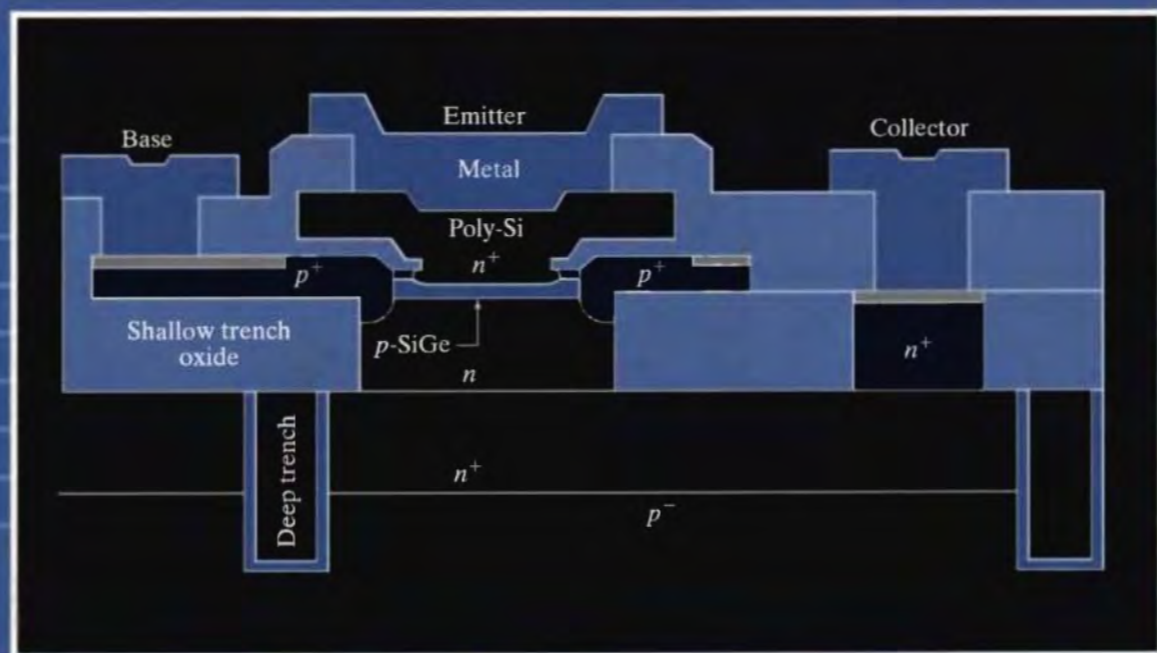


VOLUME V

INTRODUCTION TO MICROELECTRONIC FABRICATION

SECOND EDITION

RICHARD C. JAEGER



Modular Series on Solid State Devices

Gerold W. Neudeck • Robert F. Pierret, Series Editors

PHYSICAL CONSTANTS

Symbol	Name	Value
q	Magnitude of electronic charge	$1.602 \times 10^{-19} \text{C}$
m_0	Electron rest mass	$9.109 \times 10^{-31} \text{kg}$
m_p	Proton rest mass	$1.673 \times 10^{-27} \text{kg}$
c	Speed of light in vacuum	$2.998 \times 10^8 \text{m/s}$
ϵ_0	Permittivity of vacuum	$8.854 \times 10^{-12} \text{F/m}$
k	Boltzmann's constant	$1.381 \times 10^{-23} \text{J/K}$ $8.617 \times 10^{-5} \text{eV/K}$
\hbar	Planck's constant	$6.625 \times 10^{-34} \text{J-s}$ $4.135 \times 10^{-15} \text{eV-s}$
A_0	Avogadro number	$6.022 \times 10^{26} \text{molecules/kg-mole}$
kT	Thermal energy	$0.02586 \text{ eV } (T = 27^\circ \text{C})$ $0.02526 \text{ eV } (T = 20^\circ \text{C})$
E_g	Bandgap of silicon at 300K	1.12 eV
K_s	Relative permittivity of silicon	11.7
K_0	Relative permittivity of silicon dioxide	3.9
n_i	Intrinsic carrier density in silicon at 300K	$10^{10}/\text{cm}^3$

CONVERSION FACTORS

1 \AA	$= 10^{-8} \text{ cm}$	1 mil^2	$= 645.2 \text{ } \mu\text{m}^2$
	$= 10^{-10} \text{ m}$		$= 6.45 \times 10^{-6} \text{ cm}^2$
$1 \text{ } \mu\text{m}$	$= 10^{-4} \text{ cm}$	1 eV	$= 1.602 \times 10^{-19} \text{ J}$
	$= 10^{-6} \text{ m}$	λ	$= 1.24/E \text{ } \mu\text{m} \text{ (E in eV)}$
1 mil	$= 10^{-3} \text{ in}$		
	$= 25.4 \text{ } \mu\text{m}$		

MODULAR SERIES ON SOLID STATE DEVICES

Gerold W. Neudeck and Robert F. Pierret, Editors

Volume V

Introduction to Microelectronic Fabrication

Second Edition

Richard C. Jaeger

Auburn University



Prentice Hall
Upper Saddle River, New Jersey 07458

Library of Congress Cataloging-in-Publication Data

Jaeger, Richard C.

Introduction to microelectronic fabrication / Richard C. Jaeger—2nd Edition
p. cm.

(Modular series on solid state devices; v. 5)

Includes bibliographical references and index.

ISBN 0-201-44494-7

1. Integrated circuits—Very large scale integration—Design and construction—Congresses. I. Title. II. Series.

CIP Data available.

Vice President and Editorial Director, ECS: *Marcia J. Horton*

Publisher: *Tom Robbins*

Associate Editor: *Alice Dworkin*

Editorial Assistant: *Jody McDonnell*

Vice President and Director of Production and Manufacturing, ESM: *David W. Riccardi*

Executive Managing Editor: *Vince O'Brien*

Managing Editor: *David A. George*

Production Editor: *Irwin Zucker*

Director of Creative Services: *Paul Belfanti*

Manager of Electronic Composition and Digital Content: *Jim Sullivan*

Electronic Composition: *William Johnson*

Creative Director: *Carole Anson*

Art Director: *Jayne Conte*

Art Editor: *Gregory Dulles*

Manufacturing Manager: *Trudy Piscioti*

Manufacturing Buyer: *Lisa McDowell*

Marketing Manager: *Holly Stark*

Marketing Assistant: *Karen Moon*



© 2002, 1998 by Prentice Hall

Published by Prentice-Hall, Inc.

Upper Saddle River, New Jersey 07458

All rights reserved. No part of this book may be reproduced in any format or by any means, without permission in writing from the publisher.

The author and publisher of this book have used their best efforts in preparing this book. These efforts include the development, research, and testing of the theories and programs to determine their effectiveness. The author and publisher make no warranty of any kind, expressed or implied, with regard to these programs or the documentation contained in this book. The author and publisher shall not be liable in any event for incidental or consequential damages in connection with, or arising out of, the furnishing, performance, or use of these programs.

Printed in the United States of America

10 9 8 7 6 5 4 3 2

ISBN 0-201-44494-7

Pearson Education Ltd., *London*

Pearson Education Australia Pty. Ltd., *Sydney*

Pearson Education Singapore, Pte. Ltd.

Pearson Education North Asia Ltd., *Hong Kong*

Pearson Education Canada Inc., *Toronto*

Pearson Educación de México, S.A. de C.V.

Pearson Education—Japan, *Tokyo*

Pearson Education Malaysia, Pte. Ltd.

To My Family—Joan, Peter, and Stephanie

Also a special thanks to
ThePirateBay.se. I love you
guys so much for sharing my
book for free!!! Love Jesus.

Contents

PREFACE

xiii

Chapter 1 An Overview of Microelectronic Fabrication

1

- 1.1 A Historical Perspective 1
- 1.2 An Overview of Monolithic Fabrication Processes and Structures 5
- 1.3 Metal-Oxide-Semiconductor (MOS) Processes 7
 - 1.3.1 Basic NMOS Process 7
 - 1.3.2 Basic Complementary MOS (CMOS) Process 9
- 1.4 Basic Bipolar Processing 10
- 1.5 Safety
 - References 14
 - Problems 14

Chapter 2 Lithography

17

- 2.1 The Photolithographic Process 17
 - 2.1.1 Wafers and Wafer Cleaning 19
 - 2.1.2 Barrier Layer Formation 21
 - 2.1.3 Photoresist Application 21
 - 2.1.4 Soft Baking / Prebaking 22
 - 2.1.5 Mask Alignment 23
 - 2.1.6 Photoresist Exposure and Development 23
 - 2.1.7 Hard Baking 25
- 2.2 Etching Techniques 25
 - 2.2.1 Wet Chemical Etching 25
 - 2.2.2 Dry Etching Plasma Systems 26
 - 2.2.3 Photoresist Removal 26
 - 2.2.4 Metrology and Critical Dimension Control 28
- 2.3 Photomask Fabrication 28
- 2.4 Exposure Systems 28
- 2.5 Exposure Sources 34
- 2.6 Optical and Electron Microscopy 37
 - 2.6.1 Optical Microscopy 37
 - 2.6.2 Scanning Electron Microscopy 37
 - 2.6.3 Transmission Electron Microscopy 38
- 2.7 Summary 38
 - References 40
 - Further Reading 40
 - Problems 40

Chapter 3	Thermal Oxidation of Silicon	43
3.1	The Oxidation Process	43
3.2	Modeling Oxidation	44
3.3	Factors Influencing Oxidation Rate	46
3.4	Dopant Redistribution During Oxidation	51
3.5	Masking Properties of Silicon Dioxide	51
3.6	Technology of Oxidation	52
3.7	Oxide Quality	53
3.8	Selective Oxidation and Shallow Trench Formation	55
3.8.1	Trench Isolation	56
3.8.2	Chemical Mechanical Polishing (CMP)	57
3.9	Oxide Thickness Characterization	61
3.10	Process Simulation	61
	Summary	61
	References	63
	Problems	64
Chapter 4	Diffusion	67
4.1	The Diffusion Process	67
4.2	Mathematical Model for Diffusion	68
4.2.1	Constant-Source Diffusion	69
4.2.2	Limited-Source Diffusion	70
4.2.3	Two-Step Diffusion	71
4.3	The Diffusion Coefficient	72
4.4	Successive Diffusions	74
4.5	Solid-Solubility Limits	74
4.6	Junction Formation and Characterization	76
4.6.1	Vertical Diffusion and Junction Formation	76
4.6.2	Lateral Diffusion	78
4.6.3	Concentration-Dependent Diffusion	79
4.7	Sheet Resistance	81
4.7.1	Sheet-Resistance Definition	82
4.7.2	Irvin's Curves	85
4.7.3	The Four-Point Probe	88
4.7.4	Van der Pauw's Method	88
4.8	Generation-Depth and Impurity Profile Measurement	90
4.8.1	Grove-and-Stain and Angle-Lap Methods	90
4.8.2	Impurity-Profile Measurement	91
4.9	Diffusion Simulation	93
4.10	Diffusion Systems	95
4.10.1	Boron Diffusion	97
4.10.2	Phosphorus Diffusion	98
4.10.3	Arsenic Diffusion	99
4.10.4	Antimony Diffusion	100

4.11	Gettering	100	
	Summary	101	
	References	102	
	Problems	103	
Chapter 5	Ion Implantation		109
5.1	Implantation Technology	109	
5.2	Mathematical Model for Ion Implantation	111	
5.3	Selective Implantation	114	
5.4	Junction Depth and Sheet Resistance	117	
5.5	Channeling, Lattice Damage, and Annealing	118	
	5.5.1 Channeling	118	
	5.5.2 Lattice Damage and Annealing	120	
	5.5.3 Deviations from the Gaussian Theory	121	
5.6	Shallow Implantations	121	
	5.6.1 Low-Energy Implantation	122	
	5.6.2 Rapid Thermal Annealing	123	
	5.6.3 Transient Enhanced Diffusion (TED)	123	
	Summary	124	
	References	125	
	Source Listing	126	
	Problems	126	
Chapter 6	Film Deposition		129
6.1	Evaporation	129	
	6.1.1 Kinetic Gas Theory	130	
	6.1.2 Filament Evaporation	132	
	6.1.3 Electron-Beam Evaporation	132	
	6.1.4 Flash Evaporation	134	
	6.1.5 Shadowing and Step Coverage	134	
6.2	Sputtering	135	
6.3	Chemical Vapor Deposition	136	
	6.3.1 CVD Reactors	137	
	6.3.2 Polysilicon Deposition	138	
	6.3.3 Silicon Dioxide Deposition	139	
	6.3.4 Silicon Nitride Deposition	140	
	6.3.5 CVD Metal Deposition	141	
6.4	Epitaxy	141	
	6.4.1 Vapor-Phase Epitaxy	142	
	6.4.2 Doping of Epitaxial Layers	145	
	6.4.3 Buried Layers	145	
	6.4.4 Liquid-Phase and Molecular-Beam Epitaxy	148	
	Summary	148	
	References	149	
	Further Reading	149	
	Problems	149	

Chapter 7	Interconnections and Contacts	151
7.1	Interconnections in Integrated Circuits	151
7.2	Metal Interconnections and Contact Technology	153
7.2.1	Ohmic Contact Formation	153
7.2.2	Aluminum-Silicon Eutectic Behavior	154
7.2.3	Aluminum Spiking and Junction Penetration	155
7.2.4	Contact Resistance	156
7.2.5	Electromigration	157
7.3	Diffused Interconnections	158
7.4	Polysilicon Interconnections and Buried Contacts	159
7.4.1	Buried Contacts	160
7.4.2	Butted Contacts	162
7.5	Silicides and Multilayer-Contact Technology	162
7.5.1	Silicides, Polycides, and Salicides	162
7.5.2	Barrier Metals and Multilayer Contacts	164
7.6	The Liftoff Process	164
7.7	Multilevel Metallization	166
7.7.1	Basic Multilevel Metallization	166
7.7.2	Planarized Metallization	167
7.7.3	Low Dielectric Constant Interlevel Dielectrics	167
7.8	Copper Interconnects and Damascene Processes	168
7.8.1	Electroplated Copper Interconnect	168
7.8.2	Damascene Plating	168
7.8.3	Dual Damascene structures	169
	Summary	172
	References	172
	Further Reading	173
	Problems	174
Chapter 8	Packaging and Yield	177
8.1	Testing	177
8.2	Wafer Thinning and Die Separation	178
8.3	Die Attachment	178
8.3.1	Epoxy Die Attachment	179
8.3.2	Eutectic Die Attachment	179
8.4	Wire Bonding	179
8.4.1	Thermocompression Bonding	182
8.4.2	Ultrasonic Bonding	183
8.4.3	Thermosonic Bonding	184
8.5	Packages	184
8.5.1	Circular TO-Style Packages	184
8.5.2	Dual-in-Line Packages (DIPs)	184

8.5.3	Pin-Grid Arrays (PGAs)	185
8.5.4	Leadless Chip Carriers (LCCs)	186
8.5.5	Packages for Surface Mounting	186
8.6	Flip-Chip and Tape-Automated-Bonding Processes	187
8.6.1	Flip-Chip Technology	188
8.6.2	Ball Grid Array (BGA)	190
8.6.3	The Tape-Automated-Bonding (TAB) Process	191
8.6.4	Chip Scale Packages	193
8.7	Yield	194
8.7.1	Uniform Defect Densities	194
8.7.2	Nonuniform Defect Densities	195
	Summary	198
	References	198
	Further Reading	199
	Problems	199

Chapter 9 MOS Process Integration

201

9.1	Basic MOS Device Considerations	201
9.1.1	Gate-Oxide Thickness	202
9.1.2	Substrate Doping and Threshold Voltage	203
9.1.3	Junction Breakdown	204
9.1.4	Punch-through	204
9.1.5	Junction Capacitance	205
9.1.6	Threshold Adjustment	206
9.1.7	Field-Region Considerations	208
9.1.8	MOS Transistor Isolation	208
9.1.9	Lightly Doped Drain structures	210
9.1.10	MOS Transistor Scaling	210
9.2	MOS Transistor Layout and Design Rules	212
9.2.1	Metal-Gate Transistor Layout	213
9.2.2	Polysilicon-Gate Transistor Layout	217
9.2.3	More-Aggressive Design Rules	218
9.2.4	Channel Length and Width Biases	219
9.3	Complementary MOS (CMOS) Technology	221
9.3.1	<i>n</i> -Well Process	221
9.3.2	<i>p</i> -Well and Twin Well Processes	221
9.3.3	Gate Doping	222
9.3.4	CMOS Isolation	224
9.3.5	CMOS Latchup	225
9.3.6	Shallow Trench Isolation	225
9.4	Silicon on Insulator	226
	Summary	227
	References	228
	Problems	229

Chapter 10	Bipolar Process Integration	233
10.1	The Junction-Isolated Structure	233
10.2	Current Gain	235
10.3	Transit Time	236
10.4	Basewidth	237
10.5	Breakdown Voltages	239
10.5.1	Emitter-Base Breakdown Voltage	239
10.5.2	Circular Emitters	239
10.5.3	Collector-Base Breakdown Voltage	240
10.6	Other Elements In SBC Technology	242
10.6.1	Emitter Resistor	243
10.6.2	Base Resistor	244
10.6.3	Epitaxial Layer Resistor	245
10.6.4	Pinch Resistor	246
10.6.5	Substrate <i>pnp</i> Transistor	246
10.6.6	Lateral <i>pnp</i> Transistors	248
10.6.7	Schottky Diodes	249
10.7	Layout Considerations	249
10.7.1	Buried-Layer and Isolation Diffusions	249
10.7.2	Base Diffusion to Isolation Diffusion Spacing	251
10.7.3	Emitter-Diffusion Design Rules	252
10.7.4	A Layout Example	252
10.8	Advanced Bipolar Structures	253
10.8.1	Locos Isolated Self-Aligned Contact Structure	254
10.8.2	Dual Polysilicon Self-Aligned Process	254
10.8.3	The Silicon Germanium Epitaxial Base Transistor	257
10.9	Other Bipolar Isolation Techniques	259
10.9.1	Collector-Diffusion Isolation (CDI)	259
10.9.2	Dielectric Isolation	259
10.10	BICMOS	262
	Summary	263
	References	264
	Problems	265
Chapter 11	Processes for MicroElectroMechanical Systems: MEMS	269
11.1	Mechanical Properties of Silicon	270
11.2	Bulk Micromachining	271
11.2.1	Isotropic and Anisotropic Etching	271
11.2.2	Diaphragm Formation	273
11.2.3	Cantilever Beams and Released Structures	275
11.3	Silicon Etchants	277
11.3.1	Isotropic Etching	277
11.3.2	Anisotropic Etching	278

11.4	Surface Micromachining	279
11.4.1	Cantilever Beams, Bridges and Sealed Cavities	279
11.4.2	Movable In-Plane Structures	279
11.4.3	Out-of-Plane Motion	282
11.4.4	Release Problems	286
11.5	High-Aspect-Ratio Micromachining: The LIGA Molding Process	288
11.6	Silicon Wafer Bonding	289
11.6.1	Adhesive Bonding	289
11.6.2	Silicon Fusion Bonding	289
11.6.3	Anodic Bonding	291
11.7	IC Process Compatibility	292
11.7.1	Preprocessing	292
11.7.2	Postprocessing	292
11.7.3	Merged Processes	294
	Summary	295
	References	296
	Problems	298

ANSWERS TO SELECTED PROBLEMS

301

INDEX

303

Preface

The spectacular advances in the development and application of integrated circuit (IC) technology have led to the emergence of microelectronics process engineering as an independent discipline. Additionally, the pervasive use of integrated circuits requires a broad range of engineers in the electronics and allied industries to have a basic understanding of the behavior and limitations of ICs. One of the goals of this book is to address the educational needs of individuals with a wide range of backgrounds.

