, directed spotlights, surgeon headlamps, and increasingly fiber-optic lighting.

Concurrent to the design guidelines of Table 14, the lighting designer must remain cognizant of the specifications for reflectances of the walls, floor, ceiling, and other objects that occupy the design area.

TABLE 14 Suggested Illuminance Values and Other Criteria in Health-Care Facilities⁷⁹

Health Care Facility Area and Current Function	Horizontal Illuminance (lx)	Vertical Illuminance (lx)	Important Design Issues
Patient room			
Normal	200+, as home	30+, as home	Patient control, CRI >80, daylighting
Examination	1000	300	Glare reduction, CRI >80, CCT >3000 K, doctor control
Observation (night)	30	30	Red/amber CCT, no higher than 18" off floor
Nursing station	300	50	Glare reduction, CRI >80
Critical care unit			
Normal	300	50	Patient control, CRI >85, daylighting
Examination	1000	300	Glare reduction, CRI >85, CCT >3000 K, doctor control
Nursery	100	30	Glare reduction, CRI >80, lighting control, daylighting
Mental health center	As per other functions	As per other functions	Daylighting, CRI >80, CCT of 4100–5000 K with fluorescent, 3500 K otherwise, glare reduction
Operating room			
Normal	1000	500	CRI >85, matched light
Table	25000+ in 20-cm spot (adjustable)	1000	Sources for illumination, shadow, and glare reduction, doctor control, CCT of 3500–6700 K, high uniformity
Dental unit			
Normal	300	50	
Examination	24000+ in central spot	500	As operating room
Radiography unit	50, but depends on test	30, but depends on test	CRI >80, glare reduction, doctor control
Pharmacy	300	100	Glare reduction, high uniformity, CRI >80

TABLE 15 Suggested Reflectances for the Various Objects in a Health-Care Facility Room for a General Application and the Operating Room Environ⁷⁹

Surface	General Reflectance	Operating Room Reflectance
Ceiling	70–80%	90–100%
Walls	40–60%	60%
Furniture/fabrics	25–45%	0–30%
Equipment	25–45%	25–45%
Floors	20–40%	20–30%

Table 15 provides this data for two locations, general and operating room. Using these values in conjunction with the optical design process one can find suitable illumination configurations to provide the Table 14 guidelines.

Industrial Lighting The goals of industrial lighting are to provide energy efficient lighting with adequate task and ambient lighting, direction, safety, and visual comfort. Many common requirements for industrial lighting requirements are summarized in Table 16.

TABLE 16 Summary of Commonly Suggested Specifications and Guidelines for Industrial Lighting

Specifications	Goals
Maximum luminance ratios	
Task to neighboring areas	3:1 to 1:3 (in general)
Task to distant regions	10:1 to 1:10
Luminaire (including daylight sources) to adjacent surfaces	20:1
Anywhere within the FOV	40:1
Reflectivity of objects* (%)	
Ceiling	80–90
Walls	40–60
Floors	≥20
Desk/bench tops, machines, equipment	25–45
CRI	>65
Lighting schemes	Direct, indirect, and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, provide high illuminance uniformity, sufficient illuminance, and glare control
Common light sources	Daylighting, LED, fluorescent, HID, and CFL
Glare sources to be eliminated or reduced	Direct glare from light sources, veiling glare, indirect glare from shiny objects
Direct view of the luminaire (deg)	>25 (preferably >45)

*Valid for matte or diffuse finish.

The ambient lighting is provided by daylighting and/or large area, overhead, wide angle luminaires (direct and semidirect). Task lighting is provided by fixed and portable direct or diffuse light sources. Backlights are used for translucent task surfaces. Direction of illumination with respect to the view angle(s) is important in tasks where surface features of the object cause shadows to enhance depth perception. Illumination at grazing incidence is used to highlight specific feature that can scatter light and render themselves visible. To emphasize specular but uneven surfaces, specific patterns of light (like bright and dark lines) can be projected onto the task surface to be reflected into the viewer's eyes. For lighting in regions where there is a high density of workstations or areas where there are several similar process, a high level of uniformity in horizontal illuminance is recommended. In such areas, variation of horizontal illuminance should be less than 1/6th the average horizontal illuminance. Otherwise, the variations in luminance across the work space must be within the luminance ratios prescribed in Table 16. The level of horizontal illuminance needed varies with task. In industrial lighting, the quality of horizontal illuminance is of special importance as efficient and safe task performance is needed. The reader is referred to IESNA published guidelines for horizontal illuminance values for a wide variety of tasks in different industries.⁸⁰

In addition to providing the task lighting and ambient lighting, it is often necessary to provide additional lighting dedicated for emergency, safety, and security within the industrial complex and its exteriors.

Visual discomfort must be avoided especially during task performance and in situation where safety could be compromised. Special attention should be given to situations causing veiling reflections and glare. The various limits placed on the luminance ratios in Table 16 between task and non-task regions help eliminate the impact of glare.

Exterior Lighting

Of course the primary aspect of exterior lighting is to provide illumination during hours of darkness. The lighting not only provides illumination for general use, it also provides safety and security; indication of direction of travel for paths, roadways, and so forth; and architectural enhancement.

External lighting is attempting to replicate the illumination of the sun; however, it can in no fashion accomplish this feat. The illumination provided by any lighting source cannot match that of the sun, which means that the sky will appear dark rather than blue; numerous sources are required to provide the necessary illumination level so glare from the many sources is a major factor; distance from the light sources, mesopic vision and even scotopic levels may be demanded; and the many sources can confuse the viewer such that objects can be difficult to discern from one another.⁸¹ Additionally, this section is rather broad in scope, including design aspects such as roadway lighting, path lighting, outdoor event lighting, and façade illumination. Thus, the reader may need to consult specific literature for a certain type of lighting. Please consult the transportation lighting section for further information on vehicular and roadway. In the next two subsections details are presented about the issues present for external lighting design and design examples, excluding the transportation applications discussed later.

External Lighting Design Issues There are several topics that a lighting designer for external spaces must consider. First and foremost is the specific application guidelines (for example, see the section on illuminating roadways), glare issues, light pollution and trespass, and perception issues. All of these factors are interrelated, and the subsections herein provide insight into them.

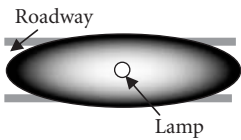

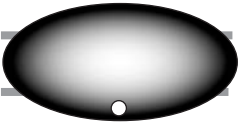
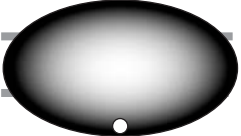
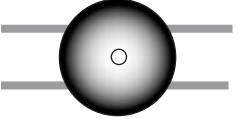
Glare A primary issue of external lighting is glare. Since users of a space in the hours of darkness are relying on nonphotopic vision, glare can blind the viewer as one approaches a bright source thus not only causing the scotopic vision receptors (e.g., rods) to naturally saturate but also the photopic vision receptors (e.g., cones). This “blinding” is due to disability glare, which will hide objects or reduce their contrast. Other levels of glare can also inhibit observation: discomfort glare, which is attributed to a large variance on and around the object under view; and annoyance glare, which is light that in the opinion of the observer should not be present, such as light trespass.

The major source of glare is directly from the source in comparison to the object under observation. Comparison between the luminaire luminance to the object’s illuminance is the factor that often defines the level of glare, from disabling to discomfort. Depending on the orientation of the object one either compares the horizontal or vertical surface luminance to the luminaire luminance. The horizontal luminance is used for horizontal surfaces such as walkways and roads, while the vertical luminance is used for vertical surfaces such as building façades, people, and structures.

Light Pollution and Light Trespass Associated to glare is light pollution and light trespass. Light pollution, also called sky glow, is light that is directed upward to the atmosphere.⁸² This sky glow hides stars from observation, and on cloudy days provides a glow, typically red in hue, which can distract from one’s appreciation of the view of their surroundings. Astronomers have been particularly vocal in the reduction of light pollution, and there is increasingly consideration of energy efficiency demands. In 2008 it was hypothesized that over \$10B U.S. was wasted through light pollution,⁸³ which means that this amount has only increased since then. Like light pollution, light trespass is unwanted light, but in this case it is light that illuminates objects outside of the intended illumination region.⁸⁴ Such stray light enters through windows, illuminates someone’s property, or can impair the vision of drivers. Light pollution and trespass arise from improperly designed luminaires. There are increasing regulations for the design of outdoor lighting to alleviate such concerns. A primary correcting factor is the use of luminaire cutoffs as presented in the section on luminaire design.

External Lighting Example Pole lighting is the most common form of outdoor illumination providing a significant amount of light, safety and security, and indication of the path of travel for walkways, roadways, and so forth. In the United States, there are five primary forms of classifications of pole-mounted luminaires, but there are numerous variations based upon the lighting requirements, the shape of the illumination region, and illumination level demands.⁸⁵ Table 17 provides descriptions of the five types, the typical application(s) of the type, and an iso-illuminance plot for each of the types. Per the iso-illuminance plots of Table 17, the larger the illumination distribution the more ground is illuminated, which means not only the roadway is illuminated but also sidewalks, shoulders, and other neighboring areas. Thus, Type IV and V luminaires tend to have more light trespass.

TABLE 17 Description, Applications, and Iso-Illuminance Plots for the Five Primary Luminaire Types Used for Pole-Mounted Outdoor Lighting⁸⁵

Type	Description	Applications	Iso-Illuminance Plot
I	Little setback to the road or mounted over the roadway; narrow oval illumination distribution	Roadways	
II	Moderate setback to pedestrian area; moderate oval illumination distribution	Pedestrian areas	
III	Large setback to the road; Moderate elliptical illumination distribution	Pedestrian areas roadways, parking lots	
IV	Great setback to the roadway; large elliptical illumination distribution	Roadways	
V	Mounted over pedestrian area rotationally symmetric illumination	Pedestrian areas parking lots	

Transportation Lighting

Transportation lighting includes the subfields of vehicular and roadway lighting. In both cases rather than just guidelines, there are many stringent governmental standards that must be met by the illumination systems. For example, there are governmental standards for traffic lights and vehicular taillights and headlights. If these standards are not met, then the lighting systems are not legal for use on our roadways.

Vehicular Lighting Vehicular lighting is the external illumination aspects of vehicles, including automobiles, motorcycles, snowmobiles, emergency vehicles, airplanes, ships, and heavy machinery including construction, industrial and agricultural. Ground vehicles are especially important due to their pervasiveness in society. For the remainder of this section, we highlight ground vehicle standards, but there are similar standards for other types of vehicles. In the United States, the U.S. Federal Government standardizes the lighting requirements through the Federal Motor Vehicle Safety Standards (FMVSS) from the offices of the National Highway Traffic Safety Administration (NHTSA) within the Department of Transportation (DOT). In Europe and a number of other countries spanning the globe, the United Nations Economic Commission for Europe (UNECE or better known as ECE) sets the standards. The FMVSS calls upon the standards delineated by the Society of Automotive Engineers (SAE) to provide the explicit lighting requirements for distinct lighting systems. While FMVSS 108 provides the framework for lighting systems on ground vehicles within the United States,⁸⁶ the SAE provides the accepted standards to design into the lighting systems.

For example, SAE standard J581 provides the upper-beam requirements to be an accepted high beam on U.S. roads,⁸⁷ while R113 is the accepted ECE standard.⁸⁸ The associated low-beam standards are SAE J582 in the U.S.⁸⁹ and ECE R112 in Europe,⁹⁰ but there is currently an active dialogue to harmonize the standards between the U.S., European, and Japanese markets. The SAE J1735 standard is currently addressing the low-beam harmonization, and in time it will also address the high-beam requirements. The end results will be, excepting the inherent difference of left-hand (e.g., United Kingdom) and right-hand (e.g., United States) traffic, harmonization has the goal of making the lighting standards the same at as many places possible across the globe. This increased standardization means that the design and fabrication costs will be reduced for manufacturers.

There are essentially two types of lighting on a ground vehicle: forward lighting for illuminating the road surface and nearby surroundings and indicator and warning lighting including turn signals, brake lights, side markers, and tail lights. Each vehicular light system is defined by its own set of regulations including the distribution of light, the color and characteristics of the protective lens, mounting requirements, and a number of enviromechanical tests including vibration, dust, moisture, corrosion, and warpage. In the United States, the test procedures and protocols are typically governed by SAE J575, while each individual standard provides the photometric requirements. In the United States, the photometric requirements are the luminous intensity distribution at a distance of 18.3 m from the lamp, while the ECE requirements are for the illuminance distribution at a distance of 25 m. The typical instrument to make this measurement is a goniometer that holds the lamp and allows it to be rotated so that a different angle is incident on a fixed photodetector at the end of an absorbing light tunnel. By rotating the lamp within the goniometer and making a series of measurements the full distribution of light can be determined; however, rather than measuring across the whole lit region, the standards dictate a series of test points and test areas that must be measured. In the remainder of this section, we highlight two applications: low-beam headlamp and stop lamp. In each subsection, the requirements for both U.S. and ECE standards are provided, a depiction of a typical distribution of light, discussion of design strategies, and additional optical requirements of the standards. Note that while we highlight these two applications and their respective standards, there are numerous other optical standards governing the external lighting systems on vehicles (see Refs. 86 to 92) for further information about these additional standards).

Design and Standards for a Low-Beam Headlamp A headlamp, as shown in Fig. 34a, must illuminate the road surface for the driver, illuminate to the sides for both the driver and any other observer, provide a low level of illumination to oncoming drivers, and in all cases remove the formation of hot spots above the horizon such that glare is not a concern to oncoming traffic. Of note in Fig. 34a is

- High-beam luminaire:
 - Rightmost recess of the lamp.
 - Faceted reflector, but other options include smooth, tailored reflector (e.g., NURBS surface as designed in CAD); faceted lens in conjunction with smooth reflector (e.g., parabolic); and projection lens in conjunction with reflector.

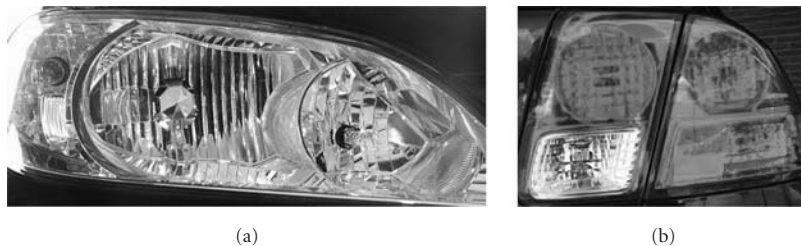


FIGURE 34 (a) A faceted headlamp including high-beam (right), low-beam (middle), and turn signal (left) luminaire. Note the yellowish tinge of the turn signal, which is due to the coating placed on the bulb used therein. (b) A faceted taillight including the following functions: tail (upper left), stop (upper right), turn signal (lower right), reflex reflector (lower middle), and backup (lower left). (See also color insert.)

- Filament source, but other options include high-intensity discharge lamp or array of white-light LEDs.
- Low-beam luminaire:
 - Middle recess of the lamp.
 - Faceted reflector, but other options are as the high-beam luminaire above.
 - Filament source, but other options are as the high-beam luminaire above.
 - Bulb shield: reflective structure in the middle of the lamp and covers direct view of the filament source. The bulb shield greatly alleviates direct light above the horizon.
- Turn signal luminaire:
 - Leftmost recess of the lamp.
 - Faceted reflector, but other options are as the high-beam luminaire above.
 - Filament source, but other options are as the high-beam luminaire above. Note that the bulb has a yellow coating placed on its glass envelope, which then provides the yellowish appearance of such lamps. This coating provides the color as specified by its respective standard.

Headlamp reflectors have a parabolic base shape in order to provide a high degree of collimation in the forward direction. The reflector, both faceted and smooth, is then deformed to provide the distribution that meets standards. Designing such deformations can be difficult, but there are software codes to assist in the process. Design guidelines include a goal of having the emitted radiation only incident once on the reflector; avoid secondary interactions with the bulb shield, bulb, or reflector shelves (i.e., the reflector sides, as shown in Fig. 34a); and angling intersegment fillers (i.e., spaces between the facets) such that they cannot be directly seen from the source emission region. In order to avoid secondary interactions with the bulb and its respective shield, one typically angles the reflected light across the luminaire, except that striking near the reflector near the source. The latter is directed away from the source. In systems with projection and faceted lenses, the reflector is there to capture bulb emission, while the lens provides the required distribution that meets standards. Due to demanding illumination requirements, it is best if most of the bulb emission is incident on the reflector; therefore, in almost all cases the filament or arc is oriented along the axis of the reflector (i.e., orthogonal to the filament as shown in Fig. 13). This axial filament orientation also ensures less light goes above the horizon in the low-beam luminaire.

Table 18 provides the SAE J582 luminous intensity requirements, while Table 19 provides the ECE R1 12 requirements. Figure 35 shows a typical SAE luminous intensity distribution that meets the requirements of Table 18. Figure 36 shows a typical ECE illuminance distribution that meets the requirements of Table 19. The SAE standards of Table 18 and Fig. 35 are in the units of candelas

TABLE 18 SAE J582 Standard Photometric Requirements for Auxiliary Low-Beam Lamps⁸⁹

Test Point (Deg, $\pm 0.25^\circ$)	Max Luminous Intensity (cd)	Min Luminous Intensity (cd)
10U to 90U	75	—
1.5U–1L to L	300	—
1.5U–1R to R	300	—
0.5U–1L to L	400	—
0.5U to 1R to 3R	400	—
0.5D–1R to 3R	25,000	2,000
0.5D–1L to L	10,000	—
5D–4R	—	3,000
5D–4L	—	3,000
1D–1R	—	10,000
3D–3R	5,000	—
4D–V	3,000	—
2.5D–15L	—	1,500
2.5D–15R	—	1,500

TABLE 19 ECE R112 Standard Photometric Requirements for Class A Passing (i.e., Low-Beam) Lamp for Right-Hand Traffic⁹⁰

Test Point Label	Horizontal Location (mm)	Vertical Location (mm)	Max. Illuminance (lx)	Min. Illuminance (lx)
B50L	1500L	250U	0.4	—
75R	500R	250D	—	6
75L	1500L	250D	12	—
50L	1500L	375D	15	—
50R	750R	375D	—	6
50V	V	375D	—	—
25L	3960L	750D	—	1.5
25R	3960R	750D	—	1.5
Zone III	See Fig. 36	See Fig. 36	0.7	—
Zone IV	2250L to 2250R	375D to 750D	—	2
Zone I	L to R	750D to D	20	—

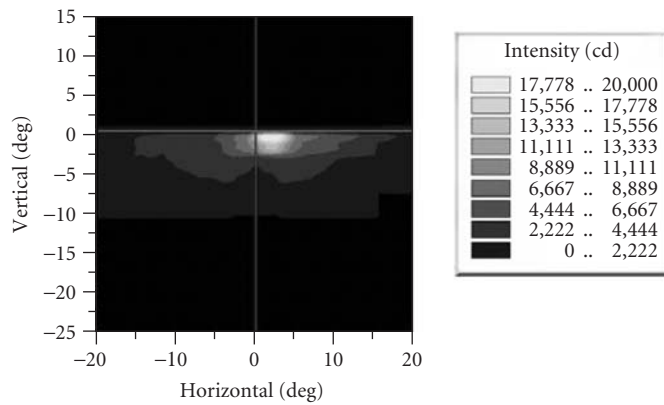


FIGURE 35 Luminous intensity (cd) distribution for the SAE low-beam requirements of Table 18. (See also color insert.)

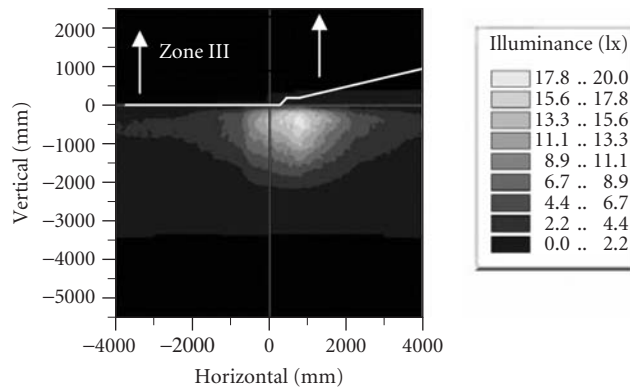


FIGURE 36 Illuminance (lx) distribution for the ECE passing/low-beam requirements of Table 19. (See also color insert.)

(cd, lumens/steradian). The test points are given in degrees and a letter designation, which mean U = up direction, D = down direction, L = left direction, and R = right direction as measured from the point (H, V), which is the center point where H = horizontal and defines the horizon, and V = vertical and defines the lane of traffic. The ECE standards of Table 19 and Fig. 36 are in the units of lux (lx, lumens/m²), and, while the letter designations still hold, the test point locations and zones are given as position coordinates in millimeters (mm).

Design and Standards for a Stop Lamp A taillight, as shown in Fig. 34*b*, is typically comprised of a number of different functions, such as a stop light, turn signal, backup lamp, and so forth. The taillight of Fig. 34*b* has the following optics in its functions.

- Tail lamp (upper left, red): for night-driving conditions or when headlamps are on. Indicates to following traffic the presence of the vehicle during reduced lighting conditions. The governmental standards are SAE J585 and ECE R7.
- Stop lamp (upper right, red): indicates when the driver has applied the brakes. This luminaire is also lit for night driving conditions but a lower output level. The governmental standards are SAE J586 and ECE R7.
- Turn Signal Lamp (lower right, red): indicates the driver is to make a turn in the designated direction. The governmental standards are SAE J588 and ECE R6.
- Reflex reflectors (middle right, red): this area, directly neighboring the turn signal is comprised of prism structures that provide retroreflection to the driver of following vehicles. It is important for dark-driving conditions to highlight the presence of this vehicle to following drivers or indicate its presence when the vehicle is not in operation. Note that this function of the taillight is passive in the sense that no internal light source is used. The governmental standards are SAE J594 and ECER3.
- Backup lamp (lower left, white): this luminaire indicates when the vehicle is in reverse. The governmental standards are SAE J593 and ECE R6.

In Fig. 34*b*, faceted reflectors are used for all the active luminaires. Typically, a transverse filament is used, as per Fig. 13, since the illumination standards for taillight functions are quite broad angularly. Thus, direct radiation from the filament is useful to filling in the required light distribution, while the reflected component of the illumination fills in the required hot spots. LEDs are replacing filament lamps to become the norm for most taillight functions. For LED sources, refractive optics that also employ total internal reflection provide a better means to meet the illumination requirements.

In the remainder of this section, we focus our attention on the design and standards of the stop lamp function, but the other taillight functions have similar requirements. The stop lamp is to inform following vehicles and pedestrians that the vehicle is slowing. In Fig. 34*b*, the stop lamp is located on the upper right of the taillight. Table 20 provides the SAE J586 luminous intensity requirements,⁹¹ while Table 21 provides the ECE R7.⁹²

Figure 37 shows a typical SAE luminous intensity distribution that meets the requirements of Table 20. Figure 38 shows a typical ECE illuminance distribution that meets the requirements of Table 21. The SAE standards of Table 20 and Fig. 37 are in the units of candelas (cd, lumens/steradian). The ECE standards of Table 21 and Fig. 38 are in the units of candelas (cd, lumens/steradian). In all cases, the letter designations as per the previous section (headlamps), but in this case the angles are with respect to rear direction of the vehicle.

Roadway Lighting There are multiple subfields that make up roadway lighting: street lighting, roadway signage, tunnel lighting, and integration of collocated pedestrian, bike, and analogous areas. The goal of all forms of roadway lighting is to reduce the potential of accidents, aid in the flow of traffic, provide a higher level of safety and security, and assist in commerce during hours of darkness. There are a multitude of light sources that affect roadway lighting, including the external lighting from vehicles (see section on vehicular lighting); direct roadway lighting, which is the subject of this section; traffic lights, as governed in the United States by the International Transportation Engineers

TABLE 20 SAE J586 Standard Photometric Requirements for a Stop Lamp*

Zone	Test Point (Deg.)	1 Lit Section	2 Lit Sections	3 Lit Sections
		Min. Luminous Intensity (cd)	Min. Luminous Intensity (cd)	Min. Luminous Intensity (cd)
I	10U-5L	9.6	11.4	13.2
	5U-20L	6	7.2	9
	5D-20L	6	7.2	9
	10D-5L	9.6	11.4	13.2
	Zone Total	52	62	74
II	5U-V	18	21	24
	H-10L	24	28.2	33
	5D-V	18	21	24
	Zone Total	100	117	135
III	5U-V	42	49.2	57
	H-5L	48	57	66
	H-V	48	57	66
	H-5R	48	57	66
	5D-V	42	49.2	57
Zone Total	380	449	520	
IV	5U-V	18	21	13.2
	H-10R	24	28.2	9
	5D-V	18	21	9
	Zone Total	100	117	13.2
V	10U-5R	9.6	11.4	74
	5U-20R	6	7.2	24
	5D-20R	6	7.2	33
	10D-5R	9.6	11.4	24
	Zone Total	52	62	135
Maximum	Any point above	300	360	420

*The stop lamp is comprised of up to three distinct lit sections over the extent of the taillight for the stop lamp function. Each zone has a number of test point minima that must be realized and the summed total of all test points within a zone. The last line of the table indicates the maximum luminous intensity that can be measured at any test point.⁹¹

TABLE 21 ECE R7 Minimum and Maximum Photometric Requirements for a Stop Lamp*

Test Point (Deg.)	1 Lamp Illumination Level		2 Lamp Illumination Levels	
	Min. Luminous Intensity (cd)	Max. Luminous Intensity (cd)	Min. Luminous Intensity (cd)	Max. Luminous Intensity (cd)
10U-5L	12	37	6	16
10U-5R	12	37	6	16
5U-20L	6	18.5	3	8
5U-10L	12	37	6	16
5U-H	42	129.5	21	56
5U-10R	12	37	6	16
5U-20R	6	18.5	3	8
V-10L	21	64.75	10.5	28
V-5L	54	166.5	27	72
V-H	60	185	30	80
V-5R	54	166.5	27	72
V-10R	21	64.75	10.5	28
5D-20L	6	18.5	3	8
5D-10L	12	37	6	16
5D-H	42	129.5	21	56
5D-10R	12	37	6	16
5D-20R	6	18.5	3	8
10D-5L	12	37	6	16
10D-5R	12	37	6	16

*The two categories are for a lamp that is either lit or not lit (e.g., 1 lamp illumination level) and for a lamp that has an unlit, partially lit, and fully lit state (e.g., 2 lamp illumination levels).⁹²

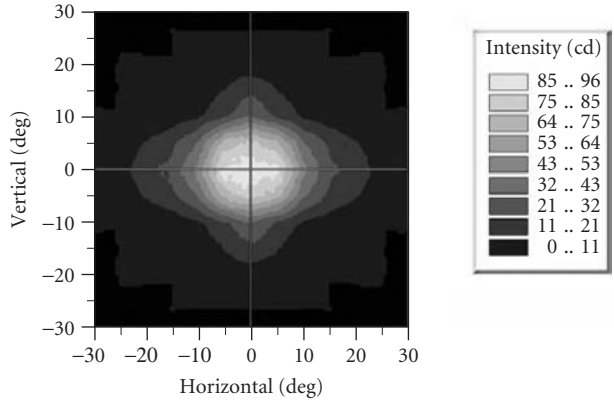


FIGURE 37 Luminous intensity (cd) distribution for the SAE stop lamp requirements of Table 20 (1 lit section). (See also color insert.)

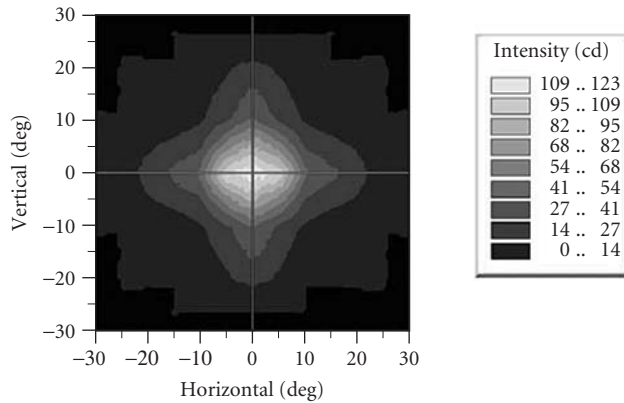


FIGURE 38 Luminous intensity (cd) distribution for the R7 stop lamp requirements of Table 21 (1 lamp illumination level). (See also color insert.)

society; and lighting from other sources such as residential, industrial, and retail. An added difficulty of roadway lighting is the harsh conditions in which they reside. There are stringent maintenance demands that include replacement of sources, cleaning of the optics, and trimming of foliage around the luminaire.⁹³

Street Lighting There are essentially three metrics for defining the lighting of a roadway: illuminance (lux, lumens/m²), luminance (nit, lumens/sr/m²), and small target visibility (STV).⁹⁴ Illuminance modeling provides the level of illumination across a surface. Luminance modeling predicts the level of light (i.e., glare), both direct and reflected, that is directed to a driver. STV is the visibility of a target array (18 × 18 cm² and 50 percent reflective) on the road. Table 22 provides the illuminance, luminance, and STV guidelines for a number of road types, four road surfaces in dry conditions, and the interaction level with pedestrians. The illuminance and luminance ratios of Table 22 provide limits such that disabling glare does not blind drivers or pedestrians. Table 23 qualifies the road types, the four road surfaces, and the pedestrian interaction levels of Table 22.

TABLE 22 Illuminance, Luminance, and STV Guidelines for the Type of Road, the Road Surface Classification, and the Interaction Level with Pedestrians⁹⁴

Road	Ped. Level	Road Surface				Illum. Uni. Ratio E_{avg}/E_{min}	Avg. Lum. I_{avg} (cd/m ²)	Avg. Lum. Uni. Ratio I_{avg}/I_{min}	Max. Lum. Uni. Ratio I_{max}/I_{min}	Veiling Lum. Ratio L_{Vmax}/L_{avg}	STV Weight Avg. VL	STV Avg. Lum. L_{avg} (cd/m ²)		STV Lum. Uni. Ratio L_{max}/L_{min}
		R1 (lx)	R2/R3 (lx)	R4 (lx)								<7.3 m	≥7.3 m	
Class A Freeway	NA	6	9	8	3	0.6	3.5	6	0.3	3.2	0.5	0.4	6	
Class B Freeway	NA	4	6	5	3	0.4	3.5	6	0.3	2.6	0.4	0.3	6	
Expressway	High	10	14	13	3	1	3	5	0.3	3.8	0.5	0.4	6	
	Medium	8	12	10	3	0.8	3	5	0.3	3.8	0.5	0.4	6	
	Low	6	9	8	3	0.6	3.5	6	0.3	3.8	0.5	0.4	6	
Major	High	12	17	15	3	1.2	3	5	0.3	4.9	1	0.8	6	
	Medium	9	13	11	3	0.9	3	5	0.3	4	0.8	0.7	6	
	Low	6	9	8	3	0.6	3.5	6	0.3	3.2	0.6	0.6	6	
Collector	High	8	12	10	4	0.8	3	5	0.4	3.8	0.6	0.5	6	
	Medium	6	9	8	4	0.6	3.5	6	0.4	3.2	0.5	0.4	6	
	Low	4	6	5	4	0.4	4	8	0.4	2.7	0.4	0.4	6	
Local	High	6	9	8	6	0.6	6	10	0.4	2.7	0.5	0.4	10	
	Medium	5	7	6	6	0.5	6	10	0.4	2.2	0.4	0.3	10	
	Low	3	4	4	6	0.3	6	10	0.4	1.6	0.3	0.3	10	

TABLE 23 Road Types and Surfaces, and Pedestrian Interaction Levels as per Table 22⁹⁴

Road Type	
Class A Freeway	Divided high traffic highways with full access control
Class B Freeway	All other divided highways with full access control
Expressway	Divided highways with limited access control
Major	Primary roadways within and leaving metropolitan areas
Collector	Service roads connecting major and local roadways
Local	Provide direct access to residential, retail, and industrial areas
Road Surface Class	
R1	Portland cement concrete and asphalt with at least 12% artificial brightener aggregates; treat as diffuse
R2	Asphalt road surface with a minimum 60% gravel aggregate; treat as both diffuse and specular
R3	Asphalt with dark aggregates; slightly specular
R4	Smooth asphalt surface; specular
Pedestrian Interaction	
High	Large amount of pedestrian traffic during hours of darkness (100 or more pedestrians in one block during an hour): retail and entertainment areas
Medium	Moderate amount of pedestrian traffic during hours of darkness (11 to 100 pedestrians in one block during an hour): office area, apartment area, older city neighborhoods
Low	Little pedestrian traffic during hours of darkness (10 or less pedestrians in one block during an hour): rural and suburban areas

The luminaire cutoff classification, as described earlier, is used extensively in the design of roadway lighting. The luminaire cutoffs suppress glare to drivers and pedestrians while also alleviating light pollution and light trespass.

There are a number of other road types (e.g., sidewalks, bike paths, and intersections), environmental conditions (e.g., fog, rain, and wet roads), and special considerations with interaction with pedestrians. The reader is encouraged to consult Ref. 94 for these additional circumstances.

Sign Lighting A number of governmental standards have been developed for the lighting of road signs, but a set of guidelines have been developed by the IESNA.⁹⁵ Light for signs can be from external sources (e.g., car headlamps) or an internal source for a transmissive sign. Externally lit signs either make use of associated lights that illuminate it or retroreflectors that reflect back to the observer. In the United States the Federal Government provides the standards that must be met for lit roadway signage.⁹⁶

Tunnel Lighting Tunnel lighting is an important factor to increase drive safety while also allowing drivers to maintain their speed. In the United States a structure is considered a tunnel when it is greater than 25 m in length.⁹⁷ Distances greater than this require additional lighting to supplement any available daylight. Tunnel lighting is broken up into a number of regions including the approach, adaptation, threshold, transition, interior, and exit zones. Each of these zones has different guidelines to ensure driver comfort; however, these guidelines are dependent on the time of the year, the average speed of traffic, and the presence of oncoming traffic in undivided tunnels. For more information please consult Ref. 97.

40.8 ACKNOWLEDGMENTS

We are grateful to Jim McGuire and Kevin Thompson for providing valuable feedback on this chapter. We would like to acknowledge the financial support on literature purchase for this chapter by Optical Research Associates and our past and present employers for allowing us the use of their software: Optical Research Associates for LightTools, Photon Engineering for FRED, Lambda Research for TracePro and Breault Research Organization for ASAP.

40.9 REFERENCES

1. M. Johnson, D. G. Stork, S. Biswas, and Y. Furuichi, "Inferring Illumination Direction Estimated from Disparate Sources in Painting: An Investigation into Jan Vermeer's Girl with a Pearl Earring," *SPIE Symposium on Electronic Imaging*, San Jose, CA, 2008.
2. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, Chapter 3.
3. IESNA 1988 Lighting for the Aged and Partially Sighted Committee, *Recommended Practice for Lighting and the Visual Environment for Senior Living*, RP-28-98, 1998.
4. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 87.
5. M. D. Egan and V. Olgay, *Architectural Lighting*, 2nd edition, McGraw Hill, New York, 2002, p. 219.
6. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, 2003, p. 141.
7. M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000, pp. 10–13.
8. Commission Internationale de l'Éclairage (CIE), *International Lighting Vocabulary*, 4th edition, Publication 17.4, 1987.
9. Commission Internationale de l'Éclairage (CIE), *Colorimetry*, 3rd edition, Publication 15, 2004.
10. Commission Internationale de l'Éclairage (CIE), *Method of Measuring and Specifying Color Rendering Properties of Light Sources*, Publication 13.3, 1995.
11. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 162.
12. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 171.
13. Commission Internationale de l'Éclairage (CIE), *CIE Equations for Disability Glare*, Publication, 146, 2002.
14. Commission Internationale de l'Éclairage (CIE), 1995c, *Recommendations for the Lighting of Roads for Motor and Pedestrian Traffic*, CIE Technical Report 115, Vienna, 1995.
15. S. K. Guth, "A Method for the Evaluation of Discomfort Glare," *Illumination Engineering*, 57:351–64 (1963).
16. Chartered Institution of Building Services Engineers (CIBSE), *Technical Memorandum 10: The Calculation of Glare Indices*, CIBSE, London, 1985.
17. D. Fischer, "The European Glare Limiting Method," *Lighting Research and Technology*, 4:97–100, 1972.
18. G. Sollner, "Glare from Luminous Ceilings," *Lichttechnik*, 24:557–560, 1972.
19. G. Sollner, "Subjective Appraisal of Discomfort Glare in High Halls, Illuminated by Luminaires for High Intensity Discharge Lamps," *Lichttechnik*, 26:169–172, 1974.
20. Commission Internationale de l'Éclairage (CIE), 1995b, *Discomfort Glare in Interior Lighting*, Publication 117, 1995.
21. Commission Internationale de l'Éclairage (CIE), *Glare from Small, Large and Complex Sources*, Publication 147, 2002.
22. R. Levin, "Position Index in VCP Calculations," *Journal of the Illuminating Engineering Society*, pp. 99–105, January 1975. (Equation adapted from M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000, p. 9–26.)
23. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, Chapter 10.
24. R. B. Gibbons and C. J., Edwards, "A Review of Disability and Discomfort Glare Research and Future Direction," *18th Biennial TRB Visibility Symposium*, College Station, TX, 2007.
25. H. H. Bjorset and E. A. Frederiksen, "A Proposal for Recommendations for the Limitation of the Contrast Reduction in Office Lighting," *Proc. CIE 19th session*, Kyoto, Japan, 1979.
26. R. Whitehead, *Residential Lighting—A Practical Guide*, John Wiley and Sons, Newark, NJ, 2004, p. 5.
27. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–17.
28. M.S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–21.
29. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–29.
30. J. C. Stover, *Optical Scattering: Measurement and Analysis*, 2nd edition, SPIE, Bellingham, WA, 1995.
31. A. Gupta, *Simulation in Light Tools*, an illumination design software by Optical Research Associates, 2008.
32. Wikipedia, en.wikipedia.org/wiki/IGES (as of April 17, 2008).