This text presents an introduction to the basic processes common to most IC technologies and provides a base for understanding more advanced processing and design courses. In order to contain the scope of the material, we deal only with material related to silicon processing and packaging. The details of many problems specifically related to VLSI/ULSI fabrication are left to texts on advanced processing, although problem areas are mentioned at various points in this text, and goals of the International Technology Roadmap for Semiconductors are discussed as appropriate.

Chapter 1 provides an overview of IC processes, and Chapters 2–6 then focus on the basic steps used in fabrication, including lithography, oxidation, diffusion, ion implantation and thin film deposition, and etching. Interconnection technology, packaging, and yield are covered in Chapters 7 and 8. It is important to understand interactions between process design, device design, and device layout. For this reason, Chapter 9 and 10 on MOS and bipolar process integration have been included. Chapter 11 provides a brief introduction to the exciting area of Microelectromechanical Systems (MEMS).

Major changes in the second edition of this text include new or expanded coverage of lithography and exposure systems, trench isolation, chemical mechanical polishing, shallow junctions, transient-enhanced diffusion, copper Damascene processes, and process simulation. The chapters on MOS and bipolar process integration have been substantially modified, and the chapter on MEMS is entirely new. The problem sets have been expanded, and additional information on measurement techniques has been included.

The text evolved from notes originally developed for a course introducing seniors and beginning graduate students to the fabrication of solid-state devices and integrated circuits. A basic knowledge of the material properties of silicon is needed, and we use Volume I of this Series as a companion text. An introductory knowledge of electronic components such as resistors, diodes, and MOS and bipolar transistors is also useful.

The material in the book is designed to be covered in one semester. In our case, the microelectronics fabrication course is accompanied by a corequisite laboratory. The students design a simple device or circuit based upon their individual capability, and the designs are combined on a multiproject polysilicon gate NMOS chip. Design, fabrication, and testing are completed within the semester. Students from a variety of disciplines, including electrical, mechanical, chemical, and materials engineering; computer science; and physics, are routinely enrolled in the fabrication classes.

Before closing, I must recognize a number of other books that have influenced the preparation of this text. These include *The Theory and Practice of Microelectronics* and *VLSI Fabrication Principles* by S. K. Ghandi, *Basic Integrated Circuit Engineering* by D. J. Hamilton and W. G. Howard, *Integrated Circuit Engineering* by A. H. Glaser and G. E. Subak-Sharpe, *Microelectronic Processing and Device Design* by R. A. Colclaser, *Semiconductor Devices—Physics and Technology* by S. M. Sze, *Semiconductor Integrated Circuit Processing Technology* by W. R. Runyon and K. E. Bean, and *The Science and Engineering of Microelectronic Fabrication* by Stephen A. Campbell.

Thanks also go to the many colleagues who have provided suggestions and encouragement for the new edition and especially to our laboratory manager Charles Ellis who has been instrumental in molding the laboratory sections of our course.

RICHARD C. JAEGER
Auburn, Alabama

CHAPTER 1

An Overview of Microelectronic Fabrication

1.1 A HISTORICAL PERSPECTIVE

In this volume, we will develop an understanding of the basic processes used in monolithic integrated-circuit (IC) fabrication. Silicon is the dominant material used throughout the IC industry today, and in order to conserve space, only silicon processing will be discussed in this book. However, all of the basic processes discussed here are applicable to the fabrication of compound semiconductor integrated circuits (ICs) such as gallium arsenide or indium phosphide, as well as thick- and thin-film hybrid ICs.

Germanium was one of the first materials to receive wide attention for use in semiconductor device fabrication, but it was rapidly replaced by silicon during the early 1960s. Silicon emerged as the dominant material, because it was found to have two major processing advantages. Silicon can easily be oxidized to form a high-quality electrical insulator, and this oxide layer also provides an excellent barrier layer for the selective diffusion steps needed in integrated-circuit fabrication.

Silicon was also shown to have a number of ancillary advantages. It is a very abundant element in nature, providing the possibility of a low-cost starting material. It has a wider bandgap than germanium and can therefore operate at higher temperatures than germanium. In retrospect, it appears that the processing advantages were the dominant reasons for the emergence of silicon over other semiconductor materials.

The first successful fabrication techniques produced single transistors on a rectangular silicon die 1–2 mm on a side. The first integrated circuits, fabricated at Texas Instruments and Fairchild Semiconductor in the early 1960s, included several transistors and resistors to make simple logic gates and amplifier circuits. From this modest beginning, the level of integration has been doubling every one to two years, and we have now reached integration levels of billions of components on a 20-mm × 20-mm die [1–3]. For example, one-gigabit dynamic random-access memory (DRAM) chips have more than 10^9 transistors and more than 10^9 capacitors in the memory array, as

well as millions of additional transistors in the access and decoding circuitry. One-gigabit RAMs are currently being produced with photographic features measuring between 0.13 and 0.18 micron (μm). MOS transistors with dimensions below $0.05 \mu\text{m}$ have been fabricated successfully in research laboratories, and these devices continue to behave as predicted by macroscopic models. So we still have significant increases in integrated-circuit density yet to come, provided that manufacturable fabrication processes can be developed for deep submicron dimensions.

The larger the diameter of the wafer, the more integrated-circuit dice can be produced at one time. Many wafers are processed at the same time, and the same silicon chip is replicated as many times as possible on a wafer of a given size. The size of silicon wafers has steadily increased from 1-, 2-, 3-, 4-, 5-, and 6-in. diameters to the point where 8-in. (200-mm) wafers are now in production. (See Fig. 1.1(a).) Wafers with 300-mm diameters will be in full production in the near future, and 450 mm wafers are projected to be in use by the end of the decade. Wafer thicknesses range from approximately 350 to 1250 microns. Large-diameter wafers must be thicker in order to maintain structural integrity and planarity during the wide range of processing steps encountered during IC fabrication.

Figure 1.1(c) shows the approximate number of 10-mm \times 10-mm dice that fit on a wafer of given diameter. For a given wafer processing cost, the more dice per wafer, the lower the individual die cost becomes. Thus, there are strong economic forces driving the IC industry to continually move to larger and larger wafer sizes.

The dramatic progress of IC miniaturization is depicted graphically in Fig. 1.2 [1–3] on pages 4 and 5. The complexities of memory chips and microprocessors have both grown exponentially with time. In the three decades since 1965, memory density has grown by a factor of more than 10 million from the 64-bit chip to the 1-Gb memory chip, as indicated in Fig. 1.2(a). Similarly, the number of transistors on a microprocessor chip has increased by a factor of more than five thousand since 1970 (Fig. 1.2 (b).)

Since the commercial introduction of the integrated circuit, these increases in density have been achieved through a continued reduction in the minimum line width, or minimum feature size, that can be defined on the surface of the integrated circuit, as shown in Fig. 1.3 on page 6. Today, most corporate semiconductor laboratories around the world are actively working on deep submicron processes with feature sizes less than $0.1 \mu\text{m}$, less than one one-thousandth the diameter of a human hair!

These trends and future projections are summarized in Table 1.1 on page 6, which is abstracted from the International Technology Road map for Semiconductors (ITRS) generated by the Semiconductor Industry Association [4]. The ITRS is updated every three years; the projections are mind-boggling, even for those of us who have worked in the industry for many years. By the year 2011, MOS transistor gate lengths are projected to reach 30 nm ($0.030 \mu\text{m}$), multigigabit DRAM chips will be commonplace, and microprocessors will have a billion transistors on die exceeding 25 mm (one inch) on an edge. It remains to be seen whether the industry meets these projections. However, progress will be impressive, even if only a fraction of the projections are achieved.*

Historically there has been a problem with the units of measure used to describe integrated circuits. Horizontal dimensions were originally specified in mils (1 mil = 0.001 in.), whereas specification of the shallower vertical dimensions commonly made

*In the past several years, the IC industry has actually managed to exceed the ITRS goals.

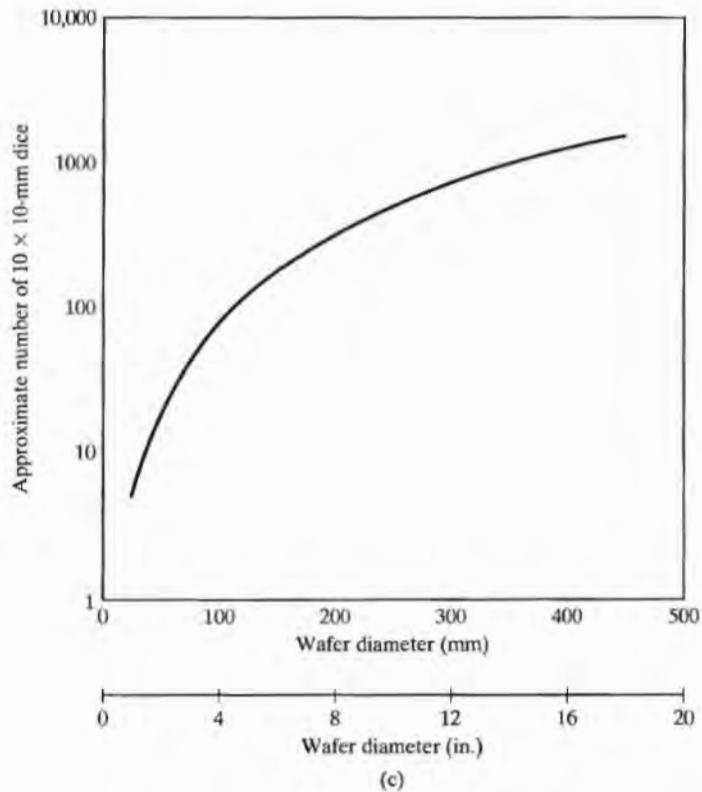
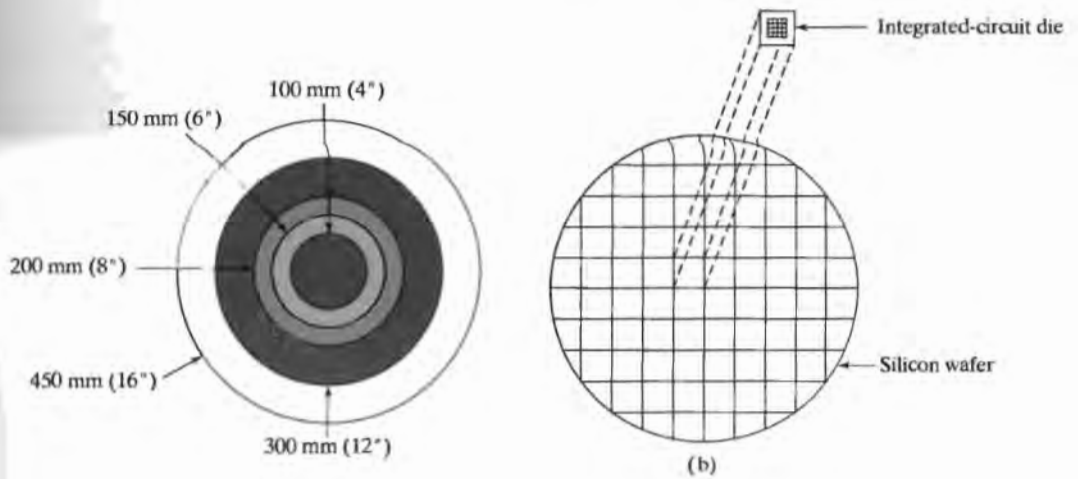


FIGURE 1.1

(a) Relative size of wafers with diameters ranging from 100 to 450 mm; (b) The same integrated circuit die is replicated hundreds of times on a typical silicon wafer; (c) the graph gives the approximate number of 10×10 mm dice that can be fabricated on wafers of different diameters.

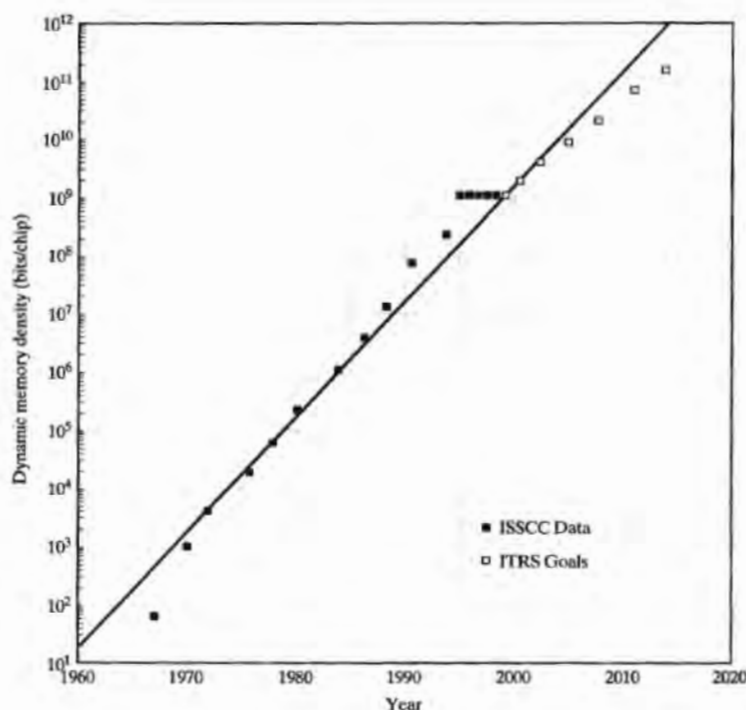


FIGURE 1.2

(a) Dynamic memory density versus year since 1960.

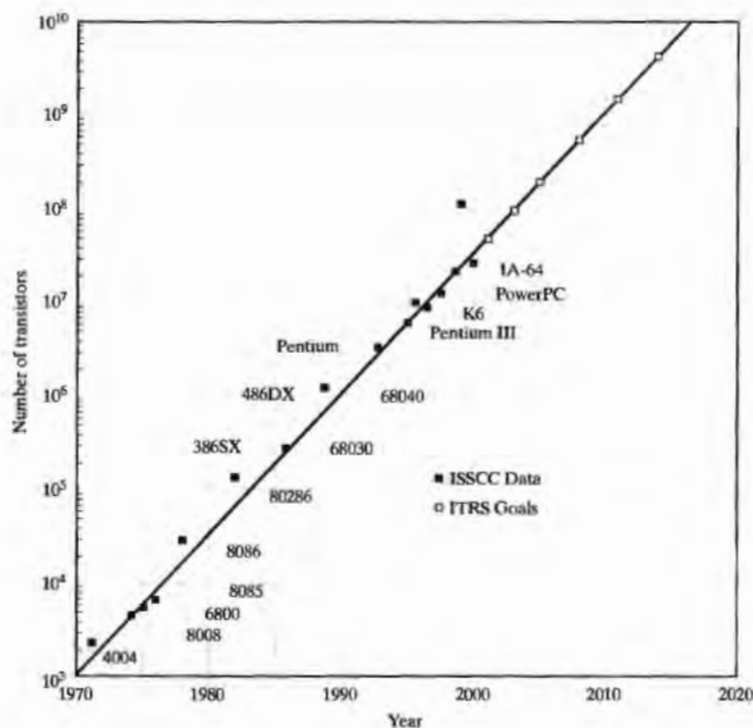


FIGURE 1.2

(b) Number of transistors in a microprocessor versus year.

use of the metric system. Today, most of the dimensions are specified using the metric system, although Imperial units are occasionally still used. Throughout the rest of this book, we will attempt to make consistent use of metric units.

1.2 AN OVERVIEW OF MONOLITHIC FABRICATION PROCESSES AND STRUCTURES

Monolithic IC fabrication can be illustrated by studying the basic cross sections of MOS and bipolar transistors in Figs. 1.4 (on page 7) and 1.5 (on page 8). The n -channel MOS transistor is formed in a p -type substrate. Source/drain regions are formed by selectively converting shallow regions at the surface to n -type material. Thin and thick silicon-dioxide regions on the surface form the gate insulator of the transistor and serve to isolate one device from another. A thin film of polysilicon is used to form the gate of the transistor, and a metal such as aluminum is used to make contact to the source and drain. Interconnections between devices can be made using the diffusions and the layers of polysilicon and metal.

The bipolar transistor in Fig. 1.5 has alternating n - and p -type regions selectively fabricated on a p -type substrate. Silicon dioxide is again used as an insulator, and aluminum is used to make electrical contact to the emitter, base, and collector of the transistor.

Both the MOS and bipolar structures are fabricated through the repeated application of a number of basic processing steps:

- Oxidation
- Photolithography
- Etching
- Diffusion
- Evaporation or sputtering
- Chemical vapor deposition (CVD)
- Ion implantation
- Epitaxy
- Annealing

Silicon dioxide can be formed by heating a silicon wafer to a high temperature (1000 to 1200 °C) in the presence of oxygen. This process is called *oxidation*. Metal films can be deposited through evaporation by heating the metal to its melting point in a vacuum. Thin films of silicon nitride, silicon dioxide, polysilicon, and metals can all be formed through a process known as *chemical vapor deposition* (CVD), in which the material is deposited out of a gaseous mixture onto the surface of the wafer. Metals and insulators may also be deposited by a process called *sputtering*.

Shallow n - and p -type layers are formed by high-temperature (1000 to 1200 °C) *diffusion* of donor or acceptor impurities into silicon or by *ion implantation*, in which the wafer is bombarded with high-energy donor or acceptor ions generated in a high-voltage particle accelerator.

In order to build devices and circuits, the n - and p -type regions must be formed selectively in the surface of the wafer. Silicon dioxide, silicon nitride, polysilicon, photo resist, and other materials can all be used to mask areas of the wafer surface to prevent

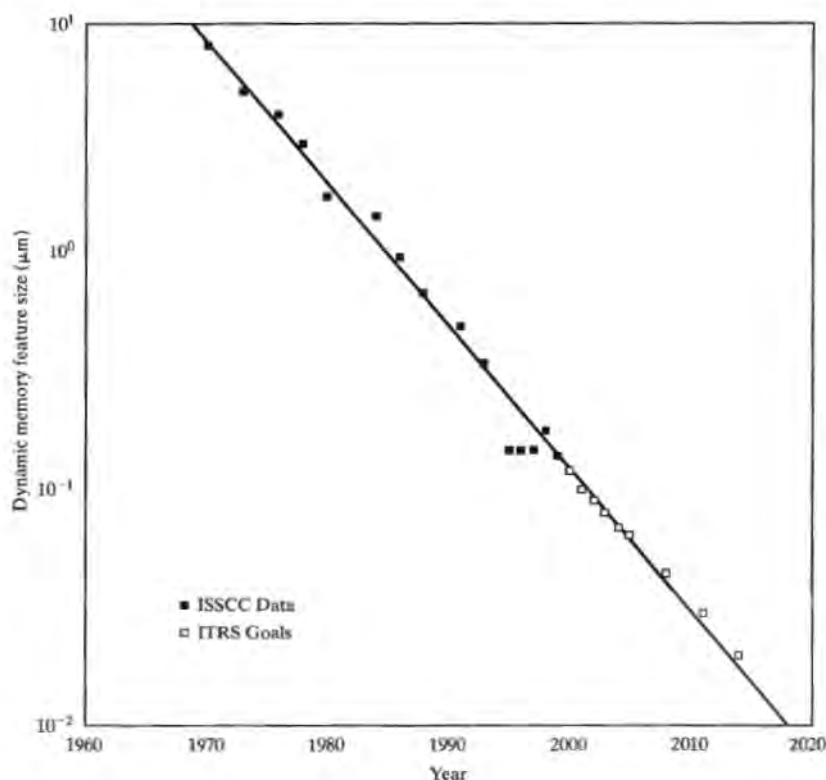


FIGURE 1.3

Feature size used in fabrication of dynamic memory as a function of time.

TABLE 1.1 International Technology Road Map for Semiconductors (ITRS) [4]

Year of First Product Shipment	Selected Projections					
	2001	2003	2005	2008	2011	2014
DRAM Metal Line Half-Pitch (nm)	150	120	100	70	50	35
Microprocessor Gate Widths (nm)	100	80	65	45	30	20
DRAM (G-bits/chip)	2.2	4.3	8.6	24	68	190
Microprocessor (M-transistors/chip)	48	95	190	540	1500	4300
DRAM Chip Area: Year of Introduction (mm ²)	400	480	526	600	690	790
DRAM Chip Area: Production (mm ²)	130	160	170	200	230	260
MPU Chip Size at Introduction (mm ²)	340	370	400	470	540	620
MPU Chip Area: Second “shrink” (mm ²)	180	210	230	270	310	350
Wafer Size (mm)	300	300	300	450	450	450

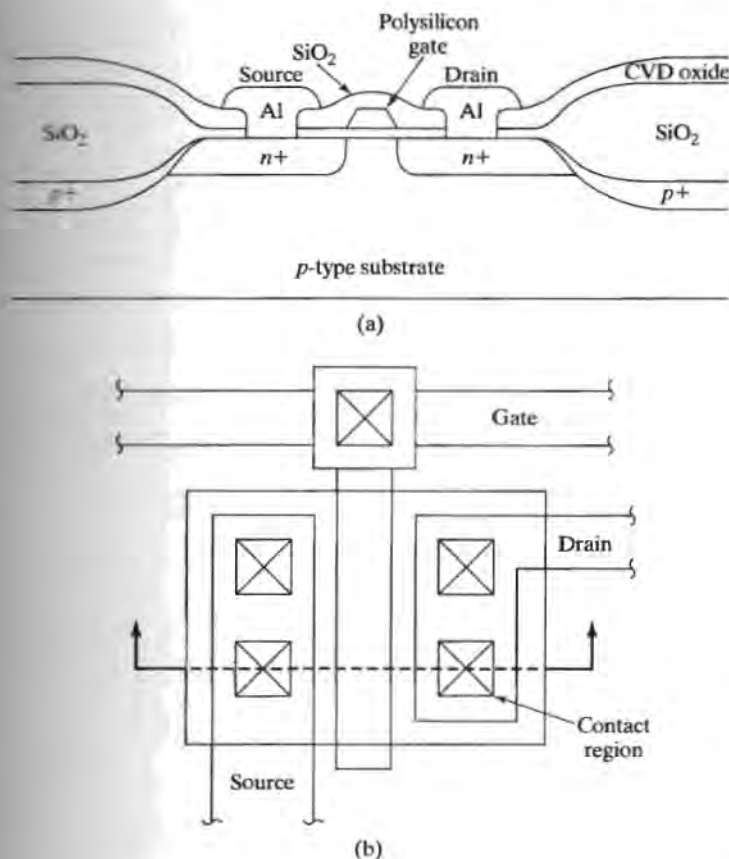


FIGURE 1.4

The basic structure of an n -channel metal-oxide-semiconductor (NMOS) transistor structure. (a) The vertical cross section through the transistor; (b) a composite top view of the masks used to fabricate the transistor in (a). The transistor uses heavily doped polysilicon as the gate "metal."

penetration of impurities during ion implantation or diffusion. Windows are cut in the masking material by etching with acids or in a plasma. Window patterns are transferred to the wafer surface from a mask through the use of optical techniques. The masks are also produced using photographic reduction techniques.

Photolithography includes the overall process of mask fabrication, as well as the process of transferring patterns from the masks to the surface of the wafer. The photolithographic process is critical to the production of integrated circuits, and the number of mask steps is often used as a measure of complexity when comparing fabrication processes.

1.3 METAL-OXIDE-SEMICONDUCTOR (MOS) PROCESSES

1.3.1 Basic NMOS Process

A possible process flow for a basic n -channel MOS process (NMOS) is shown in Fig. 1.6 on page 9 and Fig. 1.7 on page 10. The starting wafer is first oxidized to form a thin pad oxide layer of silicon dioxide (SiO_2) that protects the silicon surface. Silicon nitride is then deposited by a low-pressure chemical vapor deposition (LPCVD) process. Mask #1 defines the active transistor areas. The nitride/oxide sandwich is etched away

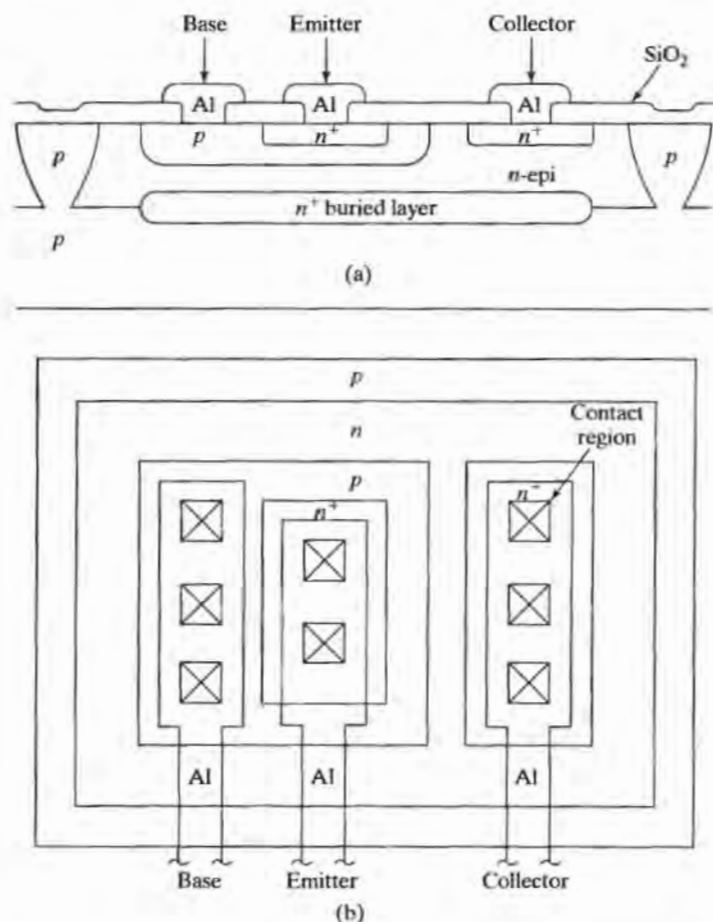


FIGURE 1.5

The basic structure of a junction-isolated bipolar transistor. (a) The vertical cross section through the transistor; (b) a composite top view of the masks used to fabricate the transistor in (a).

everywhere except where transistors are to be formed. A boron implantation is performed and followed by an oxidation step. The nitride serves as both an implantation mask and an oxidation mask. After the nitride and thin oxide padding layers are removed, a new thin layer of oxide is grown to serve as the gate oxide for the MOS transistors. Following gate-oxide growth, a boron implantation is commonly used to adjust the threshold voltage to the desired value.

Polysilicon is deposited over the complete wafer using a CVD process. The second mask defines the polysilicon gate region of the transistor. Polysilicon is etched away everywhere except over the gate regions and the areas used for interconnection. Next, the source/drain regions are implanted through the thin oxide regions. The implanted impurity may be driven in deeper with a high-temperature diffusion step. More oxide is deposited on the surface, and contact openings are defined by the third mask step. Metal is deposited over the wafer surface by evaporation or sputtering. The fourth mask step is used to define the interconnection pattern that will be etched in the metal. A passivation layer of phosphosilicate glass or silicon nitride (not shown in Fig. 1.6) is deposited on the wafer surface, and the final mask (#5) is used to define windows so that bonding wires can be attached to pads on the periphery of the IC die.

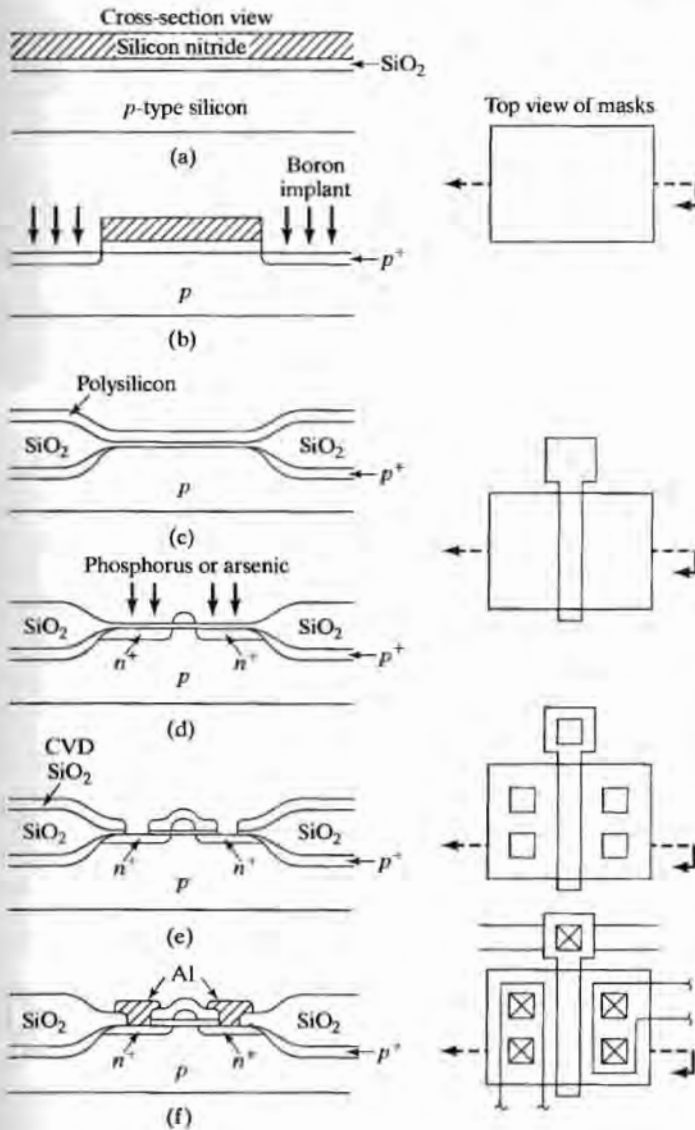


FIGURE 1.6

Process sequence for a semirecessed oxide NMOS process. (a) Silicon wafer covered with silicon nitride over a thin padding layer of silicon dioxide; (b) etched wafer after first mask step. A boron implant is used to help control field oxide threshold; (c) structure following oxidation, nitride removal, and polysilicon deposition; (d) wafer after second mask step and etching of polysilicon; (e) the third mask has been used to open contact windows following silicon dioxide deposition; (f) final structure following metal deposition and patterning with fourth mask.

This simple process requires five mask steps. Note that these mask steps use subtractive processes. The entire surface of the wafer is first coated with a desired material, and then most of the material is removed by wet chemical or dry plasma etching.

1.3.2 Basic Complementary MOS (CMOS) Process

Figure 1.8 shows the mask sequence for a basic complementary MOS (CMOS) process. One new mask, beyond that of the NMOS process, is used to define the “*n*-well,” or “*n*-tub,” which serves as the substrate for the *p*-channel devices. A second new mask step is used to define the source/drain regions of the *p*-channel transistors.

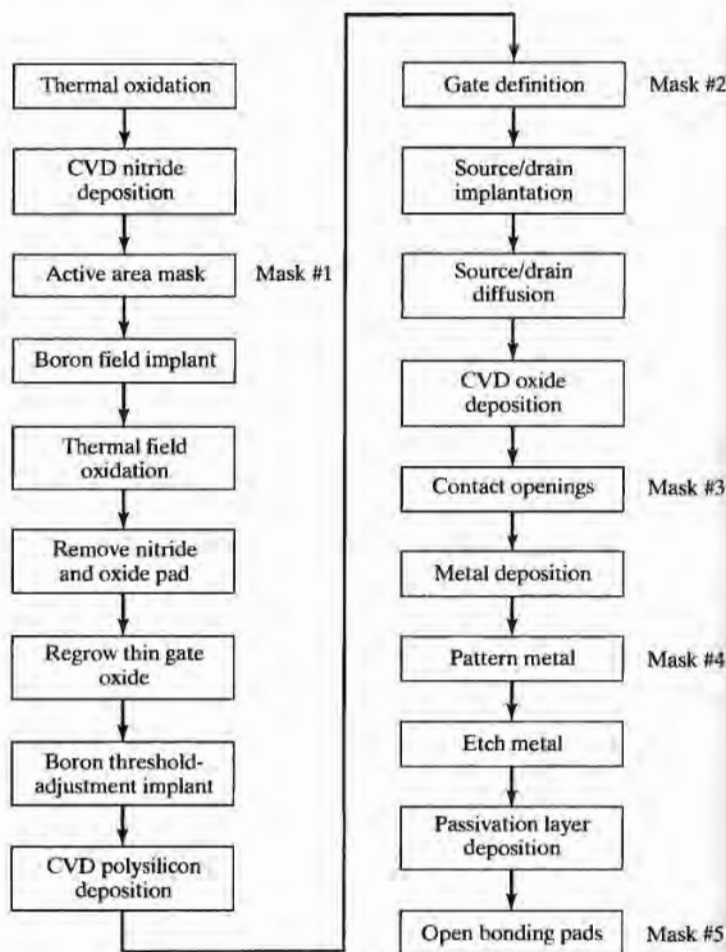


FIGURE 1.7

Basic NMOS process flowchart.