33. U. S. Product Data Association, *Initial Graphics Exchange Specification IGES 5.3*, Charlestown, SC (1996), downloadable PDF version available at www.uspro.org/documents/IGES5-3 for Download. pdf (1997, as of April 17, 2008).
34. ISO TC 184/SC4, *STEP Application Handbook ISO 10303*, Version 3, Prepared by SCRA (North Charleston, SC, 2006); downloadable PDF version available at www.tc184_sc4.org/SC4_Open/SC4_Standards_Developers_Info/Files/STEP_application_handbook_63006.pdf (2006, as of April 17, 2008).
35. Wikipedia, en.wikipedia.org/wiki/ISO_10303 (as of April 17, 2008).
36. M. S. Kaminski, K. J. Garcia, M. A. Stevenson, M. Frate, and R. J. Koshel, "Advanced Topics in Source Modeling," *SPIE Proc. of Source Modeling I* **4775**, 46 (Society of Photo-Optical Instrumentation Engineers, Bellingham, WA, 2002).
37. M. A. Stevenson, M. S. Kaminski, M. Frate, and R. J. Koshel, "Modeling Filament-Based Sources for System Tolerancing," *SPIE Proc. of Modeling and Characterization of Light Sources* **4775**, 67, Bellingham, WA, 2000.
38. T. McReynolds and D. Blythe, *Advanced Graphics Programming Using OpenGL*, Elsevier Morgan Kaufmann, San Francisco, CA (2005).
39. F. X. Sillion and C. Puech, *Radiosity and Global Illumination*, Morgan Kaufmann, San Francisco, CA (1994).
40. R. J. Koshel, "Lit Appearance Modeling of Illumination Systems," *SPIE Proc. of Novel Optical Systems Design and Optimization V* **4768**:65 (Society of Photo-Optical Instrumentation Engineers, Bellingham, WA, 2002).
41. D. Scott and R. J. Koshel, Rendering in "ASAP," optical illumination software by Breault Research Organization.
42. R. J. Koshel, Rendering in "LucidShape" and "LucidDrive," vehicular illumination software by Brandenburg, GMBH.
43. M. Zollers, Photorealistic rendering in "LightTools," optical illumination software by Optical Research Associates.
44. D. Scott, Photorealistic rendering in "FRED," optical illumination software by Photon Engineering.
45. R. E. Smallwood, "Refractory Metals and Their Industrial Applications: Asymposium," ASTM committee B-10 on Reactive and Refractory Metals and Alloys, 1982.
46. S. Shionoya and W. M. Yen, *Phosphor Handbook*, CRC Press LLC, Boca Raton, FL, 1999.
47. A. M., Srivastava and C. R. Ronda, "Phosphors," *The Electrochemical Society Interface*, summer 2003.
48. N. Tesla, "Experiments with Alternate Currents of Very High Frequency and Their Application to Methods of Artificial Illumination," delivered before the American Institute of Electrical Engineers, Columbia College, NY, May 20, 1891.
49. W.J. Cassarly and T.L.R. Davenport, "Non-Rotationally Symmetric Mixing Rods," *Proceedings of SPIE*, 6342, July 2006.
50. W.J. Cassarly, "High-Brightness LEDs," *Optics and Photonics News*, 19–23, January 2008.
51. "LEDs Move into the Ultraviolet," www.physicsworld.com, May 17, 2006
52. W.J. Cassarly and T.L.R. Davenport, "Non-rotationally Symmetric Mixing Rods," *Proceedings of SPIE*, 6342, July 2006.
53. H. Chen, C. Hsu, H. Hong, "InGaN-CdSe-ZnSe Quantum Dots White LEDs," *Photonics Technology Letters*, IEEE **18**(1):193–195 (Jan. 1, 2006).
54. "Joint Venture to Make ZnSe White LEDs," <http://optics.org/cws/article/research/16534>, accessed on Oct. 13, 2008.
55. S. J. Smith, E. M. Purcell, "Visible Light from Localized Charges Moving across a Grating," *Physical Review* **92**(4):1069 (1953).
56. E. Yablanovich, "Inhibited Spontaneous Emission in Solid-State Physics and Electronics," *Physical Review Letters*, **55**:20, 2059, 1987.
57. J. J. Wierer, M. R. Krames, J. E. Epler, N. F. Gardner, J. R. Wendt, M.I. M. Sigalas, S. R. J. Brueck, D. Li, and M. Shagam, "III-nitride LEDs with Photonic Crystal Structures," *Proceedings of SPIE*, **5739**:102–107, 2005.
58. D. L. Barton and A. J. Fischer, "Photonic Crystals Improve LED Efficiency," *SPIE Newsroom*, 10.1117/2.1200603.0160, 2006.
59. C. W. Tang and S. A. VanSlyke, "Organic Electroluminescent Diodes," *Applied Physics Letters*, **51**:913, 1987.
60. K. Müllen and U. Scherf, *Organic Light Emitting Devices: Synthesis, Properties and Applications*, 1st edition, Wiley-VCH, Weinheim, March 2006.

61. E. Ne'eman, and R. G. Hopkinson, "Critical Minimum Acceptable Window Size: A Study of Window Design and Provision of a View," *Lighting Research and Technology*, **2:1**, 17–27, 1970.
62. E. C. Keighley, "Visual Requirements and Reduced Fenestration in Offices: A Study of Multiple Apertures and Window Area," *Building Science*, **8:4**, 321–331, 1973.
63. A. M. Ludlow, "The Functions of Windows in Buildings," *Lighting Research and Technology*, **8:2**, 57–68, 1976.
64. R. Winston, J. C. Miñano, and P. Benítez, *Nonimaging Optics*, Elsevier Academic Press, Burlington, MA, 2005.
65. J. Chaves, *Introduction to Nonimaging Optics*, CRC Press, Taylor and Francis Group, Boca Raton, FL, 2008.
66. *American National Standard Practice for Roadway Lighting*, The Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, ANSI/IESNA RP-8-001, p. 6 IESNA, New York, NY, 1999, reaffirmed 2005.
67. *Luminaire Classification System for Outdoor Luminaires*, Luminaire Classification Task Group of IESNA, IESNA TM-15-07 (revised), IESNA, New York, NY, 2007.
68. W. J. Cassarly, D. Jenkins, A. Gupta, R. J. Koschel, "Hidden Devices That Light Our World Lightpipes," *Optics and Photonics News*, 34–39, August 2001.
69. J. K. Holton, "Daylighting of Buildings," *US National Bureau of Standards*, NBSIR 76-1098, October 1976.
70. W. M. C. Lam, *Sunlight as Formgiver of Architecture*, Von Nostrand Reinhold, New York, 1986.
71. M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000.
72. M. D. Egan and V. Olgyay, *Architectural Lighting*, 2nd edition, McGraw Hill, New York, NY, 2002.
73. Commission Internationale de l'Éclairage, *Methods of Characterizing Illuminance Meters and Luminance Meters: Performance, Characteristics and Specifications*, Publication 69, Bureau Central de la CIE, Vienna, 1987.
74. A. V. Arecchi, T. Messadi, and R. J. Koschel, *Field Guide to Illumination*, SPIE Press, Bellingham, WA (2007), pp. 90–92.
75. M. S. Kaminski and R. J. Koschel, "Methods of Tolerancing Injection-Molded Parts for Illumination Systems," *Proceedings of SPIE, Design of Efficient Illumination Systems*, 5186, 61, Bellingham, WA, 2003.
76. R. J. Koschel, "Illumination System Tolerancing," *Proceedings of SPIE, Optical System Alignment and Tolerancing*, 6676, 667604, Bellingham, WA, 2007.
77. *Recommended Practice for Lighting Merchandise Areas (A Store Lighting Guide)*, The IESNA Merchandise Lighting Committee, IESNA RP-2-01, IESNA, New York, NY, 2001. p. 2.
78. *Ibid*, p. 5.
79. *Lighting for Hospitals and Health Care Facilities*, The IESNA Committee for Health Care Facilities, ANSI/IESNA RP-29-06, IESNA, New York, NY, 2006, p. 4.
80. *Recommended Practice for Lighting Industrial Facilities*, Appendix A2, by the IESNA Industrial Lighting Committee, New York, NY, 2001.
81. *Lighting for Exterior Environments and IESNA Recommended Practice*, The IESNA Outdoor Environment Lighting Committee, IESNA RP-33-99, IESNA, New York, NY, 1998.
82. *Addressing Obtrusive Light (Urban Sky Glow and Light Trespass) in Conjunction with Roadway Lighting*, the Obtrusive Light Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-10-00, IESNA, New York, NY, 2000.
83. International Dark Sky Association, U.S. House of Representatives Briefing by Lee Cooper, June 20, 2008.
84. *Light Trespass: Research Results and Recommendations*, The Obtrusive Light Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-11-00, IESNA, New York, NY, 2000.
85. *A Discussion of Appendix E—"Classification of Luminance Light Distributions"*, "Roadway Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-3-95, IESNA, New York, NY, 1995.
86. *Code of Federal Regulations*, Title 49 Transportation, Volume 6, Chapter 5, Part 571, Federal Motor Vehicle Safety Standards, USA Federal Government (Washington, D.C., 2007), pp. 279–353; online at http://edocket.access.gpo.gov/cfr_2007/octqtr/pdf/49cfr571.108.pdf. Accessed 13 October 2008.
87. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 8.
88. *UNECE Standards*, ECE Transaction 505, Revision 2, Addendum 112, Regulation No. R1 13, online at <http://www.unece.org/trans/main/wp29/wp29regs/r113e.pdf>, 2001. Accessed 13 October 2008.

89. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 9.
90. *UNECE Standards*, ECE Transaction 505, Revision 2, Addendum 111, Regulation No. R112, online at <http://www.unece.org/trans/main/wp29/wp29regs/r112e.pdf>, 2001. Accessed 13 October 2008.
91. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 134.
92. *UNECE Standards*, ECE Transaction 505, Revision 4, Addendum 6, Regulation No. R7, online at <http://www.unece.org/trans/main/wp29/wp29regs/r007r4e.pdf>, 2006. Accessed 13 October 2008.
93. *Design Guide for Roadway Lighting Maintenance*, The Subcommittee on Lighting Maintenance and Light Sources of the IESNA Roadway Lighting Committee, IESNA DG-4-03, IESNA, New York, NY, 2003.
94. *American National Standard Practice for Roadway Lighting*, The Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, ANSI/IESNA RP-8-001, IESNA (New York, NY, 1999, reaffirmed 2005), p. 2.
95. *IESNA Recommended Practice for Roadway Sign Lighting*, The Sign Lighting Subcommittee of the IESNA Roadway Lighting Committee, IESNA RP-19-01, IESNA, New York, NY, 2001.
96. *Manual on Uniform Traffic Control Devices*, Sections 2A-11 to 2A-15, Federal Highway Administration, U.S. Department of Transportation, Washington, D.C., 1988.
97. *American National Standard Practice for Tunnel Lighting*, The Tunnel Lighting Subcommittee of the IESNA Roadway Lighting Committee, IESNA RP-22-05, IESNA, New York, NY, 2005.

This page intentionally left blank.

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe illuminated eyepieces, **12.12**, **12.12f**
Abbe illumination system, **39.23**, **39.23f**,
 39.34, **39.35f**
Abbe's sine condition, **34.19**
Aberration coefficients, **3.10–3.11**
Aberration curves (in lens design), **2.1–2.6**
 considerations for, **2.5–2.6**
 field plots of, **2.4–2.5**
 transverse ray plots of, **2.2–2.4**
Aberrations:
 balancing of, **11.30**, **11.35–11.36**, **11.36t**
 evaluation of, **3.9–3.11**
 (See also *specific aberrations*, e.g.: Axial chromatic aberration)
Absolute detectors, **34.27–34.30**
 electrical substitution radiometers, **34.27–34.29**
 photoionization devices, **34.29**
 predictable quantum efficiency devices,
 34.29–34.30
Absolute measurements, **34.20–34.37**
 accuracy and traceability of, **34.21**
 error propagation in, **34.22**
 error types in, **34.21–34.23**
 relative vs., **34.20–34.21**
 and uncertainty estimates, **34.21–34.23**
Absolute responsivity units (A/W), **34.31**
Absolute sources (of radiation), **34.23–34.27**
 blackbody radiator, **34.23–34.24**
 blackbody simulators, **34.24–34.26**
 synchrotron radiation, **34.26–34.27**
Absorbing media:
 in photodetectors, **26.4f**, **26.5**
 radiant power transfer through, **34.13**
Absorbing substrate (AS) chips, **17.7**, **17.7t**
Absorptance:
 defined, **35.4**
 measurement of, **35.10**
 in thermal detectors, **28.2**
 and transmittance/reflectance, **35.7**, **35.8**, **35.8t**
Absorption:
 defined, **35.4**
 quantum resonance, **22.16**, **22.17**
 stimulated, **16.7–16.8**, **16.8f**
Absorption coefficient, **32.2–32.4**, **32.3f**
 of *pin* photodiodes, **25.8**, **25.9f**
 of visible array detectors, **32.2–32.3**, **32.3f**
Absorption rate, **23.8**
Ac lamps, **15.32f**
Accent lighting, **40.14**, **40.14f**
Accuracy:
 of absolute measurements, **34.21**
 of CGHs, **14.6–14.7**, **14.6f**, **14.7f**
 as measure of systemic errors, **12.2**
Achromatic lenses, athermalized, **1.16**, **1.16f**
Achromatism, **1.14–1.15**, **1.15f**
Actinic effects (of radiation), **34.6**, **34.7**
Actinic ultraviolet action spectrum, **36.17**
Actinometry:
 conversions between radiometry/photometry
 and, **34.12**, **34.12t**
 defined, **34.7**, **34.11**
Activated-phosphor sources (of radiation), **15.49**
Active athermalization, **6.24**, **6.24f**
Active imaging, **31.29–31.30**
Active mechanical athermalization, **8.11**, **8.11f**
Acutance:
 defined, **30.2**
 of photographic systems, **29.17–29.19**,
 29.18t, **29.19f**
Adaptation, in vision, **40.9**
Additive damping, **3.18**
Adiabatic approximation, **23.21** (See also
 Markov approximation)
Advanced Photo System (APS), **30.21**,
 30.26, **30.27t**
Afocal systems:
 as attachments, **1.8**, **1.9f**
 first-order layout for, **1.7**, **1.7f**

- Agfachrome, **29.14**
- AHU pelloid, **30.4**
- Air-spaced doublet lens, **6.7**
- Air-spaced triplet lens, **6.21, 6.22f**
- Airway beacon lamps, **15.11**
- Allan Deviation, **22.2–22.4**
- Allan Variance method, **22.2, 22.3**
- Alloy disordering, **19.24**
- Alphanumeric displays, LED, **17.31–17.32**
- Aluminized phosphor-screen/window assembly, **31.14, 31.15f**
- Aluminum, diamond turning and, **10.4**
- Aluminum gallium arsenide (AlGaAs) emitters, **17.32**
- Aluminum gallium arsenide (AlGaAs) LEDs, **17.17, 17.17t, 17.28f**
- Aluminum gallium arsenide (AlGaAs) quantum well photodetectors, **25.16–25.17, 25.16f, 25.17f**
- Aluminum gallium arsenide (AlGaAs) substrate, **17.22**
- Aluminum gallium nitride (AlGaN) alloy photovoltaic detectors, **24.46**
- Aluminum gallium nitride (AlGaN) substrate, **17.22**
- Aluminum indium gallium nitride (AlInGaN) material systems, **18.1–18.2, 18.2f, 18.4**
- Aluminum indium gallium phosphide (AlInGaP) LEDs, **17.18, 17.19f**
- Aluminum indium gallium phosphide (AlInGaP) material systems, **18.1, 18.5**
- Aluminum indium gallium phosphide (AlInGaP) substrate, **17.22**
- Aluminum mirrors, mounting of, **6.20f**
- Ambient lighting, **40.12, 40.13f, 40.15f**
- Ambient temperature electrical substitution radiometers, **34.27**
- American Institute of Physics (AIP), **36.3**
- American National Standards Institute (ANSI), **4.11, 36.2**
- American Society for Testing and Materials (ASTM), **37.11**
- AMI (amplified MOS imager) MOS readout, **32.21**
- Amorphous silicon photoconductors, **32.4f, 32.31, 32.32**
- Amplification, **16.9**
- Amplifier strategies, for PZTs, **22.18**
- Amplifiers, **27.2**
 - properties of, **16.3**
 - selection of, **27.10–27.12**
 - transconductance, **27.11–27.12, 27.11f**
 - voltage, **27.10–27.11**
- Amplitude gating, **21.7**
- Amplitude modulation (AM), **19.36**
- Amplitude response, frequency vs., **22.6–22.7**
- Analytical density, **29.14**
- ance (suffix), **35.3**
- Angle measurement, **12.10–12.17**
 - autocollimeters for, **12.11–12.12, 12.11f, 12.12f**
 - interferometric methods of, **12.14**
 - levels (tools) for, **12.13–12.14, 12.13f, 12.14f**
 - mechanical methods of, **12.10–12.11, 12.11f**
 - in prisms, **12.14–12.16, 12.15f–12.17f**
 - theodolites for, **12.13**
- Angle solves, **3.6**
- Angular dilution, **39.6**
- Angular uniformity, **39.31**
- Annular flanges, **6.3–6.4, 6.4f**
- Annular polynomials:
 - for defocus, **11.38f, 11.39**
 - Zernike, **11.13–11.21, 11.14f, 11.16f, 11.17t–11.21t**
- Annunciator assemblies, **17.30**
- Anodes, in photomultipliers, **27.6, 27.7, 27.7f**
- Anomalous reflection colors, **30.17**
- Antiblooming, **32.9, 32.10f**
- Antihalation undercoat (AHU) layers, **30.4**
- AOM transducers, **22.20**
- APART (stray light analysis program), **7.11**
- Aperture(s):
 - data about, **3.4**
 - nonideal, **34.35–34.36, 34.35f**
 - numerical, **34.20, 39.1**
 - in optical design software, **3.6**
- Aperture flash mode, **39.7**
- Aperture placement (in stray light suppression), **7.5–7.10**
 - aperture stops, **7.6–7.7, 7.7f, 7.8f**
 - field stops, **7.7, 7.8f, 7.9f, 34.18–34.19, 34.19f**
 - Lyot stops, **7.8–7.10, 7.8f–7.11f**
- Aperture stops, **1.4, 3.4, 7.6–7.7, 7.7f, 7.8f, 34.18, 34.19f**
- Aphakic hazard, **36.17**
- Aplanatic optical systems, **34.19–34.20**
- Aplanats, **39.8, 39.9f**
- Apodization, **39.7**

- Apostilb (unit), 34.43, 36.7, 36.8*t*
 Apovortex surfaces, 39.11
 Approximate transfer factor (ATF), 4.1, 4.4
 Arc radiation sources (*see specific arcs, e.g.:*
 Argon arcs)
 Area image sensor arrays, 32.24–32.32
 about, 32.2
 CCD
 frame transfer, 32.26–32.28, 32.27*f*, 32.28*f*
 interline transfer, 32.28–32.32,
 32.29*f*–32.31*f*
 performance of, 32.32
 image area dimensions for, 32.25*t*
 MOS, 32.25–32.26, 32.26*f*
 Area-solid angle product, 39.5
 Argon arcs, 15.12, 15.13
 Argon ion lasers, 16.14, 16.15*f*, 16.30
 Array-mode selection, semiconductor, 19.27
 Array-mode stability, semiconductor, 19.27
 Artificial sources (of radiation), 15.3–15.53
 about, 15.3–15.4
 commercial, 15.13–15.53
 activated-phosphor sources, 15.49
 blackbody simulators, 15.14, 15.15*f*, 15.16*f*
 carbon arc sources, 15.21–15.24, 15.23*f*,
 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
 concentrated arc lamps, 15.47–15.48,
 15.48*f*, 15.49*f*
 glow modulator tubes, 15.49, 15.50*f*,
 15.51*f*, 15.52*t*
 high-energy sources, 15.40
 high-pressure enclosed arc, 15.24,
 15.28–15.34, 15.29*f*–15.35*f*
 hydrogen and deuterium arc lamps,
 15.49, 15.53*f*
 incandescent nongaseous sources,
 15.15–15.21, 15.17*f*–15.22*f*
 low-pressure enclosed arc, 15.35–15.47,
 15.36*f*, 15.36*t*–15.43*t*, 15.44*f*–15.47*f*,
 15.46*t*, 15.47*t*
 special-purpose sources, 15.53
 luminaire optics for, 40.45–40.47, 40.45*f*,
 40.46*f*
 and radiation law, 15.4–15.7, 15.5*f*, 15.5*t*, 15.6*f*
 standardized laboratory sources, 15.7–15.13
 baseline standard of radiation, 15.9, 15.9*f*,
 15.10*f*, 15.12*f*
 blackbody cavity theory, 15.7–15.9, 15.8*f*
 working standards of radiation,
 15.9–15.13, 15.10*f*, 15.12*f*, 15.13*f*
 ASAP (optical software), 7.25
 Aspheric lenses, 39.8, 39.9, 39.9*f*
 Aspheric measuring system, 10.12*f*
 Aspheric surfaces, 3.5
 Aspherical optics fabrication, 9.7–9.8, 9.7*f*
 Aspherical wavefront measurement,
 13.23–13.27
 holographic compensators, 13.25, 13.25*f*,
 13.26*f*
 infrared interferometry, 13.25
 Moiré tests, 13.26–13.27
 refractive or reflective compensators, 13.24,
 13.24*f*, 13.25
 sub-Nyquist interferometry, 13.27
 two-wavelength interferometry, 13.25, 13.26
 wavefront stitching, 13.27, 13.27*f*
 Assembly tolerances, 5.8
 Astigmatism, 2.3, 2.3*f*
 Astronomical telescopes, 1.7*f*
 Athermal laser beam expanders, 8.13–8.14
 Athermalization, 1.15–1.16, 1.16*f*, 6.22–6.24
 active, 6.24, 6.24*f*
 intrinsic, 8.7–8.8, 8.7*f*
 mechanical, 8.8–8.12
 active, 8.11, 8.11*f*
 by image processing, 8.12
 part active, part passive, 8.11–8.12, 8.12*f*
 passive, 8.8–8.10, 8.8*f*–8.10*f*
 optical, 8.12–8.15, 8.13*t*
 about, 8.12–8.13
 athermal laser beam expanders, 8.13–8.14
 diffractive optics usage, 8.15
 of separated components, 8.14, 8.15
 three-material solutions, 8.14, 8.14*t*, 8.15*t*
 passive, 6.22, 6.23*f*, 6.24
 single material design, 6.22, 6.23*f*
 Athermalized achromatic lenses, 1.16, 1.16*f*
 Atomic (gain) noise, 23.34–23.35
 Attosecond optics, 21.1–21.9
 about, 21.2
 driving lasers in, 21.4–21.6
 carrier-envelope offset frequency, 21.5
 carrier-envelope phase, 21.5*f*, 21.6
 carrier-envelope phasemeter, 21.6
 chirped pulse amplification, 21.5
 chirped pulse amplifiers, 21.6
 single-shot *f*-to-2*f* interferometer, 21.6
 high-harmonic generation, 21.2, 21.2*f*
 phase-matching in, 21.4
 ponderomotive potential in, 21.3

- Attosecond optics (*Cont.*):
 pulse characterization, 21.8–21.9, 21.8f
 attosecond streak camera, 21.9
 FROG-CRAB, 21.9
 RABITT, 21.9
 second-order autocorrelator, 21.9
 pulse generation, 21.6–21.8, 21.7f
 amplitude gating, 21.7
 attosecond pulse train, 21.6, 21.7
 double optical gating, 21.8
 polarization gating, 21.7–21.8
 two-color gating, 21.7
 quantum trajectories in, 21.3–21.4
 semiclassical model of, 21.3
 single isolated pulses in, 21.4
 strong field approximation in, 21.3
- Attosecond pulse, 21.8–21.9, 21.8f
 attosecond streak camera, 21.9
 FROG, 21.9
 FROG-CRAB, 21.9
 generation of, 21.6–21.8, 21.7f
 amplitude gating, 21.7
 attosecond pulse train, 21.6, 21.7
 double optical gating, 21.8
 polarization gating, 21.7–21.8
 two-color gating, 21.7
 RABITT, 21.9
- Attosecond pulse train, 21.6, 21.7
 Attosecond streak camera, 21.9
 Augur recombination, 19.17, 19.17f
- Autocollimeters:
 angle measurement with, 12.11–12.12,
 12.11f, 12.12f
 curvature measurement with, 12.19–12.20,
 12.20f
 defined, 12.11–12.12
- Automatic brightness control (ABC), 31.18
 Automatic spherometers, 12.19
 Autotest levels (tools), 12.14, 12.14f
- Avalanche multiplication, 25.9
 Avalanche photodetectors (APDs):
 high-speed, 26.17–26.20, 26.18f, 26.20f, 26.21f
 improvements in, 26.3
- Avalanche photodiodes, 24.62–24.70,
 24.63f–24.70f, 24.72–24.73, 24.72f, 24.73f,
 25.8–25.10, 25.9f
 defined, 24.10
 germanium, 24.70f, 24.72–24.73, 24.72f, 24.73f
 InGaAs, 24.66–24.70, 24.66f–24.69f
 silicon, 24.62–24.65, 24.63f–24.66f
- Avogadro's number, 34.11
 Axial chromatic aberration, 1.14, 2.2, 2.3f
 Axial gap prevention, 6.21, 6.22f
 Axial rays, 1.4, 1.11f, 1.12
 Azimuth angle, 35.5
 Azomethine dyes, 30.10f
 excited state properties of, 30.11–30.12,
 30.12f
 formation of, 30.10
 photochemistry of, 30.11
- Back light, 40.43, 40.44f
 Background temperature, 24.10
 Background-limited performance (BLIP), of
 infrared detector arrays, 33.24
 Backing, film, 29.4
 Backlighting, 40.1, 40.12, 40.47, 40.47f, 40.48f
 Backscattering, 20.13–20.15, 20.15f
 Backward trace, 39.7
 Baffles:
 cone-shaped secondary, 7.3–7.4, 7.3f, 7.4f
 with integrating cavities, 39.26
 in lighting design, 40.41, 40.45f, 40.46
 shields for, 7.9–7.10, 7.9f, 7.10f
 in stray light suppression, 7.10, 7.11
 two-stage, 7.10
- Balanced spherical aberrations, 11.30
 Ballasts:
 in fluorescent lamps, 40.32–40.33
 in HID lamps, 40.36
- Band pass, 38.8
 Bandwidth:
 of amplifiers, 27.10
 gain-bandwidth, 26.17
 normalization of, 36.14–36.16, 36.15t, 36.16f
 of photomultipliers, 27.7
- “Bang-bang” zoom, 1.12
 Banker lamps, 40.12, 40.46, 40.46f
 Bar spherometers, 12.19, 12.19f
 Bar-code reading, 17.34
 Bare source light, 40.43
 Barium strontium titanate (BST), 28.11, 28.12
 Baryta layer, 30.5
 Batwing lenses, 40.12
 Baud rate, 17.33
 Beacon lamps, airway, 15.11
 Beam splitters, 13.7, 34.32
 Beam transformers, 39.18, 39.18f
 Beam-forming illumination systems,
 39.22, 39.39

- Beam-smearing faceted reflectors, 39.39
- Beer-Lambert law (Beer's law), 16.9, 34.35, 38.5
- Bell-clamping (edging fabrication step), 9.6
- Bellcore, 19.39, 19.41
- Bevel gauges, 12.10, 12.11, 12.11f
- Bevel placement (on vanes), 7.13, 7.14f
- Bias angle, 31.13
- Biased *pin* photodetectors, 26.6f
- Biconical reflectance, 35.5t, 35.6f, 35.6t
- Bidirectional reflectance, 35.5t, 35.6f, 35.6t
- Bidirectional reflectance distribution function (BRDF), 7.1, 7.18–7.19, 7.22, 35.5, 35.13, 37.9
- Bidirectional scattering distribution function (BSDF), 7.2, 7.23, 7.24f, 7.25f, 35.13
- Bidirectional transmittance distribution function (BTDF), 35.3, 35.13
- Bihemispherical reflectance, 35.5t, 35.6f, 35.6t
- Binning, 38.10
- Bipolar transistors, 27.11
- Black-and-white (B&W) film, 29.4, 30.24–30.25, 30.25t
- Blackbody cavity theory, 15.7–15.9
- Blackbody D star, 24.10
- Blackbody detectivity, 24.10
- Blackbody noise-equivalent power, 24.10
- Blackbody radiation, 15.4–15.6, 15.5t
emittance of, 34.25–34.26
sources of, 15.14, 15.15f, 15.16f, 34.23–34.24
temperature vs., 36.12, 36.12f, 36.14, 36.14f
working standards for, 15.14, 15.16f
- Blackbody responsivity, 24.10
- Blackbody simulators, 15.14, 15.15f, 15.16f, 34.24–34.26
- Black-light fluorescent lamps, 15.35, 15.36t
- Bleaching, in film development, 29.14
- Blindness, flash, 40.9
- Blip detector (blip condition), 24.10
- Blocked impurity band (BIB), 33.7
- Blocking contacts, 26.3
- Blocking filters, 38.8
- Blocking gas analysis, 17.34
- Blooming:
antiblooming, 32.9, 32.10f
in image sensors, 32.6, 32.9
- Blue emitters, in LED technology, 17.18, 17.19
- Blue light, color film and, 29.13, 29.13f, 30.3–30.4
- Blue semiconductor lasers, 19.7
- Blue-enhanced photodiodes, 24.55f, 24.61–24.62, 24.61f, 24.62f
- Blur filters, 32.34, 32.34f
- “Boat grown” technique, 17.21
- Bode representation, of servo system, 22.5–22.6, 22.6f
- Bolometers, 24.5, 28.3–28.5, 28.4f
about, 28.1
carbon, 28.5
detectivity of perfect, 24.17, 24.18f
germanium low-temperature, 24.31–24.32, 24.32f, 24.33f
indium antimonide hot-electron, 24.29, 24.30, 24.30f, 24.31f
as infrared detectors, 33.9–33.10
metal, 28.4, 28.7t
properties of, 28.7t
resistive arrays of, 28.10–28.11, 28.10f
semiconductor, 28.4–28.5
superconducting, 28.5
thermistor, 24.24–24.25, 24.24f, 24.25f, 28.7t
- Bonded mountings, 6.13–6.15, 6.15f, 6.16f
- Boresight tolerances, 5.8
- Boron-doped silicon (Si:B) detectors, 24.95f, 24.96
- Bose-Einstein condensation (BEC), 23.39
- Bose-Einstein statistics, 23.9
- Bouillotte lamps, 40.12, 40.46, 40.46f
- Boules, glass, 9.3
- Boundary conditions (of optics):
defined, 3.17
methods for handling, 3.18–3.19
specification of, 4.12
- Boxcar averaging, of modulated signal sources, 27.13, 27.13f, 27.15
- Bragg reflectors, 19.41
- Bridgeman technique, 17.21
- Brightness:
of carbon arcs, 15.23f
luminance vs., 34.40
perception of, 40.4, 40.4f
of scene, 31.4–31.5
- British Glare Index (CIBSE), 40.10
- Broad bandwidth solid-state lasers, 16.34–16.35, 16.34f
- Bromine, in light bulbs, 40.30
- Brunning distance-measuring interferometers, 12.8–12.9, 12.8f
- Buffered direct injection (BDI), 33.19t, 33.20f, 33.21–33.22

- Build-and-test evaluation (for stray light suppression), 7.28
- Built-in potential, 25.6
- Bulb blackening, 40.30
- Bulb shield, 40.45*f*, 40.46
- Bulk material photodetectors, 26.4*f*, 26.5
- Bulk-grown materials, 17.8
- Buried crescent lasers, 19.24, 19.25*f*
- Buried heterostructure (BH) lasers, 19.8, 19.9*f*, 19.20*t*, 19.24, 19.36*f*
- Buried TRS (BTRS) lasers, 19.19, 19.20*t*, 19.21*f*
- Buried V-groove-substrate inner stripe (BVSIS), 19.19, 19.20*t*
- Buried-channel CCDs (BCCDs), 32.14, 33.13
- Buried-channel MOS capacitors, 32.4*f*, 32.7–32.8
- Burnished mounting, 6.3, 6.3*f*

- Cable TV (CATV), 25.11–25.12
- Cadmium selenide (CdSe) photoconductors, 24.49–24.52, 24.52*f*
- Cadmium sulfide (CdS) photoconductors, 24.49–24.52, 24.51*f*–24.53*f*
- Cadmium telluride (CdTe) detectors, 24.52, 24.54, 24.54*f*
- Cadmium zinc telluride (CdZnTe) detectors, 24.52
- Calibration:
 - artificial sources of radiation for (*see* Artificial sources (of radiation))
 - legal traceability of, 34.21
 - photometric, 34.42–34.43
 - radiometric, 34.31–34.32
 - self-calibration, 34.29
 - spectroradiometric, 38.11–38.13, 38.11*t*, 38.12*f*
- Calibration transfer devices, 34.31–34.32
- Callier coefficient, 29.7
- Candela (unit), 34.37, 34.39, 36.4, 36.5, 37.3, 37.4
- Candle power, 37.3
- Capacitive bolometers, 33.10
- Capacitive transimpedance amplifier (CTIA), 33.19*t*, 33.20*f*, 33.22–33.23
- Capacitors, MOS, 32.4*f*, 32.7–32.8
- Capillary mercury-arc lamps, 15.30–15.31, 15.31*f*
- Carbon arc light sources, 40.40
- Carbon arc sources (of radiation), 15.21–15.24, 15.23*f*, 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
- Carbon bolometers, 28.5, 28.7*t*
- Carbon-dioxide lasers, 16.16, 16.16*f*, 16.30
- Carey Lea silver (CLS), 29.13, 30.4
- Carrier confinement, 17.12, 17.12*f*–17.14*f*, 17.13, 17.17
- Carrier density, 19.30–19.33
- Carrier transit time, 26.6–26.7, 26.6*f*
- Carrier trapping, 26.9, 26.9*f*
- Carrier-envelope offset, 20.4
- Carrier-envelope offset frequency, 21.5
- Carrier-envelope (CE) phase:
 - of chirped pulse amplifiers, 21.6
 - of lasers, 21.5, 21.5*f*
- Carrier-envelope phasemeters, 21.6
- Cassegrain design, 7.3*f*, 7.11, 7.14, 7.14*f*, 7.16, 7.16*f*, 7.19, 7.20*f*
- Cathodes:
 - photo-, 27.6, 27.7*f*
 - shielding of, 27.10
- Cavity(-ies):
 - distributed feedback lasers, 16.29
 - integrating (*see* Integrating cavities, of nonimaging optics)
 - mode-locking, 16.27–16.29, 16.28*f*
 - modifying output distribution of, 39.27
 - properties of, 16.3
 - Q-switching, 16.26–16.27, 16.27*f*
 - ring lasers, 16.29
 - stability of, 16.23–16.25, 16.24*f*, 16.25*f*
 - unstable resonators, 16.25–16.26, 16.26*f*
- Cavity dumping, 16.27
- Cavity losses, 23.18
- Cavity-shaped radiometers, 34.28
- Ceilings, illuminated, 40.13*f*
- Cellulose acetate film, 29.4
- Cenco Company, 15.47
- Center for Optics Manufacturing, 9.4
- Center-of-mass motion of atoms, 23.45
- Centrally obscured system (*see* Cassegrain design)
- Centration, of spherical lenses, 9.8
- Channel stop region, 32.7
- Channeled substrate planar (CSP) lasers, 19.20*t*, 19.21*f*, 19.22, 19.36*f*
- Charge integration matrix (CIM), 33.10–33.11, 33.11*f*, 33.12*f*
- Charge pumping, 32.9*n*
- Charge sweep devices (CSDs), 33.12*f*, 33.13