Additional masks may be used to adjust the threshold voltage of the MOS transistors and are very common in state-of-the-art NMOS and CMOS processes.

Older CMOS processes use a p -well instead of an n -well. Twin-well processes have also been developed recently. Both a p -well and an n -well are formed in a lightly doped substrate, and the n - and p -channel devices can each be optimized for highest performance. Twin-well very large-scale integration (VLSI) processes use lightly doped layers grown on heavily doped substrates to suppress a CMOS failure mode called *latchup*.

1.4 BASIC BIPOLAR PROCESSING

Basic bipolar fabrication is somewhat more complex than single-channel MOS processing, as indicated in Figs. 1.9 on page 12 and 1.10 on page 13. A p -type silicon wafer is oxidized, and the first mask is used to define a diffused region called the *buried layer*, or *subcollector*. This diffusion is used to reduce the collector resistance of the bipolar transistor. Following the buried-layer diffusion, a process called *epitaxy* is used to grow single-crystal n -type silicon on top of the silicon wafer. The epitaxial growth process results in a high-

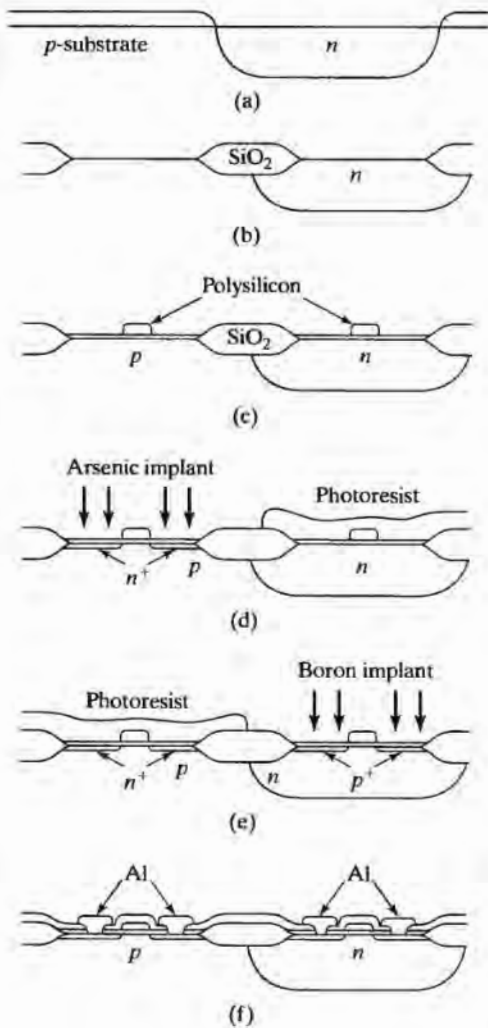


FIGURE 1.8

Cross-sectional views at major steps in a basic CMOS process. (a) Following n -well diffusion, (b) following selective oxidation, and (c) following gate oxidation and polysilicon gate definition; (d) NMOS source/drain implantation; (e) PMOS source/drain implantation; (f) structure following contact and metal mask steps.

quality silicon layer with the same crystal structure as the original silicon wafer. An oxide layer is then grown on the wafer. Mask two is used to open windows for a deep p -diffusion, which is used to isolate one bipolar transistor from another. Another oxidation follows the isolation diffusion. Mask three opens windows in the oxide for the p -type base diffusion. The wafer is usually oxidized during the base diffusion, and mask four is used to open windows for the emitter diffusion. The same diffusion step places an n^+ region under the collector contact to ensure that a good ohmic contact will be formed during subsequent metallization. Masks five, six, and seven are used to open contact windows, pattern the metallization layer, and open windows in the passivation layer just as in the NMOS process described in Section 1.3. Thus, the basic bipolar process requires seven mask levels compared with five for the basic NMOS process.

After the MOS or bipolar process is completed, each die on the wafer is tested, and bad dice are marked with ink. The wafer is then sawed apart. Good dice are mounted in various packages for final testing and subsequent sale or use.

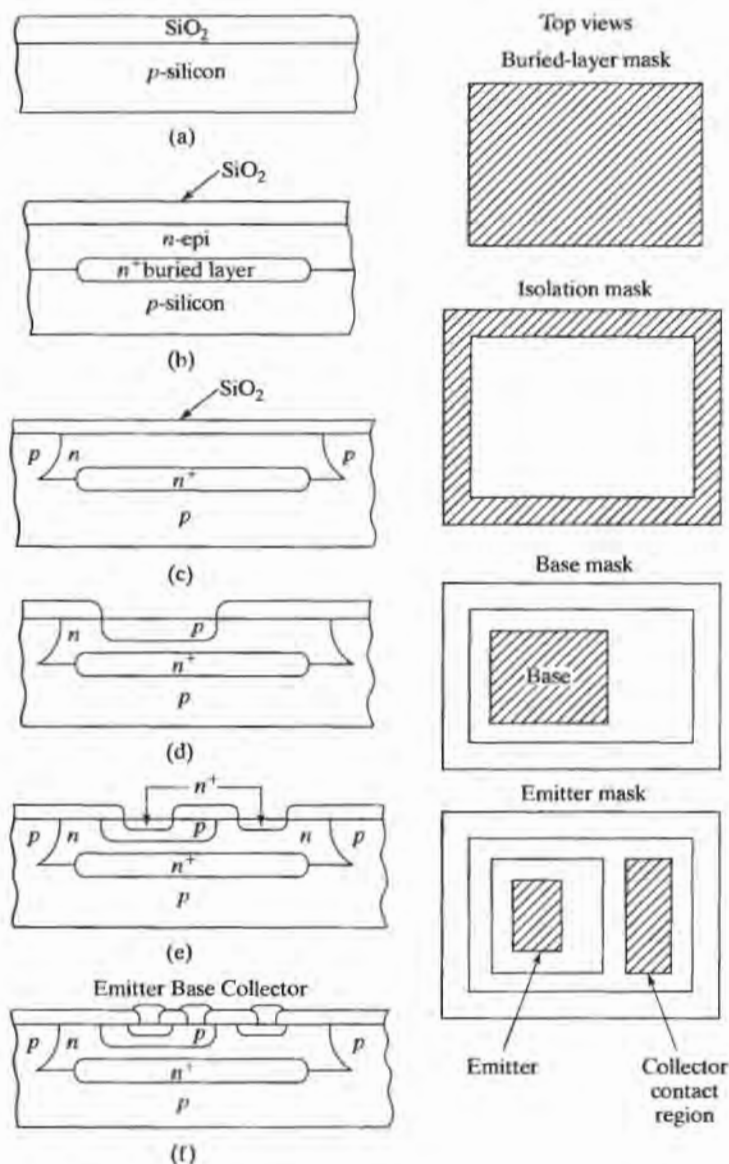


FIGURE 1.9

Cross-sectional view of the major steps in a basic bipolar process. (a) Wafer with silicon dioxide layer; (b) following buried-layer diffusion using first mask, and subsequent epitaxial layer growth and oxidation; (c) following deep-isolation diffusion using second mask; (d) following boron-base diffusion using third mask; (e) fourth mask defines emitter and collector contact regions; (f) final structure following contact and metal mask steps.

1.5 SAFETY

In the course of IC fabrication processes described throughout the rest of this text, we shall encounter a wide variety of acids, highly corrosive bases, organic and inorganic solvents, and materials with carcinogenic properties, as well as extremely toxic gases, and this represents a good opportunity to stress the need to exercise a high degree of caution before proceeding with any semiconductor processing. Because of the dangers, most laboratories require individuals to pass a safety test before they are permitted to work in the laboratory.

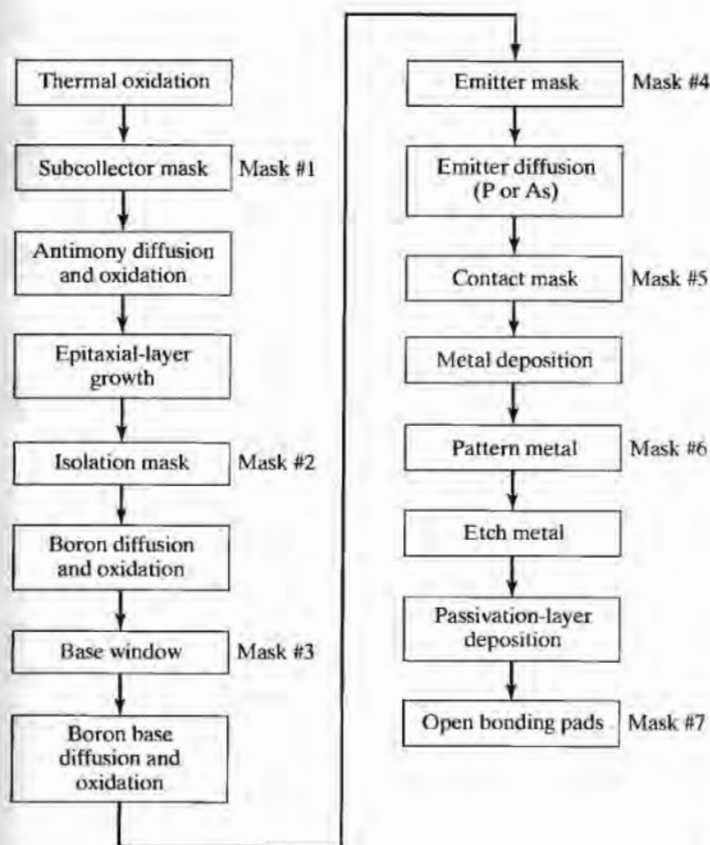


FIGURE 1.10

Basic bipolar process flowchart.

Wet processes, particularly those used for cleaning and etching, involve the use of a wide variety of acids and bases. Both can produce serious burns if they contact the skin, and even the fumes can produce irritation to the skin or serious eye damage. Rubber gloves, an apron, and eye protection should always be worn when handling these materials. However, gloves should not be relied upon to protect one during immersion in liquids, because of pinhole formation in the gloves.

Care must be exercised in handling, mixing, and disposing these liquids. Environmental standards often restrict the methods that can be used for disposal. Acids and bases must not be combined during disposal. When diluting bases or acids, concentrated chemicals should be added to water, not the reverse. Many acids will discolor the skin or give a burning sensation upon contact. However, hydrofluoric acid (HF) is much more insidious. Although a weak acid, HF readily penetrates the skin to produce deep and painful burns that are not detected until after the damage is done. Immediate medical attention is required for such burns.

Ion implantation, low-pressure chemical vapor deposition, and epitaxial growth represent just a few of the processes that may involve extremely toxic or explosive gases. For example, extreme caution must be exercised with the delivery of arsine,

phosphine, germane, silane, and anhydrous ammonia, to name just a few. In addition, many of the systems, such as ion implanters, plasma reactors, and electron-beam evaporation systems, involve lethal voltages.

As a general rule before dealing with any new chemical, one must research and study its overall properties to understand toxicity, safe handling practices, and any unusual reactions that may occur or reaction products that can be produced.

The rest of this book concentrates on the basic processes used in the fabrication of monolithic integrated circuits. Chapters 2 through 8 discuss mask making and pattern definition, oxidation, diffusion, ion implantation, film deposition, interconnections and contacts, and packaging and yield. The last three chapters introduce the integration of process, layout, and device design for MOS, bipolar, and MEMS technologies.

REFERENCES

- [1] *Digest of the IEEE International Solid-State Circuits Conference*, held in February of each year. (<http://www.sscs.org/isscc>)
- [2] *Digest of the IEEE International Electron Devices Meeting*, held in December of each year. (<http://www.ieee.org/conference/iedm>)
- [3] *Digests of the International VLSI Technology and Circuits Symposia*, co-sponsored by the IEEE and JSAP, held in June of each year. (<http://www.vlsisymposium.org>)
- [4] *The International Technology Roadmap for Semiconductors*, The Semiconductor Industry Association (SIA), San Jose, CA, 1999. (<http://www.semichips.org>)

PROBLEMS

- 1.1 Make a list of two dozen items in your everyday environment that you believe contain IC chips. A PC and its peripherals are considered to be one item. (Do not confuse electro-mechanical timers, common in clothes dryers or the switch in a simple thermostat or coffee maker, with electronic circuits.)
- 1.2 (a) Make a table comparing the areas of wafers with the following diameters: 25, 50, 75, 100, 125, 150, 200, 300 and 450 mm.
 (b) Approximately how many 1-mm \times 1-mm dice are on a 450-mm wafer?
 (c) How many 25-mm \times 25-mm dice?
- 1.3 (a) Calculate an estimate of the number of 20-mm \times 20-mm dice on a 300-mm diameter wafer, in terms of the total wafer and die areas.
 (b) Calculate the exact number of 20-mm \times 20-mm dice that actually fit on the 300-mm wafer. (It may help to draw a picture.)
- 1.4 The straight line in Fig. 1.2(a) is described by $B = 19.97 \times 10^{0.1977(Y - 1960)}$ bits/chip. If a straight-line projection is made using this equation, what will be the number of memory bits/chip in the year 2020?
- 1.5 The straight line in Fig. 1.2(b) is described by $N = 1027 \times 10^{0.1505(Y - 1970)}$ transistors. Based upon a straight-line projection of this figure, what will be the number of transistors in a microprocessor in the year 2020?

- 1.6 (a) How many years does it take for memory chip density to increase by a factor of two, based upon the equation in Problem 1.4?
 (b) How about by a factor of 10?
- 1.7 (a) How many years does it take for microprocessor circuit density to increase by a factor of two, based upon the equation in Problem 1.5?
 (b) How about by a factor of 10?
- 1.8 If you make a straight-line projection from Fig. 1.3, what will be the minimum feature size in integrated circuits in the year 2020? The curve can be described by $F = 8.214 \times 10^{-0.06079(Y-1970)} \mu\text{m}$. Do you think this is possible? Why or why not?
- 1.9 The filament of a small vacuum tube uses a power of approximately 0.5 W. Suppose that approximately 300 million of these tubes are used to build the equivalent of a 256-Mb memory. How much power is required for this memory? If this power is supplied from a 220-V ac source, what is the current required by this memory?
- 1.10 An 18-mm \times 25-mm die is covered by an array of 0.25- μm metal lines separated by 0.25- μm -wide spaces.
 (a) What is the total length of wire on this die?
 (b) How about 0.1- μm lines and spaces.
- 1.11 The curve in Fig. 1.1(b) represents the approximate number of chips on a wafer of a given diameter. Determine the exact number of 10 \times 10 mm dice that will fit on a wafer with a diameter of 200 mm. (The number indicated on the curve is 314.)
- 1.12 The cost of processing a wafer in a particular process is \$1,000. Assume that 35% of the fabricated dice are good. Find the number of dice, using Fig. 1.1(b).
 (a) Determine the cost per good die for a 150 mm wafer.
 (b) Repeat for a 200 mm wafer.
- 1.13 A certain silicon-gate NMOS transistor occupies an area of $25 \lambda^2$, where λ is the minimum lithographic feature size.
 (a) How many MOS transistors can fit on a 5 \times 5 mm die if $\lambda = 1 \mu\text{m}$?
 (b) 0.25 μm ?
 (c) 0.10 μm ?
- 1.14 A simple *pn* junction diode is shown in cross section in Fig. P1.14. Make a possible process flowchart for fabrication of this structure, including mask steps.

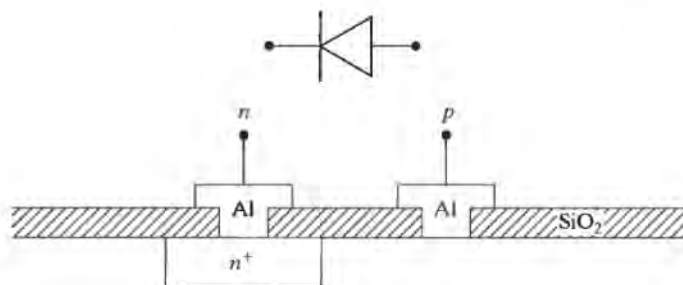


FIGURE P1.14

- 1.15** Draw a set of contact and metal masks for the bipolar transistor of Fig. 1.9. Use square contact windows with one contact to the emitter and two contacts to the base and collector regions.

CHAPTER 2

Lithography

In order to produce an integrated circuit, thin films of various materials are used as barriers to the diffusion or implantation of impurity atoms or as insulators between conductive materials and the silicon substrate. Holes, or windows, are cut through this barrier material wherever impurity penetration or contact is desired.