- Charge-coupled detector area image sensor arrays:
 frame transfer, 32.26–32.28, 32.27*f*, 32.28*f*
 interline transfer, 32.28–32.32, 32.29*f*–32.31*f*
 performance of, 32.32
- Charge-coupled detectors (CCDs), 25.10, 25.11, 25.11*f*, 31.1, 32.12–32.20, 38.9–38.10, 38.10*t*
 characteristics of, 32.17–32.20, 32.17*f*, 32.18*f*
 electronics of, 38.10
 image sensing with, 32.8
 linear arrays of, 32.21–32.24, 32.22*f*, 32.23*f*
 MIS photogate FPAs for, 33.10–33.11, 33.11*f*, 33.12*f*
 multilinear arrays of, 32.21, 32.23*f*, 32.24
 operation of, 32.12–32.14, 32.13*f*
 output of, 32.14–32.15, 32.15*f*
 performance of, 32.32
 readout from, 32.12–32.21, 32.13*f*
 types of, 32.15–32.17, 32.16*f*
- Charge-injection devices (CIDs), 31.1, 32.20, 33.10–33.11, 33.11*f*, 33.12*f*
- Chartered Institution of Building Services Engineers (CIBSE), 40.2
- Chemical beam epitaxy (CBE), 19.7
- Chemical-assisted ion beam etching (CAIBE), 19.39
- Chirped pulse amplification (CPA), 21.5, 21.5*f*, 21.6
- Chopper-stabilized amplifiers, 27.11
- Chopper-stabilized BDI, 33.19*t*, 33.20*f*, 33.21–33.22
- Chromatic aberrations, 2.2–2.4, 2.3*f*, 2.4*f*
- Chromium lasers, 16.34, 16.35
- Chromogenic film, 29.14
- Circle polynomials, 11.36*t*, 11.39
 isometric plots/interferograms/PSFs for defocus, 11.38*f*
 and noncircular pupils, 11.37, 11.39
 radial, 11.7, 11.9*f*–11.10*f*
 Zernike, 11.4, 11.6–11.12, 11.8*t*–11.9*t*, 11.9*f*–11.11*f*, 11.12*t*
- Circular discs, projected area of, 36.3*t*
- Cladding layers, 19.4
- Clarity, perception of visual, 40.5
- Clip test (of photographic film), 30.23
- Clipped Lambertian distribution, 39.3–39.4
- Clock generation, 33.16
- Closed-loop performance (in servo systems), 22.8
- Closed-loop stability issues (in servo systems), 22.8–22.12, 22.9*f*
 PID controller vs. notch filters, 22.10–22.11, 22.10*f*, 22.11*f*
 rule-of-thumb PID design for system with transducer resonance, 22.11–22.12
- Coarse-grained derivative, 23.21
- Coatings:
 lens specifications for, 4.10
 of photographic film, 29.4
- Coblentz-type thermopiles, 24.23
- Coherent states, 23.12
- Coiling, of light bulb filament, 40.30
- Cold cathode fluorescent lamps (CCFLs), 40.32
- Collares-Pereira, M., 39.17
- Collector power (in stray light suppression), 7.2
- Collectors (*see* Concentrators, nonimaging)
- Collimators:
 autocollimators, 12.12
 conic, 39.8, 39.9*f*
- Collisional broadening, emission-line, 16.5
- Colloidal silver, 29.13
- Color(s):
 anomalous reflection, 30.17
 in LEDs, 40.37
 and lighting design, 40.7–40.9
 mixing of, 40.8
 science of, 30.15–30.18, 30.16*f*, 30.17*f*
- Color aliasing, 32.34
- Color density, 29.7–29.8
- Color filter arrays (CFAs), 32.32–32.34, 32.33*f*, 32.34*f*
- Color imaging architectures, 32.32–32.34
 integral filter arrays, 32.32–32.34, 32.33*f*, 32.34*f*
 sequential, 32.32, 32.33*f*
 three-chip, 32.32, 32.33*f*
- Color negative films, 30.25–30.28, 30.27*t*
- Color photographic films:
 about, 30.2
 coating of, 29.4
 negative, 30.25–30.28, 30.27*t*
 reversal, 30.22–30.24, 30.23*t*
 structure of, 29.12–29.15, 29.13*f*, 29.14*f*, 30.3–30.5, 30.3*f*
- Color photographic paper, 30.5
- Color records, 30.4
- Color rendering, 40.8–40.9
- Color rendering index (CRI), 40.8

- Color reversal films, **30.2**, **30.22–30.24**, **30.23t**
 Color sequential systems, **32.32**, **32.33f**
 Color slide films (*see* Color reversal films)
 Color space calculations, **38.4–38.5**
 Color temperature, **34.44**, **37.4t**, **37.6–37.7**,
38.5, **40.8**
 Color transparency films (*see* Color reversal
 films)
 Color-center lasers, **16.35**
 Colorimetry, **37.11**
 Coma, with spherical aberration, **2.4**, **2.4f**
 Combined recombination, **17.3**
 Combined servo transducers, **22.19**
 Commercial sources (of radiation), **15.13–15.53**
 activated phosphor, **15.49**
 blackbody simulators, **15.14**, **15.15f**
 carbon arcs, **15.21–15.24**, **15.23f**, **15.24f**,
15.25t–15.27t, **15.28f**
 concentrated arcs, **15.47–15.49**, **15.48f**, **15.49f**
 glow modulator tubes, **15.49**, **15.50f**, **15.51f**,
15.52t
 high-energy, **15.40**
 high-pressure enclosed arc, **15.24**, **15.28–15.34**
 compact-source arcs, **15.31–15.34**,
15.32f–15.35f
 Lucalox lamps, **15.30**, **15.31f**
 mercury arcs, **15.29–15.31**, **15.30f**, **15.31f**
 multivapor arcs, **15.29**, **15.31f**
 Uviarc, **15.28–15.29**, **15.29f**, **15.30f**
 hydrogen and deuterium arcs, **15.49**, **15.53f**
 incandescent nongaseous, **15.15–15.21**
 comparisons, **15.19**, **15.19f**
 gas mantle, **15.17**, **15.18**, **15.19f**
 globar, **15.17**, **15.18f**
 Nernst glower, **15.14**, **15.15**, **15.17**, **15.17f**
 quartz-envelope lamps, **15.20**, **15.21**
 tungsten-filament lamps, **15.19**, **15.20**, **15.**
20f–15.22f
 low-pressure enclosed arc, **15.35–15.47**
 black-light fluorescent lamps, **15.35**, **15.36t**
 electrodeless discharge lamps, **15.36**, **15.44**
 germicidal lamps, **15.35**
 hollow cathode lamps, **15.35**,
15.37t–15.43t, **15.44f**
 Pluecker spectrum tubes, **15.47**, **15.47f**,
15.47t
 spectral lamps, **15.44**, **15.45**, **15.45f**,
15.46f, **15.46t**
 Sterilamps, **15.35**, **15.36f**
 special-purpose, **15.53**
- Commission Internationale de l'Eclairage
 (CIE), **40.2**
 publications from, **37.11**
 standard photometric observer, **37.2**
 Common path interferometers, **13.9**, **13.11f**
 Compact fluorescent lights (CFLs), **40.25t**,
40.26t, **40.28f**, **40.31**
 Compact-source arcs, **15.31–15.34**,
15.32f–15.35f
 Compensators:
 Dall, **13.24**, **13.24f**
 holographic, **13.25**
 Offner, **13.24**, **13.24f**
 reflective, **13.24**, **13.24f**, **13.25**
 refractive, **13.24**, **13.24f**, **13.25**
 and tolerances, **3.21**, **5.7**
 Complete monolithic FPAs, **33.10**
 Compound elliptical collectors (CECs), **39.14**,
39.15f, **39.27**, **39.37**
 Compound hyperbolic collectors (CHCs),
39.15, **39.15f**, **39.16f**, **39.37**
 Compound lens, thermal defocus of, **8.4**, **8.5f**
 Compound parabolic collectors (CPCs),
39.13–39.14, **39.13f**, **39.14f**, **39.18**, **39.19**,
39.19f
 Compressively strained QW lasers, **19.16f**, **19.17**
 Compton effect, **23.9**
 Computer graphics software, **40.21–40.23**,
40.22f–40.24f
 Computer numeric control (CNC) systems, **9.4**
 Computer-aided design (CAD) software, **40.19**
 Computer-generated holograms (CGHs),
14.1–14.9
 about, **14.1–14.3**
 accuracy limitations of, **14.6–14.7**, **14.6f**, **14.7f**
 discussion of, **14.9**
 experimental results from, **14.7–14.9**, **14.7f**,
14.8f
 interferometers using, **14.4–14.5**, **14.4f**, **14.5f**
 plotting of, **14.3–14.4**, **14.3f**
 sample, **14.2f**
 Concave facets, **39.40**, **39.40f**
 Concentrated arc lamps, **15.47–15.49**, **15.48f**,
15.49f
 Concentration:
 of radiation, **39.1**, **39.5**, **39.6**
 of solution, **38.5**
 Concentrators, nonimaging, **39.12–39.22**
 calculation of, **39.5**, **39.6**
 compound elliptical collectors, **39.14**, **39.15f**

- Concentrators, nonimaging (*Cont.*):
 compound hyperbolic collectors, 39.15, 39.15f, 39.16f
 compound parabolic collectors, 39.13–39.14, 39.13f, 39.14f
 dielectric compound parabolic collectors, 39.15, 39.16, 39.16f
 edge rays, 39.22
 geometrical vector flux, 39.21–39.22
 inhomogeneous media, 39.22
 integrating cavities with, 39.26, 39.27
 and lenses, 39.16–39.17
 and mirrors, 39.17
 multiple surface concentrators, 39.16–39.17, 39.17f
 restricted exit angle concentrators with lenses, 39.18, 39.18f
 RX, 39.17, 39.17f
 RXI, 39.17, 39.17f
 star, 39.20, 39.21
 tapered lightpipes, 39.12–39.13, 39.13f
 θ_1/θ_2 concentrators, 39.18–39.20, 39.19f
 2D vs. 3D, 39.20–39.21, 39.20f, 39.21f
- Condensers, first-order layout for, 1.10–1.11, 1.11f
- Condition of detailed balance (term), 23.23
- Conduction band, 17.4, 17.4f, 17.5f
- Conduction bandgap, 25.3, 25.3f
- Cones (eye receptors), 30.15, 30.16f, 34.37, 34.38, 36.8, 36.8f, 36.9f
- Cone-shaped secondary baffle, 7.3–7.4, 7.3f, 7.4f
- Confidence interval (CI), 34.22
- Configuration factor algebra, 34.14
- Confocal cavity technique, 12.20, 12.20t
- Confocal parameter, 16.23
- Conic collimators, 39.8, 39.9f
- Conic reflectors, 39.11, 39.11f
- Conic surfaces, 3.5
- Conical (term), 35.5
- Conical-directional reflectance, 35.5t, 35.6f, 35.6t
- Conical-hemispherical reflectance, 35.5t, 35.6f, 35.6t
- Conservation, of radiant power transfer, 34.13f
- “Conservation of complexity,” 3.7
- Constraints:
 defined, 3.17
 methods for handling, 3.18–3.19
- Constricted double-heterostructure large optical cavity (CDH-LOC), 19.19, 19.20t, 19.21f
- Consultative Committee on Photometry and Radiometry (CCPR), 36.2
- Contact scanners, 32.21, 32.22f
- Contact stresses, 6.21
- Contacting, in wafer processing, 17.24
- Contamination levels (in stray light suppression), 7.18–7.19, 7.18t, 7.19f–7.21f
- Continuous polishers (CPs), 9.7
- Continuous ring flanges, 6.4f, 6.11
- Continuous wave (cw) lasers, 23.18, 34.32
- Continuous wave (cw) power, 19.19, 19.22, 19.22f
- Convergent reflectors, 39.38–39.40, 39.38f, 39.39f
- Conversion factors:
 for English and SI units, 37.7t
 for photometric and radiometric quantities, 36.11–36.14, 36.12f–36.14f
- Convex surfaces, testing of, 14.5, 14.5f
- Coordinate measurement machines (CMMs), 9.6
- Coordinate measurement method (CMM), 40.53, 40.54
- Copper, 17.28
- Copper vapor lasers (CVLs), 16.12, 16.13f, 16.30
- Copper-doped germanium (Ge:Cu) detectors, 24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f
- Corner cube prisms, 12.16
- Cornice lighting, 40.13f
- Corning ULE, 6.18
- Correlated color temperature (CCT), 34.44, 37.7, 38.5, 40.8
- Correlated double sampling (CDS), 33.13
- Correlated emission lasers (CELs), 23.42–23.43
- Cosine law, 37.8, 37.8f
- Cosine-to-the-fourth approximation, 34.16
- Coupled cavity lasers, 19.37f, 19.38
- Coupling:
 of circulating pulses, 20.12–20.15, 20.12f, 20.15f
 étendue and source, 40.41–40.42
 in film development, 29.14
 gain, 19.29
 interface, 20.14
 output, 16.13
 phase-conjugated, 20.14
 repetition-rate, 20.14–20.15, 20.15f

- Coupling noise, resistive, 27.5, 27.6f
- Cove lighting, 40.13f, 40.16f
- Critical illumination (*see* Abbe illumination system)
- Critical objects (in stray light suppression), 7.2 imaged, 7.4, 7.5f
real-space, 7.2–7.4, 7.3f, 7.4f
- Crossed reflectors, 39.38f
- Crossed string relationship, 39.4, 39.4f
- Cryogenic electrical substitution radiometers, 34.28
- Crystalline optics, 9.8
- Current density, 19.12–19.13, 19.12f, 19.13f
- Current-confined constricted double-heterostructure large optical cavity (CC-CDH-LOC), 19.19, 19.20t, 19.21f
- Curvature measurement, 12.17–12.25
mechanical methods of, 12.17–12.19, 12.18f, 12.19f, 12.19t
optical methods of, 12.19–12.21, 12.20f, 12.20t, 12.21f
- Cusp surface, of diamond-turned optics, 10.10, 10.10f
- Cutoff wavelength, 24.10
- Cyanine dyes, 30.13, 30.13f
- Dall compensators, 13.24, 13.24f
- Damped least-squares (DLS) method, 3.17–3.19
- Damping:
additive, 3.18
of field by reservoir, 23.33–23.34
- Damping factor, 3.18
- Dark counts (of photomultipliers), 27.8
- Dark current:
absorption coefficient, 25.8, 25.9f
in CCDs, 32.20
correction of, 34.33
defined, 24.10
diffusion current, 25.7
generation-recombination current, 25.7–25.8
histogram of, 32.12n
in photosensing elements, 24.19–24.20, 32.10–32.12, 32.11f
in *pin* photodiodes, 25.7–25.8
quantum efficiency, 25.8
responsivity, 25.8
tunneling current, 25.8
- Dark regions, 17.28
- Dark signal correction, 34.33
- Dark-line defects, 17.28
- Dashed rays, 1.12, 1.12f
- Daylight:
as natural light source, 40.40–40.41
simulation of, 40.17
spectrum of, 40.40f
- Daylighting schemes, luminaires for, 40.47–40.50, 40.49f–40.51f
- Day/night cameras, 31.28–31.29
- Dazzle, 40.9
- dc carbon arcs, 15.25t
- dc lamps, 15.32f
- Decay time, 16.4
- Decorative lighting, 40.14, 40.16f
- Deep-diffused stripe (DDS) lasers, 19.23
- Defocus:
aberrations of, 11.30
annular polynomials for, 11.38f, 11.39
thermal, 8.4, 8.5f
- Density (of photographic films), 29.6–29.8, 29.7f
- Density of states, 19.9–19.10, 19.10f
- Density-operator approach, to quantum theory of lasers, 23.14–23.33
derivation of Scully-Lamb master equation, 23.17–23.22
cavity losses, 23.18
laser master equation, 23.19–23.20, 23.19f
micromaser master equation, 23.20–23.22
photon statistics, 23.22–23.27
laser, 23.22–23.26, 23.23f, 23.25f
micromaser, 23.26–23.27, 23.27f
spectrum, 23.28–23.33
laser field, 23.28–23.31, 23.30f
micromaser field, 23.31–23.33
time evolution of the field in Jaynes-Cummings model, 23.15–23.17, 23.15f
- Depletion, of charge, 25.6
- Depletion layer generation current, 32.10–32.11, 32.11f
- Depth of focus, 4.7–4.8, 4.8f
- Design software, optical (*see* Optical design software)
- Detailed balance, condition of, 23.23
- Detective quantum efficiency (DQE), 24.10, 29.1, 29.23
- Detective time constant, 24.10
- Detectivity:
of film, 29.23
of infrared detector arrays, 33.23–33.24
normalized, 38.9

- Directional total absorptance, **35.8t**
 Direct-reading autocollimators, **12.12**
 Disability glare, **40.9–40.10**
 Discomfort, visual, **40.9–40.12**
 Discomfort glare, **40.9–40.12, 40.11t**
 Discrete energy levels, **16.4**
 Disk PZT transducers, **22.17–22.18**
 Dislocation reduction, **18.2, 18.2f**
 Dispersion (of light), **38.8**
 Dispersivity, **30.9**
 Displacement current, **26.7**
 Displays (term), **40.1** (*See also specific displays, e.g.: Monolithic LED displays*)
 Distance measurement (*see* Length measurements)
 Distortion plot, **2.4, 2.4f**
 Distortion tolerances, **5.8**
 Distracting glare, **40.9**
 Distributed backscattering, **20.14**
 Distributed Bragg reflector (DBR) lasers, **19.38, 19.40**
 Distributed feedback (DFB) lasers, **16.29, 19.36, 19.38**
 Distributed grating surface-emitting lasers, **19.40–19.41, 19.40f**
 Distribution temperature, **34.43–34.44, 37.7**
 Divergent reflectors, **39.38–39.40, 39.38f, 39.39f**
 D-log H curve (for photographic films), **29.8–29.10, 29.8f**
 Domes, mounting of, **6.11, 6.12f**
 Doped extrinsic silicon, **33.7, 33.8f**
 Doping, substrate, **17.20**
 Doppler broadening, **16.5, 16.6, 16.6f, 16.9**
 Doppler linewidth (*see* Full width at half maximum)
 Double heterojunction (DH) LEDs, **17.13, 17.13f–17.15f**
 Double heterostructure (DH) lasers, **19.4, 19.5f, 19.7, 19.12–19.15, 19.18–19.19, 19.19f**
 Double heterostructure *pin* photodiodes, **26.13**
 Double monochromators, **35.9, 38.15f**
 Double optical gating, **21.8**
 Double sampling, correlated, **33.13**
 Double-beam spectrophotometers, **35.8–35.9**
 Double-channel planar buried heterostructure (DC-PBH) lasers, **19.24, 19.25f, 19.34f**
 Double-pass photodetectors, **26.4f**
 Doublet lens, air-spaced, **6.7**
 Downconversion, parametric, **23.14**
 Downhill optimizer, **3.17**
 Draft angle, **39.10**
 Drag-wiping (cleaning), **10.9**
 Drift:
 in CCDs, **32.17**
 frequency vs. time, **22.2**
 low offset, **27.11**
 photogenerated charge collection by, **32.5**
 thermocouple junctions as source of, **27.6, 27.6f**
 Driving lasers, in attosecond optics, **21.4–21.6, 21.5f**
 Drop-in assembly, **6.6, 6.6f**
 Dual beam detection, **22.13**
 Dual in-line octocouplers, **17.32, 17.32f**
 Dumet (alloy), **40.29**
 Dye lasers, **16.31–16.32, 16.32f, 20.15–20.16**
 Dye-forming reaction, in film development, **29.14**
 Dyes:
 azomethine, **30.10f, 30.11–30.12, 30.12f**
 cyanine, **30.13, 30.13f**
 light-absorbing, **30.7**
 photographic, **30.10–30.13, 30.10f, 30.12f**
 yellow filter, **30.4**
 Dynodes, in photomultipliers, **27.6–27.9, 27.7f**
 Eberhard effects, **30.3**
 Eccentric pupil design (*see* Z-system)
 ECE (United Nations Economic Commission for Europe), **40.63–40.64**
 Edge lit backlight, **40.47, 40.47f**
 Edge rays, **39.22, 39.38**
 Edge-absorbing photodetectors, **26.4f**
 Edge-emitting lasers (ELASERs), **25.15**
 Edge-emitting LEDs (ELEDs), **25.15**
 Edge-illuminated photodetectors, **26.4f, 26.5**
 Edging step (of optics fabrication), **9.6**
 Einstein (unit), **34.11**
 Einstein's light quanta, **23.6–23.9, 23.8f**
 Einstein's particle hypothesis, **23.7**
 Elastomeric mountings, **6.4, 6.4f, 6.5, 6.12**
 Electrical contact (light bulb), **40.29f**
 Electrical parasitics (laser), **19.34–19.35, 19.34f**
 Electrical substitution radiometers, **34.27–34.29**
 Electrical transfer function, with series inductance, **26.13**
 Electrodeless discharge lamps, **15.36, 15.44**
 Electrodeless fluorescent lamps, **40.36–40.37**
 Electrodeless lamps, **40.25t, 40.26t, 40.36–40.37**

- Electrodeless sulfur lamps (ESLs), **40.36–40.37**
- Electroluminescent light sources, **40.37–40.39**,
40.38f, **40.38t**, **40.39f**
- Electron bombardment (EB), **31.23**
- Electron current, **26.7**
- Electron lenses, **31.8**, **31.8f**
- Electron-bombarded SSAs (EBSSAs),
31.23–31.27
- digital cameras, **31.24–31.26**
- modulation transfer function and limiting
resolution of, **31.26–31.27**
- proximity-focused, **31.23–31.24**, **31.23f**,
31.24f
- Electronically scanned staring FPAs,
33.16–33.17
- Electro-optic modulators (EOMs), **22.14**, **22.20**
- Electro-Optical Industries, Inc. (EOI), **15.14**,
15.15f, **15.16f**
- Elliptical polynomials, **11.21**, **11.25–11.27**,
11.26t–11.27t, **11.36t**, **11.38f**
- Emission, stimulated (*see* Stimulated emission)
- Emission lasers, correlated, **23.42–23.43**
- Emission linewidth (of radiation), **16.4–16.7**,
16.6f, **16.7f**
- Emission-line broadening, **16.4–16.7**, **16.6f**, **16.7f**
- Emissivity:
of blackbody cavity, **15.7**, **15.8f**
tungsten, **40.28f**
- Emittance, **39.2**
- and absorptance, **35.8**
- calculating, **34.25–34.26**
- defined, **35.7**
- measurement of, **35.14–35.16**, **35.15f**
- Emitted photon wavelength, **17.4–17.5**
- Emitters:
AlGaAs, **17.32**
blue, **17.18**, **17.19**
GaAsP, **17.32**
- Emulsions, photographic, **24.100**, **24.101f**, **29.4**,
30.7
- Enclosed arcs, **15.24**, **15.28–15.47**
- high-pressure, **15.24**, **15.28–15.34**
- capillary mercury-arc lamps, **15.30–15.31**,
15.31f
- compact-source arcs, **15.31–15.34**,
15.32f–15.35f
- Lucalox lamps, **15.30**, **15.31f**
- mercury arcs, **15.29**, **15.30f**
- multivapor arcs, **15.29**, **15.31f**
- Uviarc, **15.28–15.29**, **15.29f**, **15.30f**
- Enclosed arcs (*Cont.*):
low-pressure, **15.35–15.47**
- black-light fluorescent lamps, **15.35**, **15.36t**
- electrodeless discharge lamps, **15.36**, **15.44**
- germicidal lamps, **15.35**
- hollow cathode lamps, **15.35**, **15.37–15.43t**,
15.44f
- Pluecker spectrum tubes, **15.47**, **15.47f**,
15.47t
- spectral lamps, **15.44**, **15.45**, **15.45f**,
15.46f, **15.46t**
- Sterilamps, **15.35**, **15.36f**
- End loss, **19.6**
- Energy:
levels of, **16.4**, **16.7**
- luminous, **37.4t**, **37.6**
- measurement of, **34.32**
- nomenclature for, **36.4**, **36.5**
- radiant, **34.7**, **37.4t**, **37.6**
- units of, **34.5–34.6**
- Energy band structure, **17.3–17.6**, **17.3f–17.5f**
- Energy bandgap, **25.3**, **25.3f**
- English units, and SI units, **37.7**, **37.7t**
- Entrance pupil, **34.18**, **34.19f**
- Entrance window, **34.19**, **34.19f**
- Environmental specifications, optical, **4.10**
- Epitaxial growth, **17.8**, **17.21**
- Epitaxial technology (for LEDs), **17.21–17.23**
- Epoxy, in indicator lamps, **17.29**
- Equivalent neutral density (END), **29.15**
- Equivalent noise input (ENI), **24.11**
- Equivalent veiling luminance (EVL), **40.10**
- Error functions, in optical design software,
3.17, **3.19–3.20**
- Error types, in absolute measurements,
34.21–34.23
- Étendue, **34.15**
- defined, **40.42**
- geometrical, **38.8**
- in nonimaging optics, **39.2**, **39.3**,
39.4f, **39.5**
- and source coupling, **40.41–40.42**
- Étendue loss, **39.6**
- Evaluation function (of optical software),
3.8–3.16
- of aberrations, **3.9–3.11**
- paraxial ray-trace, **3.8–3.9**, **3.9f**
- ray-trace, **3.11–3.13**, **3.12f**
- by spot-diagram analysis, **3.13–3.16**
- Event-driven programs (optical software), **3.7**

- Exact rays (term), 3.3, 3.11–3.12
- Excimer lasers, 16.30–16.31, 16.31*f*
- Excimers, in fluorescent lamps, 40.31
- Excited state, of azomethine dyes, 30.11–30.12, 30.12*f*
- Exciton recombination, 17.6
- Exit pupil, 34.18, 34.19*f*
- Exitance, 39.2
- defined, 34.8
 - luminous, 37.4*t*, 37.5, 37.5*f*
 - radiant, 15.4–15.6, 15.5*t*, 15.6*f*; 37.4*t*, 37.5, 37.5*f*
- Exposure:
- luminous, 37.4*t*, 37.6
 - of photographic films, 29.5–29.6
 - radiant, 37.4*t*, 37.6
- Extended baffle shields, 7.9–7.10, 7.9*f*, 7.10*f*
- Extended wavelength photodetectors, 25.10, 25.10*t*
- Exterior lighting, 40.61–40.62, 40.63*t*
- External cavity diode lasers (ECDLs), 22.21–22.23, 22.22*f*
- Extreme infrared (IR) light, 25.2
- Extrinsic photoconductors, 25.5, 25.5*f*
- Extrinsic photodetectors, 24.7, 24.7*f*
- Extrinsic semiconductor transition, 24.11
- Eye, human (*see* Human eye)
- Fabrication, optical, 9.3–9.9
- about, 9.3
 - aspherical, 9.7–9.8, 9.7*f*
 - crystalline, 9.8
 - by diamond turning (*see* Diamond turning)
 - diamond turning vs. traditional, 10.6
 - guide to methods of, 10.3*t*
 - material formation for, 9.3–9.4
 - methods of, 10.3*t*
 - plano, 9.7
 - spherical, 9.4–9.6
- Fabry-Perot interferometers, 16.19*f*
- Faceted reflectors, 39.10*f*, 39.39–39.41, 39.39*f*, 39.40*f*
- Failures per 10⁹ hours (FITS), 17.25
- False-colored infrared film, 30.22
- Far infrared (FIR) radiation, 24.3, 25.2
- Far ultraviolet radiation, 15.12, 15.13
- Fatigue, thermal, 17.25
- Federal Motor Vehicle Safety Standards (FMVSS), 40.63
- Feedback:
- optical, 16.2
 - resonant optical, 19.38, 19.38*f*
 - stabilization of, 34.32
- Femtoseconds, 20.1
- Fermi occupation functions, 19.11
- Ferroelectric bolometer arrays, 28.11, 28.12, 28.12*f*
- Ferroelectric detectors, 33.10
- Fiber lasers, 16.34
- Fiber optics:
- fiber alignment for, 17.33
 - focal-length measurement with, 12.25
 - LED considerations with, 17.33–17.34
 - in nonimaging optics, 39.21
- Fiber-optic octocouplers, 17.33–17.34
- Fiberoptic-coupled (FO) II SSAs, 31.20–31.22, 31.20*f*, 31.21*f*, 31.21*t*
- Field curvature plot, 2.4, 2.4*f*, 2.5
- Field effect transistors (FETs), 27.10
- Field lenses, first-order layout for, 1.8, 1.10, 1.10*f*
- Field of view (FOV), 3.4, 7.11, 24.11, 31.1
- Field patch trace, 39.7–39.8
- Field plots, of aberration curves, 2.4–2.5
- Field stops, 7.7, 7.8*f*, 7.9*f*, 34.18–34.19, 34.19*f*
- Field-enhanced pyroelectric arrays (*see* Ferroelectric bolometer arrays)
- Figures of merit (FOM), 24.13, 33.23–33.28
- for infrared photodetectors, 25.12
 - minimum resolvable temperature (MRT), 33.27–33.28, 33.27*f*
 - NE ΔT , 33.25*f*, 33.26*f*
 - in spectroradiometry, 38.5–38.6
- Filament notching, 40.30
- Filaments:
- lamp, 15.20*f*
 - light bulb, 40.25, 40.27, 40.29–40.30, 40.29*f*
- Fill gases, light bulb, 40.29*f*, 40.30
- Film, photographic (*see* Photographic films)
- Filters:
- blocking, 38.8
 - blur, 32.34, 32.34*f*
 - in II SSA cameras, 31.7
 - infrared, 40.12
 - interference, 34.36
 - neutral density, 40.52
 - notch, 22.10–22.11, 22.10*f*, 22.11*f*
 - UV, 40.12
- Finitely distant objects, systems with, 1.6, 1.6*f*

- First-order layout techniques, 1.3–1.16
 achromatism, 1.14–1.15, 1.15*f*
 afocal attachments, 1.8, 1.9*f*
 afocal systems, 1.7, 1.7*f*
 athermalization, 1.15–1.16, 1.16*f*
 axial/principal rays, 1.12
 component power minimization, 1.13, 1.13*f*
 condensers, 1.10–1.11, 1.11*f*
 defined, 1.4
 field lenses, 1.8, 1.10, 1.10*f*
 magnifiers and microscopes, 1.8
 ray-tracing, 1.4–1.5
 reasonableness of layout, 1.13–1.14
 two-component systems, 1.5–1.7
 zoom or varifocal systems, 1.11–1.12
 “Fitness for use,” 17.25
 Five-axis machining, 10.7*f*
 5 × 7 matrix LED displays, 17.31–17.32
 Fixed interferogram evaluation, 13.14–13.15
 Fixed-orientation mirrors, 6.17
 Fixer, film, 29.5
 Fizeau interferometers, 12.14, 13.8–13.9, 13.9*f*,
 13.10*f*, 13.18, 14.4, 14.5*f*
 Flaming arcs, 15.23, 15.24*f*, 15.26*t*–15.27*t*
 Flanges:
 annular, 6.3–6.4, 6.4*f*
 continuous ring, 6.4*f*, 6.11
 Flash blindness, 40.9
 Flash lamps, 16.16–16.17, 16.17*f*
 Flex-Pivots (flexures), 6.19
 Flexure mountings, 6.5, 6.5*f*, 6.15–6.17, 6.16*f*
 Flicker:
 in fluorescent lamps, 40.32–40.33
 impact of, 40.12
 in incandescent lamps, 40.30
 Flicker floor, 22.3
 Flicker noise, 24.11
 Flight control, 19.3
 Flip chip packaging, 18.6, 18.6*f*
 Flip-and-fold approach, 39.29*f*
 Floating diffusion, 32.14, 32.15*f*
 Floating gate output amplifiers, 32.15
 Fluorescence, 34.13, 40.30
 Fluorescent lamps, 40.30–40.33
 applications for, 40.26*t*
 characteristics of, 40.25*t*
 construction of, 40.33*f*, 40.34*f*
 elements of, 40.31*f*
 emission spectrum of, 40.32*f*, 40.35*f*
 types of, 40.28*f*
 Flux:
 luminous, 37.4, 37.4*t*, 37.6
 radiant, 37.3, 37.4*t*
 total luminous, 37.4*t*, 37.6
 total radiant, 37.4*t*, 37.6
 Flux budget, for fiber optics, 17.33
 Flux density (*see* Irradiance)
 Fly-by-light (FBL), 19.3
 FM spectroscopy, 22.13–22.14
 F-number, 34.20
 Focal depth, 4.7–4.8, 4.8*f*
 Focal length, 1.5–1.7, 12.21–12.25
 fiber optics, 12.25
 focimeters, 12.22–12.23, 12.22*f*, 12.23*f*
 Fourier transforms, 12.24
 microlenses, 12.24
 Moiré deflectometry, 12.23, 12.24*f*
 nodal slide bench, 12.22, 12.22*f*
 Talbot autoimages, 12.23, 12.24
 Focal plane arrays (FPAs), 33.3
 hybrid, 33.14–33.23
 microbolometer, 33.13–33.14
 MIS photogate, 33.10–33.11, 33.11*f*, 33.12*f*
 monolithic, 33.10–33.14
 Focal ratio, 34.20
 Focimeters, 12.22–12.23, 12.22*f*, 12.23*f*
 Focus athermalization techniques, 6.22–6.24
 active athermalization, 6.24, 6.24*f*
 passive athermalization, 6.22, 6.23*f*, 6.24
 single material designs, 6.22, 6.23*f*
 Focus distance, 1.5–1.7
 Focus shift, thermal, 8.2–8.4, 8.3*t*, 8.4*t*
 Fog, film, 29.9
 Fokker-Planck equation, 23.37
 Foot-candle (unit), 34.43, 36.7, 36.7*t*, 37.7*t*
 Foot-lambert (unit), 34.43, 36.7, 36.8*t*, 37.7,
 37.7*t*
 Forbes method, 3.20
 Fore-optics, 38.7
 Forward light, 40.43, 40.44*f*
 Forward looking infrared (FLIR), 33.4
 Foucault test, 12.19, 13.2–13.3, 13.2*f*, 13.3*f*
 Fourier analysis, of interferograms,
 13.16–13.17, 13.17*f*
 Fourier transform spectrophotometers, 35.9
 Fourier transforms, for focal-length
 determination, 12.24
 Four-phase CCDs, 32.13–32.14, 32.13*f*, 32.16*f*
 Frame interline transfer (FIT) CCDs,
 32.27*f*, 32.29

- Frame transfer (FT) CCD image sensors, 32.26–32.28, 32.27f, 32.28f, 32.32
- Frame transfer (FT) CCD TDI FPAs, 33.11, 33.12f
- Free-electron lasers (FELs), 16.36–16.37, 16.37f, 23.43–23.45
- Free-electron-laser (FEL) lamps, 15.11, 15.12, 15.13f
- Frequency:
- and drift, 22.2
 - phase and amplitude responses vs., 22.6–22.7, 22.6f, 22.7f
 - stability of, 22.2
- Frequency comb, 20.1, 20.2, 20.7–20.9
- Frequency discriminators:
- for laser locking, 22.12–22.14
 - optical cavity-based, 22.14–22.16, 22.17f
- Frequency modulation (FM), 19.36, 19.36f, 22.4
- Frequency shift keying (FSK), 19.36
- Frequency-controlled lasers, 22.7, 22.7f
- Frequency-resolved optical gating (FROG), 21.8, 21.9
- Frequency-resolved optical gating for complete reconstruction of attosecond bursts (FROG-CRAB), 21.9
- Frequency-selective-feedback lasers, 19.37f, 19.38
- Fresnel lenses, 39.9–39.10, 39.10f, 40.45f, 40.46
- Full width at half maximum (FWHM), 16.5, 16.6, 20.3, 21.2
- Functional specifications (optical design), 4.2
- Fundamental array mode, 19.27
- Furniture-integrated lighting system, 40.13f
- Gain:
- defined, 16.9
 - in photomultipliers, 27.7
- Gain coefficient, 16.9–16.10
- Gain coupling, 19.29
- Gain medium, 16.3
- Gain saturation, 16.10
- Gain stability margin, 22.9–22.10
- Gain-bandwidth (GB), 26.17
- Gain-coupled arrays, 19.27
- Gain-guided phased array, 19.28f
- Galilean telescopes, 1.7f
- Gallium aluminum arsenide (GaAlAs) LEDs, 17.12, 17.12f, 17.13f
- Gallium arsenide (GaAs) lasers, 19.7
- Gallium arsenide (GaAs) LEDs, 17.8, 17.9, 17.9f
- Gallium arsenide (GaAs) semiconductor diode lasers, 16.18–16.19, 16.18f
- Gallium arsenide phosphide (GaAsP) emitters, 17.32
- Gallium arsenide phosphide (GaAsP) LEDs, 17.9–17.10, 17.10f, 17.15–17.17
- energy band diagram for, 17.5f
 - homojunction in, 17.10, 17.11f
 - light degradation in, 17.27f
 - performance summary of chips in, 17.16t
- Gallium arsenide phosphide (GaAsP) photodiodes, 24.49, 24.49f, 24.50f
- Gallium arsenide phosphide (GaAsP) substrate, 17.22
- Gallium arsenide (GaAs) quantum well photodetectors, 25.16–25.17, 25.16f, 25.17f
- Gallium nitride (GaN) photovoltaic detectors, 24.42, 24.43, 24.45f, 24.46, 24.46f, 24.47
- Gallium nitride (GaN) substrate, 17.22
- Gallium phosphide (GaP), 17.16, 17.21–17.22
- Gallium phosphide (GaP) dynodes, 24.42, 24.44f
- Gallium phosphide (GaP) photodiodes, 24.47–24.49, 24.48f
- Gallium phosphide (GaP) substrate, 17.20–17.22
- Gallium-doped germanium (Ge:Ga) infrared detectors, 24.100
- Gallium-doped silicon (Si:Ga) infrared detectors, 24.95, 24.95f, 24.96, 24.96f
- Galvo-driven Brewster plates, 22.18–22.19
- Gas chromatography-mass spectroscopy (GC-MS), 33.4
- Gas lights, 40.40
- Gas mantle, 15.17–15.19, 15.19f
- Gaseous laser gain media, 16.30–16.31, 16.31f
- Gas-filled lamps, 34.31
- Gate modulation, 33.19t, 33.20f, 33.22
- Gated integration, 27.12–27.13, 27.13f, 27.15
- Gating:
- amplitude, 21.7
 - double optical, 21.8
 - frequency resolved, 21.8, 21.9
 - FROG, 21.9
 - FROG-CRAB, 21.9
 - polarization, 21.7–21.8
 - two-color, 21.7
- Gauss illuminated eyepieces, 12.12

- Gaussian intensity distribution, **39.28, 39.29f**
 Gaussian line shape, **16.6f**
 Gaussian mode, **16.21**
 Gaussian parameters (optical design), **4.5–4.6, 4.6t**
 Gaussian-shaped beam, **16.22**
 General Conference on Weights and Measures (CGPM), **36.2**
 General Electric, **15.29, 15.30, 15.48, 19.29t**
 General system data (optics), **3.3, 3.4**
 Generating step (of optics fabrication), **9.4**
 Generation noise, **24.11**
 Generation-recombination (GR) current, **25.7–25.8**
 Generation-recombination (GR) noise, **24.11**
 Geometrical configuration factor (GCF), **7.1, 7.2, 7.22**
 Geometrical etendue, **38.8**
 Geometrical optical transfer function (GOTF), **3.16**
 Geometrical optics, **3.16**
 Geometrical vector flux, **39.21–39.22**
 Geometry-controlled lasers, **19.37f, 19.38**
 Germanium (Ge) avalanche photodiodes, **24.70f, 24.72–24.73, 24.72f, 24.73f**
 Germanium (Ge) bolometers, **28.5, 28.7t**
 Germanium (Ge) detectors:
 copper-doped, **24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f**
 gallium-doped infrared, **24.100**
 gold-doped, **24.83–24.85, 24.84f–24.86f**
 intrinsic photodetectors, **24.70–24.73, 24.70f–24.73f**
 mercury-doped, **24.84f, 24.92–24.95, 24.93f–24.95f**
 pn and *pin*, **24.70–24.71, 24.70f–24.72f**
 zinc-doped, **24.84f, 24.98–24.100, 24.99f**
 Germanium gallium arsenide (GeGaAs) photodiodes, **34.31**
 Germanium (Ge) intrinsic photodetectors, **24.70–24.73, 24.70f–24.73f**
 Germanium (Ge) low-temperature bolometers, **24.31–24.32, 24.32f, 24.33f**
 Germanium photodiodes, **38.9, 38.9t**
 Germicidal lamps, **15.35**
 Glare:
 and exterior lighting, **40.62**
 limiting of, **40.10, 40.41**
 and visual discomfort, **40.9–40.12, 40.11t**
 and windows, **40.41**
 Glare stops (*see* Lyot stops)
 Glass:
 formation of optical, **9.3–9.4**
 optical, **5.9**
 as photographic film emulsion, **29.4**
 tolerances for, **5.9**
 Glass envelope, light bulb, **40.29, 40.29f**
 Glazing, window, **40.41**
 Global, **15.17, 15.18f, 15.19, 15.19f**
 Glossiness, **40.5**
 Glow lamps, **40.39**
 Glow modulator tubes, **15.49, 15.50f, 15.51f, 15.52t**
 Gobs, glass, **9.4**
 Golay cell detectors, **28.2, 28.6, 28.7t**
 Gold, diamond turning and, **10.5**
 Gold-doped germanium (Ge:Au) detectors, **24.83–24.85, 24.84f–24.86f**
 Gold-germanium (Au-Ge) alloys, **17.24**
 Goldpoint blackbody, **15.9**
 Gold-zinc (Au-Zn) alloys, **17.24**
 Goniometers (goniophotometers), **12.10, 40.52–40.53, 40.53f, 40.54f**
 Gouffé method, **15.7–15.9, 15.8f**
 Graded-index separate-confinement heterostructure (GRIN SCH) quantum lasers, **19.14, 19.14f**
 Gradient-freeze technique, **17.21**
 Grains and graininess:
 of photographic films, **29.5**
 of photographic images, **29.18t, 29.19–29.22, 29.21f**
 of silver halide crystals, **29.4**
 Granularity, photographic film speed and, **30.19**
 Grating equation, **38.7–38.8**
 Grating surface-emitting laser array, **19.40–19.41, 19.40f**
 Gray gel, **30.4**
 Graybody, **35.7**
 Green light, color film and, **29.13, 29.13f**
 Green-emitting AlInGaP devices, **17.18**
 Grinding, aspheric, **9.8**
 Ground loop noise, **27.5, 27.6f**
 Ground state, **16.4**
 Grown homojunctions, LED, **17.8, 17.9, 17.9f**
 Growth techniques:
 for epitaxial layers, **17.21–17.23**
 substrate, **17.20, 17.21**
 Guard ring, **24.11**

- Guide to the Expression of Uncertainty in Measurement* (IS), 38.6
- Gurney-Mott mechanism, 29.5
- H-aggregates, 30.13
- Haidinger interferometers, 12.14
- Halation, 30.4
- Halogen lamps, 15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Halon, 38.12–38.13
- Halophosphates, 40.31
- Hamiltonian rays, 3.12
- Hanbury-Brown-Twiss (HB&T) effect, 23.13, 23.13*f*, 23.14
- Hard mounting, of optics, 6.1–6.4, 6.3*f*, 6.4*f*
- Hartmann test, 13.4–13.6, 13.5*f*
- Hartmann-Shack test, 13.6–13.7, 13.6*f*
- H&D curve (*see* D-log H curve)
- Headlamps:
 - design of, 40.21, 40.23, 40.23*f*
 - low-beam, 40.64–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*
- Health-care facility lighting, 40.58–40.60, 40.60*t*
- Heat-pipe blackbody furnace, 15.9*f*
- Height solves, 3.6
- Heisenberg-Langevin approach, to quantum theory of lasers, 23.33–23.35
- Helium-cadmium (He-Cd) lasers, 16.6, 16.15, 16.15*f*, 16.30
- Helium-neon (He-Ne) lasers, 16.15, 16.15*f*, 16.30
- Hemispherical emittance, 35.15
- Hemispherical total absorptance, 35.8*t*
- Hemispherical-conical reflectance, 35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-directional reflectance, 35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-spectral absorptance, 35.8*t*
- Heterodyne interferometers, 13.22
- Heterojunction lasers, 19.4
- Heterojunctions, 17.12, 17.12*f*–17.15*f*, 17.13, 17.17, 26.9
- Hewlett-Packard double-frequency distance-measuring interferometer, 12.9–12.10, 12.9*f*
- Hexagonal polynomials, 11.21, 11.22*t*–11.25*t*, 11.36*t*, 11.38*f*, 11.39
- High-accuracy spectrophotometers, 35.9
- High-brightness visible LEDs (HB-LEDs), 18.1–18.6
 - about, 18.1
 - epitaxial growth of, 18.3
 - packaging of, 18.5–18.6, 18.5*f*, 18.6*f*
 - processing of, 18.3–18.4, 18.3*f*
 - semiconductor material systems for, 18.1–18.2
 - solid-state lighting with, 18.4–18.5, 18.4*f*
 - structure for modern InGaN, 18.2*f*
 - substrates for, 18.2–18.3, 18.2*f*
- High-dye-yield yellow couplers, 30.6
- High-energy radiation, 15.40, 30.19–30.20
- High-gain oscillators, 20.10–20.12, 20.11*f*
- High-intensity carbon arc lamps, 15.21–15.23, 15.24*f*
- High-intensity discharge (HID) lamps, 40.33–40.36
 - applications for, 40.26*t*
 - characteristics of, 40.25*t*
 - CMH, 40.25*t*, 40.26*t*, 40.33, 40.36
 - construction of, 40.34*f*
 - emission spectrum of, 40.35*f*
 - Hg, 40.25*t*, 40.26*t*, 40.33, 40.36
 - HPS, 40.25*t*, 40.26*t*, 40.33
 - MH, 40.25*t*, 40.26*t*, 40.33, 40.35, 40.35*f*, 40.36
- High-intensity reciprocity failure, of photographic films, 29.12
- High-order harmonic generation, 21.2, 21.2*f*
- High-power diode lasers, 19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
- High-power laser arrays, 19.26–19.30, 19.28*f*, 19.28*t*, 19.29*t*, 19.30*f*
- High-power lasers, 19.24, 19.25*f*, 19.25*t*
- High-power semiconductor lasers, 19.18–19.30
 - arrays in, 19.26–19.29, 19.28*f*, 19.28*t*, 19.29–19.30, 19.29*t*, 19.30*f*
 - commercial diode, 19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
 - future directions for, 19.23–19.26, 19.25*f*, 19.25*t*, 19.26*t*, 19.27*f*
 - mode-stabilized lasers with reduced facet intensity, 19.18–19.19, 19.19*f*
- High-power strained QW lasers, 19.26, 19.26*t*
- High-pressure, short-arc xenon lamps, 15.35*f*
- High-pressure enclosed arcs, 15.24, 15.28–15.34
 - capillary mercury-arc lamps, 15.30–15.31, 15.31*f*

- High-pressure enclosed arcs (*Cont.*):
 compact-source arcs, 15.31–15.34,
 15.32f–15.35f
 Lucalox lamps, 15.30, 15.31f
 mercury arcs, 15.29, 15.30f
 multivapor arcs, 15.29, 15.31f
 Uviarc, 15.28–15.29, 15.29f, 15.30f
- High-pressure mercury-arc lamps,
 15.29f, 15.30f
- High-speed modulation, of semiconductor
 lasers, 19.30–19.36, 19.31f–19.36f
- High-speed optical recording systems, 19.3
- High-speed photoconductors, 26.20–26.23,
 26.21f–26.23f
- High-speed photodetectors, 26.1–26.24
 about, 26.3
 avalanche photodetectors, 26.17–26.20,
 26.18f, 26.20f, 26.21f
 photoconductors, 26.20–26.23, 26.21f–26.23f
pin photodiodes, 26.10, 26.12–26.15
 resonant, 26.15, 26.15f
 vertically illuminated, 26.3, 26.4f, 26.5,
 26.10, 26.12–26.13, 26.12f
 waveguide, 26.13–26.14, 26.14f
- Schottky photodiodes, 26.16, 26.16f, 26.17f
- speed limitations on, 26.5–26.10
 carrier transit time, 26.6–26.7, 26.6f
 carrier trapping, 26.9, 26.9f
 diffusion current, 26.8, 26.8f, 26.9
 packaging, 26.9–26.10, 26.10f, 26.11f
 RC time constant, 26.7–26.8, 26.7f
 structures of, 26.3–26.5, 26.4f
- High-speed photographic films, 30.18–30.20
- High-voltage power supply (HVPS), 31.1,
 31.9, 31.10f
- Hindle mounts, 6.19, 6.19f
- Hole current, 26.7
- Hole-accumulated photodiodes (HADs), 32.4f,
 32.8
- Hollow cathode lamps, 15.35, 15.37t–15.43t,
 15.44f
- Hollow lightpipes, 39.30–39.31
- Holograms, computer-generated (*see* Computer-
 generated holograms)
- Holographic compensators, 13.25
- Homogeneous broadening, emission-line, 16.5,
 16.6, 16.6f, 16.9
- Homogeneous reflectors, 39.39
- Homogeneous temperature change, 8.2–8.6,
 8.3t, 8.4t, 8.5f, 8.6f
- Homojunction lasers, 19.4
- Homojunctions, LED, 17.8–17.10, 17.9f–17.11f
- Horizontal illuminance, 40.7, 40.18f
- Horizontally illuminated photodetectors,
 26.4f, 26.5
- Hot cathode fluorescent lamps, 40.32
- Hot-electron bolometers, 24.29, 24.30, 24.30f,
 24.31f
- Hottel strings, 39.4, 39.14
- Hub mounting, 6.17, 6.18f
- Hubble telescope, 11.4, 13.24
- Hue, 40.5
- Human eye, 30.15–30.16, 30.16f, 34.6
 cones in, 36.8, 36.8f, 36.9f
 rods in, 36.8–36.10, 36.8f, 36.9f
 wavelengths detectable by the, 36.8–36.10,
 36.8f, 36.9f
- Humidity specifications, for lenses, 4.10
- Hybrid arrays, pyroelectric, 28.11–28.12,
 28.11f, 28.12f
- Hybrid FPAs:
 direct readout architectures of,
 33.15–33.18
 electronically scanned staring FPAs,
 33.16–33.17
 output circuits, 33.18
 TDI scanning FPAs, 33.17, 33.17f
 X-Y addressing and clock generation,
 33.16
- input circuits of, 33.18–33.23, 33.19t, 33.20f
 buffered direct injection, 33.19t, 33.20f,
 33.21
 capacitive transimpedance amplifier,
 33.19t, 33.20f, 33.22–33.23
 chopper-stabilized BDI, 33.19t, 33.20f,
 33.21–33.22
 direct detector integration, 33.18, 33.19t,
 33.20f
 direct injection, 33.18–33.21, 33.19t,
 33.20f, 33.21f
 gate modulation, 33.19t, 33.20f, 33.22
 thermal expansion match in, 33.14
- Hybrid reflectors (*see* Faceted reflectors)
- Hyde maxim, 3.22
- Hydrogen arc lamps, 15.49, 15.53f
- Hypo (film fixer), 29.5
- Ideal mode-locked lasers, 20.7
- Ideal thermal detectors, 28.2–28.3, 28.3f
- Ilford Photo Corporation, 29.25

- Illuminance, *37.4t*, *37.5*, *37.5f*, *39.2t*
 defined, *34.11*, *34.40*, *40.1*
 guidelines on levels of, *40.7t*
 and lighting design, *40.7*
 and luminance, *37.9*, *37.9f*
 retinal, *34.40–34.42*
 uniformity of, *40.7*
 unit conversions for, *36.7t*, *36.8t*
 units of, *34.43*
- Illuminance meters, *34.42*, *40.51*, *40.52f*
- Illuminated ceilings, *40.13f*
- Illuminated eyepieces, *12.12*, *12.12f*
- Illuminated objects (in stray light suppression),
 7.5, *7.5f*, *7.6f*
- Illuminating Engineering Society (IES),
 40.19
- Illumination:
 guidelines on, *40.7t*
 in nonimaging objects, *39.1*
 (See also Uniform illumination, of
 nonimaging optics)
- Illumination Engineering Society of North
 America (IESNA), *36.2*, *36.3*, *37.11*, *40.2*,
 40.7t
- Illumination subsystem, of nonimaging optics,
 39.22
- Image dissectors, *39.21*, *39.21f*
- Image height, *1.4*
- Image intensifiers (IIs), *31.7–31.18*,
 31.8f–31.10f
 defined, *31.8*
 input window/photocathode assemblies for,
 31.10–31.12, *31.11f*, *31.12f*
 MCP IIs, *31.7*, *31.9*, *31.9f*, *31.10f*
 and microchannel plates, *31.12–31.14*,
 31.12f, *31.13f*, *31.13t*
 phosphor screen assemblies for, *31.14–31.16*,
 31.14t, *31.15f*
 proximity-focused MCP IIs, *31.16–31.18*,
 31.17t, *31.18f*, *31.19f*
- Image irradiance, *4.7*
- Image lag, *32.6*
- Image processing, *8.12*
- Image quality, *4.6–4.7*
- Image sensors, *32.2–32.12*, *32.3f*, *32.21–32.34*
 antiblooming in, *32.9*, *32.10f*
 area arrays of, *32.24–32.32*, *32.25t*
 CCD performance, *32.32*
 frame transfer CCDs, *32.26–32.28*,
 32.27f, *32.28f*
- Image sensors, area arrays of (*Cont.*):
 interline transfer CCDs, *32.28–32.32*,
 32.29f–32.31f
 MOS, *32.25–32.26*, *32.26f*
 color imaging with, *32.32–32.34*, *32.33f*, *32.34f*
 dark current in, *32.10–32.12*, *32.11f*
 junction photodiodes, *32.3–32.6*, *32.4f*, *32.6f*
 linear arrays of, *32.21–32.24*, *32.22f*, *32.23f*
 MOS capacitors, *32.7–32.8*
 photoconductors, *32.8–32.9*
 pinned photodiodes, *32.8*
- Image size, *1.4*
- Image specifications, for lenses, *4.3*, *4.6–4.8*, *4.8f*
- Image structure, of photographic systems, *29.17*
- Imaged critical objects, *7.4*, *7.5f*
- Image-intensified (II) electronic imaging,
 31.1–31.30
 about, *31.2–31.3*, *31.3f*
 applications for, *31.27–31.30*
 active imaging, *31.29–31.30*
 day/night cameras, *31.28–31.29*
 mosaic II SSA cameras, *31.29*
 optical multichannel analyzers, *31.27–31.28*
 range gating and LADAR, *31.28*
 image-intensifier modules of, *31.7–31.18*,
 31.8f–31.10f
 input window/photocathode assemblies,
 31.10–31.12, *31.11f*, *31.12f*
 microchannel plates, *31.12–31.14*, *31.12f*,
 31.13f, *31.13t*
 phosphor screens, *31.14–31.16*, *31.14t*,
 31.15f
 proximity-focused MCP IIs, *31.16–31.18*,
 31.17t, *31.18f*, *31.19f*
 optical interface of, *31.3–31.7*
 considerations, *31.6–31.7*, *31.6f*
 photometry and camera lens, *31.5–31.6*
 quantum limited imaging conditions,
 31.3–31.4
 radiometry, *31.4–31.5*
 self-scanned arrays, *31.19–31.27*, *31.19f*
 electron-bombarded, *31.23–31.27*, *31.23f*,
 31.24f
 fiberoptic-coupled, *31.20–31.22*, *31.20f*,
 31.21f, *31.21t*
 lens-coupled, *31.22–31.23*, *31.22f*
 Image-intensified self-scanned arrays (II SSAs),
 31.19–31.27, *31.19f*
 for active imaging, *31.29–31.30*
 camera for, *31.5–31.6*

- Image-intensified self-scanned arrays (II SSAs)
(*Cont.*):
 electron-bombarded, 31.23–31.27, 31.23f, 31.24f
 fiberoptic-coupled, 31.20–31.22, 31.20f, 31.21f, 31.21t
 lens-coupled, 31.22–31.23, 31.22f
- Impedance:
 in amplifiers, 27.10–27.11
 in photodetectors, 24.19–24.20
- Impurity band conduction (IBC), 33.7
- Incandescence, 40.25
- Incandescent sources (of radiation), 40.25, 40.27–40.30
 calibration of, 34.31
 characteristics of, 40.25t
 elements of, 40.29f
 nongaseous, 15.15–15.21
 comparisons of, 15.19, 15.19f
 gas mantle, 15.17, 15.18, 15.19f
 glowbar, 15.17, 15.18f
 Nernst glower, 15.14, 15.15, 15.17, 15.17f
 quartz-envelope lamps, 15.20, 15.21
 tungsten-filament lamps, 15.19, 15.20, 15.20f–15.22f
 tungsten emissivity in, 40.28f
- Index-guided lasers, 19.8, 19.27, 19.28f
- Indicator lamps, LED, 17.27f, 17.29–17.30, 17.29f
- Indirect glare, 40.9
- Indirect lighting, 40.14, 40.15, 40.15f, 40.16f, 40.46f
- Indirect semiconductors, 17.4, 17.4f, 17.5f, 17.6
- Indium antimonide (InSb) hot-electron bolometers, 24.29, 24.30, 24.30f, 24.31f
- Indium antimonide (InSb) intrinsic photovoltaic detectors, 24.80–24.83, 24.82f, 24.83f
- Indium arsenide (InAs) photovoltaic detectors, 24.75, 24.77–24.78, 24.77f–24.79f
- Indium gallium arsenic phosphide (InGaAsP) laser material system, 19.7
- Indium gallium arsenide (InGaAs) detectors, 24.65–24.70, 24.66f–24.69f
- Indium gallium arsenide (InGaAs) photodetectors, 25.10, 25.10t
- Indium gallium arsenide (InGaAs) photodiodes, 24.66–24.70, 24.66f–24.69f, 34.31
- Indium gallium nitride (InGaN) HB-LEDs, 18.2f
- Indium phosphide (InP) laser material system, 19.7
- Induction lamps (ILs), 40.36–40.37
- Inductive pickup noise, 27.5, 27.6f
- Inductively coupled plasma (ICP), 18.3
- Industrial lighting, 40.60–40.61, 40.61f
- “Infant mortality period,” 17.26, 17.26f
- Infectious film developers, 29.5
- Infinitely distant objects, systems with, 1.5–1.6, 1.6f
- Information capacity, of photographic systems, 29.24
- Infrared detector arrays, 33.1–33.31
 about, 33.3–33.4
 applications for, 33.4
 current status of, 33.28–33.30, 33.28f, 33.29f, 33.29t
 future trends and technology directions of, 33.30–33.31, 33.30f, 33.31f
 hybrid FPAs, 33.14–33.23
 detector interface input circuit, 33.18–33.23, 33.19t, 33.20f, 33.21f
 hybrid readout, 33.15–33.23
 readout, 33.17f
 thermal expansion match in, 33.14
 monolithic FPAs, 33.10–33.14
 direct-charge-injection silicon FPAs, 33.11f, 33.13
 microbolometer FPAs, 33.13–33.14
 MIS photogate FPAs, 33.10–33.11, 33.11f, 33.12f
 scanning and staring, 33.14
 silicon FPAs, 33.11–33.13, 33.11f, 33.12f
 operating principles of, 33.7–33.10, 33.8f, 33.9f
 performance of, 33.23–33.28
 detectivity, 33.23–33.24
 minimum resolvable temperature, 33.27–33.28, 33.27f
 NE ΔT , 33.24–33.27, 33.25f, 33.26f
 percentage of BLIP, 33.24
 scanning and staring, 33.6–33.7, 33.6f
 spectral bands for, 33.4–33.5, 33.6f
- Infrared detectors:
 gallium-doped germanium, 24.100
 gallium-doped silicon, 24.95, 24.95f, 24.96, 24.96f
- Infrared film, 30.22
- Infrared filters, 40.12
- Infrared interferometry, 13.25
- Infrared LED chips, 17.8, 17.9, 17.9f
- Infrared photodetectors, 25.12, 25.15

- Infrared (IR) radiation, **34.6, 40.41**
 extreme, **25.2**
 far, **24.3, 25.2**
 forward looking, **33.4**
 long-wavelength, **24.3, 33.3–33.5, 33.6f**
 medium-wavelength, **24.3, 25.2, 33.3, 33.5, 33.6f**
 near, **24.3, 25.2**
 short-wavelength, **24.3, 33.3, 33.5**
 very long-wavelength, **24.3**
- Infrared radiometry, standards for, **15.11–15.12, 15.12f**
- Inhomogeneous broadening, emission-line, **16.5, 16.6f**
- Inhomogeneous media, **39.22**
- Inhomogeneous reflectors, **39.39**
- Injection-locked lasers, **19.37f, 19.38**
- Input circuits, of hybrid FPAs, **33.18–33.23, 33.19t, 33.20f**
 buffered direct injection, **33.19t, 33.20f, 33.21**
 capacitive transimpedance amplifier, **33.19t, 33.20f, 33.22–33.23**
 chopper-stabilized BDI, **33.19t, 33.20f, 33.21–33.22**
 direct detector integration, **33.18, 33.19t, 33.20f**
 direct injection, **33.18–33.21, 33.19t, 33.20f, 33.21f**
 gate modulation, **33.19t, 33.20f, 33.22**
- Input optics, **38.7**
- Input windows, of image intensifiers, **31.9–31.12, 31.9f, 31.11f, 31.12f**
- Institute for Electrical and Electronic Engineering (IEEE), **36.3**
- Insulators, light bulb, **40.29f**
- Integral color filter arrays (CFAs), **32.32–32.34, 32.33f, 32.34f**
- Integral density, **29.14**
- Integrated lasers, with 45° mirror, **19.39–19.40, 19.39f**
- Integrated transmittance, **35.3**
- Integrating cavities, of nonimaging optics, **39.24–39.27**
 efficiency vs. luminance, **39.26, 39.27f**
 modifying cavity output distribution, **39.27**
 with nonimaging concentrator/collectors, **39.26, 39.27**
 nonuniformities with spherical, **39.24–39.26, 39.24f, 39.25f**
- Integrating spheres (devices), **35.9, 35.11–35.13, 35.11f, 37.9–37.10, 37.10f**
- Integrating-bucket phase shifting, **13.21, 13.21f**
- Intensity:
 defined, **34.9, 40.1**
 luminous, **39.2t**
 nomenclature for, **36.4**
 radiant, **39.2t**
- Interface coupling, **20.14**
- Interference filters, **34.36**
- Interferograms, **13.14–13.18**
 from direct interferometry, **13.17–13.18**
 fixed, **13.14–13.15**
 Fourier analysis of, **13.16–13.17, 13.17f**
 interpolation of, **13.15–13.16**
- Interferometers, **13.7–13.12**
 Brunning distance-measuring, **12.8–12.9, 12.8f**
 common-path, **13.9, 13.11f**
 computer-generated holograms for, **14.4–14.5, 14.4f, 14.5f**
 distance-measuring, **12.7–12.10, 12.7f–12.9f**
 Fabry-Perot, **16.19f**
 Fizeau, **12.14, 13.8–13.9, 13.9f, 13.10f, 13.18, 14.4, 14.5f**
 Haidinger, **12.14**
 heterodyne, **13.22**
 lateral-shearing, **12.14, 13.9–13.12, 13.11f, 13.12f**
 Michelson, **12.5, 12.6, 12.14**
 microinterferometers, **10.13, 10.13f**
 multiple-pass, **13.13**
 multiple-reflection, **13.13**
 nonreacting, **12.7**
 point diffraction, **13.11f**
 radial-shearing, **13.12, 13.13f**
 reversing-shearing, **13.12, 13.13f**
 rotational-shearing, **13.12, 13.13f**
 sensitivity of, **13.13–13.14, 13.14f**
 single-shot f -to- $2f$, **21.6**
 Twyman-Green, **13.7–13.8, 13.7f, 13.8f, 13.18**
 Zernike phase-contrast method applied to, **13.13–13.14, 13.14f**
 (*See also specific interferometers, e.g.: Lateral-shearing interferometers*)
- Interferometric plots, for orthonormal aberrations, **11.36–11.37, 11.37f, 11.38f**
- Interferometry:
 of angles, **12.14**
 direct, **13.17–13.18**

- Interferometry (*Cont.*):
 infrared, 13.25
 of medium distances, 12.6–12.10, 12.7f–12.9f
 phase-shifting, 13.18–13.23, 13.18f–13.20f
 heterodyne interferometer, 13.22
 integrating bucket method, 13.21, 13.21f
 phase errors, 13.22
 phase stepping method of, 13.20, 13.20f
 phase-lock method, 13.23, 13.23f
 simultaneous measurement, 13.22
 two steps plus one method, 13.21, 13.22
 pulse-train, 20.12, 20.12f
 of small distances, 12.5, 12.6
 sub-Nyquist, 13.27
 two-wavelength, 13.25, 13.26
- Interior lighting, 40.55–40.61
 for health-care facilities, 40.58–40.60, 40.60t
 for industry, 40.60–40.61, 40.61f
 for offices, 40.55, 40.56t
 for residences, 40.57, 40.58, 40.59t
 for retail, 40.55–40.57, 40.56t–40.58t
- Interlayer interimage effects (IIEs), 30.19
- Interline transfer (IT) CCD image sensors, 32.28–32.32, 32.29f–32.31f
- Interline-transfer (IT) CCD FPAs, 33.11–33.13, 33.12f
- Internally processed (IP) photocathodes, 31.24
- International Astronomical Union (IAU), 36.3
- International Bureau of Weights and Measures (BIPM), 36.2
- International candle (unit), 37.3
- International Commission on Illumination (CIE), 36.2
- International Committee for Weights and Measures (CIPM), 36.2, 38.6
- International Graphics Exchange Specification (IGES), 40.19
- International Standards Organization (ISO), 4.10, 4.11, 36.2, 40.19
- International System of Units (*see* SI units)
- International Union of Pure and Applied Physics (IUPAP), 36.2
- Intervalence band absorption (IVBA), 19.17
- Interwoven pulse trains, 20.13
- Intrinsic athermalization, 8.7–8.8, 8.7f
- Intrinsic infrared detectors, 33.7, 33.8f
- Intrinsic photoconductors, 25.5, 25.5f
- Intrinsic photodetectors, 24.7, 24.7f
 germanium, 24.70–24.73, 24.70f–24.73f
 indium antimonide photovoltaic, 24.80–24.83, 24.82f, 24.83f
- Intrinsic semiconductor transition, 24.11
- Invariants, 1.11
- Inverse square law, 34.14, 37.8
- Inversion layer, in MOS transistors, 25.11, 25.11f
- Inverted channel substrate planar (ICSP) lasers, 19.20t, 19.23
- Involute reflectors, 39.11–39.12, 39.12f
- Ion bombardment strip lasers, 19.8, 19.9f
- Ion current measurement devices, 34.29
- Irradiance:
 defined, 34.8, 36.5, 37.4t, 37.5, 39.2t
 excitance and emittance vs., 39.2
 image specifications for, 4.7
 spectral, 36.14, 38.1–38.2, 38.11t, 38.13–38.16, 38.13f–38.16f
- Irradiance response units, 34.31
- Isoelectronic dopants, 17.16
- Isoelectronic trap, 17.6, 17.6f
- Isolated pulses, in attosecond optics, 21.4
- Isometric plots, for orthonormal aberrations, 11.36–11.37, 11.37f, 11.38f
- Isotope broadening, 16.6
- Isotropic (term), 36.4, 39.3
- Isotropic point source, 36.4
- Iterated rays, 3.12
- J-aggregates, 30.13, 30.14
- Jaynes-Cummings model, 23.15–23.17, 23.15f
- Jet polishing, 9.6
- Johnson noise (*see* Thermal noise)
- Johnson noise power density, 28.3
- Jones (unit), 24.11, 24.13
- Joule (unit), 34.5–34.6, 37.6
- Judd-Vos modified function, 36.10
- Junction photodiodes, 32.3–32.6, 32.4f, 32.6f
- Kaleidoscope effect, 39.28
- Keck telescope, 11.4
- Kick operator, 23.17
- Kirchhoff's law, 34.25, 35.7, 35.8t
- Knife-edge test (*see* Foucault test)
- Kodachrome film, 29.14, 30.23
- Kodacolor, 29.14
- Kodak Gold film, 30.25
- Kodak Royal Gold film, 30.25

- Kodak Technical Pan Film, **29.19t**
- Kohler illumination, **1.11, 1.11f, 39.23–39.24, 39.23f, 39.34, 39.35f**
- Laboratory sources (of radiation), **15.7–15.13**
 baseline standard for, **15.9, 15.9f, 15.10f, 15.12f**
 blackbody cavity theory, **15.7–15.9, 15.8f**
 working standards for, **15.9–15.13, 15.10f, 15.12f, 15.13f**
- Labsphere, **38.12**
- Lagrangian rays, **3.12**
- Lamb shift, **23.13**
- Lambda Research Corporation, **7.27**
- Lambert (unit), **34.43, 36.7, 36.8t**
- Lambertian approximation (of radiant flux transfer), **34.14–34.18**
 and lambertian sources, **34.14–34.17, 34.15f, 34.16f**
 radiant flux transfer through lambertian reflecting sphere, **34.17–34.18**
- Lambertian surface, **37.8**
- Lambert's cosine law, **37.8, 37.8f**
- Lamps:
 configurations of, **15.20f**
 modeling of, **40.17**
 standards for, **15.11**
 (See also *specific types of lamps, e.g.: Airway beacon lamps*)
- Lamps for Scientific Purposes* (G. M. B. H. Osram), **15.20**
- Land (term), **34.35, 34.35f**
- Lapping step (of optics fabrication), **9.5**
- Large-area detectors, **25.12**
- Laser(s), **12.7, 16.1–16.37**
 about, **16.2–16.3**
 diagram of, **16.2f**
 electromagnetic spectrum involving, **16.2, 16.3, 16.3f**
 and laser gain medium, **16.4–16.19**
 emission linewidth and line broadening of radiating species, **16.4–16.7, 16.6f, 16.7f**
 energy levels and radiation, **16.4**
 gain saturation, **16.10**
 optimization of output coupling from laser cavity, **16.13, 16.14**
 population inversions, **16.8–16.10, 16.12–16.13, 16.13f, 16.14f**
 pumping techniques to produce inversions, **16.14–16.19**
- Laser(s), and laser gain medium (*Cont.*):
 stimulated absorption and emission, **16.7–16.8, 16.8f**
 threshold conditions, **16.10–16.12, 16.11f**
 as light sources, **40.39**
 in measurement, **12.2, 12.6**
 and optical cavities or resonators, **16.19–16.25, 16.19f**
 configurations and cavity stability, **16.23–16.25, 16.24f, 16.25f**
 longitudinal laser modes, **16.20, 16.20f**
 transverse laser modes, **16.21–16.23, 16.21f–16.23f**
 probability flow diagram for, **23.23f**
 quantum theory of (*see* Quantum theory of lasers)
 as radiometric characterization tool, **34.32**
 semiconductor arrays of, **19.26–19.29, 19.28f, 19.28t**
 special laser cavities, **16.25–16.29**
 distributed feedback lasers, **16.29**
 mode-locking, **16.27–16.29, 16.28f**
 Q-switching, **16.26–16.27, 16.27f**
 ring lasers, **16.29**
 unstable resonators, **16.25–16.26, 16.26f**
 two-dimensional high-power arrays of, **19.29–19.30, 19.29t, 19.30f**
 as two-level system, **20.23–20.24, 20.25t**
 types of, **16.29–16.37, 16.31f, 16.32f, 16.34f, 16.35f, 16.37f**
 (See also *related topics, e.g.: Laser stabilization*)
- Laser beam expanders, athermal, **8.13–8.14**
- Laser cavities, **16.25–16.29**
 distributed feedback lasers, **16.29**
 mode-locking, **16.27–16.29, 16.28f**
 Q-switching, **16.26–16.27, 16.27f**
 ring lasers, **16.29**
 unstable resonators, **16.25–16.26, 16.26f**
- Laser detection and ranging (LADAR), **31.28, 31.30**
- Laser field, spectral properties of the, **23.28–23.31**
- Laser gain media, **16.4–16.19**
 dielectric solid-state, **16.32–16.34**
 gaseous, **16.30–16.31, 16.31f**
 liquid, **16.31–16.32, 16.32f**
 properties associated with, **16.4–16.19**
 emission linewidth and line broadening of radiating species, **16.4–16.7, 16.6f, 16.7f**
 energy levels and radiation, **16.4**

- Laser gain media, properties associated with (*Cont.*):
 gain saturation, 16.10
 optimization of output coupling from laser cavity, 16.13, 16.14
 population inversions, 16.8–16.10, 16.12–16.13, 16.13*f*, 16.14*f*
 pumping techniques to produce inversions, 16.14–16.19
 stimulated absorption and emission, 16.7–16.8, 16.8*f*
 threshold conditions with mirrors, 16.10–16.12, 16.11*f*
 in vacuum, 16.36–16.37, 16.37*f*
- Laser linewidth, 23.34–23.35
- Laser locking, frequency discriminators for, 22.12–22.14
- Laser master equation, 23.19–23.20, 23.19*f*
- Laser phase-transition analogy, 23.35–23.40, 23.37*t*, 23.38*f*, 23.39*f*
- Laser photon statistics, 23.22–23.26, 23.23*f*, 23.25*f*
- Laser power measurement, 34.32
- Laser resonators, 16.20*f*, 16.23–16.26, 16.23*f*–16.26*f*
- Laser scanning, 40.54
- Laser scribing, 17.24–17.25
- Laser stabilization, 22.1–22.24
 about, 22.1
 Allan Deviation, 22.2–22.3
 and frequency discriminators for laser locking, 22.12–22.14
 frequency vs. time drift, 22.2
 future directions for, 22.23–22.24
 and optical cavity-based frequency discriminators, 22.14–22.16, 22.17*f*
 quantifying frequency stability, 22.2
 and quantum resonance absorption, 22.16, 22.17
 representative/example designs of, 22.20–22.23, 22.22*f*
 and servos, 22.5–22.12
 Bode representation of servos, 22.5–22.6, 22.6*f*
 closed-loop performance, 22.8
 closed-loop stability issues, 22.8–22.12, 22.9*f*–22.1*f*, 22.10–22.12
 measurement noise, 22.7–22.8
 phase and amplitude responses, 22.6–22.7, 22.6*f*, 22.7*f*
- Laser stabilization (*Cont.*):
 spectral noise density, 22.3–22.5
 and transducers, 22.17–22.20
- Laser stripe structures, 19.8, 19.9*f*
- Lasers and Optronics Buying Guide*, 15.14
- Lasing, without inversion, 23.40–23.42, 23.41*f*
- Latent image (LI), 30.7
- Latent-image speck, 29.5
- Lateral aberration curves (*see* Transverse ray plots)
- Lateral antiblooming, 32.9, 32.10*f*
- Lateral color, 1.14, 2.5, 2.5*f*
- Lateral *pin* photodetectors, 25.15–25.16, 25.15*f*
- Lateral-collection photodetectors, 26.4*f*
- Lateral-shearing interferometers, 12.14, 13.9–13.12, 13.11*f*, 13.12*f*
- Lathe assembly technique, 6.7–6.8, 6.8*f*
- Lead salt lasers, 19.7–19.8
- Lead selenide (PbSe) detectors, 24.76*f*, 24.78, 24.79, 24.79*f*–24.82*f*
- Lead sulfide (PbS) photoconductors, 24.73–24.74, 24.74*f*–24.77*f*, 24.74*t*
- Lead tin telluride (PbSnTe) photovoltaic detectors, 24.92, 24.93*f*
- Lead-in wires, light bulb, 40.29, 40.29*f*
- Leakage current, 24.19–24.20, 32.10–32.12, 32.11*f*
- Leaky-mode arrays, 19.29
- Least-squares method, 3.17–3.19
- Legacy films, 30.23–30.25
- Legal traceability (of calibration), 34.21
- Legendre polynomials, 11.5, 11.30
- Length measurements, 12.2–12.10
 interferometric measurement, 12.5–12.10, 12.7*f*–12.9*f*
 stadia and range finders, 12.2–12.4, 12.3*f*, 12.4*f*
 standards for, 12.2
 time-based and optical radar, 12.4, 12.5, 12.6*f*
- Lens(es), 39.32–39.37, 39.33*f*
 assembly adjustment of, 6.8–6.11, 6.10*f*
 coating specifications for, 4.10
 component specifications for, 4.2, 4.3
 and concentrators, 39.16–39.18, 39.18*f*
 data entry for, 3.2–3.8
 humidity specifications for, 4.10
 for II electronic imaging, 31.5–31.6
 image specifications for, 4.3, 4.6–4.8, 4.8*f*
 literature on nonimaging, 39.32, 39.33