Masks contain the patterns of windows that are transferred to the surface of the silicon wafer using a process called *photolithography*. Photolithography makes use of a highly refined version of the photoengraving process. The patterns are first transferred from the mask to a light-sensitive material called *photoresist*. Chemical or plasma etching is then used to transfer the pattern from the photoresist to the barrier material on the surface of the wafer. Each mask step requires successful completion of numerous processing steps, and the complexity of an IC process is often measured by the number of photographic masks used during fabrication. This chapter will explore the lithographic process, including mask fabrication, photoresist processes, and etching.

2.1 THE PHOTOLITHOGRAPHIC PROCESS

Photolithography encompasses all the steps involved in transferring a pattern from a mask to the surface of the silicon wafer. The various steps of the basic photolithographic process given in Figs. 2.1 and 2.2 will each be discussed in detail next.

Ultraclean conditions must be maintained during the lithography process. Any dust particles on the original substrate or that fall on the substrate during processing can result in defects in the final resist coating. Even if defects occur in only 10% of the chip sites at each mask step, fewer than 50% of the chips will be functional after a seven-mask process is completed. Vertical laminar-flow hoods in clean rooms are used to prevent particulate contamination throughout the fabrication process. Clean rooms use filtration to remove particles from the air and are rated by the maximum number of particles per cubic foot or cubic meter of air, as shown in Table 2.1. Clean rooms have evolved from Class 100 to the Class 1 facilities now being used for VLSI/ULSI processing. For comparison, each cubic foot of ordinary room air has several million dust particles exceeding a size of 0.5 μm .

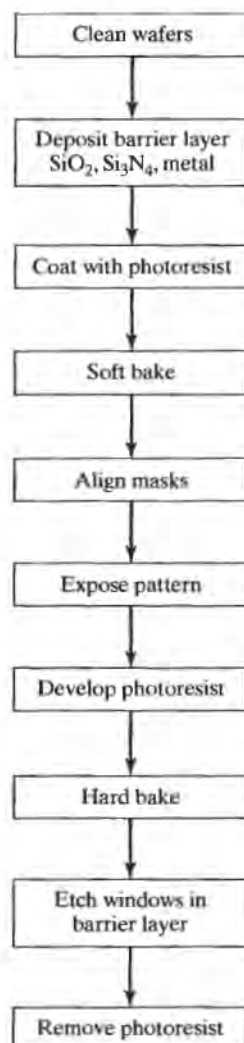


FIGURE 2.1
Steps of the photolithographic process.

TABLE 2.1 Ratings by Class of Effectiveness of Filtration in Clean Rooms

Class	Number of 0.5- μm particles per ft ³ (m ³)	Number of 5- μm particles per ft ³ (m ³)
10,000	10,000 (350,000)	65 (23,000)
1,000	1,000 (35,000)	6.5 (2,300)*
100	100 (3,500)	0.65 (230)*
10	10 (350)	0.065 (23)*
1	1 (35)*	0.0065 (2.3)*

*It is very difficult to measure particulate counts below 10 per ft³.

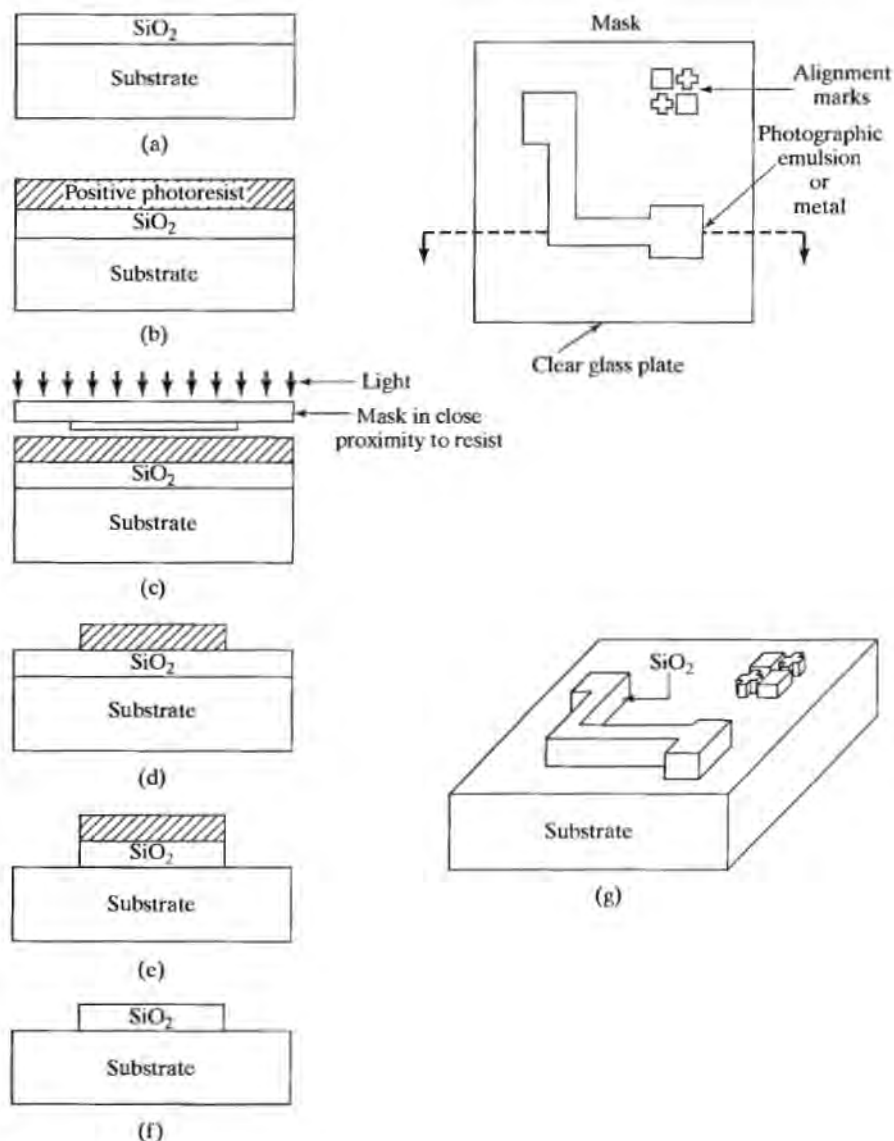


FIGURE 2.2

Drawings of a wafer through the various steps of the photolithographic process. (a) Substrate covered with silicon dioxide barrier layer; (b) positive photoresist applied to the surface of the wafer; (c) mask in close proximity to the surface of the resist-covered wafer; (d) substrate following resist exposure and development; (e) substrate following etching of the silicon dioxide layer; (f) oxide barrier on wafer surface after resist removal; (g) view of substrate with silicon dioxide pattern on the surface.

2.1.1 Wafer and Wafer Cleaning

IC fabrication starts with *n*- or *p*-type silicon wafers supplied with a specified resistivity. The wafers range in thickness from 250 to 500 μm . Two-hundred-mm (eight-inch)

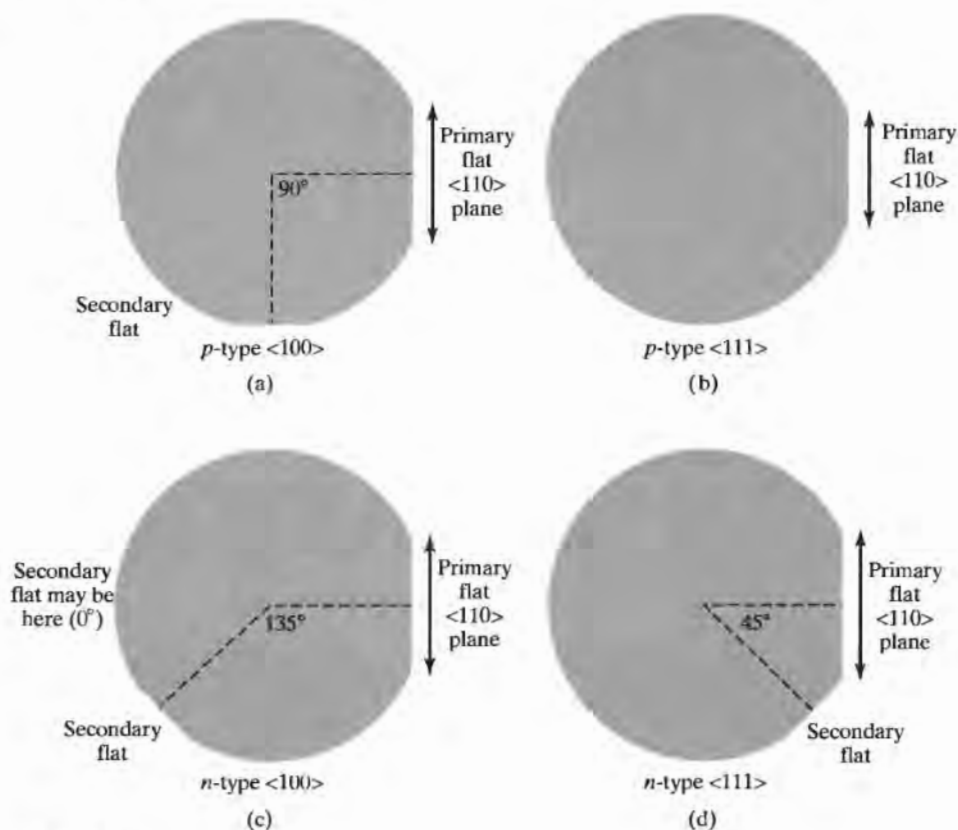


FIGURE 2.3

Illustration of wafer flat standard used to identify 100 mm wafers.

diameter wafers are widely used at the time of this writing, and processing equipment for 300-mm (12-inch) wafers is already becoming available. Wafers with diameters of 1, 1.5, 2, 3, 4, 5, and 6 in. have all been used at various stages in history, and Table 1.1 indicates that 450-mm wafers will be in use by the end of the decade.

The silicon wafers are identified by a standard system¹ of straight edges or wafer flats. These flats are ground into the silicon ingot before it is sliced into wafers and are used to indicate the wafer type (*n*-type or *p*-type) and the surface orientation (<100> or <111>) as indicated in Fig. 2.3. The primary wafer flat identifies the <110> crystal plane. (It is important to note, however, that this identification system is not always utilized, particularly with wafer sizes of 150 mm and above.)

Prior to use, wafers are chemically cleaned to remove particulate matter on the surface, as well as any trace of organic, ionic, and metallic impurities. A cleaning step utilizing a solution of hydrofluoric acid removes any oxide that may have formed on the wafer surface. A typical cleaning process is presented in Table 2.2.

One very important chemical used in wafer cleaning and throughout microelectronic fabrication processes is deionized (DI) water. DI water is highly purified and filtered to remove all traces of ionic, particulate, and bacterial contamination. The

¹A SEMI (Semiconductor Equipment and Materials International) standard.

TABLE 2.2 Silicon Wafer Cleaning Procedure[7, 8]

-
- | | |
|----|--|
| A. | Solvent Removal |
| 1. | Immerse in boiling trichloroethylene (TCE) for 3 min. |
| 2. | Immerse in boiling acetone for 3 min. |
| 3. | Immerse in boiling methyl alcohol for 3 min. |
| 4. | Wash in DI water for 3 min. |
| B. | Removal of Residual Organic/Ionic Contamination |
| 1. | Immerse in a (5:1:1) solution of $H_2O-NH_4OH-H_2O_2$; heat solution to 75–80 °C and hold for 10 min. |
| 2. | Quench the solution under running DI water for 1 min. |
| 3. | Wash in DI water for 5 min. |
| C. | Hydrous Oxide Removal |
| 1. | Immerse in a (1:50) solution of $HF-H_2O$ for 15 sec. |
| 2. | Wash in running DI water with agitation for 30 sec. |
| D. | Heavy Metal Clean |
| 1. | Immerse in a (6:1:1) solution of $H_2O-HCl-H_2O_2$ for 10 min at a temperature of 75–80 °C. |
| 2. | Quench the solution under running DI water for 1 min. |
| 3. | Wash in running DI water for 20 min. |
-

theoretical resistivity of pure water at 25 °C is 18.3 Mohm-cm. Basic DI water systems achieve resistivities of 18 Mohm-cm with fewer than 1.2 colonies of bacteria per milliliter and with no particles larger than 0.25 μm .

Cleanliness and contamination control is such an overarching concern in VLSI/ULSI fabrication that an NSF Industry University Cooperative Research Center in Micro Contamination Control was established at the University of Arizona.

2.1.2 Barrier Layer Formation

After cleaning, the silicon wafer is covered with the material that will serve as a barrier layer. The most common material is silicon dioxide (SiO_2), so we will use it as an example here. Silicon nitride (Si_3N_4), polysilicon, photoresist, and metals are also routinely used as barrier materials at different points in a given process flow. In subsequent chapters, we will discuss thermal oxidation, chemical vapor deposition, sputtering, and vacuum evaporation processes, all of which are used to produce thin layers of these materials.

The original silicon wafer has a metallic gray appearance. Once an SiO_2 layer is formed on the silicon wafer, the surface will have a color that depends on the SiO_2 thickness. The finished wafer will have regions with many different thicknesses. Each region will produce a different color, resulting in beautiful, multicolored IC images, photographs of which appear in many books and magazines.

2.1.3 Photoresist Application

After the SiO_2 layer is formed, the surface of the wafer is coated with a light-sensitive material called *photoresist*. The surface must be clean and dry to ensure good photoresist adhesion. Freshly oxidized wafers may be directly coated, but if the wafers have been stored, they should be carefully cleaned and dried prior to application of the resist.

Lack of adhesion of photoresist to many film surfaces is a commonly encountered problem in silicon processing. In order to promote adhesion, the wafer surface is



FIGURE 2.4

Rite Track 88e wafer processing system. (Courtesy of Rite Track Equipment Services, Inc.)

treated with an adhesion promoter such as hexamethyldisilazane (HMDS) prior to photoresist application. This treatment provides good photoresist adhesion to a variety of films, including silicon dioxide (SiO_2), silicon dioxide containing phosphorous, polycrystalline silicon, silicon nitride (Si_3N_4), and aluminum.

Photoresist is typically applied in liquid form. The wafer is held on a vacuum chuck and then spun at high speed for 30 to 60 sec to produce a thin uniform layer. Speeds of 1,000 to 5,000 rpm result in layers ranging from 2.5 to 0.5 μm , respectively. The actual thickness of the resist depends on its viscosity and is inversely proportional to the square root of the spinning speed.

Figure 2.4 shows an automated cassette-to-cassette wafer track system that automatically dispenses a tightly controlled amount of photoresist onto each wafer and controls the spinning profile to achieve highly reproducible resist film thicknesses. Each cassette typically contains 25 wafers.

2.1.4 Soft Baking or Prebaking

A drying step called *soft baking*, or *prebaking*, is used to improve adhesion and remove solvent from the photoresist. Times range from 5 to 30 min in an oven at 60 to 100°C in an air or nitrogen atmosphere. The soft-baking process is specified on the resist manufacturer's data sheet and should be followed closely. After soft baking, the photoresist is ready for mask alignment and exposure.

2.1.5 Mask Alignment

The complex pattern from a photo mask (or just mask), a square glass plate with a patterned emulsion or metal film on one side, must be transferred to the surface of the wafer. Each mask following the first² must be carefully aligned to the previous pattern

²In many MEMS (see Chapter 11 for an introduction) and sensor applications, the first mask must also be carefully aligned to the crystallographic axes. Double-sided alignment is also used with infrared systems that can "see" through the wafer or with specialized optics systems that permit simultaneous viewing of both sides of the wafer.

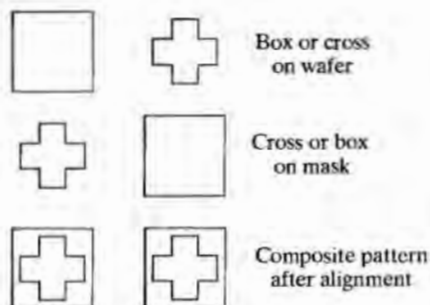


FIGURE 2.5

A simple set of alignment marks. At some steps a cross may be aligned within a box. At others, a box may be placed around the cross. The choice depends on the type of resist being used at a given mask step.

on the wafer. Manual operation of alignment and exposure equipment was used in early fabrication systems. However, VLSI/ULSI designs require extremely small geometrical features (minimum line width or space) and tight alignment tolerances. For example, 100 nm (0.1 μm) lithography will require a worst-case alignment error of 35 nm (mean + 3σ), and computer-controlled alignment systems are required to achieve these required levels of alignment precision.

With basic manual alignment equipment, the wafer is held on a vacuum chuck and carefully moved into position below the mask using an adjustable x - y stage. The mask is spaced 25 to 125 μm above the surface of the wafer during alignment. If contact printing is being used, the mask is brought into contact with the wafer after alignment.

Alignment marks are introduced on each mask and transferred to the wafer as part of the IC pattern. The marks are used to align each new mask level to one of the previous levels. A sample set of alignment marks is shown in Fig. 2.5. For certain mask levels, the cross on the mask is placed in a box on the wafer. For other mask levels, the box on the mask is placed over a cross on the wafer. The choice depends on the type of resist used during a given photolithographic step. Split-field optics are used to simultaneously align two well-separated areas of the wafer.