- Lens(es) (*Cont.*):
- multicomponent assemblies of (*see* Multicomponent lens assemblies)
 - optical parameters for, 4.9
 - perfect, 4.4
 - reflector/lens-array combinations, 39.34–39.37, 39.36f, 39.37f
 - single-lens arrays, 39.33–39.34, 39.34f
 - tandem-lens arrays, 39.34, 39.35f–39.37f
 - thermal defocus of, 8.4, 8.5f
 - thermal focus shift of, 8.2–8.4, 8.3t, 8.4t (*See also specific lenses, e.g.: Air-spaced doublet lens*)
- Lens design:
- aberration curves in, 2.1–2.6
 - software for, 40.20
 - and tolerancing calculations, 5.10
- Lens setup routine (in optical software), 3.7
- Lens-coupled II SSAs, 31.22–31.23, 31.22f
- Levels (tools), 12.13–12.14, 12.13f, 12.14f
- Lever mechanism mountings, 6.19, 6.19f
- Lifetime classification, of photodetector materials, 26.5
- Light:
- spectrum of, 25.2
 - speed of, 12.2
- Light amplification by stimulated emission of radiation (*see* Laser(s))
- Light detection and ranging (LIDAR), 25.12
- Light distribution, 40.5
- Light extraction, 17.6–17.8, 17.7t
- Light loss factor (LLF), 40.17
- Light piping, 30.6–30.7, 30.6f
- Light pollution, 40.43, 40.62
- Light quanta, 23.6–23.9, 23.8f
- Light scattering, 30.5–30.7, 30.6f
- Light shelves, 40.48, 40.50f
- Light sources, 40.24–40.41, 40.28f
- applications for, 40.26t
 - carbon arc sources, 40.40
 - characteristics of, 40.25t
 - daylight, 40.40–40.41, 40.40f
 - electrodeless lamps, 40.36–40.37
 - electroluminescent sources, 40.37–40.39, 40.38f, 40.38t, 40.39f
 - fluorescent lamps, 40.30–40.33, 40.31f, 40.32f
 - glow lamps, 40.39
 - high-intensity discharge lamps, 40.33–40.36, 40.33f–40.35f
- Light sources (*Cont.*):
- incandescent sources, 40.25, 40.27–40.30, 40.28f, 40.29f
 - low-pressure sodium lamps, 40.33–40.36, 40.33f–40.35f
 - nuclear sources, 40.39
 - pure Xe arc lamps, 40.39
 - short arc sources, 40.39
 - types of, 40.27f
- Light stability, 30.10
- Light stabilization, 30.12–30.13
- Light trespass, 40.43, 40.62
- Light-absorbing dye, 30.7
- Light bulbs, 40.27–40.30
- base of, 40.29, 40.29f
 - CFL/fluorescent/minature, 40.28f
 - elements of, 40.29f
 - shapes of, 15.20f, 15.30f
 - sizes of, 15.30f
 - types of, 40.27t
- Light-emitting diodes (LEDs), 17.1–17.34
- conversion of, luminous intensity to radiant intensity, 36.13, 36.13f
 - device structures of, 17.8–17.15
 - diffused homojunctions, 17.9–17.10, 17.10f, 17.11f
 - double heterojunctions, 17.13, 17.13f–17.15f
 - grown homojunctions, 17.8, 17.9, 17.9f
 - single heterojunctions, 17.12, 17.12f
 - epitaxial technology for, 17.21–17.23
 - lamps with, 40.37–40.39
 - applications for, 40.26t
 - characteristics of, 40.25t
 - materials/emitted colors of, 40.38t
 - photonic crystal, 40.39f
 - structure of, 40.38f
 - and LED-based products, 17.29–17.32, 17.29f–17.31f
 - and light extraction, 17.6–17.8, 17.7t
 - and light-generation processes, 17.2–17.6, 17.3f–17.6f
 - material systems for, 17.15–17.19
 - AlInGaP system, 17.18, 17.19f
 - Al_xGa_{1-x}As system, 17.17, 17.17t
 - blue LED technology, 17.18, 17.19
 - GaAs_{1-x}P_x system, 17.15–17.17, 17.16t
 - octocouplers in, 17.32–17.34, 17.32f
 - production levels for, 17.2
 - quality/reliability of, 17.25–17.28

- Light-emitting diodes (LEDs) (*Cont.*):
 substrate technology for, 17.20–17.21, 17.20*t*
 wafer processing for, 17.23–17.25
 (*See also specific light-emitting diodes, e.g.:*
 High-brightness visible LEDs)
- Lighting, 40.1–40.71
 about, 40.1–40.3
 exterior, 40.61–40.62, 40.63*t*
 functions of, 40.12–40.14, 40.13*f*–40.16*f*
 insufficient, 40.9
 interior, 40.55–40.61
 health-care facility lighting, 40.58–40.60,
 40.60*t*
 industrial lighting, 40.60–40.61, 40.61*f*
 office lighting, 40.55, 40.56*t*
 residential lighting, 40.57, 40.58, 40.59*t*
 retail lighting, 40.55–40.57, 40.56*t*–40.58*t*
 perception of, 40.4–40.6, 40.4*f*
 for transportation, 40.63–40.71
 roadway lighting, 40.67, 40.69–40.71,
 40.70*t*, 40.71*t*
 vehicular lighting, 40.63–40.67, 40.64*f*,
 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
 vision biology, 40.3–40.6
 (*See also Luminaires*)
- Lighting design, 40.6–40.23
 and color, 40.7–40.9
 and context, 40.6
 and functions of lighting, 40.12–40.14,
 40.13*f*–40.16*f*
 geometries in, 40.13*f*, 40.14–40.15, 40.15*f*,
 40.16*f*
 goals of, 40.6
 and illuminance, 40.7, 40.7*t*
 and properties of objects and impact, 40.16
 system layout and simulation in, 40.16–
 40.23, 40.18*f*, 40.22*f*–40.24*f*
 and visual discomfort, 40.9–40.12, 40.11*t*
- Lighting Design and Application (IESNA),
 39.8
- Lighting design software, 40.21
- Lighting geometries, 40.13*f*, 40.14–40.15,
 40.15*f*, 40.16*f*
- Lighting Handbook (IESNA), 36.7, 40.17
- Lighting measurement, 40.51–40.54
 goniometers, 40.52–40.53, 40.53*t*, 40.54*f*
 illuminance meters, 40.51, 40.52*f*
 luminance meters, 40.52, 40.52*f*
 reflectometers, 40.52
 surface measurement systems, 40.53, 40.54
- Lighting system layout and simulation,
 40.16–40.23, 40.18*f*
 computer graphics software for, 40.21–40.23,
 40.22*f*, 40.23*f*
 IGES standard for, 40.19
 optical analysis and design software for,
 40.20
 optical design and analysis software for,
 40.20–40.21
 software tools for, 40.18–40.23, 40.22*f*–40.24*f*
 source modeling software for, 40.19–40.20
 STEP standard for, 40.19
- Lightness (term), 40.4
- Light-output degradation, 17.26–17.28, 17.28*f*
- Lightpipes, 39.27–39.32
 angular uniformity of, 39.31
 applications for, 39.32
 length of, 39.30
 periodic distributions of, 39.30
 shapes of, 39.27–39.30, 39.28*f*–39.30*f*
 solid vs. hollow, 39.30–39.31
 tapered, 39.12–39.13, 39.13*f*, 39.31–39.32,
 39.31*f*
- LightTools (optical software), 7.26
- Light-trap silicon photodiodes, 34.30
- Limiting resolution, of EBSSAs, 31.27
- Linear image sensor arrays, 32.2, 32.21–32.24,
 32.22*f*, 32.23*f*
- Linear lasers, 20.18–20.19, 20.19*f*
- Linear variable differential transformers
 (LVDTs), 10.12
- Linearity, of photoemissive detectors, 24.41
- Line-scanned imaging systems, 31.29
- Linewidth, spectral density and, 22.4–22.5
- Lippmann emulsions, 29.4
- Liquid Encapsulated Czochralski (LEC)
 technique, 17.20, 17.21
- Liquid laser gain media, 16.31–16.32, 16.32*f*
- Liquid phase epitaxy (LPE), 17.21–17.22, 19.6,
 19.19
- Lit-appearance modeling, 40.21–40.23, 40.23*f*,
 40.24*f*
- Lithographic etching, 18.3–18.4, 18.3*f*
- Lithographic projection lens, 6.10*f*
- LLL sensitivity, 31.19, 31.20
- Local area networks, 19.3
- Localized avalanche breakdown, 17.28
- Lock-in amplifiers, 27.14, 27.14*f*, 27.15, 38.10
- Long wavelength lasers, 19.8

- Long wavelength QW lasers, 19.17–19.18, 19.17*f*
- Longitudinal laser modes, 16.19*f*, 16.20, 16.20*f*
- Long-wavelength infrared (LWIR), 24.3, 33.3–33.5, 33.6*f*
- Lorentzian line shape, 16.6*f*
- Louvers, 40.41, 40.45*f*, 40.46
- Low-beam headlamps, 40.64–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*
- Low-intensity carbon arc lamps, 15.21–15.24, 15.24*f*, 15.28*f*
- Low-intensity reciprocity failure, of photographic films, 29.12
- Low-level-light (LLL) TV imaging, 31.1–31.4
- Low-pressure enclosed arcs, 15.35–15.47
 - black-light fluorescent lamps, 15.35, 15.36*t*
 - electrodeless discharge lamps, 15.36, 15.44
 - germicidal lamps, 15.35
 - hollow cathode lamps, 15.35, 15.37*t*–15.43*t*, 15.44*f*
 - Pluecker spectrum tubes, 15.47, 15.47*f*, 15.47*t*
 - spectral lamps, 15.44, 15.45, 15.45*f*, 15.46*f*, 15.46*t*
 - Sterilamps, 15.35, 15.36*f*
- Low-pressure lamps, 15.36*f*
- Low-pressure sodium (LPS) lamps, 40.33–40.36
 - applications for, 40.26*t*
 - characteristics of, 40.25*t*
 - construction of, 40.34*f*
 - emission spectrum of, 40.35*f*
- Low-speed photographic films, 30.18–30.20
- Low-temperature bolometers, 24.31–24.32, 24.32*f*, 24.33*f*, 28.5
- Low-temperature (LT) grown photoconductors, 26.23
- Lucalox lamps, 15.30, 15.31*f*
- Lumen (unit), 36.6, 37.6, 39.2*t*
- Lumen lighting simulation, 40.17
- Luminaires, 40.24–40.50
 - applications for, 40.26*t*
 - calculation of needed, 40.16–40.17
 - characteristics of, 40.25*t*
 - classification system for, 40.43–40.45, 40.43*t*, 40.44*f*
 - defined, 40.1
 - design of, 40.41–40.50
 - conics shapes and intensity distribution, 40.43*t*
 - etendue and source coupling, 40.41–40.42
 - design of (*Cont.*):
 - luminaire classification system, 40.43–40.45, 40.43*t*, 40.44*f*
 - methods, 40.42–40.43, 40.43*t*
 - light sources for, 40.24–40.41, 40.27*f*, 40.28*f*
 - carbon arc sources, 40.40
 - daylight, 40.40–40.41, 40.40*f*
 - electrodeless lamps, 40.36–40.37
 - electroluminescent sources, 40.37–40.39, 40.38*f*, 40.38*t*, 40.39*f*
 - fluorescent lamps, 40.30–40.33, 40.31*f*, 40.32*f*
 - glow lamps, 40.39
 - high intensity discharge lamps, 40.33–40.36, 40.33*f*–40.35*f*
 - incandescent sources, 40.25, 40.27–40.30, 40.28*f*, 40.29*f*
 - low-pressure sodium lamps, 40.33–40.36, 40.33*f*–40.35*f*
 - nuclear sources, 40.39
 - pure Xe arc lamps, 40.39
 - short arc sources, 40.39
 - optics of, 40.45–40.50
 - for artificial sources, 40.45–40.47, 40.45*f*, 40.46*f*
 - backlighting, 40.47, 40.47*f*, 40.48*f*
 - for daylight sources, 40.47–40.50, 40.49*f*–40.51*f*
- Luminance, 37.4*t*, 37.5, 37.5*f*, 39.2*t*
 - calibration of, 34.42
 - defined, 34.11, 34.40, 36.7, 40.1
 - and illuminance, 37.9
 - of integrating cavities, 39.26
 - in nonimaging optics, 39.2*t*, 39.3
 - uniformity of, 40.7
 - units of, 34.43
- Luminance contrast, 40.6, 40.10
- Luminance meters, 40.52, 40.52*f*
- Luminance ratio, 40.7
- Luminous efficacy, 18.5
- Luminous efficiency, 17.15
- Luminous energy, 37.4*t*, 37.6
- Luminous exitance, 37.4*t*, 37.5, 37.5*f*
- Luminous exposure, 37.4*t*, 37.6
- Luminous flux, 15.11, 34.10–34.11, 34.39, 34.42, 37.4, 37.4*t*, 37.6, 38.2
- Luminous flux density, 34.11
- Luminous intensity, 15.11, 34.11, 34.40, 37.4, 37.4*t*, 39.2*t*
- Lux (unit), 34.43, 36.7, 36.7*t*

- Luxmeters, 40.52*f*
Lux-second, 29.6
Lyot stops, 7.8–7.10, 7.8*f*–7.11*f*
- Macrofocal reflectors, 39.11
Magnesium, as *p*-type impurity, 17.23
Magnetic shielding, 27.10
Magneto-rheological finishing (MRF), 9.5
Magnification, 1.4
Magnifiers, first-order layout for, 1.8
Main event loop, 3.7
Maksutov sphere, 14.7, 14.8*f*
Maksutov test, 14.7–14.9, 14.8*f*
Mandel Q_M parameter, 23.26
Markov approximation, 23.21
Martin Black coating, 7.14–7.17, 7.14*f*–7.16*f*,
7.23, 7.25*f*
Masers, micro- (*see* Micromasers)
Master-oscillator power amplifier (MOPA),
19.41
Materials:
 formation of, for optics, 9.3–9.4
 specifications for, 4.9
 tolerancing and properties of, 5.9
Matte, 30.3
Maximized D star, 24.11
Maximum spectral luminous efficiency
 (of radiation), 37.2
Maxwell's principle, 30.16
Mazars (microwave amplification by z-motion-
induced emission of radiation), 23.45
Mean-square-spot size (MSS), 3.15
Mean-time-to-failure (MTTF), 25.13, 25.14
Measurement(s), 35.8–35.16
 of absorptance, 35.10
 of emittance, 35.14–35.16, 35.15*f*
 of lighting, 40.51–40.54, 40.52*f*, 40.53*t*, 40.54*f*
 of reflectance, 35.10–35.13, 35.10*f*–35.12*f*,
35.14*t*
 terminology, 35.2–35.3, 35.2*f*
 of transmittance, 35.8–35.10, 35.9*f*
 (*See also specific types of measurement, e.g.:*
 Nonlinearity measurement)
Measurement noise, 22.7–22.8
Measurements Assurance Program (MAP), 35.13
Mechanical athermalization, 8.8–8.12
 active, 8.11, 8.11*f*
 by image processing, 8.12
 part active, part passive, 8.11–8.12, 8.12*f*
 passive, 8.8–8.10, 8.8*f*–8.10*f*
- Mechanical scribing, 17.24
Mechanical specifications:
 for lenses, 4.9
 optical vs., 4.2
Mechanical tolerances, 5.2
Mechanical vibrations, 27.5, 27.6*f*
Mechanically clamped mountings, 6.12, 6.13*f*
Medium-wavelength infrared (MWIR) radia-
tion, 24.3, 25.2, 33.3, 33.5, 33.6*f*
Memory colors, 30.21
Mercury arc lamps, 15.29, 15.29*f*–15.31*f*, 15.34*f*
Mercury cadmium telluride (HgCdTe)
 detectors, 24.86–24.92, 24.87*f*
 infrared, 33.7, 33.8*f*
 photoconductors, 24.86*f*, 24.88–24.90,
24.89*f*–24.91*f*
 photodetectors, 25.10, 25.10*t*
 photovoltaic, 24.86*f*, 24.88*f*, 24.90–24.92,
24.91*f*, 24.92*f*
Mercury-doped germanium detectors, 24.84*f*,
24.92–24.95, 24.93*f*–24.95*f*
Mercury-free fluorescent lamps, 40.31
Mercury-halide fluorescent lamps, 40.33*f*
Mercury-vapor fluorescent lamps, 40.30,
40.31, 40.33*f*
Mercury-xenon lamps, 15.34*f*
Meridional rays, 3.3
Merit function, in optical design software, 3.17
Mesa etching, 18.3–18.4, 18.3*f*
Mesa photodiodes, 25.14, 25.15*f*
Mesopic vision, 34.37, 36.9, 37.2
Metal insulator semiconductor (MIS) photo-
gate FPAs, 33.10–33.11, 33.11*f*, 33.12*f*
Metal insulator semiconductors (MISs), 33.4
Metallic mirrors, mounting of, 6.19–6.20, 6.20*f*
Metalorganic chemical vapor deposition
(MOCVD), 19.6–19.7, 19.20*t*, 19.23
Metalorganic vapor phase epitaxy, 17.21, 17.22
Metal-oxide semiconductors (MOSs):
 linear arrays of, 32.21–32.24, 32.22*f*, 32.23*f*
 readouts from, 32.20–32.21
Metal-oxide-semiconductor (MOS) area array
image sensors, 32.25–32.26, 32.26*f*
Metal-oxide-semiconductor (MOS) capacitors,
32.4*f*, 32.7–32.8
Metal-oxide-semiconductor (MOS) detectors,
25.11, 25.11*f*
Metal-semiconductor-metal (MSM)
photodetectors, 26.3, 26.4*f*
Meter (unit), 12.2

- Meter, 1875 Treaty of the, **34.20**, **36.2**
 Meter candle (unit), **34.43**, **36.7**
 Metrology, optical, **12.1–12.25**
 angle measurements, **12.10–12.17**
 autocollimeters, **12.11–12.12**, **12.11f**, **12.12f**
 interferometric methods, **12.14**
 levels (tools), **12.13–12.14**, **12.13f**, **12.14f**
 mechanical methods, **12.10–12.11**, **12.11f**
 in prisms, **12.14–12.16**, **12.15f–12.17f**
 theodolites, **12.13**
 curvature measurements, **12.17–12.25**
 mechanical methods, **12.17–12.19**, **12.18f**,
 12.19f, **12.19t**
 optical methods, **12.19–12.21**, **12.20f**,
 12.20t, **12.21f**
 of diamond-turned optics, **10.12–10.13**,
 10.12f, **10.13f**
 focal length measurements, **12.21–12.25**,
 12.22f–12.24f
 length measurements, **12.2–12.10**
 interferometers, **12.5–12.10**, **12.7f–12.9f**
 stadia and range finders, **12.2–12.4**, **12.3f**,
 12.4f
 time-based and optical radar, **12.4**, **12.5**,
 12.6f
 straightness measurements, **12.10**
 terminology, **12.2**
 Michelson interferometers, **12.5**, **12.6**, **12.14**
 Microbolometer FPAs, **33.13–33.14**
 Microchannel plate tubes (MCPTs), **24.32**,
 24.33f, **24.40**
 Microchannel plates (MCPs), **31.1**, **31.9**, **31.9f**,
 31.12–31.14, **31.12f**, **31.13f**, **31.13t**
 Microchannel-plate image intensifiers
 (MCP IIs), **31.7**
 high-voltage power supply for, **31.9**, **31.10f**
 proximity-focused, **31.9**, **31.9f**, **31.16–31.18**,
 31.17t, **31.18f**, **31.19f**
 Microdensitometers, **29.6**, **29.15–29.16**, **29.15f**
 Microinterferometers, **10.13**, **10.13f**
 Microlens arrays, **32.30**, **32.31**, **32.31f**
 Microlenses, **12.24**
 Micromaser master equation, **23.20–23.22**
 Micromasers, **23.26–23.27**, **23.27f**, **23.45**
 Microminiature lamps, **15.53**
 Microplasma formation, **17.28**
 Microscopes:
 first-order layout for, **1.8**
 Nomarski, **10.11**, **10.11f**
 traveling, **12.20**, **12.21**, **12.21f**
 Microwave powered lamps (*see* Electrodeless
 sulfur lamps)
 Military Sensing Information Analysis Center
 (SENSIAC), **15.6**
 Millilambert (unit), **36.7**
 Milliphot (unit), **36.7**, **36.7t**
 Miniature lamps, **15.53**, **40.26t**, **40.28f**
 Minimum resolvable temperature (MRT)
 (of infrared detector arrays), **33.2**,
 33.27–33.28, **33.27f**
 Minkowitz, S., **12.7**
 Minkowitz distance-measuring interferometers,
 12.7, **12.8**, **12.8f**
 Minority-carrier recombination, **17.2**
 Mirror scatter relationships (stray light), **7.18**
 Mirrored tiling, **39.28**, **39.30**, **39.30f**
 Mirrors:
 on amplifiers, **16.3**
 and concentrators, **39.17**
 mounting of [*see* Mounting (of optical
 components)]
 nonabsorbing, **19.23–19.24**
 threshold conditions with, **16.10–16.12**,
 16.11f
 Mixing rods (*see* Lightpipes)
 Modal dispersion, **17.34**
 Mode-locked lasers, **16.27–16.29**, **16.28f**, **20.7**
 Mode-stabilized lasers, **19.18–19.19**, **19.19f**
 Modulation, as laser stabilization technique,
 22.13–22.14
 Modulation noise, **24.11**
 Modulation transfer function (MTF), **4.1**, **4.2**,
 4.4, **4.5f**, **4.8f**, **29.18**, **29.18t**, **29.19f**, **30.5**,
 31.26–31.27, **33.2**, **33.7**
 Modulators, electro-optic, **22.14**, **22.20**
 MOFSET bucket brigades, **33.17**
 Moiré deflectometry, **12.23**, **12.24f**
 Moiré tests, of spherical aberrations, **1**
 3.26–13.27
 Molecular beam epitaxy (MBE), **17.21–17.23**,
 19.6, **19.7**, **25.16**
 Moments normalization (technique),
 36.15–36.16, **36.15t**, **36.16f**
 Monochromators, **35.9**, **38.7–38.8**, **38.14–38.16**,
 38.14f–38.16f
 Monolithic built-up mirror substrate, **6.18**, **6.18f**
 Monolithic FPAs, **33.10–33.14**, **33.11f**, **33.12f**
 Monolithic LED displays, **17.30**
 Monolithic silicon bolometers, **28.10–28.11**,
 28.10f

- Monolithic two-dimensional (2D) laser arrays, **19.39**
- Monte Carlo simulations, **39.7**
- Mordants, **30.4**
- Mosaic II SSA cameras, **31.29**
- Mounting (of optical components), **6.1–6.24**, **6.17f**, **6.18f**
 and contact stresses, **6.21**
 of domes, **6.11**, **6.12f**
 hard, **6.2–6.4**
 of individual rotationally symmetric optics, **6.2–6.5**
 lever-mechanism, **6.19**, **6.19f**
 of moderate-sized mirrors, **6.17–6.20**, **6.20f**
 in multicomponent lens assemblies, **6.5–6.11**
 drop-in assembly, **6.6**, **6.6f**
 lathe assembly, **6.7–6.8**, **6.8f**
 lens adjustments at assembly, **6.8–6.11**, **6.10f**
 “poker chip” assembly, **6.8**, **6.9f**
 tightly toleranced assembly, **6.7**, **6.7f**
 of small mirrors/prisms, **6.11–6.17**
 bonded mountings, **6.13–6.15**, **6.15f**
 elastomeric mountings for mirrors, **6.12**
 flexure mountings for small mirrors/prisms, **6.15–6.17**, **6.16f**
 mechanically clamped mountings, **6.12**, **6.13f**
 spring-loaded mountings, **6.13**, **6.14f**
 soft, **6.4–6.5**
 and temperature effects, **6.21–6.24**, **6.22f**
 of windows, **6.11**, **6.11f**
- Multichannel detectors, **38.9**
- Multicomponent lens assemblies, **6.5–6.11**
 drop-in, **6.6**, **6.6f**
 lathe, **6.7–6.8**, **6.8f**
 and lens adjustments, **6.8–6.11**, **6.10f**
 “poker chip,” **6.8**, **6.9f**
 tightly toleranced, **6.7**, **6.7f**
- Multiphase reflectors, **39.39**
- Multiple quantum well (MQW) LEDs, **18.1**, **18.2f**
- Multiple surface concentrators, **39.16–39.17**, **39.17f**
- Multiple wafer transducers, **22.18–22.19**
- Multiple-pass interferometers, **13.13**
- Multiple-reflection interferometers, **13.13**
- Multiple-segment LED displays, **17.31**
- Multiplier photodiodes, **24.11**
- Multiplier phototubes (*see* Photomultiplier tubes (PMTs))
- Multiplier tubes, **24.6**
- Multipop (film exposure technique), **30.22**
- Multiquantum well (MQW) lasers, **19.14**, **19.15**, **19.15f**, **19.24**, **19.25t**
- Multiquantum-well buried heterostructure (MQW BH) lasers, **19.25t**
- Multispeed choppers, **15.14**
- Multivapor arcs, **15.29**, **15.30f**, **15.31f**
- Murty’s lateral shear interferometer, **12.14**, **13.12f**
- Mylar film, **29.4**
- National Physical Laboratory (NPL), **15.4**
- National Search Engine for Standards, **4.11**
- Natural broadening, emission-line, **16.5**, **16.6**
- Natural color (NC) film, **30.27**
- Natural linewidth (of transition), **16.5**
- Near infrared (NIR) radiation, **24.3**, **25.2**
- Negative color photographic films, **30.25–30.28**, **30.27t**
- Neodymium (Nd) glass lasers, **16.32–16.33**
- Neodymium-yttrium vanadate (Nd:VO₄) lasers, **16.33**
- Neodymium-yttrium-aluminum-garnet (Nd:YAG) lasers, **16.32–16.33**
- Neodymium-yttrium-lithide-fluoride (Nd:YLF) lasers, **16.33**
- Neon signs, **40.39**
- Nernst glower, **15.14**, **15.15**, **15.17**, **15.17f**, **15.19**, **15.19f**
- Nesonian illumination (*see* Abbe illumination system)
- Neutral density filters, **40.52**
- Nit (unit), **34.43**, **36.7**, **36.8t**
- Nitride LEDs, **17.19**, **18.3f**
- Nitrogen doping, **17.16**
- Nitrogen-doped GaAsP, **17.16**, **17.21–17.22**
- Nodal slide bench, **12.22**, **12.22f**
- Noise, **27.3–27.6**
 1/*f*, **27.4**
 atomic, **23.34–23.35**
 in CCDs, **32.20**, **32.32**
 excess, **24.11**
 generation, **24.11**
 generation-recombination, **24.11**
 ground loop, **27.5**, **27.6f**
 inductive pickup, **27.5**, **27.6f**
 measurement, **22.7–22.8**
 modulation, **24.11**
 nonessential, **27.4–27.6**, **27.5f**, **27.6f**

- Noise (*Cont.*):
 pattern, 32.12
 in photodetectors, 24.19–24.20, 24.20f
 in photoemissive detectors, 24.39, 24.40, 24.41f
 resistive coupling, 27.5, 27.6f
 RMS, 24.12
 shot, 24.12, 27.3, 27.3f, 32.12
 and signal detection, 27.1
 spatial, 33.26–33.27, 33.26f
 stray capacitance, 27.5, 27.6f
 temperature, 24.12
 thermal, 24.13, 27.4, 32.20
- Noise equivalent irradiance (NEI), 24.11
- Noise equivalent power (NEP), 24.10, 24.12, 24.14, 25.12, 28.2, 38.9, 38.10, 38.10t
- Noise equivalent quanta (NEQ), 29.1, 29.23
- Noise equivalent temperature difference (NETD), 28.8–28.9, 33.2, 33.24–33.27, 33.25f, 33.26f
- Noise spectral density (NSD), 33.2
- Noise spectrum, 24.11
- Nomarski microscope, 10.11, 10.11f
- Nonabsorbing mirrors (NAMs), 19.23–19.24
- Nonblackbody radiation source, 15.10
- Nonchromogenic film, 29.14
- Noncircular pupils, 11.4, 11.37, 11.39
- Non-diffraction-limited optics, 8.6
- Nonequilibrium errors, 34.25
- Nonessential noise, 27.4–27.6, 27.5f, 27.6f
- Nonessential ray tracing, 40.20–40.21
- Nonideal aperture, 34.35–34.36, 34.35f
- Nonimaging concentrators (*see* Concentrators, nonimaging)
- Nonimaging optics, 39.1–39.41
 about, 39.1–39.2
 aspheric lenses in, 39.8, 39.9, 39.9f
 calculations for, 39.2–39.6
 clipped Lambertian distribution, 39.3–39.4
 concentration, 39.5, 39.6
 dilution, 39.6
 etendue, 39.2, 39.3, 39.4f
 Hottel strings, 39.4, 39.4f
 Lambertian, 39.3
 luminance, 39.3
 projected solid angle, 39.5, 39.5f
 solid angle, 39.5, 39.5f
 concentration of, 39.12–39.22
 2D vs. 3D geometries, 39.20–39.21, 39.20f, 39.21f
 calculation, 39.5, 39.6
- Nonimaging optics, concentration of (*Cont.*):
 compound elliptical collectors, 39.14, 39.15f
 compound hyperbolic collectors, 39.15, 39.15f, 39.16f
 compound parabolic collectors, 39.13–39.14, 39.13f, 39.14f
 dielectric compound parabolic collectors, 39.15, 39.16, 39.16f
 edge rays, 39.22
 geometrical vector flux, 39.21–39.22
 inhomogeneous media, 39.22
 multiple surface concentrators, 39.16–39.17, 39.17f
 restricted exit angle concentrators with lenses, 39.18, 39.18f
 tapered lightpipes, 39.12–39.13, 39.13f
 θ_1/θ_2 concentrators, 39.18–39.20, 39.19f
 conic reflectors in, 39.11, 39.11f
 Fresnel lenses in, 39.9–39.10, 39.10f
 involute reflectors in, 39.11–39.12, 39.12f
 macrofocal reflectors in, 39.11
 software modeling of, 39.6–39.8
 spherical lenses in, 39.8, 39.9f
 terminology, 39.2, 39.2t
 uniform illumination of, 39.22–39.41
 classic projection system uniformity, 39.23–39.24, 39.23f
 faceted structures, 39.39–39.41, 39.39f, 39.40f
 integrating cavities, 39.24–39.27, 39.24f, 39.25f, 39.27f
 lens arrays, 39.32–39.37, 39.33f–39.37f
 lightpipes, 39.27–39.32, 39.28f–39.31f
 tailored reflectors, 39.37–39.39, 39.38f
- Nonimaging software modeling, 39.6–39.8
- Noninterferometric optical testing, 13.1–13.7
 Foucault test, 13.2–13.3, 13.2f, 13.3f
 Hartmann test, 13.4–13.6, 13.5f
 Hartmann-Shack test, 13.6–13.7, 13.6f
 Ronchi test, 13.3–13.4, 13.3f, 13.4f
- Nonlinearity correction factor, 34.33
- Nonlinearity measurement, 34.34–34.35
- Nonneutral (color) density, of photographic films, 29.7–29.8
- Nonradiative recombination, 17.2
- Nonreacting interferometers, 12.7
- Nonsequential ray tracing, 39.6
- Nonsequential surfaces data, 3.6–3.7
- Nonuniform rational B-splines (NURBS), 39.6, 40.41

- Nonuniformity, of radiation distribution, 34.25, 34.28, 34.35
- Normal equations, 3.18
- Normalized detectivity, 38.9
- Normalized detector irradiance (NDI), 7.22
- Notch filters, PID controller vs., 22.10–22.11, 22.10f, 22.11f
- np* silicon photodiodes, 34.30
- Nuclear light sources, 40.39
- Numeric displays, LED, 17.30–17.31, 17.30f, 17.31f
- Numerical aperture (NA), 34.20, 39.1
- Numerically controlled machines (CNCs), 12.11
- Nutting's law, 29.6
- Nyquist condition, 13.27
- Nyquist noise (*see* Thermal noise)
- Object counting, reflexive sensors for, 17.34
- Objective tone reproduction, 29.16–29.17, 29.17f
- Octocouplers, 17.32–17.34, 17.32f
- Off-axis angles, 7.23
- Off-axis chromatic aberrations, 2.2–2.4, 2.3f, 2.4f
- Off-axis irradiance, 34.16, 34.16f
- Off-axis rejection (OAR), 7.23
- Office lighting, 40.55, 40.56t
- Offner compensators, 13.24, 13.24f
- Offset, thermocouple junctions as source of, 27.6, 27.6f
- Offset drift, 27.11
- Offset subtraction error, 34.32–34.33
- Ohmic contact, 17.13
- 1/*f* noise, 25.12, 27.4
- 110 photographic film, 30.21, 30.25
- Open arcs, 15.22
- Open-loop gain function, 22.9–22.10, 22.9f
- Open-tube diffusion, 17.24
- Optical analysis software, 40.20
- Optical athermalization, 8.12–8.15, 8.13t–8.15t
- Optical cavities, 19.18, 19.19f
- Optical cavity technique, 12.20, 12.20f
- Optical cavity-based frequency discriminators, 22.14–22.16, 22.17f
- Optical choppers, 27.14
- Optical communication systems, 19.3
- Optical components:
 purchasing of, 9.9
 specifications for systems vs., 4.3
- Optical confinement factor, 19.11
- Optical density, of photographic films, 29.6–29.8, 29.7f
- Optical design software, 3.1–3.24, 40.20–40.21
 about, 3.2
 and computing environment, 3.21
 data entry for, 3.2–3.8
 design process flowchart, 3.3f
 evaluation function of, 3.8–3.16
 aberrations, 3.9–3.11
 paraxial analysis, 3.8–3.9, 3.9f
 ray tracing, 3.11–3.13, 3.12f
 spot-diagram analysis, 3.13–3.16
 global optimization with, 3.21
 optimization function of, 3.16–3.21
 programming considerations for, 3.7–3.8
 purchasing of, 3.22–3.24
 setup routine in, 3.7
 simulation with, 3.21
- Optical image transformers, 39.21
- Optical multichannel analyzers (OMAs), 31.27–31.28
- Optical parametric oscillators (OPOs), 20.20–20.22, 20.20f, 20.21f
- Optical path difference (OPD), 2.1, 2.6, 3.12–3.13, 3.12f, 8.1, 8.7, 13.14–13.15
- Optical power (*see* Radiant flux (power))
- Optical pulse(s), 20.2–20.15
 coupling of circulating, 20.12–20.15, 20.12f, 20.15f
 in high gain oscillators, 20.10–20.12, 20.11f
 in ideal cavity, 20.6–20.7
 and pulse train, 20.2–20.9
 single, 20.2–20.3, 20.3f
 toward steady-state, 20.9–20.12
- Optical pumping, 16.16–16.19, 16.16f–16.18f
- Optical radar, 12.4, 12.5
- Optical software (for stray light suppression), 7.24–7.27
 advantages/disadvantages of, 7.29, 7.29f
 ASAP, 7.25
 CODE V, 7.26
 FRED, 7.25–7.26
 LightTools, 7.26
 SPEOS, 7.27
 TracePro, 7.27
 ZEMAX, 7.26–7.27
- Optical specifications, 4.1–4.12
 about, 4.1–4.2
 element description, 4.8–4.10
 environmental, 4.10

- Optical specifications (*Cont.*):
 image, 4.3, 4.6–4.8, 4.8f
 mechanical vs., 4.2
 preparing, 4.5–4.6, 4.6t
 presentation of, 4.10–4.11
 problems with writing, 4.11–4.12
 for systems vs. components, 4.3
 wavefront, 4.3–4.5, 4.5t
- Optical tolerances (*see* Tolerances)
- Optical transfer function (OTF), 29.17
- Optimization function (of optical software),
 3.16–3.21
 by damped least-squares method, 3.17–3.19
 and error functions, 3.19–3.20
 global, 3.21
 multiconfiguration, 3.20
 by orthonormalization, 3.19
 by simulated annealing, 3.19
 and tolerancing, 3.20–3.21
- OPTIS (simulation software), 7.27
- Optoelectronic integrated circuit (OEIC) chip,
 25.15
- Optronic Laboratories, Incorporated, 15.49
- Ordinary rays, 3.12
- Organic dye lasers, 16.31–16.32, 16.32f
- Organic LEDs (OLEDs), 40.37–40.39
- Organometallic vapor phase deposition
 (OMVPE), 19.6–19.7, 19.20t, 19.23
- Orthonormal polynomials, 11.3–11.40
 and aberration balancing, 11.30, 11.35–
 11.36, 11.36t
 about, 11.4–11.5
 circle, for noncircular pupils, 11.37, 11.39
 defined, 11.5–11.6
 discussion of, 11.39–11.40
 elliptical, 11.21, 11.25–11.27, 11.26t–11.27t
 hexagonal, 11.21, 11.22t–11.25t
 isometric, interferometric, and PSF plots for
 orthonormal aberrations, 11.36–11.37,
 11.37f, 11.38f
 rectangular, 11.27–11.28, 11.28t, 11.29t
 slit, 11.30, 11.35t
 square, 11.30, 11.31t–11.34t
 Zernike annular, 11.13–11.21, 11.14f, 11.16f,
 11.17t–11.21t
 Zernike circle, 11.6–11.12, 11.8t–11.9t,
 11.9f–11.11f, 11.12t
- Orthonormalization, 3.19
- Oscillation, relaxation, 16.12, 19.31–19.34,
 19.31f
- Oscillators:
 high-gain, 20.10–20.12, 20.11f
 optical-parametric, 20.20–20.22,
 20.20f, 20.21f
- Out-of-band radiation errors, 34.36
- Output amplifier noise, in CCDs, 32.20
- Output circuits, direct readout architectures,
 33.18
- Output coupling mirror, 16.11
- Output gate (OG), 32.14
- Output windows, in proximity-focused MCP
 IIs, 31.9, 31.9f
- Overhead lighting, 40.12, 40.13f, 40.14, 40.46f
- Overloosening and overtightening (tolerancing
 problems), 5.11
- Oxide stripe lasers, 19.8, 19.9f
- Packages and packaging:
 HB-LED, 18.5–18.6, 18.5f, 18.6f
 of photodetectors, 26.9–26.10, 26.10f, 26.11f
 reliability of LED, 17.25–17.26
- Paint modeling, 40.17
- Paralyzing glare, 40.9
- Parametric downconversion, 23.14
- Paraxial ray tracing, 3.5, 3.8–3.9, 3.9f
- Paraxial rays, 3.3
- PART (stray light analysis program), 7.11
- Part active, part passive athermalization,
 8.11–8.12, 8.12f
- Particle hypothesis, Einstein's, 23.7
- Particle pumping, 16.14–16.16, 16.15f, 16.16f
- Parts per million (ppm), 17.25
- Passivation, in wafer processing, 17.23
- Passive athermalization, 6.22, 6.23f, 6.24,
 8.8–8.10, 8.8f–8.10f
- Pattern effect, of light pulses, 19.32, 19.32f
- Pattern noise, 32.12
- Penalty function method, 3.18, 3.19
- Pen-Ray, 15.36f
- Pentaprisms, 6.14f, 6.15f, 12.12, 12.12f
- Perception, of lit environment, 40.1–40.2,
 40.4–40.6, 40.4f
- Perceptual constancy, 40.5
- Perfect lens, 4.4
- Perfect reflecting diffusers, 37.8
- Perfect transmitting diffusers, 37.8
- Periscopes, 1.10, 1.10f
- Petzval surface, 2.4, 2.5
- Phase conjugated coupling, 20.14
- Phase diffusion coefficient, 23.29

- Phase errors, phase-shifting interferometry and, 13.22
- Phase modulation index, 22.4
- Phase response, frequency vs., 22.6–22.7, 22.6f, 22.7f
- Phase space, 39.3
- Phase stability margin, 22.9
- Phase-interruption broadening, emission-line, 16.5
- Phase-lock phase shifting, 13.23
- Phase-locked laser arrays, 19.26, 19.27, 19.28f, 19.28t
- Phase-matching, in attosecond optics, 21.4
- Phase-shifting interferometry, 13.18–13.23, 13.18f–13.20f
- heterodyne interferometer, 13.22
 - integrating bucket method, 13.21, 13.21f
 - phase errors, 13.22
 - phase-lock method, 13.23, 13.23f
 - phase-stepping method, 13.20, 13.20f
 - simultaneous measurement, 13.22
 - two-steps-plus-one method, 13.21, 13.22
- Phase-stepping phase shifting, 13.20, 13.20f
- Phonon broadening, emission-line, 16.5
- Phonon coupling, 17.16
- Phosphor salts, 40.31
- Phosphor screens:
 - of image intensifiers, 31.14–31.16, 31.14t, 31.15f
 - in proximity-focused MCP IIs, 31.9, 31.9f
- Phosphor-type designation system, 31.14t
- Phot (unit), 34.43, 36.7, 36.7t
- Photo cell (*see* Photodiodes)
- Photo excitation, 25.5
- Photo gain, in photoconductors, 25.5–25.6, 25.5f
- Photocapacitors, MOS (*see* Metal-oxide-semiconductor capacitors)
- Photocathodes, 27.6, 27.7f, 31.9, 31.9f
- assemblies of, 31.10–31.12, 31.11f, 31.12f
 - internally and remotely processed, 31.10, 31.24
- Photochemical blue-light hazard, 36.17
- Photoconductive (PC) arrays, 33.4
- Photoconductive gain, 24.11
- Photoconductors, 24.7–24.8, 24.7f
- electronics of, 38.10
 - fabrication of, 32.2
 - image sensing with, 32.8–32.9
 - operating principles of, 25.4–25.6, 25.4f, 25.5f
- Photoconductors (*Cont.*):
- photo gain in, 25.5–25.6, 25.5f
 - types of, 25.5, 25.5f
 - (*See also specific photoconductors, e.g.:*
 - Amorphous silicon photoconductors)
- Photodetection, 25.1–25.17
- about, 25.2–25.3
 - applications for, 25.11–25.12
 - future directions in, 25.15–25.17, 25.15f–25.17f
 - materials for, 25.3t
 - operating principles of, 25.3–25.11, 25.3f, 25.4f
 - extended wavelength photodetectors, 25.10, 25.10t
 - photoconductors, 25.4–25.6, 25.4f, 25.5f
 - photogate, 25.10, 25.11, 25.11f
 - pin* photodiodes, 25.6–25.10, 25.7f, 25.9f
 - reliability of, 25.13–25.14, 25.13t
- Photodetectors, 24.3–24.101
- AlGaN alloy photovoltaic detectors, 24.46
 - CdS photoconductors, 24.49–24.52, 24.51f–24.53f
 - CdSe photoconductors, 24.49–24.52, 24.52f
 - CdTe detectors, 24.52, 24.54, 24.54f
 - CdZnTe detectors, 24.52
 - GaAsP photodiodes, 24.49, 24.49f, 24.50f
 - GaN photovoltaic detectors, 24.42, 24.43, 24.45f, 24.46, 24.46f, 24.47
 - GaP photodiodes, 24.47–24.49, 24.48f
 - Ge intrinsic photodetectors, 24.70–24.73, 24.70f–24.73f
 - Ge low-temperature bolometers, 24.31–24.32, 24.32f, 24.33f
 - Ge:Au detectors, 24.83–24.85, 24.84f–24.86f
 - Ge:Cu detectors, 24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f
 - Ge:Ga infrared detectors, 24.100
 - Ge:Hg detectors, 24.84f, 24.92–24.95, 24.93f–24.95f
 - Ge:Zn detectors, 24.84f, 24.98–24.100, 24.99f
 - HgCdTe detectors, 24.86–24.92, 24.86f–24.92f
 - InAs photovoltaic detectors, 24.75, 24.77–24.78, 24.77f–24.79f
 - InGaAs detectors, 24.65–24.70, 24.66f–24.69f
 - InGaAs photodiodes, 34.31
 - InSb hot-electron bolometers, 24.29, 24.30, 24.30f, 24.31f
 - InSb intrinsic photovoltaic detectors, 24.80–24.83, 24.82f, 24.83f
 - life test for, 25.13, 25.14, 25.14f, 25.15f

- Photodetectors (*Cont.*):
 manufacturers' specifications for,
 24.21–24.101, 24.21*t*
 noise, impedance, dark and leakage current
 in, 24.19–24.20, 24.20*f*
 PbS photoconductors, 24.73–24.74,
 24.74*f*–24.77*f*, 24.74*t*
 PbSe detectors, 24.76*f*, 24.78, 24.79,
 24.79*f*–24.82*f*
 PbSnTe photovoltaic detectors, 24.92, 24.93*f*
 performance/sensitivity of, 24.13–24.18
 background-limited case, 24.14–24.17,
 24.16*f*, 24.16*t*, 24.17*f*
 strong-signal case, 24.14
 thermal detectors, 24.17, 24.18*f*
 photoemissive detectors, 24.32–24.42, 24.33*f*,
 24.35*f*–24.41*f*, 24.43*f*, 24.44*f*
 photographic emulsions, 24.100, 24.101*f*
 planar, 25.14, 25.15*f*
 pyroelectric detectors, 24.26–24.29,
 24.26*f*–24.29*f*
 quantum, 24.6–24.10, 24.7*f*–24.9*f*
 quantum well, 25.16–25.17, 25.16*f*, 25.17*f*
 quantum well infrared, 25.15
 responsivity and quantum efficiency of,
 24.18, 24.19
 Si photovoltaic detectors, 24.52*f*, 24.54–
 24.65, 24.55*f*–24.66*f*
 Si:B detectors, 24.95*f*, 24.96
 SiC UV detectors, 24.47, 24.47*f*
 Si:Ga infrared detectors, 24.95, 24.95*f*,
 24.96, 24.96*f*
 and signal detection, 27.2
 spectral response of, 24.18, 24.19*f*
 speed of, 24.20, 24.21
 stability of, 24.21, 24.21*f*
 terminology, 24.10–24.13
 thermal detectors, 24.4–24.6, 24.4*f*–24.6*f*
 thermistor bolometers, 24.24–24.25, 24.24*f*,
 24.25*f*, 28.7*t*
 thermocouples, 24.22–24.23, 24.22*f*
 thermopiles, 24.23–24.24, 24.23*f*
 TiO₂ UV detectors, 24.47, 24.48*f*
 types of, 25.3, 25.4*f*
 uniformity of, 24.20
 (*See also specific photodetectors, e.g.:*
 Avalanche photodetectors)
- Photodiode CCDs (PD-CCDs), 33.11*f*, 33.12
 Photodiode linear arrays (PDAs), 38.9,
 38.10, 38.10*t*
- Photodiode MOSs (PD-MOSs), 33.11*f*, 33.13
 Photodiodes (PDs):
 CCD, 33.11*f*, 33.12
 defined, 24.11
 electronics of, 38.10
 GaAsP, 24.49, 24.49*f*, 24.50*f*
 GaP, 24.47–24.49, 24.48*f*
 Ge avalanche, 24.70*f*, 24.72–24.73, 24.72*f*,
 24.73*f*
 GeGaAs, 34.31
 InGaAs, 34.31
 InGaAs avalanche, 24.66–24.70, 24.66*f*–24.69*f*
 junction, 32.3–32.6, 32.4*f*, 32.6*f*
 MOS, 33.11*f*, 33.13
pin (*see pin photodiodes*)
p⁺np, 32.8
 silicon, 34.29, 34.30
 silicon avalanche, 24.62–24.65, 24.63*f*–24.66*f*
 silicon *pn*, 24.52*f*, 24.55–24.58, 24.55*f*–24.59*f*
 UV-enhanced, 24.55*f*, 24.61–24.62, 24.61*f*,
 24.62*f*
 (*See also specific photodiodes, e.g.:* Avalanche
 photodiodes)
- Photoelectromagnetic (PEM) detectors,
 24.9, 24.9*f*
- Photoemissive detectors, 24.6, 24.7*f*
 gallium phosphide dynodes, 24.42, 24.44*f*
 linearity of, 24.41
 manufacturers of, 24.42
 noise from, 24.39, 24.40, 24.41*f*
 operating temperature of, 24.40
 photon counting for, 24.42, 24.43*f*, 24.44*f*
 quantum efficiency of, 24.35*f*–24.38*f*,
 24.36–24.38
 recommended circuit for, 24.42, 24.43*f*
 response time of, 24.40
 responsivity of, 24.35*f*, 24.38
 sensitive area of, 24.41
 sensitivity of, 24.34, 24.35*f*–24.39*f*
 sensitivity profile of, 24.41
 short-wavelength considerations for,
 24.34, 24.40*f*
 specifications for, 24.32–24.42, 24.33*f*
 stability of, 24.41
- Photogates, 25.10, 25.11, 25.11*f*
 Photographic detectors, 24.9–24.10
 Photographic dyes, 30.10–30.13,
 30.10*f*, 30.12*f*
 Photographic emulsions, 24.100, 24.101*f*,
 29.4, 30.7

- Photographic film speed, 29.9–29.10
 in Advanced Photo System, 30.26
 in color negative film, 30.25
 and granularity, 30.19
 high vs. low, 30.18–30.20
 and sensitivity to high-energy radiation,
 30.19–30.20
- Photographic films, 29.3–29.16, 30.18–30.28
 about, 30.2–30.3
 black-and-white (B&W) film, 30.24–30.25,
 30.25*t*
 color, 29.12–29.15, 29.13*f*, 29.14*f*
 color negative film, 30.25–30.28, 30.27*t*
 color reversal film, 30.22–30.24, 30.23*t*
 development effects on, 29.12, 29.13*f*
 D-log H curve for, 29.8–29.10, 29.8*f*, 29.10*f*
 exposure of, 29.5–29.6
 grain element of, 29.5
 granularity of, 30.19
 high-speed vs. low-speed, 30.18–30.20
 and light scattering by silver halide crystals,
 30.5–30.7, 30.6*f*
 microdensitometers, 29.15–29.16, 29.15*f*
 optical density of, 29.6–29.8, 29.7*f*
 processing of, 29.5
 professional vs. amateur film, 30.20–30.22
 reciprocity failure of, 29.11–29.12
 spectral sensitivity of, 29.11, 29.11*f*
 speed of, 30.18–30.20
 structure of color, 30.3–30.5, 30.3*f*
 structure of silver halide photographic layers
 in, 29.4
- Photographic materials, 30.1–30.28
 about, 30.1–30.2
 dyes, 30.10–30.13
 about, 30.10, 30.10*f*
 excited state properties, 30.11–30.12, 30.12*f*
 light stabilization methods, 30.12–30.13
 photochemistry of azomethine dyes, 30.11
 films, 30.18–30.28
 black-and-white film, 30.24–30.25, 30.25*t*
 color negative film, 30.25–30.28, 30.27*t*
 color reversal film, 30.22–30.24, 30.23*t*
 professional vs. amateur film,
 30.20–30.22
 speed, 30.18–30.20
 optics of, 30.2–30.7
 about, 30.2–30.3
 light scatter by silver halide crystals,
 30.5–30.7, 30.6*f*
- Photographic materials, optics of (*Cont.*):
 structure of color films, 30.3–30.5, 30.3*f*
 structure of color papers, 30.5
 and photographic spectral sensitizers,
 30.13–30.18
 about, 30.13–30.14, 30.14*f*
 color science, 30.15–30.18, 30.16*f*, 30.17*f*
 photophysics of spectral sensitizers on silver
 halide surfaces, 30.14–30.15, 30.15*f*
 silver halide light detectors, 30.7–30.9, 30.8*f*
- Photographic papers:
 about, 30.2–30.3
 and light scattering by silver halide crystals,
 30.5–30.7, 30.6*f*
 structure of color, 30.5
- Photographic spectral sensitizers (*see* Spectral
 sensitizers, photographic)
- Photographic systems, 29.16–29.25
 acutance of, 29.17–29.19, 29.18*t*, 29.19*f*
 detective quantum efficiency of, 29.23
 graininess in, 29.19–29.22, 29.21*f*
 image structure of, 29.17
 information capacity of, 29.24
 manufacturers of, 29.25
 performance of, 29.16–29.17, 29.17*f*, 29.18*t*
 resolving power of, 29.24
 sharpness in, 29.22
 signal-to-noise ratio of, 29.22–29.23
- Photoionization detectors, 24.10
- Photoionization devices, 34.29
- Photoionization yield, 34.29
- Photometry, 34.37–34.44, 36.1–36.17, 36.3*f*,
 36.3*t*
 about, 36.2–36.4
 approximations, 36.10, 36.10*f*, 36.11*f*
 basis of physical, 37.1–37.2, 37.2*f*
 calibrations in, 34.42–34.43
 concepts/terminology of, 34.10–34.11,
 34.38*t*, 34.39–34.40, 34.43–34.44, 39.2,
 39.2*t*
 conversion between radiometric and
 photometric quantities, 34.12*t*,
 36.11–36.14, 36.12*f*–36.14*f*
 defined, 34.37, 37.1
 and human eye, 36.8–36.10, 36.8*f*, 36.9*f*
 illuminance-luminance relationship, 37.9,
 37.9*f*
 integrating sphere device, 37.9–37.10, 37.10*f*
 inverse square law, 37.8
 Lambert's cosine law, 37.8, 37.8*f*

- Photometry (*Cont.*):
 normalization, 36.14–36.17, 36.15*t*, 36.16*f*
 and photopic/scotopic/mesopic vision,
 34.37–34.39, 34.38*t*
 practice in, 37.11
 quantities and units in, 37.4–37.8, 37.4*t*,
 37.5*f*, 37.7*t*
 radiometry vs., 34.6
 retinal illuminance, 34.40–34.42
 symbols/nomenclature of, 36.5–36.10,
 36.6*t*–36.8*t*
 weighting functions, 36.17
- Photomicrographic lamps, 15.47–15.49,
 15.48*f*, 15.49*f*
- Photomultiplier tubes (PMTs), 24.11,
 24.32–24.34, 24.33*f*, 24.38–24.42, 24.38*f*,
 24.39*f*, 24.43*f*, 38.9, 38.9*t*
 applications for, 27.6–27.10, 27.7*f*
 base design of, 27.8–27.10
 electronics of, 38.10
- Photon(s), 23.6–23.14, 25.2, 34.30, 36.4
 Einstein's light quanta, 23.6–23.9, 23.8*f*
 photon-photon correlations, 23.13–23.14,
 23.13*f*
 quantum electrodynamics, 23.9–23.13
- Photon counting, 27.15
 defined, 24.12
 of modulated signal sources, 27.14
 in photoemissive detectors, 24.42, 24.43*f*,
 24.44*f*
- Photon density, 19.30–19.35
- Photon detectors:
 background-limited case of, 24.14–24.17,
 24.16*f*, 24.16*t*, 24.17*f*
 strong-signal case of, 24.14
- Photon dose, 34.6, 34.11
- Photon Engineering, 7.25
- Photon flux, 34.5, 34.11
- Photon infrared detectors, 33.7, 33.8*f*
- Photonic excitation, 25.3, 25.3*f*
- Photonics Directory of Optical Industries*, 15.14
- Photopic vision, 34.37–34.39, 34.38*t*,
 36.8*f*–36.10*f*, 37.2, 40.3
- Photovoltaic (PV) arrays, 33.4
- Photovoltaic detectors, 24.8, 24.8*f*, 24.9
 aluminum gallium nitride, 24.46
 gallium nitride, 24.42, 24.43, 24.45*f*, 24.46,
 24.46*f*, 24.47
 indium antimonide, 24.80–24.83, 24.82*f*,
 24.83*f*
- Photovoltaic detectors (*Cont.*):
 indium arsenide, 24.75, 24.77–24.78,
 24.77*f*–24.79*f*
 lead tin telluride, 24.92, 24.93*f*
 mercury cadmium telluride, 24.86*f*, 24.88*f*,
 24.90–24.92, 24.91*f*, 24.92*f*
 silicon, 24.52*f*, 24.54–24.65, 24.55*f*–24.66*f*
- Photovoltaic Schottky barrier detectors (SBDs),
 33.7, 33.8*f*
- Physical optics, 3.16
- Physical photometry, 34.37, 36.4, 37.2
- Pickups, 3.6
- Piezoelectric transducers (PZTs):
 amplifier strategies for, 22.18
 disk vs. tube, 22.17–22.18
- Piezoelectric-based (PZT-based) systems, 22.1
- pin* junctions, 26.3
- pin* photodetectors:
 biased, 26.6*f*
 high-speed, 26.10, 26.12–26.15, 26.12*f*,
 26.14*f*, 26.15*f*
 lateral, 25.15–25.16, 25.15*f*
- pin* photodiodes, 24.54, 24.55*f*–24.60*f*,
 24.58–24.61, 25.4*f*, 25.6–25.10, 25.7*f*, 32.4*f*
 absorption coefficient of, 25.8, 25.9*f*
 avalanche photodiodes, 25.8–25.10, 25.9*f*
 dark current in, 25.7–25.8
 diffusion current of, 25.7
 equivalent circuit of, 26.7, 26.7*f*
 generation-recombination current of,
 25.7–25.8
 germanium, 24.70–24.71, 24.70*f*–24.72*f*
 InGaAs, 24.66, 24.67*f*
 operating principles of, 25.6–25.10,
 25.7*f*, 25.9*f*
 quantum efficiency of, 25.8
 resonant, 26.15, 26.15*f*
 responsivity of, 25.8
 silicon, 24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*,
 24.60*f*
 tunneling current of, 25.8
 vertically illuminated, 26.3, 26.4*f*, 26.5,
 26.10, 26.12–26.13, 26.12*f*
 waveguide, 26.13–26.14, 26.14*f*
- Pitch-based edging (fabrication step), 9.6
- Pixel summation, 38.10
- Planar buried heterostructure (PBH) lasers,
 double-channel, 19.24, 19.25*f*, 19.34*f*
- Planar diffused silicon photodiodes, 24.56*f*
- Planar photodetectors, 25.14, 25.15*f*

- Planckian radiation (*see* Blackbody radiation)
- Planck's formula, 23.6, 23.7
- Planck's law, 34.23, 34.24, 37.10–37.11
- Plane diffusers, 38.14, 38.14f, 38.15f
- Plano optics fabrication, 9.7
- Plastic, as photographic film emulsion, 29.4
- Plastic-packaged LEDs, 17.25–17.26
- Platings, diamond turning of, 10.5
- Platinum silicon (PtSi) infrared detectors, 33.7, 33.8f, 33.29, 33.29t
- Pluecker spectrum tubes, 15.47, 15.47f, 15.47t
- pn* junctions, 24.8, 24.8f
- pn* photodetectors, 24.70–24.71, 24.70f–24.72f
- pn* photodiodes, 24.52f, 24.54–24.58, 24.55f–24.59f, 34.30
- p⁺np* photodiodes, 32.8
- Point diffraction interferometers, 13.11f
- Point source irradiance transmittance (PSIT), 7.23
- Point source normalized irradiance transmittance (PSNIT), 7.22–7.23
- Point source power transmittance (PSPT), 7.23
- Point source transmittance (PST), 7.5, 7.6f, 7.22–7.23
- Point spread function (PSF), 3.16, 11.36–11.37, 11.37f, 11.38f
- Point-to-point approximation (of radiant flux transfer), 34.14
- “Poker chip” assembly, 6.8, 6.9f
- Polar angles, 35.5
- Polarization gating, 21.7–21.8
- Polarization-dependent systems, 34.33
- Polaroid Corporation, 29.12, 29.25
- Polaroid Instant Color Film, 29.14
- Polaroid “One Film,” 29.14
- Pole-mounted luminaires, 40.62, 40.63t
- Polishers, continuous, 9.7
- Polishing step (of optics fabrication), 9.5–9.6, 9.8
- Polychromatic radiation, 34.9–34.10
- Polyethylene terephthalate film, 29.4
- Polynomial numbering, 11.39
- Polynomial-ordering number, 11.7
- Polynomials, orthonormal (*see* Orthonormal polynomials)
- Polytetrafluoroethylene (PTFE), 35.13, 38.12–38.13
- Ponderomotive potential, 21.3
- Population inversions, 16.8–16.10, 16.12–16.13
described, 16.8–16.10
mechanism for achieving, 16.12–16.13, 16.13f, 16.14f
optical pumping for, 16.16–16.19, 16.16f–16.18f
particle pumping for, 16.14–16.16, 16.15f, 16.16f
semiconductor diode laser pumping for, 16.19
- Positive image (in photography), 29.9
- Potassium dihydrogen phosphate (KDP), 10.2
- Power measurement, for lasers, 34.32
- Power spectrum of granularity, 29.21
- Power supply, high-voltage, 31.9, 31.10f
- Predictable quantum efficiency (PQE) devices, 34.29–34.30
- Preloads, 6.2
- Primaries (colors stimuli), 30.16
- Principal rays, 1.4, 1.11f, 1.12
- Prisms, 40.45f, 40.46
angle measurement in, 12.14–12.16, 12.15f–12.17f
mounting of, 6.11–6.17
bonded mountings, 6.13–6.15, 6.15f, 6.16f
flexure mountings, 6.15–6.17, 6.16f
mechanically clamped mountings, 6.12, 6.13f
spring-loaded mountings, 6.13, 6.14f
penta, 6.14f, 6.15f, 12.12, 12.12f
right-angle, 12.15, 12.16, 12.16f, 12.17f
Risley, 12.4, 12.4f
rotating glass block, 12.4, 12.4f
sliding, 12.4, 12.4f
Zerodur, 6.16, 6.16f
- Procedural programs (optical software), 3.7
- Projected area, in photometry/radiometry, 36.3–36.4, 36.3f, 36.3t
- Projected solid angle (PSA), 39.5, 39.5f
- Projection density, 29.7, 29.7f
- Projection lenses, 6.6f
- Projection systems, 39.23–39.24, 39.23f
- Proportional integral derivative (PID) controllers, 22.10–22.12, 22.10f, 22.11f
- Proportional-integral (PI) amplifier circuit, 22.5f
- Proton stripe lasers, 19.35f
- Proton-bombardment-defined lasers, 19.41
- Proximity-focus electronic lens, 31.8, 31.23–31.24, 31.23f, 31.24f

- Proximity-focused MCP IIs, **31.9, 31.9f, 31.16–31.18, 31.17t, 31.18f, 31.19f**
- Psychophysical photometry, **34.37**
- P*-type impurities, **17.23**
- Pulse height, of photomultipliers, **27.7–27.8**
- Pulse train interferometry, **20.12, 20.12f**
- Pulse trains:
 about, **20.3–20.5, 20.4f, 20.5f**
 attosecond, **21.6, 21.7**
 and backscattering, **20.13–20.15, 20.15f**
 soliton solution and steady-state, **20.5–20.9**
- Pulsed lasers, **23.18**
- Pumping (for population inversions), **16.14–16.19**
 optical, **16.16–16.19, 16.16f–16.18f**
 particle, **16.14–16.16, 16.15f, 16.16f**
 semiconductor diode laser, **16.19**
- Punch through (in color films), **30.4**
- Pupils, noncircular, **11.4**
- Push processing (of film), **30.22**
- Pyramidal error, **12.14, 12.15, 12.15f**
- Pyroelectric detectors, **24.6, 24.6f, 24.26–24.29, 24.26f–24.29f, 28.2, 28.6, 28.6f, 28.7, 28.7t, 33.10**
- Pyroelectric electrical substitution radiometers, **34.27–34.28**
- Pyroelectric hybrid arrays, **28.11–28.12, 28.11f, 28.12f**
- Q*-switched lasers, **16.26–16.27, 16.27f**
- Quality, image, **4.6–4.7**
- Quality Assurance of Ultraviolet Measurements in Europe (QASUME), **38.5**
- Quantized center-of-mass motion, of atoms, **23.45**
- Quantum box, **19.18**
- Quantum cascade lasers, **16.36**
- Quantum dot, **16.7, 19.18, 26.4f, 26.5**
- Quantum efficiency (QE), **25.3, 25.4, 34.29**
 defined, **24.12**
 detective, **29.23**
 of photodetectors, **24.18, 24.19**
 of photoemissive detectors, **24.35f–24.38f, 24.36–24.38**
 of photomultipliers, **27.7**
 of *pin* photodiodes, **25.8**
- Quantum electrodynamics (QED), **9.6, 23.9–23.13**
- Quantum limited imaging (QLI), **31.3–31.4**
- Quantum photodetectors, **24.6–24.10, 24.7f–24.9f**
- Quantum resonance absorption, **22.16, 22.17**
- Quantum sensitivity, **30.9**
- Quantum theory of lasers, **23.14–23.35**
 about, **23.5–23.6**
 density-operator approach to, **23.14–23.33**
 derivation of Scully-Lamb master equation, **23.17–23.22, 23.19f**
 photon statistics, **23.22–23.27, 23.23f, 23.25f, 23.27f**
 spectral properties, **23.28–23.33, 23.30f**
 spectrum, **23.28–23.33**
 time evolution of the field in Jaynes-Cummings model, **23.15–23.17, 23.15f**
 Heisenberg-Langevin approach to, **23.33–23.35**
- Quantum trajectories, in attosecond optics, **21.3–21.4**
- Quantum well (QW) detectors, **26.4f, 26.5**
- Quantum well infrared photodetectors (QWIPs), **25.15–25.17, 25.16f, 25.17f, 33.9**
- Quantum well (QW) lasers, **16.7, 19.9–19.18, 19.20t**
 GRIN SCH single, **19.14, 19.14f**
 long wavelength, **19.17–19.18, 19.17f**
 schematic of, **19.10f**
 strained, **19.15–19.17, 19.16f**
 threshold modal gain, **19.12–19.15, 19.12f, 19.13f, 19.15f**
- Quantum well (QW) photodetectors, **25.16–25.17, 25.16f, 25.17f**
- Quantum wire, **16.7, 19.18, 26.4f, 26.5**
- Quartz-envelope lamps, **15.20, 15.21**
- QW ridge (QWR) waveguide lasers, **19.19, 19.20t**
- Rabi cycles and Rabi cycling, **23.21**
 defined, **20.24**
 off resonance, **20.27**
 on resonance, **20.26–20.27, 20.26f, 20.27f**
- Rack-stack laser arrays, **19.29, 19.30f**
- Radial circle polynomial, **11.7, 11.9f–11.10f**
- Radial shearing interferometers, **13.12, 13.13f**
- Radian (rad), **36.3**
- Radiance, **34.9, 34.9f, 37.4t, 37.5, 38.2, 38.11t, 38.13–38.16, 38.13f–38.16f, 39.2t**
- Radiance conservation theorem, **34.12–34.13**
- Radiance temperature (unit), **37.4t, 37.6**
- Radiance units, **34.24**

- Radiant energy, 34.7, 37.4t, 37.6
- Radiant exitance (emittance), 15.4–15.6, 15.5t, 15.6f, 35.3, 37.4t, 37.5, 37.5f
- Radiant exposure, 37.4t, 37.6
- Radiant flux (power), 34.7, 34.11, 34.17–34.18, 36.4, 36.6t, 37.3, 37.4t, 37.6, 38.2
- Radiant incidence (*see* Irradiance)
- Radiant intensity, 36.4, 37.4, 37.4t, 39.2t
- Radiant power transfer, 34.12–34.13, 34.13f
- Radiant transfer approximations, 34.13–34.20
- approximate radiance at an image, 34.19–34.20
 - lambertian, 34.14–34.18, 34.15f, 34.16f
 - point-to-point, 34.14
 - radiometric effect of stops and vignetting, 34.18–34.19, 34.19f
- Radiation, 34.23–34.27
- actinic effects of, 34.6, 34.7
 - artificial sources of (*see* Artificial sources (of radiation))
 - baseline standard of, 15.9, 15.9f, 15.10f, 15.12f
 - from blackbodies, 34.23–34.24
 - from blackbody simulators, 34.24–34.26
 - between circular source and detector, 34.15–34.16, 34.15f
 - commercial sources of (*see* Commercial sources (of radiation))
 - incandescent sources of (*see* Incandescent sources (of radiation))
 - infrared (*see* Infrared (IR) radiation)
 - and lasers, 16.4
 - photographic film speed and sensitivity to high-energy, 30.19–30.20
 - polychromatic, 34.9–34.10
 - from synchrotrons, 34.26–34.27
 - through absorbing media, 34.13
 - transfer of, 7.21–7.22
 - ultraviolet (*see* Ultraviolet (UV) radiation)
 - working standards of, 15.9–15.13, 15.10f, 15.12f, 15.13f
- Radiation law, 15.4–15.7, 15.5f, 15.5t, 15.6f
- Radiative lifetimes, 16.4, 17.4
- Radiative recombination, 17.2
- Radiators, blackbody, 34.23–34.24
- Radio frequency (rf) modulation, 22.14
- Radiometers and radiometry, 34.3–34.37, 36.1–36.17, 36.3f, 36.3t
- about, 34.5–34.7, 36.2–36.4
- Radiometers and radiometry (*Cont.*):
- approximate (*see* Radiant transfer approximations)
 - cavity-shaped, 34.28
 - concepts/terminology of, 34.7, 39.2, 39.2t
 - conversion between radiometric and photometric quantities, 34.12t, 36.11–36.14, 36.12f–36.14f
 - defined, 37.1
 - electrical substitution, 34.27–34.29
 - geometrical concepts of, 34.8–34.9, 34.9f
 - of II electronic imaging, 31.4–31.5
 - illuminance-luminance relationship, 37.9f
 - integrating sphere device, 37.9–37.10, 37.10f
 - laser as characterization tool for, 34.32
 - normalization, 36.14–36.17, 36.15t, 36.16f
 - photometry vs., 34.6
 - Planck's law, 37.10–37.11
 - quantities and units in, 37.3–37.7, 37.4t, 37.5f, 37.7t
 - spectral dependence of, 34.9–34.10
 - Stefan-Boltzmann's law, 37.11
 - symbols/units/nomenclature of, 36.4–36.5
 - thermopile-based, 34.27
 - weighting functions, 36.17
 - Wien's displacement law, 37.11
- Range finders, 12.3–12.4, 12.3f, 12.4f
- Range gating, 31.28–31.30
- Ray displacement, 3.12
- Ray intercept curves, 2.2–2.4, 3.13
- Ray sets, 3.20
- Ray tracing, 1.4–1.5, 1.4f, 3.11–3.13, 3.12f
- in lighting simulation, 40.20
 - nonsequential, 39.6
 - in optical design software, 3.11–3.13
 - paraxial, 3.5, 3.8–3.9, 3.9f
- Rays:
- axial rays, 1.4, 1.11f, 1.12
 - dashed rays, 1.12, 1.12f
 - edge rays, 39.22, 39.38
 - exact, 3.3, 3.11–3.12
 - hamiltonian rays, 3.12
 - iterated rays, 3.12
 - lagrangian rays, 3.12
 - meridional rays, 3.3
 - ordinary, 3.12
 - paraxial, 3.3
 - principal rays, 1.4, 1.11f, 1.12
- RC time constant, 26.7–26.8, 26.7f
- Reactive ion etching (RIE), 18.3, 19.39

- Readouts (of visible array detectors),
 32.12–32.21
 CCD, 32.12–32.20, 32.13f, 32.15f–32.18f
 MOS, 32.20–32.21
- Real-space critical objects, 7.2–7.4, 7.3f, 7.4f
- Reasonableness, of layout, 1.13–1.14
- Recessed lighting, 40.13f
- Reciprocity failure, of photographic films,
 29.11–29.12
- Recombination:
 combined, 17.3
 exciton, 17.6
 in GaAs, 17.8, 17.9, 17.9f
 minority-carrier, 17.2
 nonradiative, 17.2
 radiative, 17.2
- Reconstruction of attosecond beating by
 interference of two-photon transition
 (RABITT), 21.9
- Rectangular polynomials, 11.27–11.28, 11.28t,
 11.29t, 11.36t
- Red light, and color film, 29.13, 29.13f
- Reduction scanners, in linear sensors,
 32.21, 32.22f
- Reflectance:
 classification of materials by, 35.4t
 defined, 35.4–35.5
 geometrical definitions of, 35.6f
 and illuminance/luminance, 37.9
 measurement of, 35.10–35.13, 35.10f–35.12f
 nomenclature for, 35.5t, 35.6t
 spectral, 38.2, 38.17–38.18, 38.18f
 standards of, 35.14t
 and transmittance/absorptance, 35.7, 35.8,
 35.8t
- Reflecting telescopes, 11.4
- Reflection(s):
 actual/idealized, 35.2f
 defined, 35.4
 veiling, 40.12
- Reflection density, of photographic films,
 29.8
- Reflective compensators, for spherical
 aberrations, 13.24, 13.24f, 13.25
- Reflective-refractive (RX) concentrators, 39.17,
 39.17f
- Reflectometers, 35.10, 40.52
- Reflectors:
 conic, 39.11, 39.11f
 convergent, 39.38–39.40, 39.38f, 39.39f
- Reflectors (*Cont.*):
 CPC-type (*see* Compound parabolic
 collectors)
 divergent, 39.8f, 39.9f, 39.38–39.40
 faceted, 39.39–39.41, 39.39f, 39.40f
 headlamp, 40.64
 homogeneous/inhomogenous, 39.39
 involute, 39.11–39.12, 39.12f
 and lens array combinations, 39.34–39.37,
 39.36f, 39.37f
 luminaire, 40.45, 40.45f
 macrofocal, 39.11
 tailored, 39.37–39.39, 39.38f
- Reflexive sensors, 17.34
- Refraction index, 17.34
- Refractive compensators, for spherical
 aberrations, 13.24, 13.24f, 13.25
- Refractive index, 3.6, 34.13
- Relative measurements, absolute vs.,
 34.20–34.21
- Relative visual performance (RVP) model,
 40.5–40.6
- Relaxation oscillation, 16.12, 19.31–19.34,
 19.31f
- Rem jet, 30.4
- Remotely processed (RP) photocathodes,
 31.10, 31.24
- Repetition rate coupling, 20.14–20.15, 20.15f
- Rescattering model, semiclassical, 21.3
- Reset gate (RG), 32.14
- Residential lighting, 40.57, 40.58, 40.59t
- Residual amplitude modulation (RAM),
 22.14
- Resistive bolometers, 28.10–28.11, 28.10f, 33.9
- Resistive coupling noise, 27.5, 27.6f
- Resolution, of optical system, 4.6
- Resolving power, photographic, 29.24
- Resonant optical feedback, 19.38, 19.38f
- Resonant photodetectors, 26.4f
- Resonant *pin* photodiodes, 26.15, 26.15f
- Response time:
 defined, 24.12
 of photodetectors, 25.4
 of photoemissive detectors, 24.40
- Responsive quantum efficiency (*see* Quantum
 efficiency)
- Responsivity:
 blackbody, 24.10
 of photodetectors, 24.18, 24.19, 25.4
 of photoemissive detectors, 24.35f, 24.38