2.1.6 Photoresist Exposure and Development

Following alignment, the photoresist is exposed through the mask with high-intensity ultraviolet light. Resist is exposed wherever silicon dioxide is to be removed. The photoresist is developed with a process very similar to that used for developing ordinary photographic film, using a developer supplied by the photoresist manufacturer. Any resist that has been exposed to ultraviolet light is washed away, leaving bare silicon dioxide in the exposed areas of Fig. 2.6(d). A photoresist acting in the manner just described is called a *positive resist*, and the mask contains a copy of the pattern that will remain on the surface of the wafer. Windows are opened wherever the exposing light passes through the mask.

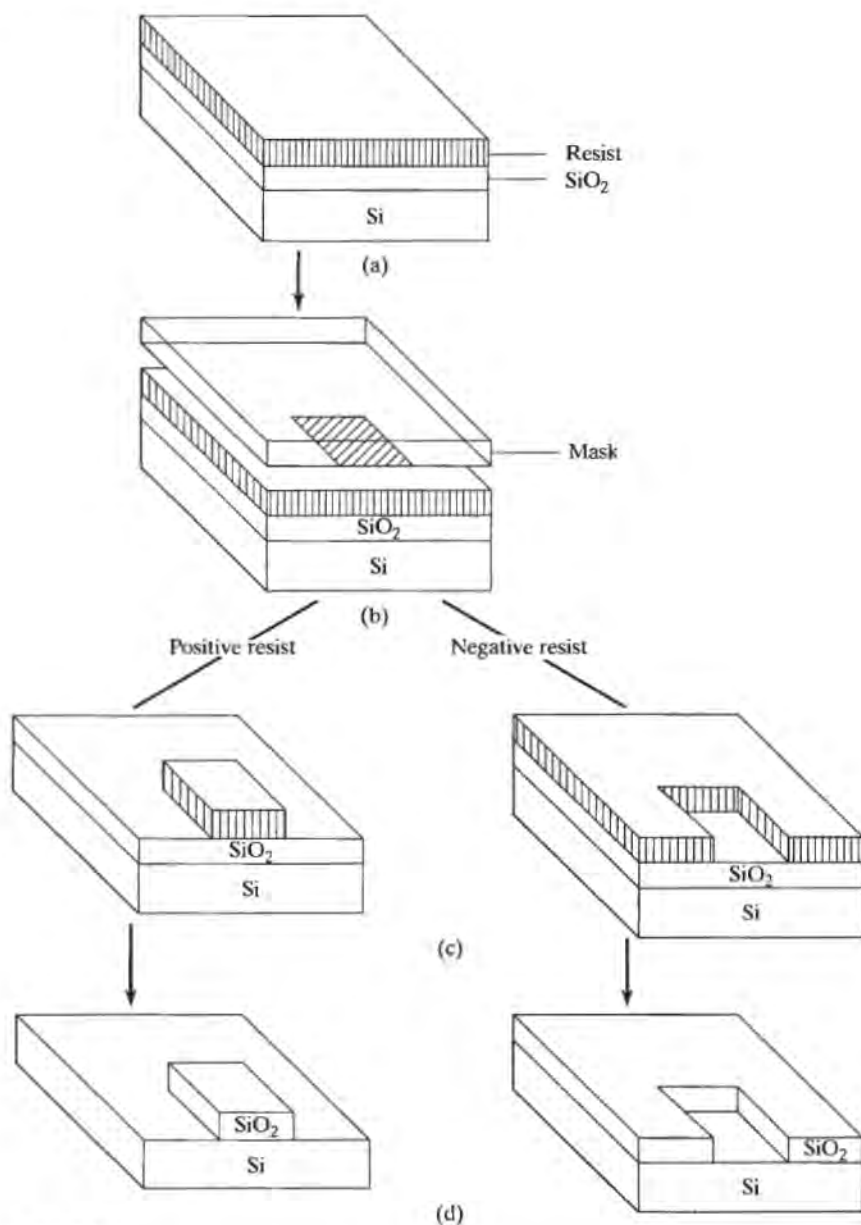


FIGURE 2.6

Resist and silicon dioxide patterns following photolithography with positive and negative resists.

Negative photoresists can also be used. A negative resist remains on the surface wherever it is exposed. Figure 2.6 shows simple examples of the patterns transferred to a silicon dioxide barrier layer using positive and negative photoresists with the same mask. Negative resists were widely used in early IC processing. However, positive resist yields better process control in small-geometry structures and is now the main type of resist used in VLSI processes.

2.1.7 Hard Baking

Following exposure and development, a baking step is used to harden the photoresist and improve adhesion to the substrate. A typical process involves baking in an oven for 20 to 30 min at 120 to 180°C. Details of this step are again specified on the manufacturer's photoresist data sheets.

2.2 ETCHING TECHNIQUES

Chemical etching in liquid or gaseous form is used to remove any barrier material not protected by hardened photoresist. The choice of chemicals depends on the material to be etched. A high degree of selectivity is required so that the etchant will remove the unprotected barrier layer much more rapidly than it attacks the photoresist layer.

2.2.1 Wet Chemical Etching

A buffered oxide etch (BOE, or BHF) is commonly used to etch windows in silicon dioxide layers. BOE is a solution containing hydrofluoric acid (HF), and etching is performed by immersing the wafers in the solution. At room temperature, HF etches silicon dioxide much more rapidly than it etches photoresist or silicon. The etch rate in BOE ranges from 10 to 100 nm/min at 25 °C, depending on the density of the silicon dioxide film. Etch rate is temperature dependent, and temperature is carefully monitored during the etch process. In addition, etch rates depend on the type of oxide present. Oxides grown in dry oxygen etch more slowly than those grown in the presence of water vapor. A high concentration of phosphorus in the oxide enhances the etch rate, whereas a reduced etch rate occurs when a high concentration of boron is present. High concentrations of these elements convert the SiO_2 layer to a phosphosilicate or borosilicate glass.

HF and water both wet silicon dioxide, but do not wet silicon. The length of the etch process may be controlled by visually monitoring test wafers that are etched along with the actual IC wafers. Occurrence of a hydrophobic condition on the control wafer signals completion of the etch step.

Wet chemical etching tends to be an isotropic process, etching equally in all directions. Figure 2.7(a) shows the result of isotropic etching of a narrow line in silicon dioxide. The etching process has etched under the resist by a distance equal to the thickness of the film. This "etch bias" becomes a serious problem in processes requiring line widths with dimensions similar to the thickness of the film.

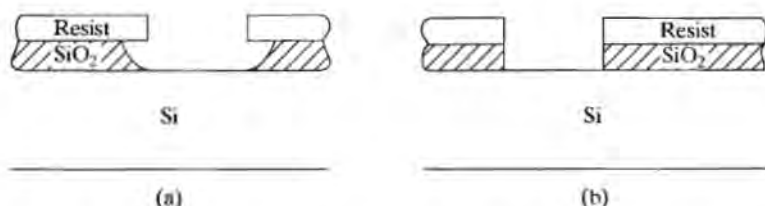


FIGURE 2.7

Etching profiles obtained with (a) isotropic wet chemical etching and (b) dry anisotropic etching in a plasma or reactive-ion etching system.

2.2.2 Dry Etching Plasma Systems

Dry plasma etching processes are widely used in VLSI fabrication. Highly anisotropic etching profiles can be obtained as shown in Fig. 2.7(b), avoiding the undercutting problem of Fig. 2.7(a) characteristic of wet processes. Dry processes require only small amounts of reactant gases, whereas wet etching requires disposal of relatively large amounts of liquid chemical wastes.

Plasma systems use RF excitation to ionize a variety of source gases in a vacuum system. The RF power source typically operates at a frequency of 13.56 MHz, which is set aside by the Federal Communications Commission (FCC) for industrial and scientific purposes. However, plasma systems can also operate at frequencies as low as a few hundred kilohertz, and microwave excitation is in use in certain systems.

The mode of operation of the plasma system depends upon the operating pressure, as shown in Table 2.3, as well as the structure of the reaction chamber. Standard plasma etching corresponds to the highest of the three pressure regimes, and a conceptual drawing for a parallel-plate plasma-etching system is shown in Fig. 2.8(a). In this case, the electrode structure is symmetric, and the wafer to be etched is placed upon the grounded electrode. Free radicals such as fluorine or chlorine are created in the plasma and react at the wafer surface to etch silicon, silicon dioxide, silicon nitride, organic materials, and metals. A sample of possible source gases used to etch these materials appears in Table 2.4, but a much broader range of choices is available [3]. The basic plasma-etching process is isotropic, and additional atomic species such as argon, hydrogen, and oxygen are often introduced to improve etch rate or selectivity.

Ion milling uses energetic noble gas ions such as Ar^+ to bombard the wafer surface. Etching occurs by physically knocking atoms off the surface of the wafer. Highly anisotropic etching can be obtained, but selectivity is often poor. Metals can be used as barrier materials to protect the wafer from etching. Ion milling operates in the lowest of the three pressure ranges given in Table 2.3. In this case, ions are accelerated toward the surface by a strong electric field that can be introduced by adding a variable external dc-bias voltage between the electrodes.

Reactive-ion etching (RIE) combines the plasma and sputter etching processes. Plasma systems are used to ionize reactive gases, and the ions are accelerated to bombard the surface. Etching occurs through a combination of the chemical reaction and momentum transfer from the etching species and is highly anisotropic. The voltage required to accelerate ions from the plasma toward the wafer surface can be developed by introducing an asymmetry into the structure of the plasma chamber as indicated in Fig. 2.8(b). In this drawing, the surface area of the upper electrode is made larger than that of the lower electrode, the upper electrode is now grounded, and the wafer is placed on the electrode driven by the RF source. The physical asymmetry of the system produces a self-bias between the electrodes that provides the acceleration potential required to direct the ions toward the wafer surface.

2.2.3 Photoresist Removal

After windows are etched through the SiO_2 layer, the photoresist is stripped from the surface, leaving a window in the silicon dioxide. Photoresist removal typically uses proprietary-liquid resist strippers, which cause the resist to swell and lose adhesion to the substrate. Dry processing may also be used to remove resist by oxidizing (burning) it in an oxygen plasma system, a process often called *resist ashing*.

TABLE 2.3 Etching Pressure Ranges

Etching Mode	Pressure (Torr)
Ion Milling	10^{-4} – 10^{-3}
Reactive Ion Etching/Ion Milling	10^{-3} – 10^{-1}
Plasma Etching	10^{-1} –5

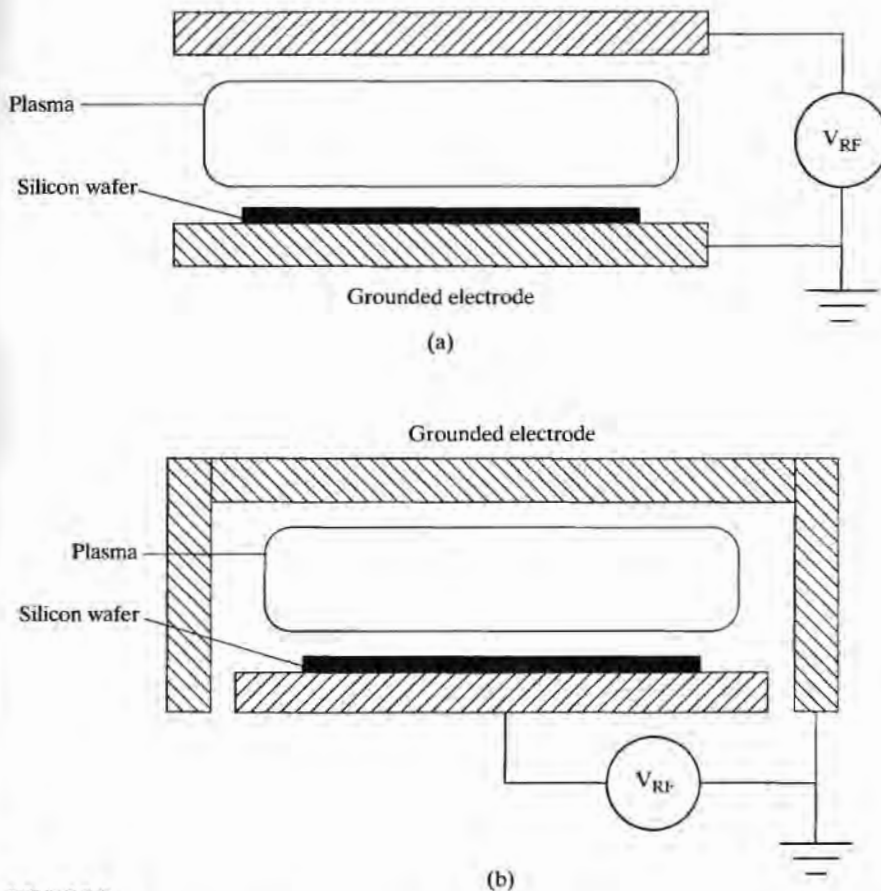


FIGURE 2.8

(a) Concept of a parallel plate plasma etcher (b) Asymmetrical reactive ion etching (RIE) system

TABLE 2.4 Plasma-Etching Sources

Material	Source Gases
Organic Materials	O_2 , SF_6 , CF_4
Polysilicon	CCl_4 , CF_4 , NF_3 , SF_6
Silicon Dioxide	CF_4 , C_2F_6 , C_3F_8 , CHF_3
Silicon Nitride	CF_4 , C_2F_6 , CHF_3 , SF_6
Aluminum	CCl_4 , Cl_2 , BCl_3
Titanium	$C_2Cl_2F_4$, CF_4
Tungsten	Cl_2

2.2.4 Metrology and Critical Dimension Control

It is extremely important to be able to maintain accurate control of critical dimensions (CDs) through photolithography and etching processes, as well as subsequent process steps. The ITRS contains projections of the required levels of CD control. The ability to reliably measure the fabricated features with the required accuracy and repeatability is itself a major problem, and semiconductor process metrology has emerged as a separate discipline of its own that concentrates on the development of the test structures and instrumentation required to support high-yield manufacturing.

2.3 PHOTOMASK FABRICATION

Photomask fabrication involves a series of photographic processes outlined in Fig. 2.9. An IC mask begins with a large-scale drawing of each mask. Early photomasks were cut by hand in a material called *rubylith*, a sandwich of a clear backing layer and a thin red layer of Mylar. The red layer was cut with a stylus and peeled off, leaving the desired pattern in red. The original rubylith copy of the mask was 100 to 1,000 times larger than the final integrated circuit and was photographically reduced to form a reticle for use in a step-and-repeat camera, as described later.

Today, computer graphics systems and optical or electron beam pattern generators have supplanted the use of rubylith. An image of the desired mask is created on a computer graphics system. Once the image is complete, files containing the commands needed to drive a pattern generator are created. An optical pattern generator uses a flash lamp to expose the series of rectangles composing the mask image directly onto a photographic plate called the *reticle*. An electron beam system draws the pattern directly in an electron-sensitive material.

Reticle images range from 1 to 10 times final size. A step-and-repeat camera is used to reduce the reticle image to its final size and to expose a two-dimensional array of images on a master copy of the final mask. On a 200-mm wafer, it is possible to get approximately 1,200 copies of a 5-mm \times 5-mm IC chip! Figure 2.10 shows examples of a computer graphics plot, a reticle, and a final mask for a simple integrated circuit.

A final master copy of the mask is usually made in a thin film of metal, such as chrome, on a glass plate. The mask image is transferred to photoresist, which is used as an etch mask for the chrome. Working emulsion masks are then produced from the chrome master. Each time a mask is brought into contact with the surface of the silicon wafer, the pattern can be damaged. Therefore, emulsion masks can only be used for a few exposures before they are thrown away.

2.4 EXPOSURE SYSTEMS

Because contact printing can damage the surfaces of both the mask and the wafer, manufacturing lines utilize proximity and projection printing systems, as illustrated in Fig. 2.11. However, contact printing is still used in research and prototyping situations, because it can economically achieve high-resolution pattern transfer. In proximity printing, the mask is brought in very close proximity to the wafer, but does not come in contact with the wafer during exposure, thus preventing damage to the mask. Projection printing uses a dual-lens system to project a portion of the mask image onto the wafer surface. The wafer and masks may be scanned, or the system may operate in a step-and-repeat mode. The actual mask and lenses are mounted many centimeters from the wafer surface.