- Responsivity (*Cont.*):
 of *pin* photodiodes, 25.8
 spectral, 24.12, 38.3, 38.18–38.19
 of spectroradiometers, 38.11–38.12, 38.12*f*
- Restricted exit angle concentrators,
 39.18, 39.18*f*
- Retail lighting, 40.55–40.57, 40.56*t*–40.58*t*
- Reticles and reticulation, 12.13, 28.11, 28.12
- Retina, 34.37
- Retinal damage, 40.9
- Retinal illuminance, 34.40–34.42
- Retinal thermal hazard, 36.17
- Retroreflection, measurement of, 35.13
- Reverse bias, 26.3
- Reversing shear interferometers, 13.12, 13.13*f*
- Rhenium, 40.27
- Ribbon-type tungsten filaments, 15.20*f*
- Ridge waveguide (RWG) lasers, 19.8, 19.9*f*
- Right-angle prisms, 12.15, 12.16, 12.16*f*, 12.17*f*
- Ring flanges, continuous, 6.4*f*, 6.11
- Ring lasers, 16.29
 with additional Kerr crystal, 20.17–20.18,
 20.17*f*
 dye, 20.15–20.16
 Ti:sapphire, with saturable absorber,
 20.16–20.17, 20.16*f*, 20.17*f*
 in two-level system analogy, 20.24, 20.25*f*
- Rings, aperture, 3.20
- Risley prisms, 12.4, 12.4*f*
- Ritchey-Chretien two-mirror imaging system,
 39.17
- RLM lamp, 40.46, 40.46*f*, 40.47
- RMS noise, 24.12
- RMS signal, 24.12
- rms-granularity, 29.19–29.21
- Roadway lighting, 40.67, 40.69–40.71
 and disability glare, 40.10
 and discomfort glare, 40.12
 sign lighting, 40.71
 street lighting, 40.69–40.71, 40.70*t*, 40.71*t*
 tunnel lighting, 40.71
- Robertson's correlated color temperature
 calculation, 38.5
- Rods (eye receptors), 30.15, 30.16*f*, 34.37,
 36.8–36.10, 36.8*f*, 36.9*f*
- Rome Air Development Center, 7.19
- Ronchi test, 13.3–13.4, 13.3*f*, 13.4*f*
- Roof-mirror-lens arrays, 32.21, 32.22*f*
- Room temperature vulcanizing (RTV) sealing
 compound, 6.4
- Root-mean-square (rms) wavefront error, 4.1,
 4.3, 4.7, 4.8
- Rotating glass block prisms, 12.4, 12.4*f*
- Rotational shear interferometers, 13.12, 13.13*f*
- Rotationally symmetric aspheric lenses, 9.7, 9.7*f*
- Rotationally symmetric optics:
 hard mounting of, 6.2–6.4, 6.3*f*, 6.4*f*
 soft mounting of, 6.4–6.5, 6.4*f*, 6.5*f*
- Ruby lasers, 16.12, 16.13*f*, 16.32
- Rule-of-thumb PID design, 22.11–22.12
- “Rule-of-thumb” tolerance, 6.2
- Rydberg states, 23.21
- Sapphire (Ti:Al₂O₃) lasers, titanium-doped,
 16.34, 16.34*f*
- Sapphire (Ti:Al₂O₃) ring lasers, titanium-
 doped, 20.16–20.17, 20.16*f*, 20.17*f*
- Sapphire substrate (for HB-LEDs), 18.2, 18.3,
 18.5, 18.6
- Satellite spheres, 39.26
- Saturated colors, 40.7
- Saturation, 40.5, 40.9
- Sawing (of LEDs), 17.24, 17.25
- Scanning arrays, 33.6, 33.6*f*, 33.14
- Scanning FPAs, 33.17, 33.17*f*
- Scanning white light interferometry (SWLI),
 10.13
- Scattered radiation effect, 34.33
- Scattering:
 and photographic film, 29.18
 rescattering, 21.3
 by silver halide crystals, 30.5–30.7, 30.6*f*
 spectral, 38.10
 surface, 7.23
- Scattering sensors, 17.34
- Schawlow-Townes linewidth, 23.18
- Schmidt-Cassegrain design, 7.20
- Schottky barrier detectors (SBDs), 33.7, 33.8*f*,
 33.12–33.13
- Schottky contact, 26.3
- Schottky junctions, 26.3
- Schottky photodiodes, 26.16, 26.16*f*, 26.17*f*
- Scotopic vision, 34.37–34.39, 34.38*t*, 36.8*f*,
 36.9, 36.9*f*, 36.10, 36.11*f*, 37.2, 40.3
- Scully-Lamb master equation, 23.15,
 23.17–23.22
 cavity losses, 23.18
 laser master equation, 23.19–23.20, 23.19*f*
 micromaser master equation, 23.20–23.22
- Sealed beam lights, 40.26*t*

- Sealed-ampoule diffusion, 17.23
- Secondary spectrum, of radiation, 1.14, 1.15
- Second-order autocorrelator, 21.9
- Seidel (third-order monochromatic) aberrations, 3.9–3.10
- Seidel astigmatism, 11.27, 11.30, 11.35, 11.39, 11.40
- Self-athermalized, 8.7–8.8, 8.7f
- Self-calibration, of silicon photodiodes, 34.29
- Selfoc lenses, 32.21, 32.22f
- Self-phase modulation, 20.6
- Self-scanned array, 31.2
- Selwyn coefficient, 29.20
- Selwyn's law, 29.20
- Semiconductor bolometers, 28.4–28.5
- Semiconductor laser pumping, 16.19
- Semiconductor lasers, 16.35–16.36, 16.35f, 19.1–19.43
 applications for, 19.3–19.4
 arrays of, 19.26–19.29, 19.28f, 19.28t
 fabrication and configurations of, 19.6–19.8, 19.9f
 gain mechanism of, 19.4
 high-power semiconductor lasers, 19.18–19.30
 arrays, 19.26–19.29, 19.28f, 19.28t
 commercial, 19.19–19.23, 19.20t, 19.21f, 19.22f
 future directions for, 19.23–19.26, 19.25f, 19.25t, 19.26t, 19.27f
 mode-stabilized lasers, 19.18–19.19, 19.19f
 two-dimensional, 19.29–19.30, 19.29t, 19.30f
 high-speed modulation of, 19.30–19.36, 19.31f–19.36f
 operation of, 19.4–19.6, 19.4f–19.6f
 quantum cascade lasers, 16.36
 quantum well lasers, 19.9–19.18
 GRIN SCH single, 19.14, 19.14f
 long wavelength, 19.17–19.18, 19.17f
 schematic of, 19.10f
 strained, 19.15–19.17, 19.16f
 threshold modal gain, 19.12–19.15, 19.12f, 19.13f, 19.15f
 spectral properties of, 19.36–19.39, 19.37f, 19.38f
 surface-emitting lasers, 19.39–19.41
 distributed grating, 19.40–19.41, 19.40f
 integrated laser with 45° mirror, 19.39–19.40, 19.39f
 vertical cavity, 16.36, 19.41, 19.42f, 19.43t
- Semiconductor photodetectors, 25.2, 38.9, 38.9t
- Semiconductors:
 arrays of, 16.36
 direct, 17.4, 17.5f
 indirect, 17.4, 17.4f, 17.5f, 17.6
 material systems for, 18.1–18.2
 properties of substrates for, 17.20t
 waveband structure of, 17.3–17.6, 17.3f–17.5f
- Semisealed-ampoule diffusion, 17.24
- Sensitive area, of photoemissive detectors, 24.41
- Sensitivity:
 defined, 24.12
 film spectral, 29.11, 29.11f
 film speed and, to high-energy radiation, 30.19–30.20
 of interferometers, 13.13–13.14, 13.14f
 of photoemissive detectors, 24.34, 24.35f–24.39f, 24.41
 quantum, 30.9
- Sensitometry variation, with film processing, 29.10, 29.10f
- Sensors:
 area arrays of, 32.24–32.32, 32.25t
 about, 32.2
 frame transfer CCD, 32.26–32.28, 32.27f, 32.28f
 image area dimensions for, 32.25t
 interline transfer CCD, 32.28–32.32, 32.29f–32.31f
 linear image, 32.2, 32.21–32.24, 32.22f, 32.23f
 metal-oxide-semiconductor, 32.25–32.26, 32.26f
 image, 32.2–32.12, 32.3f, 32.21–32.34
 antiblooming in, 32.9, 32.10f
 color imaging with, 32.32–32.34, 32.33f, 32.34f
 dark current in, 32.10–32.12, 32.11f
 junction photodiodes, 32.3–32.6, 32.4f, 32.6f
 linear arrays of, 32.21–32.24, 32.22f, 32.23f
 MOS capacitors, 32.7–32.8
 photoconductors, 32.8–32.9
 pinned photodiodes, 32.8
 LED detectors in, 17.34
 staggered linear CCD, 32.23f, 32.24
 time-delay-and-integrate linear, 32.23f, 32.24

- Separate confinement heterostructure waveguide, **19.24**
- Separated absorption, grading, and multiplication layer APDs (SAGM APDs), **26.17, 26.18f, 26.20f**
- Separated absorption and multiplication layer APDs (SAM APDs), **26.3, 26.17**
- Servo stability, **22.8**
- Servos, **22.5–22.12**
- Bode representation of, **22.5–22.6, 22.6f**
 - closed-loop performance, **22.8**
 - closed-loop stability issues, **22.8–22.12, 22.9f–22.11f**
 - design with time delay, **22.19–22.20**
 - measurement noise not a performance limit, **22.7–22.8**
 - phase and amplitude responses vs. frequency, **22.6–22.7, 22.6f, 22.7f**
- Seven-segment LED displays, **17.10, 17.11f**
- Shapes, projected areas of common, **36.3–36.4, 36.3f, 36.3t**
- Sharpness, of photographic images, **29.18, 29.19, 29.22**
- Shells, mounting of, **6.11, 6.12f**
- Shock specifications, for lenses, **4.10**
- Short arc light sources, **15.34, 15.35f, 40.39**
- Short-wavelength infrared (SWIR), **24.3, 33.3, 33.5**
- Shot noise, **24.12, 27.3, 27.3f, 32.12**
- SI units, **34.20, 37.7, 37.7t**
- Side lighting, **40.12**
- Sidecar TDI, **33.17, 33.17f**
- Sign lighting, **40.71**
- Signal analysis, **27.12–27.15**
- boxcar averaging, **27.13, 27.13f**
 - categories of, **27.3**
 - gated integration, **27.12–27.13, 27.13af**
 - lock-in amplifiers, **27.13, 27.14, 27.14f**
 - photon counting, **27.14**
 - selection of technique, **27.14–27.15**
 - transient photon counting, **27.14**
 - of unmodulated sources, **27.12**
- Signal detection, **27.1–27.12**
- and amplifiers, **27.10–27.12, 27.11f**
 - and noise sources, **27.3–27.6, 27.3f, 27.5f, 27.6f**
 - photomultiplier applications in, **27.6–27.10, 27.7f**
 - technique selection for, **27.2–27.3, 27.2f**
- Signal-to-noise ratio (S/N), **22.12, 27.1, 27.3, 29.1, 29.22–29.23, 33.2, 38.10**
- Silicon (Si):
- and diamond turning, **10.5**
 - doped extrinsic, **33.7, 33.8f**
 - Si:Ga infrared detectors, **24.95, 24.95f, 24.96, 24.96f**
- Silicon avalanche photodiodes (APDs), **24.62–24.65, 24.63f–24.66f**
- Silicon bolometers, **28.7t**
- Silicon carbide (SiC) LED devices, **17.18**
- Silicon carbide (SiC) substrate (for HB-LEDs), **18.2, 18.3**
- Silicon carbide (SiC) UV detectors, **24.47, 24.47f**
- Silicon CCDs (SCCDs), **33.11–33.13, 33.11f, 33.12f, 33.17**
- Silicon nitride layer, **17.23**
- Silicon oxide (SiO₂) passivation, **18.4**
- Silicon oxynitride layer, **17.23**
- Silicon (Si) photoconductors, **32.4f, 32.31, 32.32**
- Silicon (Si) photodiodes, **38.9, 38.9t**
- avalanche, **24.62–24.65, 24.63f–24.66f**
 - high-quality, **34.30**
 - light-trap, **34.30**
 - np*, **34.30**
 - pin*, **24.55f–24.57f, 24.58–24.61, 24.59f, 24.60f**
 - pn*, **24.52f, 24.55–24.58, 24.55f–24.59f, 34.30**
 - self-calibration of, **34.29**
 - UV- and blue-enhanced, **24.55f, 24.61–24.62, 24.61f, 24.62f**
- Silicon (Si) photovoltaic detectors, **24.54–24.65, 24.55f, 24.56f**
- avalanche photodiodes, **24.62–24.65, 24.63f–24.66f**
 - pin* photodiodes, **24.55f–24.57f, 24.58–24.61, 24.59f, 24.60f**
 - pn* photodiodes, **24.52f, 24.55–24.58, 24.55f–24.59f**
 - UV- and blue-enhanced photodiodes, **24.55f, 24.61–24.62, 24.61f, 24.62f**
- Silicon-intensifier-target (SIT) vidicons, **31.8**
- Silver, colloidal, **29.13**
- Silver halide crystals, **30.1, 30.5–30.7, 30.6f**
- Silver halide light detectors, **30.7–30.9, 30.8f**
- Silver halide surfaces, **29.4, 30.14–30.15, 30.15f**
- Simple lens, thermal focus shift of, **8.2–8.4, 8.3t, 8.4t**

- Simulated annealing, 3.19
 Simultaneous measurement, in phase-shifting interferometry, 13.22
 Simultaneous multiple surfaces (SMSs), 39.17, 39.17*f*
 Sine plate, 12.10, 12.11*f*
 Single heterojunction LEDs, 17.12, 17.12*f*
 Single isolated pulses, 21.4
 Single material designs, 6.22, 6.23*f*
 Single monochromators, 38.15*f*, 38.16*f*
 Single optical pulse, 20.2–20.3, 20.3*f*, 20.6–20.7
 Single point diamond turning (SPDT), 6.1, 6.20
 Single quantum well (SQW) LEDs, 18.1
 Single-frequency lasers, 19.37*f*
 Single-lens arrays, 39.33–39.34, 39.34*f*
 Single-longitudinal-mode lasers, 19.38
 Single-pass photodetectors, 26.4*f*
 Single-shot f -to- $2f$ interferometers, 21.6
 Single-use cameras, 30.26
 Size-of-source effect, 34.33
 Skew ray limits, 39.20, 39.20*f*
 Skewness, 40.42
 Skot (unit), 36.7
 Skytubes, 40.50*f*
 Skywells, 40.50*f*
 Sliding prisms, 12.4, 12.4*f*
 Slit polynomials, 11.30, 11.35*t*, 11.36*t*
 Slot interrupters, 17.34
 Small-signal gain coefficient, 16.10
 Smoke detectors, 17.34
 Snubbing, 27.9–27.10
 Society of Automotive Engineers (SAE), 40.2, 40.63–40.64
 Society of Photooptical and Instrumentation Engineers (SPIE), 25.2, 39.12
 Soft mounting, 6.1, 6.4–6.5, 6.4*f*
 Software:
 for lighting simulation, 40.18–40.23
 for nonimaging modeling, 39.6–39.8
 for optical design (*see* Optical design software)
 for stray light suppression, 7.24–7.27
 Solar collection, 39.1
 Solar light pipes (SLPs), 40.49, 40.51*f*
 Solid angles, 34.9, 34.9*f*, 37.4, 39.5, 39.5*f*
 Solid lightpipes, 39.30–39.31
 Solid-state lasers, 16.12, 16.13, 16.17–16.18, 16.18*f*, 22.20–22.21
 Solid-state lighting, 18.4–18.5, 18.4*f*
 Solid-state photomultipliers (SSPM), 33.9
 Soliton solution, 20.5–20.9
 Solves (term), 3.5
 Source coupling, 40.41–40.42
 Source diameter, for fiber optics, 17.33
 Source modeling, 39.7
 Source modeling software, 40.19–40.20
 Source modulation, 27.3
 Spaciousness, perception of, 40.5
 Spatial dilution, 39.6
 Spatial noise, 33.26–33.27, 33.26*f*
 Special-purpose sources (of radiation), 15.53
 Specifications, optical (*see* Optical specifications)
 Speckle effects, 39.32
 Spectra Diode Labs, 19.29, 19.29*t*, 19.41
 Spectral (term), 35.2, 35.3
 Spectral D-double star, 24.12
 Spectral density, 22.3–22.4
 Spectral dependence (of radiometric quantities), 34.9–34.10
 Spectral detectivity, 24.12
 Spectral D-star, 24.12
 Spectral emittance, 35.7, 35.15
 Spectral errors, 34.36
 Spectral irradiance, 38.1–38.2, 38.11*t*, 38.13–38.16, 38.13*f*–38.16*f*
 Spectral irradiance calibration transfer devices, 34.31
 Spectral irradiance lamps, 15.11, 15.12, 15.13*f*
 Spectral lambertian source, 34.17
 Spectral lamps, 15.44, 15.45, 15.45*f*, 15.46*f*, 15.46*t*
 Spectral luminous efficiency, for photopic vision, 36.8, 36.8*f*, 36.9*f*, 36.16*f*, 37.2
 Spectral noise density, 22.3–22.5
 Spectral noise equivalent power, 24.12
 Spectral properties:
 of laser field, 23.28–23.31, 23.30*f*
 of micromaser field, 23.31–23.33
 of semiconductor lasers, 19.36–19.39, 19.37*f*, 19.38*f*
 Spectral radiance, 38.2, 38.11*t*, 38.13–38.16, 38.13*f*–38.16*f*
 Spectral radiance calibration transfer devices, 34.31
 Spectral radiance ribbon filament lamps, 15.11, 15.12*f*
 Spectral radiance units, 34.23–34.24

- Spectral reflectance, 35.4, 35.5, 38.2, 38.17–38.18, 38.18f
- Spectral response, of photodetectors, 24.18, 24.19f
- Spectral responsivity, 24.12, 38.3, 38.18–38.19
- Spectral scattering, 38.10
- Spectral sensitivity, of photographic films, 29.11, 29.11f
- Spectral sensitizers, photographic, 30.13–30.18
 about, 30.13–30.14, 30.14f
 color science of, 30.15–30.18, 30.16f, 30.17f
 photophysics of, on silver halide surfaces, 30.14–30.15, 30.15f
- Spectral transmittance, 35.10, 38.2, 38.17, 38.17f
- Spectrally stray radiation errors, 34.36
- Spectralon, 38.12, 38.13
- Spectrophotometers and spectrophotometry, 34.6, 35.8–35.9, 38.17, 38.17f
- Spectroradiometry, 38.1–38.19
 about, 38.1
 calculations for, 38.3–38.5
 calibration of, 38.11–38.13, 38.11t, 38.12f
 computer software for, 38.11
 detectors in, 38.8–38.10, 38.9t, 38.10t
 electronics of, 38.10
 errors in, 38.5–38.6
 figures of merit in, 38.5–38.6
 input (fore-) optics in, 38.7
 monochromators in, 38.7–38.8
 quantities used in, 38.1–38.2
 spectroradiometers, 38.18, 38.18f
 system designs in, 38.13–38.19
 spectral irradiance/radiance, 38.13–38.16, 38.13f–38.16f
 spectral reflectance, 38.17–38.18, 38.18f
 spectral responsivity, 38.18–38.19
 spectral transmittance, 38.17, 38.17f
- Specular reflectance, 35.10, 35.13
- Specular transmittance, 35.3, 35.9f
- Specular vanes, 7.17, 7.17f
- Speed:
 of LEDs, 17.33
 of photodetectors, 24.20, 24.21
 (See also Photographic film speed)
- SPEOS (optical software), 7.27
- Sphere(s):
 aberrations in, 11.30
 integrating (see Integrating spheres)
- Sphere(s) (*Cont.*):
 nonuniformities with integrating, 39.24–39.26, 39.24f, 39.25f
 projected area of, 36.3–36.4, 36.3f, 36.3t
- Spherical lambertian source, 34.17
- Spherical lenses, 39.8, 39.9f
- Spherical optics fabrication, 9.4–9.6
- Spherochromatism, 2.2
- Spherometers, 12.18–12.19, 12.18f, 12.19f, 12.19t
- Spline surfaces, 39.6
- Spokes, aperture, 3.20
- Spontaneous emission lasers (see Correlated emission lasers)
- Spontaneous emission rate, 23.8
- Spot-diagram analysis, 3.13–3.16
- Spring-loaded mountings, 6.13, 6.14f
- Square polynomials, 11.30, 11.31t–11.34t, 11.36t
- Stability:
 light, 30.10
 of photodetectors, 24.21, 24.21f
 of photoemissive detectors, 24.41
- Stabilization, light, 30.12–30.13 (See also Laser stabilization)
- Stable resonators, 16.23, 16.23f
- Stadia, 12.2–12.3, 12.3f
- Staggered linear CCD image sensor, 32.23f, 32.24
- Stagnation, 3.17
- Staircase APDs, 26.3
- Standard for the Exchange of Product model data (STEP), 40.19
- Standards:
 baseline, of radiation sources, 15.9f, 15.10f, 15.12f
 for detectors, 38.12–38.13
 for infrared radiometry, 15.11–15.12, 15.12f
 international, 4.11
 for length measurements, 12.2
 for lighting, 40.19
 for lighting system layout and simulation, 40.19
 for optical image quality, 4.6
 published, 4.10
 of reflectance, 35.14t
 search engine for, 4.11
 of spectral transmittance, 35.10
 for vehicular lighting, 40.63–40.64, 40.66f, 40.66t
 working, of radiation sources, 15.9–15.13, 15.10f, 15.12f, 15.13f
- Star concentrators, 39.20, 39.21

- Staring arrays, 33.6–33.7, 33.6f, 33.14
 Staring FPAs, 33.16–33.17, 33.29, 33.29t
 Steady-state pulse train, 20.5–20.9
 Stefan-Boltzmann law, 34.24, 37.11
 Steradian (sr), 36.3, 37.4, 37.4f
 Sterilamps, 15.35, 15.36f
 Stilb (unit), 34.43, 36.7, 36.8t
 Stimulated absorption, 16.7–16.8, 16.8f
 Stimulated emission, 16.2, 16.7–16.9, 16.8f, 23.8
 Stop lamps, 40.64f, 40.67, 40.68t, 40.69f
 Stop shifting, 2.5, 2.6f
 Stops:
 aperture, 34.18, 34.19f
 field, 34.18–34.19, 34.19f
 Straddling springs, 6.13, 6.14f
 Straightness measurement, 12.10
 Strained QW lasers, 19.15–19.17, 19.16f
 Stray capacitance noise, 27.5, 27.6f
 Stray light, 29.15–29.16, 29.15f
 Stray light suppression, 7.1–7.32
 about, 7.1–7.2
 aperture placement in, 7.5–7.10
 aperture stops, 7.6–7.7, 7.7f, 7.8f
 field stops, 7.7, 7.8f, 7.9f
 Lyot stops, 7.8–7.10, 7.8f–7.11f
 baffles in, 7.10, 7.11
 and BRDF characteristics, 7.23, 7.24f
 Cassegrain design with aperture stop at
 primary (example), 7.3f
 contamination levels in, 7.18–7.19, 7.18t,
 7.19f–7.21f
 evaluation methods for, 7.27–7.29, 7.29f
 illuminated objects in, 7.5, 7.5f, 7.6f
 imaged critical objects in, 7.4, 7.5f
 information sources on, 7.31–7.32
 issues with, 7.30–7.31
 and point source transmittance definitions,
 7.22–7.23
 radiation transfer equation for, 7.21–7.22
 real-space critical objects in, 7.2–7.4, 7.3f, 7.4f
 software for, 7.24–7.27
 and stray radiation paths, 7.22
 strut design in, 7.20, 7.21, 7.21f
 and surface scattering characteristics, 7.23
 vane spacing and depth in, 7.13–7.17
 angle considerations, 7.13–7.16, 7.14f, 7.15f
 bevel placement, 7.13, 7.14f
 depth considerations, 7.16, 7.16f, 7.17f
 specular vanes, 7.17, 7.17f
 vanes in, 7.11–7.12, 7.12f, 7.13f
 Stray radiation paths, 7.9, 7.22
 Street lighting, 40.69–40.71, 40.70t, 40.71t
 Stress tolerance, 6.3
 Stretched segment displays, 17.30–17.31,
 17.30f, 17.31f
 Strip mirror integrator (SMI), 39.40
 Strong field approximation, 21.3
 Strong VW reflectometer, 35.10f
 Strut design (in stray light suppression), 7.20,
 7.21, 7.21f
 Subjective tone reproduction, 29.16
 Submillimeter (SubMM) radiation, 24.3
 Subminiature lamps, 15.53
 Sub-Nyquist interferometry, 13.27
 Substrate(s):
 absorbing, 17.7, 17.7t
 for HB-LEDs, 18.2–18.3, 18.2f
 LED, 17.20–17.21, 17.20t
 mirror, 6.17–6.18, 6.17f, 6.18f
 transparent, 17.7, 17.7t
 Suncatchers, 40.48, 40.50f
 Superconducting bolometers, 28.5, 28.7t
 Superposition (of uniformity), 39.2, 39.32,
 39.33f
 Superposition-of-sources nonlinearity mea-
 surement, 34.33
 Supersensitizers, 30.14, 30.15, 30.15f
 Support wires, light bulb, 40.29f, 40.30
 Surface emitting lasers (SELS) (SLASERs),
 19.39–19.41, 19.39f, 19.40f, 19.42f, 19.43t,
 25.15
 Surface emitting LEDs (SLEDs), 25.15
 Surface finishing, of diamond-turned optics,
 10.9–10.11, 10.9f–10.11f
 Surface generation current, 32.10, 32.11,
 32.11f
 Surface measurement systems, 40.53, 40.54
 Surface mount device (SMD) package, 18.5f
 Surface mount LEDs (SMDs), 40.37
 Surface profilometers, 9.6
 Surface scattering, 7.23
 Surface-channel CCDs, 32.14
 Surface-channel MOS capacitors, 32.4f, 32.7
 Surfaces, modeling of, 40.17
 Suspended luminaires, 40.13f
 Synchronotron radiation, 34.26–34.27
 Synchronous pumping, 16.29
 System specifications, for lenses, 4.3
 Système International (SI), 12.2, 36.2, 37.3
 (See also SI units)

- Taillights, 40.21, 40.22*f*, 40.64*f*, 40.67, 40.68*t*, 40.69*f*
- Tailored (T) reflectors, 39.37–39.39, 39.38*f*
- Tailoring (of uniformity), 39.2
- Talbot autoimages, 12.23, 12.24
- Talbot's law, 34.33–34.34
- Tandem-lens arrays, 39.34, 39.35*f*–39.37*f*
- Tapered lightpipes, 39.12–39.13, 39.13*f*, 39.31–39.32, 39.31*f*
- Task lighting, 40.12, 40.14
- Taylor-Hobson Form TalySurf, 9.6
- Technical specifications, 4.2
- Tehis method, 40.53, 40.54
- Telecentric distribution, 39.18, 39.18*f*
- Telescope(s):
- astronomical, 1.7*f*
 - Galilean, 1.7*f*
 - Hubble, 11.4, 13.24
 - Keck, 11.4
 - reflecting, 11.4
 - unit magnification Galilean, 12.4, 12.4*f*
- Telescoping input optics, 38.7
- Temperature:
- color, 37.4*t*, 37.6–37.7
 - correlated color, 37.7, 38.5
 - distribution, 37.7
 - and mounted optics, 6.21–6.24, 6.22*f*–6.24*f*
 - radiance, 37.4*t*, 37.6
- Temperature control, of PZT transducers, 22.19
- Temperature noise, 24.12
- Temperature specifications, for lenses, 4.10
- Temperature-dependence effects, 34.36–34.37
- Templates (for curvature measurement), 12.17
- Tensile-strained QW lasers, 19.16, 19.16*f*, 19.17
- Test plates (for curvature measurement), 12.17
- Testing, 13.1–13.27
- aspherical wavefront measurement, 13.23–13.27
 - holographic compensators, 13.25, 13.25*f*, 13.26*f*
 - infrared interferometry, 13.25
 - Moiré tests, 13.26–13.27
 - refractive or reflective compensators, 13.24, 13.24*f*, 13.25
 - sub-Nyquist interferometry, 13.27
 - two-wavelength interferometry, 13.25, 13.26
 - wavefront stitching, 13.27, 13.27*f*
 - computer-generated holograms in (see Computer-generated holograms)
- Testing (*Cont.*):
- of convex surfaces, 14.5
 - interferogram evaluation, 13.14–13.18
 - direct interferometry, 13.17–13.18
 - fixed interferograms, 13.14–13.15
 - Fourier analysis of interferograms, 13.16–13.17, 13.17*f*
 - global and local interpolation of interferograms, 13.15–13.16
 - interferometric, 13.7–13.12
 - common path interferometer, 13.9, 13.11*f*
 - Fizeau interferometer, 13.8–13.9, 13.9*f*, 13.10*f*
 - lateral shearing interferometers, 13.9–13.12, 13.11*f*, 13.12*f*
 - multiple-pass interferometers, 13.13
 - multiple-reflection interferometers, 13.13
 - radial, rotational, and reversal shearing interferometers, 13.12, 13.13*f*
 - sensitivity of interferometers, 13.13–13.14, 13.14*f*
 - Twyman-Green interferometer, 13.7–13.8, 13.7*f*, 13.8*f*
 - Zernike phase-contrast method applied to interferometers, 13.13–13.14, 13.14*f*
 - noninterferometric, 13.1–13.7
 - Foucault test, 13.2–13.3, 13.2*f*, 13.3*f*
 - Hartmann test, 13.4–13.6, 13.5*f*
 - Hartmann-Shack test, 13.6–13.7, 13.6*f*
 - Ronchi test, 13.3–13.4, 13.3*f*, 13.4*f*
 - phase-shifting interferometry, 13.18–13.23, 13.18*f*–13.20*f*
 - heterodyne interferometer, 13.22
 - integrating bucket method, 13.21, 13.21*f*
 - phase errors, 13.22
 - phase stepping, 13.20, 13.20*f*
 - phase-lock method, 13.23, 13.23*f*
 - simultaneous measurement, 13.22
 - two steps plus one method, 13.21, 13.22
 - in wafer processing, 17.24
- Thef-number, 38.8
- Theodolites, 12.13
- Thermal arrays, 28.7–28.12
- about, 28.7–28.8
 - noise equivalent temperature difference in, 28.8–28.9
 - pyroelectric hybrid, 28.11–28.12, 28.11*f*, 28.12*f*
 - resistive bolometer, 28.10–28.11, 28.10*f*
 - theoretical limits of, 28.9–28.10, 28.9*f*
 - thermoelectric, 28.12, 28.12*f*

- Thermal circuit theory, 28.2
- Thermal coefficient of resistance (TCR), 33.2, 33.14
- Thermal compensation, 8.1–8.15
 about, 8.2
 and effect of thermal gradients, 8.6–8.7
 and homogeneous thermal effects, 8.2–8.5, 8.3*t*, 8.4*t*, 8.5*f*
 intrinsic athermalization, 8.7–8.8, 8.7*f*
 mechanical athermalization, 8.8–8.12, 8.8*f*–8.12*f*
 optical athermalization, 8.12–8.15, 8.13*t*–8.15*t*
 tolerable homogeneous temperature change, 8.5–8.6, 8.6*f*
- Thermal defocus, of compound lens, 8.4, 8.5*f*
- Thermal detector(s), 24.4–24.6, 24.4*f*, 28.1–28.12, 38.9, 38.9*t*
 arrays of, 28.7–28.12
 about, 28.7–28.8
 noise equivalent temperature difference, 28.8–28.9
 pyroelectric hybrid arrays, 28.11–28.12, 28.11*f*, 28.12*f*
 resistive bolometer arrays, 28.10–28.11, 28.10*f*
 theoretical limits, 28.9–28.10, 28.9*f*
 thermoelectric arrays, 28.12, 28.12*f*
- bolometer, 24.5*f*, 28.3–28.5, 28.4*f*
- Golay cell, 28.6
- ideal, 28.2–28.3, 28.3*f*
- performance/sensitivity of, 24.17, 24.18*f*
- properties of, 28.7, 28.7*t*
- pyroelectric, 24.6, 24.6*f*, 28.7
- and thermal circuit theory, 28.2
- thermistor, 24.5
- thermocouple, 28.4
- thermopile, 24.5*f*, 28.4–28.5
- Thermal expansion, 33.14
- Thermal fatigue, 17.25
- Thermal focus shift, 8.2–8.4, 8.3*t*, 8.4*t*
- Thermal gradients, effect of, 8.6–8.7
- Thermal infrared detectors, 33.7, 33.8*f*
- Thermal noise, 24.13, 27.4, 32.20
- Thermal properties, of high-power lasers, 19.26, 19.27*f*
- Thermal stability, of plastic packaging materials, 17.26
- Thermistor bolometers, 24.24–24.25, 24.24*f*, 24.25*f*, 28.7*t*
- Thermistors, 24.5
- Thermocouple junctions, noise from, 27.6, 27.6*f*
- Thermocouples, 24.5, 28.7*t*
 about, 28.1
 manufacturers' specifications for, 24.22–24.23, 24.22*f*
 as thermal detectors, 28.4
- Thermoelectric arrays, 28.12, 28.12*f*
- Thermopiles, 24.5, 28.7*t*
 defined, 24.13
 manufacturers' specifications for, 24.23–24.24, 24.23*f*
 as thermal detectors, 28.4–28.5
- θ_1/θ_2 concentrators, 39.18–39.20, 39.19*f*
- Thick window chips, 17.7, 17.7*t*
- Thin doublet, 1.15–1.16
- Thin lenses, 1.5
- Thin teflon diffusers, 38.15*f*
- Thin-disk lasers, 16.18
- 35-mm photographic films, 30.21, 30.25
- Thoria (in incandescent lights), 40.27
- Threaded retaining rings, 6.3, 6.3*f*
- 3D concentrators, 2D vs., 39.20–39.21, 39.20*f*, 39.21*f*
- Three-chip color systems, 32.32, 32.33*f*
- Three-material athermal solutions, 8.14, 8.14*t*, 8.15*t*
- Three-phase CCDs, 32.15, 32.16*f*
- Three-step rescattering model, 21.3
- Threshold carrier density, 19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold current, 19.6, 19.6*f*
- Threshold modal gain, 19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold voltage, 25.11
- Tightly toleranced assembly, 6.7, 6.7*f*
- Time delay integration (TDI), 33.4
- Time delay integration (TDI) linear sensors, 32.23*f*, 32.24
- Time delay integration (TDI) scanning FPAs, 33.17, 33.17*f*
- Time evolution of the field, 23.15–23.17, 23.15*f*
- Time-averaged color mixing, 40.8
- Time-based measurement, 12.2, 12.4, 12.5, 12.6*f*
- Time-dependent error, 34.35
- Time-of-flight distance measurement, 12.4, 12.5

- Titanium oxide (TiO_2) UV detectors, 24.47, 24.48*f*
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) lasers, 16.34, 16.34*f*
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) ring lasers, 20.16–20.17, 20.16*f*, 20.17*f*
- Tolerance budgeting, 5.3
- Tolerance verification, 5.3
- Tolerances, 5.2–5.8
- assembly, 5.8
 - basis for, 5.2–5.3
 - boresight, 5.8
 - budgeting of, 5.3
 - distortion, 5.8
 - optical vs. mechanical, 5.2
 - verification of, 5.3
 - wavefront, 5.3–5.7, 5.4*f*, 5.5*f*, 5.5*t*, 5.6*t*, 5.7*f*
- Tolerancing, 5.8–5.11
- and aberration balancing, 11.35, 11.36
 - about, 5.1–5.2
 - and material properties, 5.9
 - measurement practices for, 5.8–5.9
 - and optimization, 3.20–3.21
 - problems in, 5.11
 - procedures for, 5.9–5.10
 - shop practices for, 5.8
- Tone reproduction, 29.16–29.17, 29.17*f*
- Total flux into a hemisphere, 34.15
- Total hemispherical emittance, 35.15, 35.15*f*
- Total internal reflection (TIR), 39.12, 39.17, 40.41
- Total internal reflection (TIR) Fresnel lenses, 39.10
- Total luminous flux, 37.4*t*, 37.6
- Total radiant flux, 37.4*t*, 37.6
- Total transmittance, 35.3, 35.9*f*
- Traceability:
- of absolute measurements, 34.21
 - errors in, 34.28
- TracePro (optical software), 7.27
- Transconductance amplifiers, 27.11–27.12, 27.11*f*
- Transducer resonance, 22.8, 22.11–22.12
- Transducers, 22.17–22.20
- Transformers, in voltage amplifiers, 27.11
- Transient photon counting, 27.14
- Transmission, 4.7
- actual/idealized, 35.2*f*
 - defined, 35.3
- Transmission density, of photographic films, 29.6–29.7, 29.7*f*
- Transmissive sensors, 17.34
- Transmittance, 35.3
- measurement of, 35.8–35.10, 35.9*f*
 - and reflectance/absorptance, 35.7, 35.8, 35.8*t*
 - spectral, 38.2, 38.17, 38.17*f*
- Transmitter speed, for fiber optics, 17.33
- Transparency, 39.23, 40.5
- Transparency point, 19.5
- Transparent substrate (TS) chips, 17.7, 17.7*t*
- Transportation lighting, 40.63–40.71
- roadway lighting, 40.67, 40.69–40.71, 40.70*t*, 40.71*t*
 - vehicular lighting, 40.63–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
- Transverse electromagnetic mode (TEM), 16.21–16.23, 16.22*f*
- Transverse junction stripe (TJS) lasers, 19.8, 19.9*f*, 19.23–19.24, 19.36*f*
- Transverse laser modes, 16.21–16.23, 16.21*f*–16.23*f*
- Transverse ray plots, 2.2–2.4, 3.13
- Traveling microscopes, 12.20, 12.21, 12.21*f*
- Traveling wave photodetectors, 26.4*f*, 26.5, 26.14*f*
- Treaty of the Meter of 1875, 34.20, 36.2
- Triphosphors, 40.31, 40.32*f*
- Triplet lens, air-spaced, 6.21, 6.22*f*
- Tristimulus values, 38.3–38.4
- Troffers, fluorescent luminaire, 40.47
- Troland (unit), 34.41–34.42, 37.7, 37.8
- Trough reflectors, 40.46*f*, 40.47
- Trumpet (term), 39.15, 39.16*f*, 39.17
- Tubular PZT transducers, 22.17–22.18
- Tungsten:
- in HID lamps, 40.35
 - in incandescent lights, 40.25, 40.27, 40.29
- Tungsten lamps, 15.13, 40.26*t*, 40.28*f*
- Tungsten-arc lamps, 15.47–15.48, 15.48*f*, 15.49*f*
- Tungsten-filament lamps, 15.11, 15.12, 15.13*f*, 15.19, 15.20, 15.20*f*–15.22*f*, 34.31
- Tungsten-halogen lamps, 15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Tunnel diagram (*see* Williamson construction)
- Tunnel lighting, 40.71
- Tunneling current, 25.8
- Twin-channel lasers (TCLs), 19.27
- Twin-channel substrate mesa (TCSM) lasers, 19.20*t*, 19.21*f*, 19.23

- Twin-ridge structure (TRS) lasers, **19.19**,
19.20t, **19.21f**, **19.22–19.23**
- 2D (term), **39.4**
- 2D concentrators, 3D vs., **39.20–39.21**, **39.20f**,
39.21f
- 2D high-power laser arrays, **19.29–19.30**,
19.29t, **19.30f**
- Two-color gating, **21.7**
- Two-component systems, first-order layout for,
1.5–1.7
- Two-interference pattern distance-measuring
interferometer, **12.7**, **12.7f**
- Two-mirror imaging system, **39.17**
- Two-phase CCDs, **32.15–32.16**, **32.16f**
- Two-stage baffle, **7.10**
- Two-step rescattering model, **21.3**
- Two-steps-plus-one phase shifting,
13.21, **13.22**
- Two-wavelength interferometry, **13.25**, **13.26**
- Twyman-Green interferograms, **13.10f**, **13.18f**
- Twyman-Green interferometers, **13.7–13.8**,
13.7f, **13.8f**
- Type A errors (in absolute measurement),
34.21–34.23
- Type B errors and error sources (in absolute
measurement), **34.32–34.37**
- defined, **34.21–34.23**
- nonideal aperture, **34.35–34.36**, **34.35f**
- nonlinearity of detector, **34.34–34.35**
- nonuniformity, **34.35**
- offset subtraction, **34.32–34.33**
- polarization effects, **34.33**
- scattered radiation effect, **34.33**
- size-of-source effect, **34.33**
- spectral errors, **34.36**
- temperature-dependence effects, **34.36–34.37**
- time-dependent error, **34.35**
- Ultrashort cavity microlasers, **19.39**
- Ultrashort optics, **20.1–20.28**
- about, **20.1–20.2**
- cavities with two circulating pulses,
20.15–20.22
- linear lasers, **20.18–20.19**, **20.19f**
- optical parametric oscillators, **20.20–20.22**,
20.20f, **20.21f**
- ring dye lasers, **20.15–20.16**
- ring lasers, **20.17–20.18**, **20.17f**
- Ti:sapphire ring lasers, **20.16–20.17**,
20.16f, **20.17f**
- Ultrashort optics (*Cont.*):
- coupling of circulating pulses, **20.12–20.15**,
20.12f, **20.15f**
- optical pulses and pulse trains, **20.2–20.9**
- single optical pulse, **20.2–20.3**, **20.3f**
- soliton solution and steady-state pulse
train, **20.5–20.9**
- train of pulses, **20.3–20.5**, **20.4f**, **20.5f**
- and quantum mechanical two-level system,
20.22–20.28
- coherent interaction, **20.22–20.23**
- experimental demonstration, **20.24–20.27**,
20.25f–20.27f
- impact of analogy, **20.27–20.28**
- laser as two-level system, **20.23–20.24**,
20.25t
- Rabi cycling, **20.26–20.27**, **20.26f**, **20.27f**
- steady-state pulse, **20.9–20.12**, **20.11f**
- Ultrasonic-assisted machining, **10.5**
- Ultraviolet (UV) detectors:
- silicon carbide, **24.47**, **24.47f**
- TiO₂, **24.47**, **24.48f**
- Ultraviolet (UV) enhanced photodiodes,
24.55f, **24.61–24.62**, **24.61f**, **24.62f**
- Ultraviolet (UV) filters, **40.12**
- Ultraviolet (UV) radiation, **34.6**
- and color film, **30.3**
- far, **15.12**, **15.13**
- spectrum of, **25.2**
- vacuum, **24.3**
- Uncrossed reflectors, **39.38**, **39.38f**
- Unified Glare Rating (UGR), **40.10–40.11**, **40.11t**
- Uniform illumination, of nonimaging optics,
39.22–39.41
- with classic projection systems, **39.23–39.24**,
39.23f
- faceted structures in, **39.39–39.41**, **39.39f**,
39.40f
- integrating cavities in, **39.24–39.27**, **39.24f**,
39.25f, **39.27f**
- lens arrays in, **39.32–39.37**, **39.33f–39.37f**
- lightpipes in, **39.13f**, **39.27–39.32**,
39.28f–39.30f
- tailored reflectors, **39.37–39.39**, **39.38f**
- Uniformity:
- angular, **39.31**
- control of, **39.1–39.2**
- of luminance/illuminance, **40.7**, **40.13f**
- of photodetectors, **24.20**
- and visual discomfort, **40.9**

- Unit conversions:
 - for English and SI units, 37.7, 37.7t
 - for illuminance, 36.7t, 36.8t
 - for photometric and radiometric quantities, 36.11–36.14, 36.12f–36.14f
- Unit magnification Galilean telescope, 12.4, 12.4f
- Unlit-appearance modeling, 40.21
- Unmodulated signal sources, 27.12
- Unstable resonators, 16.25–16.26, 16.26f
- Unstrained QW lasers, 19.15–19.16, 19.16f
- Uplight, 40.43, 40.44f, 40.45
- U.S. Air Force three-bar target, 4.6
- Useful life period, of LEDs, 17.26, 17.26f
- Uviarc, 15.28–15.29, 15.29f, 15.30f

- Vacuum, laser gain media in, 16.36–16.37, 16.37f
- Vacuum lamps, 34.31
- Vacuum ultraviolet (VUV) radiation, 24.3
- Valence band, 17.3, 17.4, 17.4f
- Valence lighting, 40.13f
- Vanes (in stray light suppression), 7.11–7.17
 - defined, 7.11–7.12, 7.12f, 7.13f
 - placement design for, 7.12f
 - and scatter path, 7.13f
 - spacing and depth of, 7.13–7.17, 7.14f–7.17f
- Vapor exposure, in LED packaging, 17.26
- Vapor phase epitaxy (VPE), 17.21, 17.22
- Variable temperature blackbody, 15.10f
- Variable-orientation mirrors, 6.17
- Varifocal systems, first-order layout for, 1.11–1.12
- Vector flux, 39.21–39.22
- Vehicular lighting, 40.63–40.67, 40.64f, 40.65t, 40.66f, 40.66t, 40.68t, 40.69f
- Veiling reflections, 40.12, 40.14
- Verification (of tolerance), 5.3
- Vertical antiblooming, 32.9, 32.10f
- Vertical Bridgeman technique, 17.21
- Vertical cavity lasers, 19.41, 19.42f, 19.43t
- Vertical cavity semiconductor lasers, 16.36
- Vertical cavity surface-emitting lasers (VCSELs), 16.36
- Vertical illuminance, 40.7
- Vertically integrated photodiode (VIP) FPAs, 33.10
- Vertically illuminated *pin* photodiodes, 26.3, 26.4f, 26.5, 26.10, 26.12–26.13, 26.12f

- Very-long-wavelength infrared (VLWIR) radiation, 24.3
- Very-long-wavelength semiconductor lasers, 19.7–19.8
- Vibration specifications, for lenses, 4.10
- Vibration-resistant optical reference cavity, 22.16, 22.17f
- Vignetting, 3.4, 34.19
- Virtual phase CCDs, 32.16–32.17, 32.16f
- Visible array detectors, 32.1–32.34
 - about, 32.2
 - image sensing elements of, 32.2–32.12, 32.3f
 - antiblooming, 32.9, 32.10f
 - dark current, 32.10–32.12, 32.11f
 - junction photodiode, 32.3–32.6, 32.4f, 32.6f
 - MOS capacitor, 32.7–32.8
 - photoconductor, 32.8–32.9
 - pinned photodiode, 32.8
 - readout elements of, 32.12–32.21
 - CCD, 32.12–32.20, 32.13f, 32.15f–32.18f
 - MOS, 32.20–32.21
 - sensor architectures of, 32.21–32.34
 - area image sensor arrays, 32.24–32.32, 32.25t, 32.26f–32.31f
 - color imaging, 32.32–32.34, 32.33f, 32.34f
 - linear image sensor arrays, 32.21–32.24, 32.22f, 32.23f
- Visible light photon counters (VLPCs), 33.9
- Visible (VIS) radiation, 24.3, 25.2
- Vision, 40.3–40.6, 40.9
 - biology of, 40.3–40.4
 - and perception, 40.4–40.5
 - photopic/scotopic/mesopic, 34.37–34.39, 37.2
 - (*See also* Human eye)
- Visual clarity, perception of, 40.5
- Visual discomfort, 40.9–40.12, 40.11t
- Visual discomfort probability (VCP), 40.10
- Visual photometry, 36.4
- Visual science, 34.37
- Vivid color (VC) film, 30.27
- Voltage amplifiers, 27.10–27.11

- Wafer processing, 17.23–17.25
- Wall slot lighting, 40.13f
- Wall-grazing illumination, 40.13f
- Wall-washing illumination, 40.13f
- Watanabe, F., 39.33
- Watt (unit), 39.2t

- Wave modulation distance meter, 12.5, 12.6f
- Waveband materials, 8.3t, 8.4t
- Waveband structure of semiconductors, 17.3–17.6, 17.3f–17.5f
- Wavefront error (*W*), 4.1, 4.3, 4.7, 4.8
- Wavefront measurement, aspherical (*see* Aspherical wavefront measurement)
- Wavefront stitching, 13.27, 13.27f
- Wavefront tolerancing, 5.3–5.7, 5.4f, 5.5f, 5.5t, 5.6t, 5.7f
- Wavefronts, from lenses, 4.3–4.5, 4.5t
- Waveguide photodetectors, 26.4f, 26.5
- Waveguide *pin* photodiodes, 26.13–26.14, 26.14f
- Wavelength, in fiber optics, 17.33–17.34
- Wavelength errors, 34.36
- Wearout period, 17.26, 17.26f
- Weighting functions, 36.17
- Well capacity, 25.11
- Welsbach mantle, 15.17, 15.18
- Whiffletrees (lever mechanisms), 6.19
- White light, 18.4–18.5, 18.4f, 40.7, 40.8, 40.24
- White surfaces, reflectivity of, 17.31
- White-light LEDs, 40.37, 40.38
- WI 9 lamps, 15.21f
- WI 14 lamps, 15.21f
- WI 16/G lamps, 15.21f, 15.22f
- WI 17/G lamps, 15.22f
- WI 40/G lamps, 15.22f
- WI 41/G lamps, 15.22f
- Wiener spectrum, 29.21
- Wien's displacement law, 15.7, 34.23, 34.24, 37.11
- Williamson construction, 39.12–39.13, 39.13f, 39.28, 39.29f, 39.31, 39.32
- Window/photocathode assemblies, of image intensifiers, 31.10–31.12, 31.11f, 31.12f
- Windows:
and daylight sources, 40.41, 40.47, 40.48, 40.49f, 40.50f
mounting of optical, 6.11, 6.11f, 6.12f
- Wire-wound thermopile arrays, 24.23
- Work function (of photons), 25.2
- Xenon lamps, 15.34f, 15.35f, 40.31, 40.35f
- X-ray lasers, 16.31
- X-Y addressing, 33.16
- Y-coupled junctions, 19.27, 19.29
- Yellow filter dyes, 30.4
- Yellow light, 29.13, 29.13f
- Yttrium aluminum garnet (YAG) phosphor, 18.4
- ZEMAX (optical software), 7.26–7.27
- Zernike phase-contrast test, 13.13–13.14, 13.14f
- Zernike polynomials, 5.9
annular, 11.13–11.21, 11.14f, 11.17t–11.21t
circle, 11.4, 11.6–11.12, 11.8t–11.9t, 11.9f–11.11f, 11.12t, 11.39
- Zerodur prisms, 6.16, 6.16f
- Zinc, 17.23
- Zinc diffusion, 17.9–17.10, 17.10f
- Zinc doping, 17.20
- Zinc oxide (ZnO) doped GaP, 17.16, 17.21–17.22
- Zinc selenide (ZnSe) LED devices, 17.19
- Zinc-doped germanium (Ge:Zn) detectors, 24.84f, 24.98–24.100, 24.99f
- Zirconium arc lamps, 15.47, 15.48f
- Zonal cavity lighting simulation, 40.17
- Zoom systems, 1.11–1.12, 3.20
- Z-system (eccentric pupil design), 7.11, 7.12f, 7.15–7.17, 7.15f, 7.17f, 7.19, 7.21f

COLOR PLATES

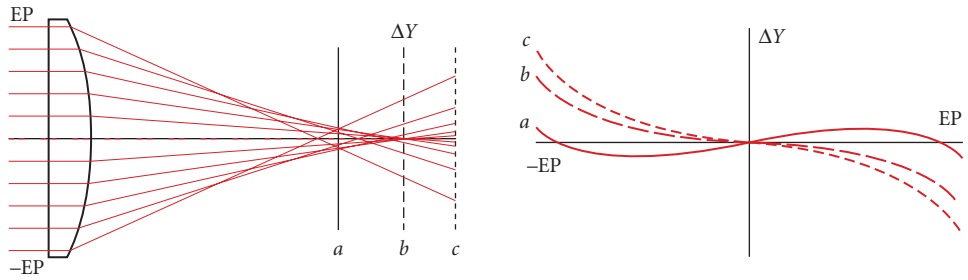


FIGURE 2.1 (Left) Rays exiting a lens are intercepted at three evaluation planes. (Right) Ray intercept curves plotted for the evaluation planes: (a) at the point of minimum ray error (circle of least confusion); (b) at the paraxial image plane; and (c) outside the paraxial image plane.

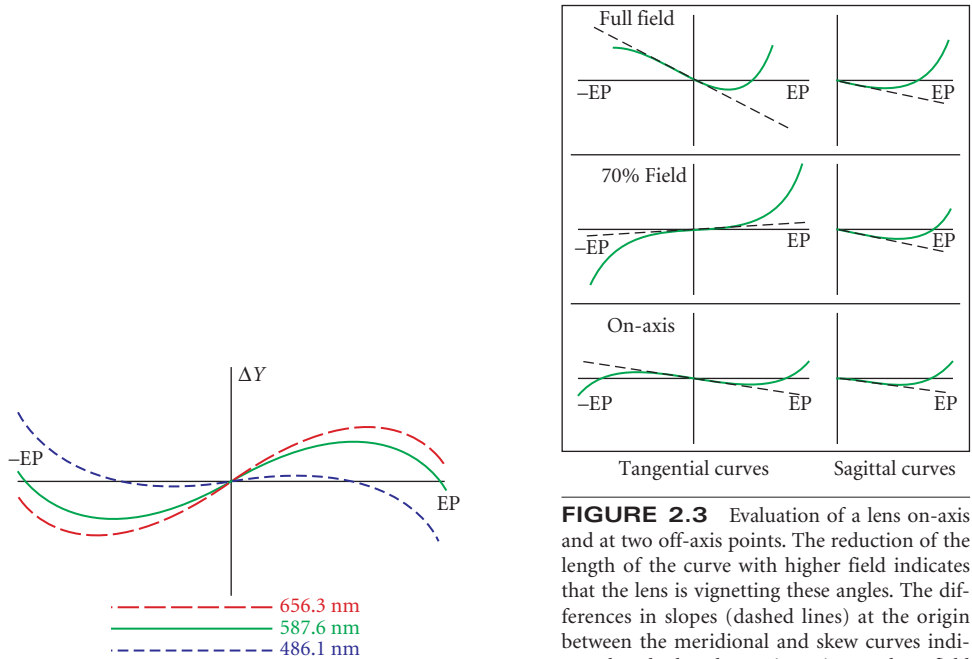


FIGURE 2.2 Meridional ray intercept curves of a lens with spherical aberration plotted for three colors.

FIGURE 2.3 Evaluation of a lens on-axis and at two off-axis points. The reduction of the length of the curve with higher field indicates that the lens is vignetting these angles. The differences in slopes (dashed lines) at the origin between the meridional and skew curves indicate that the lens has astigmatism at these field angles. The variation in the slopes with field indicates the presence of field curvature.

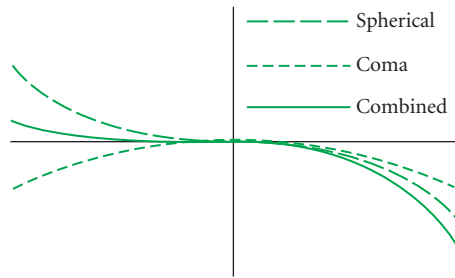


FIGURE 2.4 Ray intercept curve showing coma combined with spherical aberration.

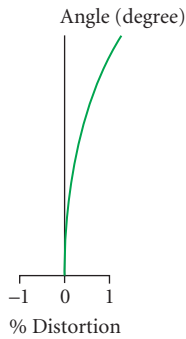


FIGURE 2.5 Field curve: distortion plot. The percentage distortion is plotted as a function of field angle. Note that the axis of the dependent variable is the horizontal axis.

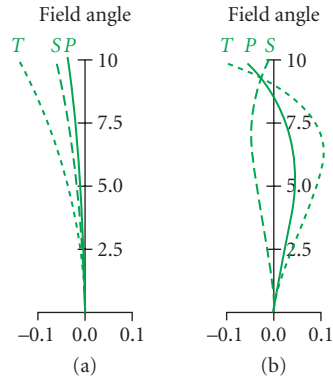


FIGURE 2.6 Field curve: field curvature plot. The locations of the tangential T and sagittal S foci are plotted for a full range of field angles. The Petzval surface P is also plotted. The tangential surface is always three times farther from the Petzval surface than from the sagittal surface: (a) an uncorrected system and (b) a corrected system.

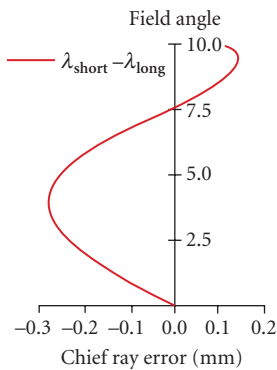


FIGURE 2.7 Field curve: lateral color plot. A plot of the transverse ray error between red and blue chief ray heights in the image plane for a full range of field angles. Here the distance along the horizontal axis is the color error in the image plane.

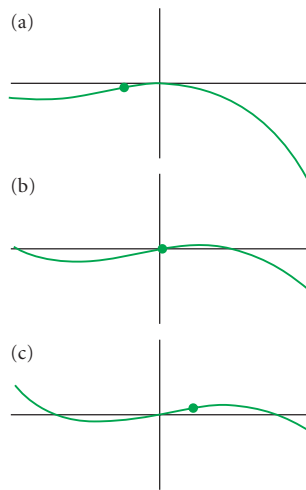


FIGURE 2.8 The effect of stop shifting on the meridional ray intercept curves of a double Gauss lens. (a) Stop located in front of the normal centrally located stop. (b) Stop at the normal stop position. (c) Stop behind the normal stop position. The dot locates the point on the curve where the origin is located for case (b).



FIGURE 40.3 Accent lighting.



FIGURE 40.4 Wall sconces for providing ambient lighting and the much needed vertical illumination in various situations.

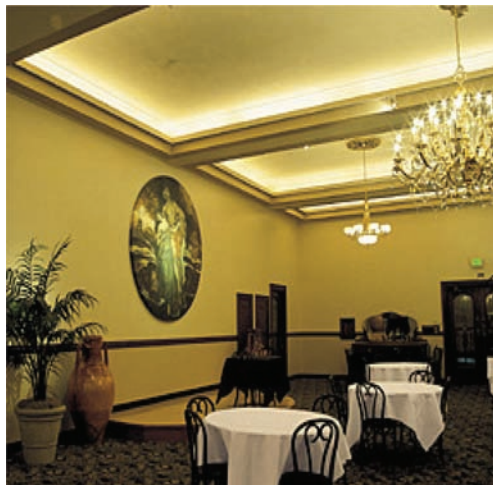


FIGURE 40.5 Indirect lighting with cove lighting in a restaurant using light strips. The chandelier provides the decorative lighting without significantly contributing to any other lighting function.

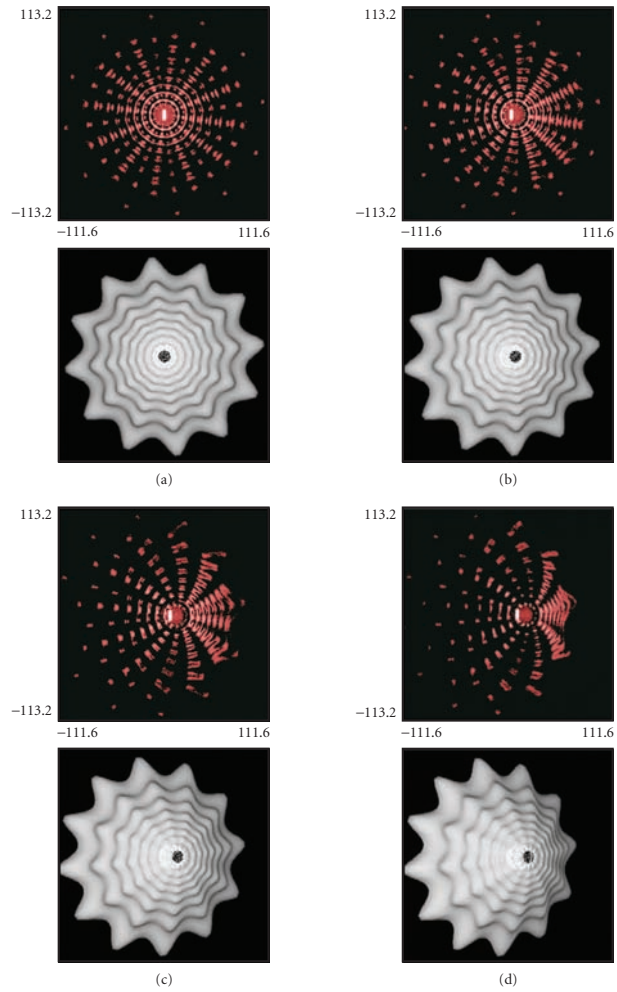


FIGURE 40.8 Views of the lit appearance (upper) of a star-shaped taillight (lower) at four horizontal angles of (a) 0°; (b) 10°; (c) 20°; and (d) 30°.

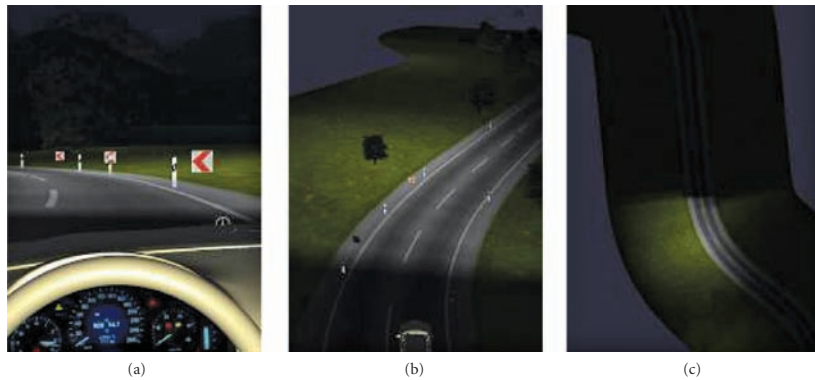


FIGURE 40.9 Three perspectives of lit-scene renderings from a low-beam headlamp: (a) driver's view; (b) 20 m above and behind automobile; and (c) bird's eye view.

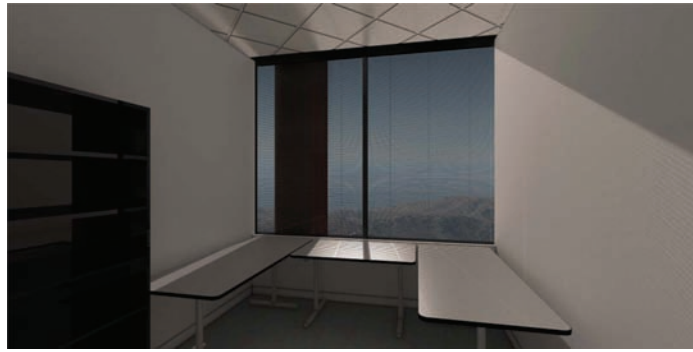


FIGURE 40.10a Rendering of a lit office room.

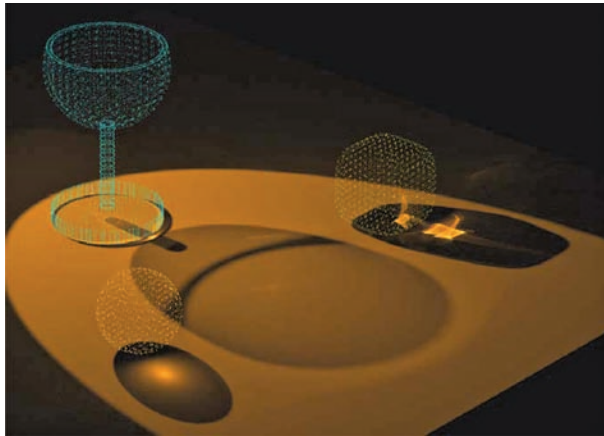


FIGURE 40.10b Rendering of a lit desk with three objects located on it (wine glass, ice cube, and crystal ball) to show both diffuse and specular effects.

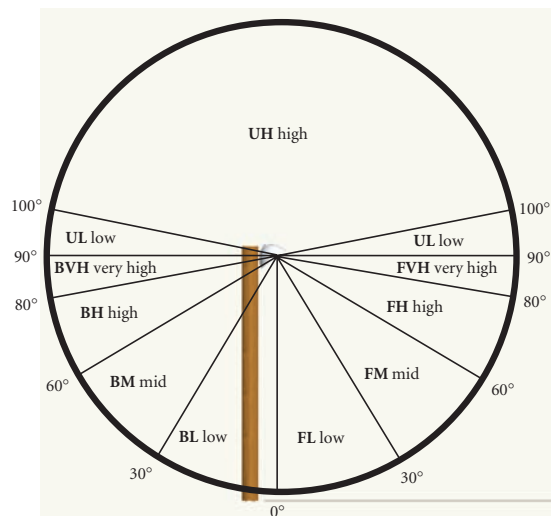


FIGURE 40.22 Layout of the light classification system subzones.

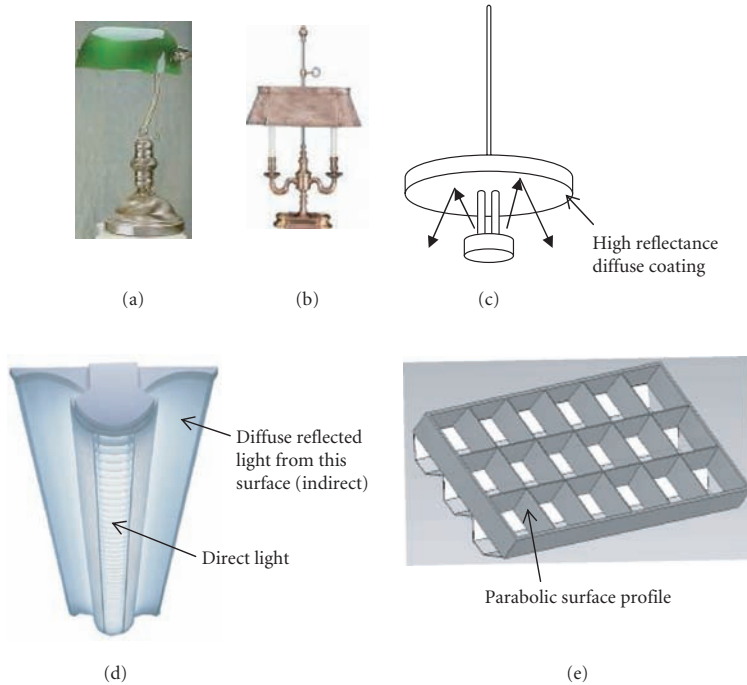


FIGURE 40.24 Depictions of luminaires: (a) Bankers lamp: multiple bounces inside the reflector create a wide angled uniform illumination; (b) Bouillotte lamp: vertical fluorescent tubes provide diffuse illumination; (c) indirect lighting with RLM fixture where the top surface reflects light into a wide angular range; (d) overhead direct-indirect lighting fixture using fluorescent tubular bulbs; and (e) parabolic louvered trough reflector for fluorescent tubes.



FIGURE 40.26 A conference room with artificial skylight made up of backlit ceiling image tiles.

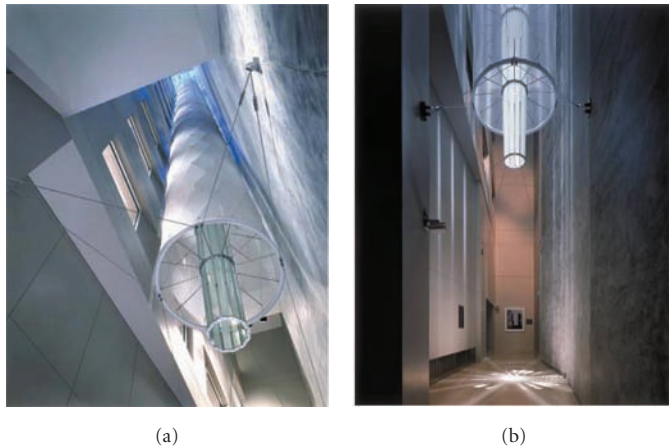


FIGURE 40.31 A Solar light pipe. (a) A 140-ft-tall light gathering and distributing device that presents daylight down into the core of a building that has no other access to daylight. (b) Light projected (10-in diameter) on the floor.

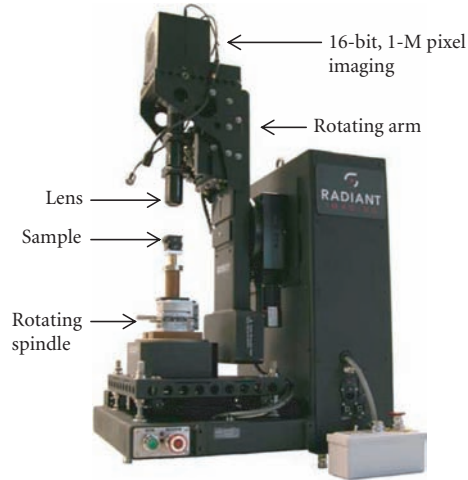


FIGURE 40.33 Photograph of a source measurement goniometer that is used to ascertain the luminance distribution of the source. The system wobble (electro-mechanical-software runout) is $15\ \mu\text{m}$ to allow for measuring small light sources like LED die.

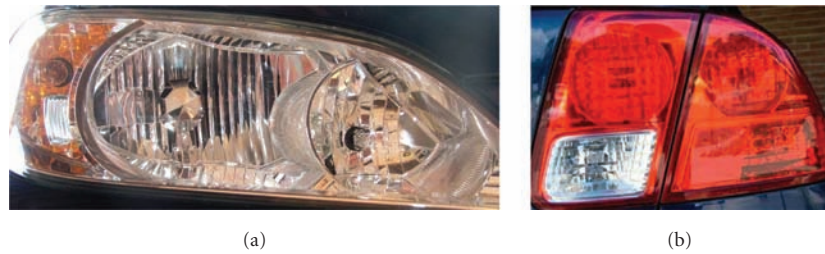


FIGURE 40.34 (a) A faceted headlamp including high-beam (right), low-beam (middle), and turn signal (left) luminaire. Note the yellowish tinge of the turn signal, which is due to the coating placed on the bulb used therein. (b) A faceted taillight including the following functions: tail (upper left), stop (upper right), turn signal (lower right), reflex reflector (lower middle), and backup (lower left).

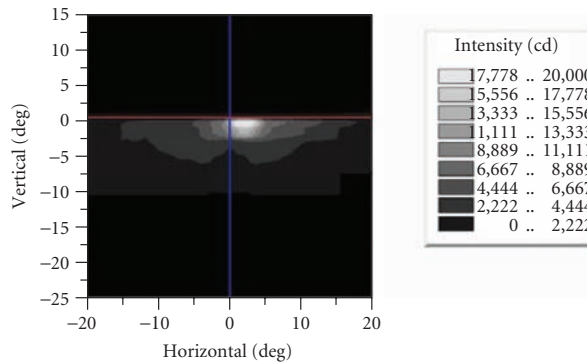


FIGURE 40.35 Luminous intensity (cd) distribution for the SAE low-beam requirements of Table 18.

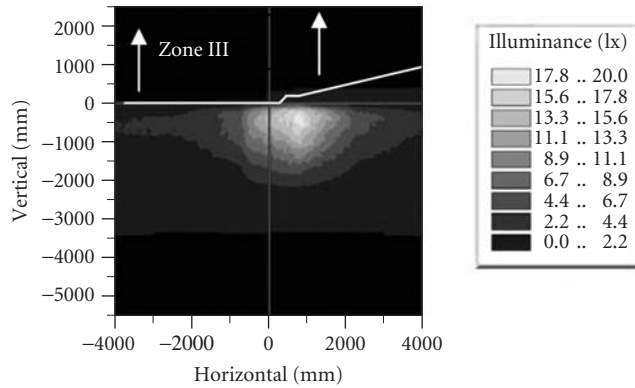


FIGURE 40.36 Illuminance (lx) distribution for the ECE passing/low-beam requirements of Table 19.

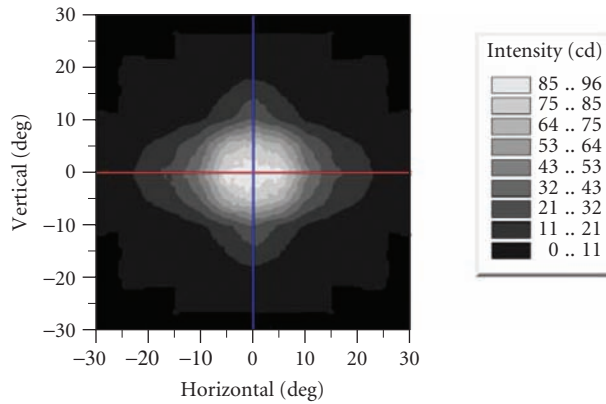


FIGURE 40.37 Luminous intensity (cd) distribution for the SAE stop lamp requirements of Table 20 (1 lit section).

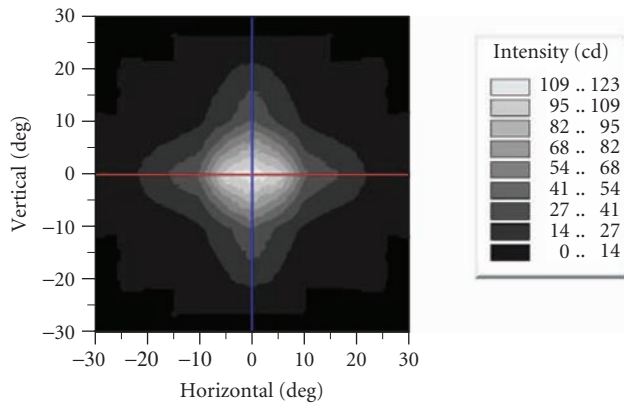


FIGURE 40.38 Luminous intensity (cd) distribution for the R7 stop lamp requirements of Table 21 (1 lamp illumination level).