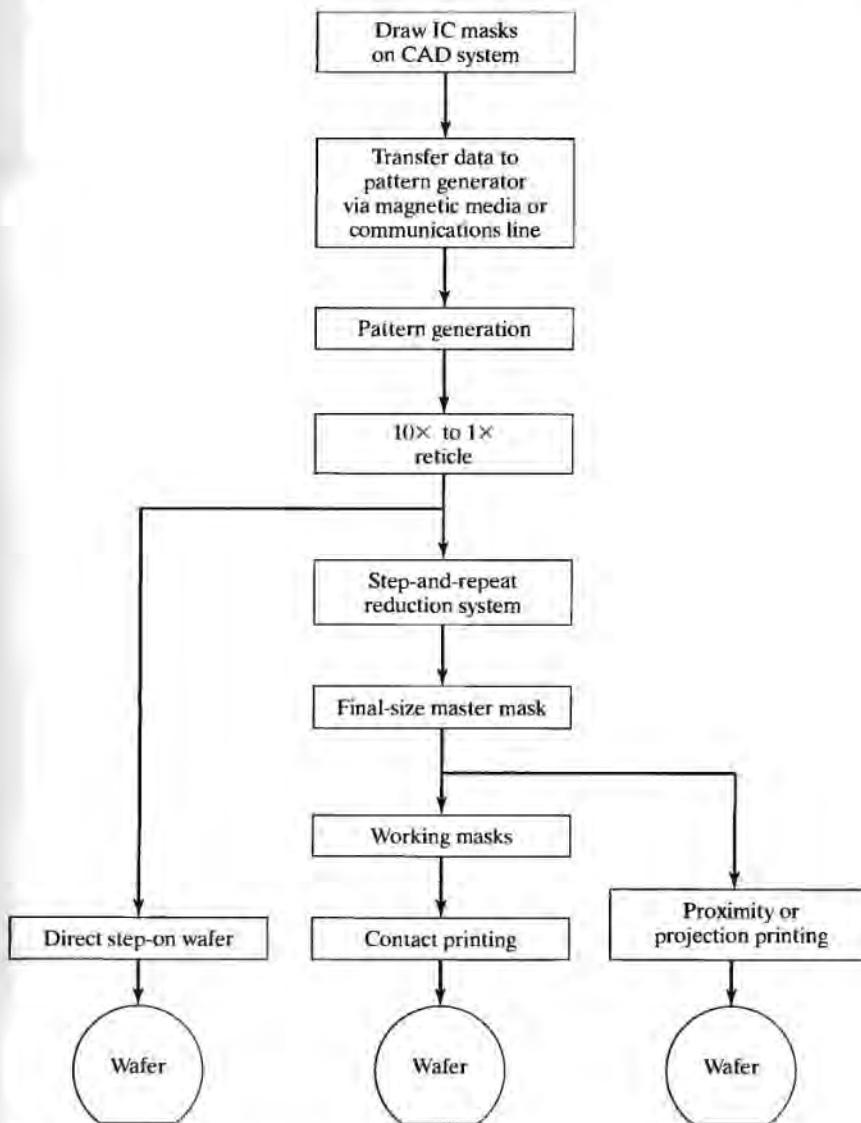


FIGURE 2.9

Outline of steps in the mask fabrication process.

For large-diameter wafers, it is impossible to achieve uniform exposure and to maintain alignment between mask levels across the complete wafer, particularly for submicron feature sizes. High-resolution VLSI lithography systems now use some form of exposure of the individual die pattern directly onto the wafer. A projection system is used with a reticle to expose the IC die pattern directly on the wafer. No step-and-repeat masks of the circuit are produced. The pattern is aligned and exposed separately at each die site.

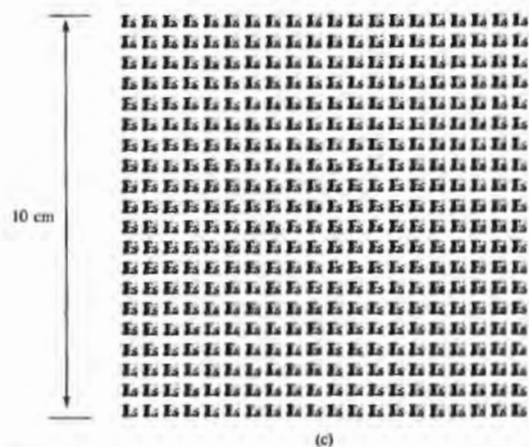
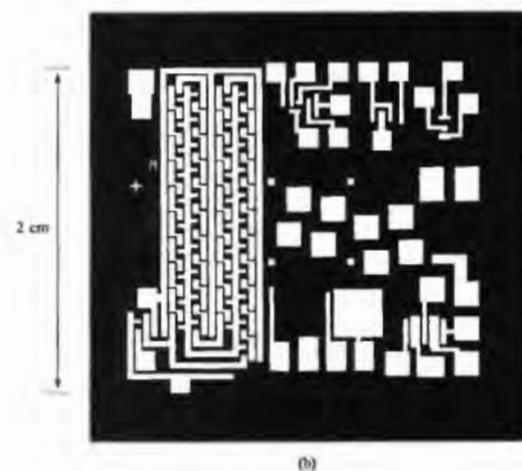
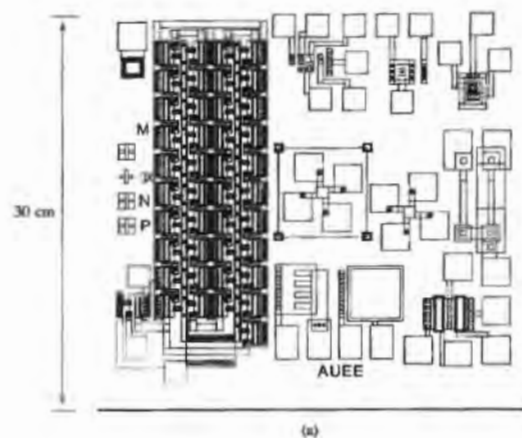


FIGURE 2.10

Mask fabrication. (a) Composite computer graphics plot of all masks for a simple integrated circuit; (b) 10X reticle of metal-level mask; (c) final-size emulsion mask with 400 copies of the metal level of the integrated circuit in (a).

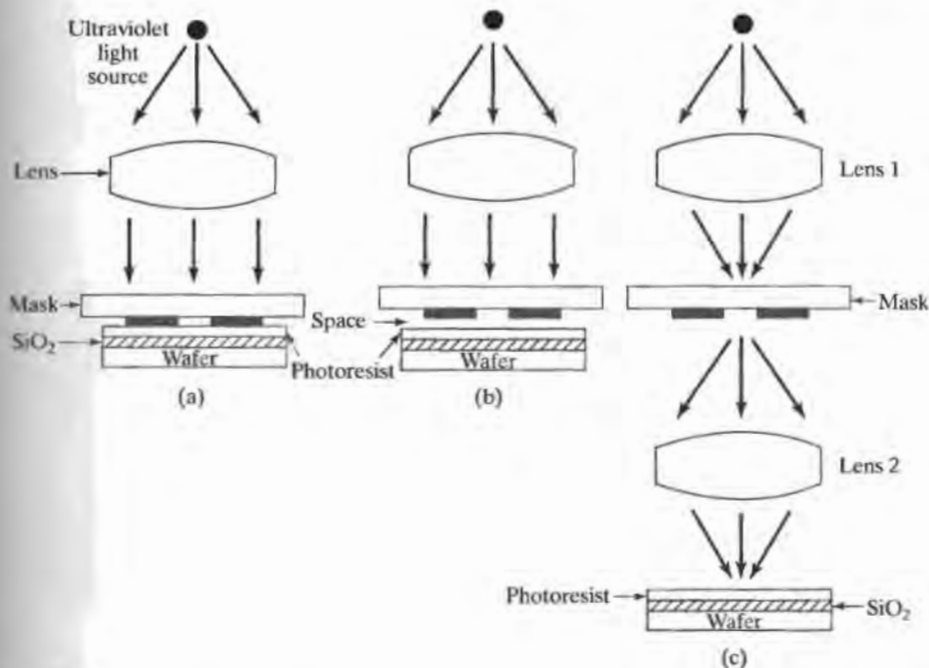


FIGURE 2.11

Artist's conception of various printing techniques. (a) Contact printing, in which wafer is in intimate contact with mask; (b) proximity printing, in which wafer and mask are in close proximity; (c) projection printing, in which light source is scanned across the mask and focused on the wafer. Copyright, 1983, Bell Telephone Laboratories, Incorporated. Reprinted by permission from Ref. [5].

Figure 2.12 on page 32 shows an artist's diagram of a direct step-on-wafer system (usually termed a "stepper"). A single die image is projected directly onto the surface of the wafer. The reticle pattern may range from 1 to 10 times the final die size. The wafer is moved (stepped) from die site to die site on the wafer, and the pattern is aligned and exposed at each individual site. The drawing in Fig. 2.12 actually hides the complexity of these systems, as can be seen from the drawing of a complete stepper system shown in Fig. 2.13 on page 33. These systems are often housed in their own environmentally controlled sections of clean rooms.

Another variation can be used when the individual die pattern becomes too large. The step-and-scan method projects only a narrow rectangular stripe of the reticle image onto the wafer. The wafer and the reticle are scanned in tandem until the complete reticle pattern is transferred to the wafer. The wafer is then indexed to the next site, and the process of alignment and scanning proceeds again. Large die images or multiple dice can be patterned using this technique.

The minimum feature size that can be reproduced using optical lithography is intimately tied to the wavelength of light used for exposure, and experts have repeatedly predicted the demise of optical lithography for many years.³ However, the technology continues to be pushed further and further into the submicron regime. A

³The author remembers attending a number of lithography panel sessions where it was predicted that optical lithography could not be used below 1–2 μm .

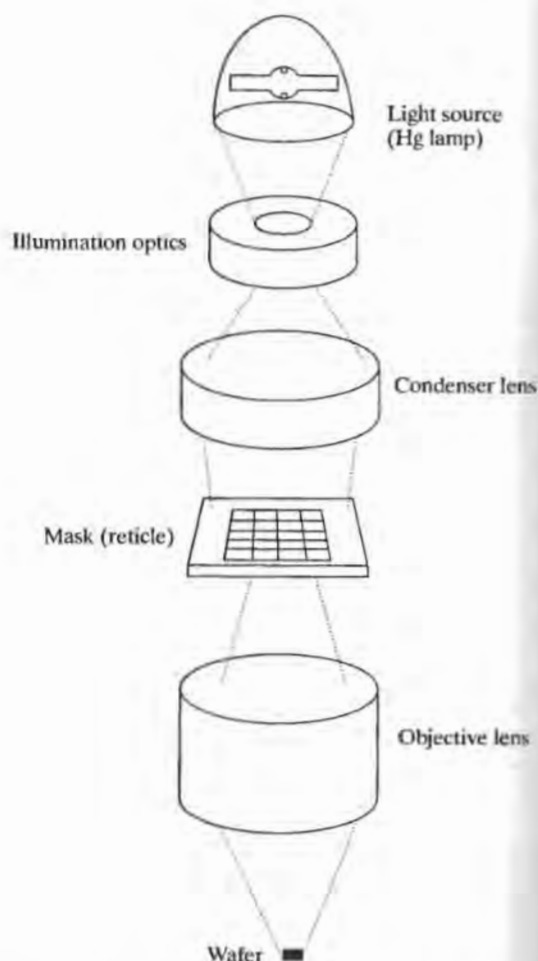


FIGURE 2.12

Concept of lens system for a wafer stepper.

rough estimate of the minimum feature size F (line or space) that can be transferred to the wafer surface is given by

$$F = 0.5 \frac{\lambda}{NA} \quad (2.1)$$

where λ is the wavelength of the illumination, and NA is the numerical aperture of the lens defined in terms of the convergence angle θ in Fig. 2.14 on page 33:

$$NA = \sin \theta. \quad (2.2)$$

For $NA = 0.5$, Eq. (2.1) predicts the minimum feature size to be approximately the same as the wavelength of the optical illumination. A second concern is the depth of field DF over which focus is maintained; an estimate for DF is

$$DF = 0.6 \frac{\lambda}{(NA)^2} \quad (2.3)$$



FIGURE 2.13

The true complexity of a wafer stepper is apparent in this system drawing. (Courtesy of ASM Lithography, Inc.)

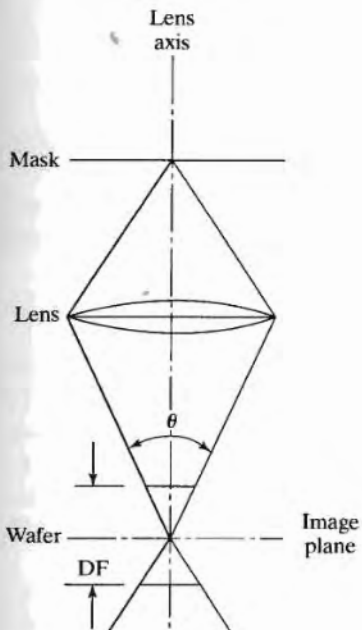


FIGURE 2.14

Optical focal plane and depth of focus

Based upon Eqs. (2.1) and (2.3) with $NA = 0.5$, one-tenth-micron technology ($F = 0.10 \mu\text{m}$) would require a 100-nm illumination source and have only a 240 nm ($0.24 \mu\text{m}$) depth of field. Thus, wafer planarity at each exposure step represents a critical issue, since focus will be maintained over a distance of only $\pm 0.12 \mu\text{m}$ from the primary focal plane.

2.5 EXPOSURE SOURCES

For many years, the source of illumination for photolithography has been high-pressure mercury (Hg) or Hg-rare gas discharge lamps. A typical emission spectra from a Hg-Xe lamp is given in Fig. 2.15. Output is relatively low in the deep ultraviolet (DUV) region (200–300 nm), but exhibits several strong peaks in the UV region between 300 and 450 nm. To minimize problems in the lens optics, the lamp output must be filtered to select one of the spectral components. The most common monochromatic selections are the 436-nm, or “g-line,” and 365-nm, or “i-line,” spectral components.

It can be observed in the NTRS lithography projections in Table 2.3, however, that DUV sources are required for lithography for 0.25- μm technology and below. Excimer lasers are the choice at these wavelengths with the KrF laser used as the 248-nm source and ArF for 193 nm. It is not clear what the lithography source will be for technology generations of below 130 nm.

Phase-shifting mask technology is representative of the inventions that have been found in the drive to squeeze the most out of optical lithography. The conceptual diagram in Fig. 2.16 compares the imaged light intensity profile of two closely spaced lines. In Fig. 2.16(a), resolution is lost, because diffraction has caused overlap of the individual line images. In Fig. 2.16(b), a 180° phase-shifting layer is applied over one of the openings on the mask, and the two individual lines appear well defined in the intensity profile at the image plane. A minimum feature size approaching one-half of the wavelength of the illumination source can be achieved using a phase-shifting mask and $NA \geq 0.5$. For highly complex IC patterns, however, designing the phase-shifting masks represents a significant challenge.

Various alternatives to optical lithography have been explored since the mid 1960s. Electron beams can be focused to spots of the order of $0.10 \mu\text{m}$ and can be used to directly write IC patterns in electron-sensitive resists. However, this process is relatively slow, since the pattern must be rewritten at each die site, and the throughput of electron-beam systems has never been sufficient for IC manufacturing. On the other hand, it is an excellent technology for producing the 1-X to 10-X reticles used in stepper systems, and “e-beam” lithography has become an extremely important technology for mask fabrication.

X-rays with energies in the 0.1–5-keV range have wavelengths that range from 10 to 0.3 nm. Thus, even the finest feature sizes in Table 2.3 would represent many wavelengths if x-rays were used for illumination. Mask generation is one of several barriers to the use of x-ray lithography in IC production. Heavy metals such as gold can be used as mask materials for x-ray lithography, but limit practical x-ray wavelengths to the 0.4–2-nm range. In addition to mask fabrication, x-ray lithography also requires a new set of illumination, resist, and alignment technologies.

Research efforts continue to expand the capabilities of electron-beam and x-ray lithography as outlined in Table 2.5 on page 36. Extreme ultraviolet [10], e-beam direct write, e-beam projection, x-ray proximity, and ion-beam projection all offer potential for future lithography systems. However, significant innovation will be required along the way to achieve the ITRS goals.

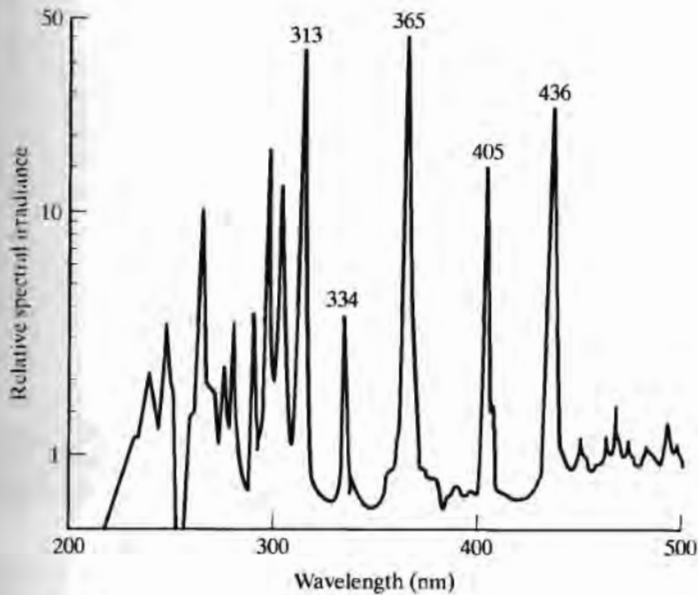


FIGURE 2.15

Spectral content of an Xe-Hg lamp (Courtesy of SVG)

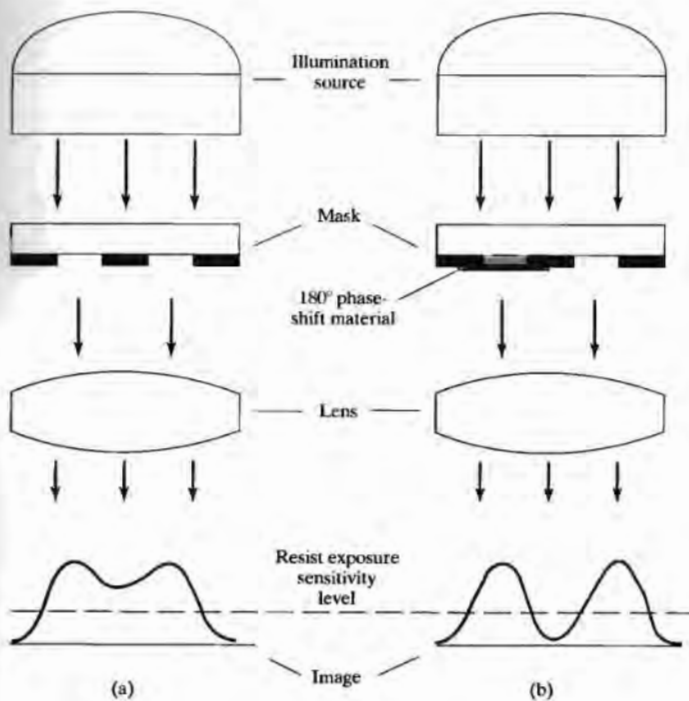


FIGURE 2.16

Pattern transfer of two closely spaced lines (a) using conventional mask technology (b) using a phase-shifting mask

TABLE 2.5 ITRS Lithography Projections

Year	2001	2003	2005	2008	2011	2014
Dense Line Half-Pitch (nm)	150	120	100	70	50	35
Worst-Case Alignment Tolerance Mean + 3 σ (nm)	52	42	35	25	20	15
Minimum Feature Size F (nm)						
Microprocessor Gate Width	100	80	65	45	30	20
Critical Dimension Control (nm)						
Mean + 3 σ Post Etching	9	8	6	4	3	2
Equivalent Oxide Thickness (nm)	1.5–1.9	1.5–1.9	1.0–1.5	0.8–1.2	0.6–0.8	0.5–0.6
Lithography Technology Options	248 nm DUV	248 nm + PSM 193 nm DUV	193 nm + PSM 157 nm E-beam projection Proximity x-ray Ion Projection	157 nm + PSM E-beam projection E-beam direct write EUV Ion Projection Proximity x-ray	EUV E-beam projection E-beam direct write Ion Projection	EUV E-beam projection E-beam direct write Ion Projection Innovation

DUV: deep ultraviolet; EUV: extreme ultraviolet; PSM: phase-shift mask.

2.6 OPTICAL AND ELECTRON MICROSCOPY

Most of us believe the old adage “a picture is worth a thousand words,” and visual inspection of masks and subsequent determination of the physical structure (morphology) of the patterns transferred to the wafer surface is extremely important in VLSI fabrication. Three methods—optical microscopy, Scanning Electron Microscopy (SEM), and Transmission Electron Microscopy (TEM)—find wide application in the visualization of VLSI morphologies and provide increasingly higher levels of magnification.

2.6.1 Optical Microscopy

Optical microscopes are a common laboratory tool familiar to most of us, and they are used to inspect and monitor the wafers throughout the fabrication process. The resolution of an optical microscope corresponds to the minimum feature size introduced in Section 2.4. Using white light for illumination with wavelengths centered on $0.5\ \mu\text{m}$, and with a numerical aperture of 0.95, the resolution is approximately $0.25\ \mu\text{m}$. On the other hand, the resolution of the human eye itself is approximately $0.25\ \mu\text{m}$. Hence, optical microscopes typically have a maximum magnification of 1,000X. The lower end of the magnification range is usually 1X–5X.

Analytical microscopes usually can operate in either the bright-field or dark-field mode. Bright-field operation is the mode that we most often encounter. The sample is illuminated by light perpendicular to the plane of the sample directly through the optics of the microscope. Light is reflected from the sample back up into same optical path in the microscope. For dark-field mode, the sample is illuminated from an oblique angle, and light that is reflected or refracted from features on the surface of the sample enters by the microscope lens system. The surface of the sample appears mostly dark with the surface features standing out in bright contrast against the dark background. In this manner, surface features that are “washed out” in bright-field mode can be clearly observed.

2.6.2 Scanning Electron Microscopy

In the Scanning Electron Microscope (SEM), the surface of the sample is bombarded with a low-energy ($0.5\text{--}40\ \text{KeV}$) beam of electrons. The incident electron beam causes low-energy ($0\text{--}50\ \text{eV}$) secondary electrons to be ejected from the inner shells of the atoms making up the surface of the sample under analysis. An image is formed by scanning the surface of the sample and recording the intensity of the secondary electron current. The magnitude of the secondary electron current depends upon the materials present and on the curvature of the surface, and significant contrast can be achieved due to varying surface morphology and materials. The SEM extends the minimum resolution limit to $20\text{--}30\ \text{\AA}$, with magnifications up to 300,000. At a magnification of 10,000, the SEM provides a depth of field of $2\text{--}4\ \mu\text{m}$, which makes it an extremely useful tool for investigating VLSI structures. As an example, the SEM image of a MEMS structure is shown in Fig. 2.17 on page 38. The image is a micro-mirror that has been raised out of the surface plane by a gear-driven linear actuator.

Some types of SEMs can suffer from electrical charge-up of the sample by the electron beam, particularly on insulating surfaces. This can be eliminated by coating the surface with a thin conducting layer of gold. However, this requires special processing of samples prior to their imaging by the SEM.

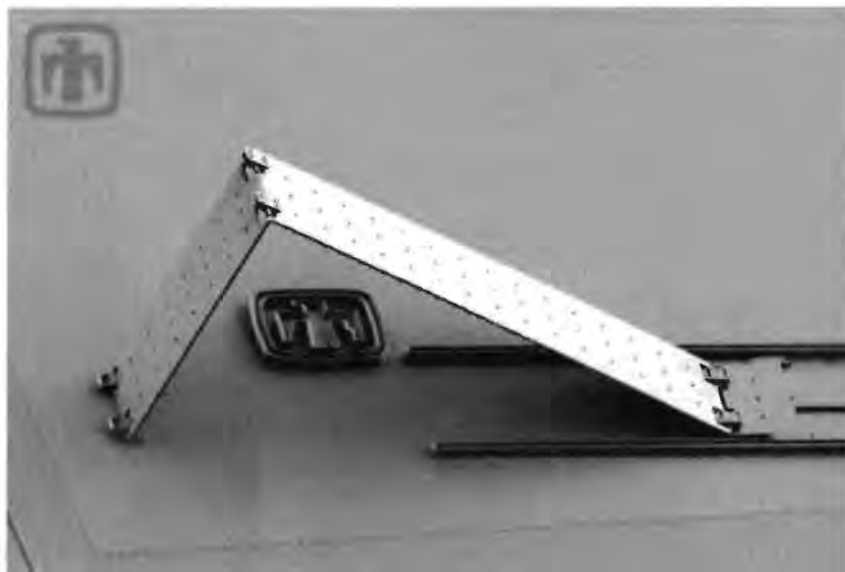


FIGURE 2.17

SEM image of a three-dimensional micromechanical system (MEMS) structure (Courtesy Sandia National Laboratories)

2.6.3 Transmission Electron Microscopy

The Transmission Electron Microscope (TEM) extends the resolution of microscopy another order of magnitude down to the 2\AA range, which corresponds to a distance below the radius of most atoms. In this instrument, a 60–400-KeV beam of electrons is used to illuminate a thin sample only $0.5\text{--}2\text{ }\mu\text{m}$ thick. The amplitude of the electron current that passes through the sample is detected, and an image is created as the beam scans the sample. In a MOS structure, for example, the TEM can display an image of the transition from the regular array of atoms in the silicon lattice to the irregular amorphous layer of the silicon dioxide gate insulator as depicted in Fig. 2.18. Although the TEM provides very high resolution, its application requires special preparation of the extremely thin samples.

SUMMARY

Photolithography is used to transfer patterns from masks to photoresist on the surface of silicon wafers. The resist protects portions of the surface while windows are etched in barrier layers such as silicon dioxide, silicon nitride, or metal. The windows may be etched using either wet- or dry-processing techniques. Wet chemical etching tends to etch under the edge of the mask, causing a loss of linewidth control at small

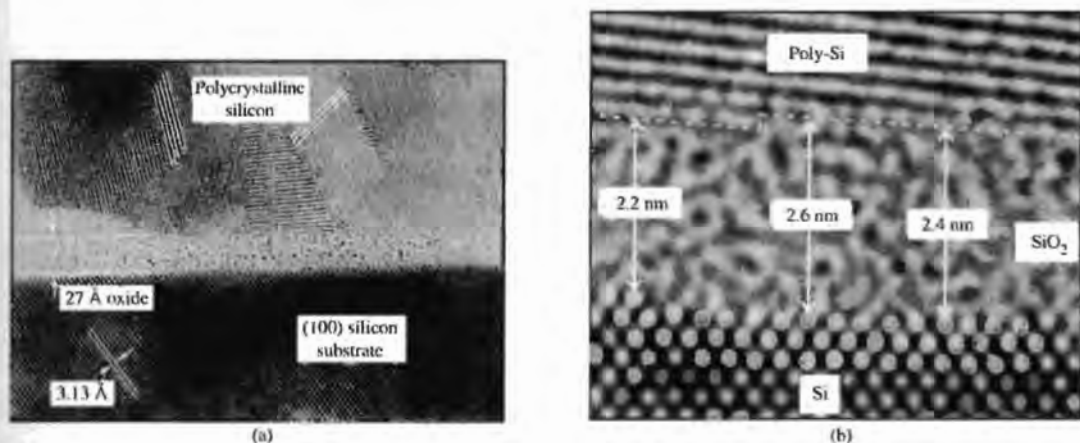


FIGURE 2.18

Cross-sectional high-resolution transmission electron microscope images for MOS structures with (a) 27-Å and (b) 24-Å image. The polysilicon grains are easily noticeable in (a); the Si/SiO₂ and poly-Si/SiO₂ interfaces are shown in part (b). On a local, atomic scale, thickness variations of 2-3 Å are found which are a direct result of atomic silicon steps at both interfaces. Copyright 1969 by International Business Machines Corporation; reprinted with permission from Ref. [9].

dimensions. Dry etching can yield highly anisotropic etching profiles and is required in most VLSI processing.

After etching, impurities can be introduced into the wafer through the windows using ion implantation or high-temperature diffusion, or metal can be deposited on the surface making contact with the silicon through the etched windows. Masking operations are performed over and over during IC processing, and the number of mask steps required is used as a basic measure of process complexity.

Mask fabrication uses computer graphics systems to draw the chip image at 100 to 2,000 times final size. Reticles 1 to 10 times final size are made from this computer image, using optical pattern generators or electron-beam systems. Step-and-repeat cameras are used to fabricate final masks from the reticles, or direct step-on-wafer systems may be used to transfer the patterns directly to the wafer.

Today, we are reaching the limits of optical lithography. Present equipment can define windows that are approximately 0.15 μm wide. (Just a few years ago, experts were predicting that 1-2 μm would be the limit!) The wavelength of light is too long to produce much smaller geometrical features, because of fringing and interference effects. Electron-beam and X-ray lithography are now being used to fabricate devices with geometrical features smaller than 0.10 μm, and lithography test structures have reproduced shapes with minimum feature sizes below 0.05 μm.

REFERENCES

- [1] *The International Technology Roadmap for Semiconductors*, The Semiconductor Industry Association (SIA), San Jose, CA: 1999, (<http://www.semichips.org>)
- [2] S. Nonogaki, T. Ueno, and T. Ito, *Microlithography Fundamentals in Semiconductor Devices and Fabrication Technology*, Marcel Dekker, Inc. New York, 1998.
- [3] W. R. Runyon and K. E. Bean, *Semiconductor Integrated Circuit Processing Technology*, Addison Wesley, Reading, MA, 1990.
- [4] L. F. Thompson, C. G. Wilson, and M. J. Bowden, Eds., *Introduction to Microlithography*, American Chemical Society, Washington, D.C., 1983.
- [5] S. M. Sze, Ed., *VLSI Technology*, McGraw-Hill, New York, 1983.
- [6] D. J. Elliot, *Integrated Circuit Fabrication Technology*, McGraw-Hill, New York, 1982.
- [7] W. Kern and D. A. Poutinen, "Cleaning Solutions Based upon Hydrogen Peroxide for Use in Silicon Semiconductor Technology," *RCA Review*, 31, 187–206 (June 1970).
- [8] W. Kern, "Purifying Si and SiO₂ Surfaces with Hydrogen Peroxide," *Semiconductor International*, pp. 94–99 (April 1984).
- [9] D. A. Buchanan et al., "Scaling the Gate Dielectric: Materials, Integration and Reliability," *IBM J. Research and Development*, vol. 43, no. 3, pp. 245–264, May 1999.
- [10] J. E. Bjorkholm, "EUV Lithography - The Successor to Optical Lithography?," *Micro Technology Journal*, vol. Q3'98, pp. 1–8, 1998.

FURTHER READING

- [1] P. Burggraaf, "Optical Lithography to 2000 and Beyond," *Solid-State Technology*, vol. 42, no. 2, pp. 31–41, February 1999.
- [2] J. A. McClay and A. S. L. McIntyre, "157 nm Optical Lithography: The Accomplishments and the Challenges," *Solid-State Technology*, vol. 42, no. 6, pp. 57–68, June 1999.
- [3] M. C. King, "Principles of Optical Lithography," in *VLSI Electronics*, vol. 1, N. G. Einspruch, Ed., Academic Press, New York, 1981.
- [4] J. H. Bruning, "A Tutorial on Optical Lithography," in *Semiconductor Technology*, D. A. Doane, D. B. Fraser, and D. W. Hess, Eds., Electrochemical Society, Pennington, NJ, 1982.
- [5] R. K. Watts and J. H. Bruning, "A Review of Fine-Line Lithographic Techniques: Present and Future," *Solid-State Technology*, 24, 99–105 (May 1981).
- [6] J. A. Reynolds, "An Overview of E-Beam Mask-Making," *Solid-State Technology*, 22, 87–94 (August 1979).
- [7] E. C. Douglas, "Advanced Process Technology for VLSI Circuits," *Solid-State Technology*, 24, 65–72 (May 1981).
- [8] L. M. Ephrath, "Etching Needs for VLSI," *Solid-State Technology*, 25, 87–92 (July 1982).
- [9] N. D. Wittels, "Fundamentals of Electron and X-ray Lithography," in *Fine Line Lithography*, R. Newman, Ed., North-Holland, Amsterdam, 1980.
- [10] D. Maydan, "X-ray Lithography for Microfabrication," *Journal of Vacuum Science and Technology*, 17, 1164–1168 (September/October 1980).

PROBLEMS

- 2.1 A complex CMOS fabrication process requires 25 masks. (a) What fraction of the dice must be good (i.e., what yield must be obtained) during each mask step if we require 30% of the final dice to be good? (b) How about if we require 70% to be good?
- 2.2 The mask set for a simple rectangular *pn* junction diode is shown in Fig. P2.2. The diode is formed in a *p*-type substrate. Draw a picture of the horizontal layout for the diode that

results when a worst-case misalignment of $3\text{ }\mu\text{m}$ occurs in both the x - and y -directions on each mask level.

- (a) Assume that both the contact and metal levels are aligned to the diffusion level.
- (b) Assume that the contact level is aligned to the diffusion level and the metal level is aligned to the contact level.

2.3 Figure P2.3 shows a resist pattern on top of a silicon dioxide film $1\text{ }\mu\text{m}$ thick. Draw the silicon-silicon dioxide structure after etching and removal of the photoresist for:

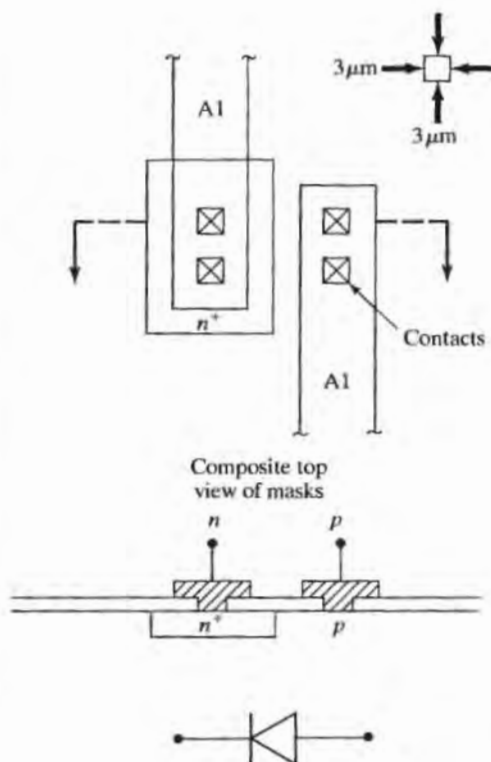


FIGURE P2.2

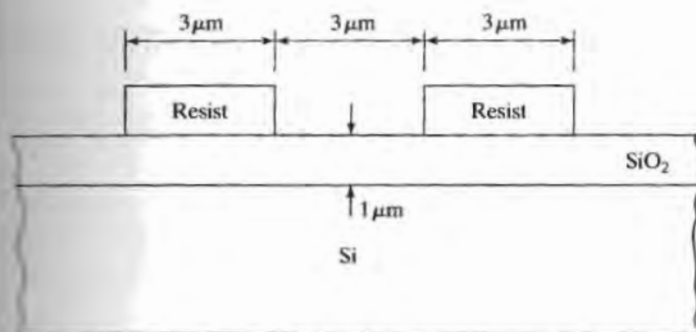


FIGURE P2.3

- (a) Isotropic wet chemical etching.
 - (b) Anisotropic dry etching with no undercutting.
- 2.4 (a)** What type of photoresist must be used with each of the three mask levels (n -diffusion window, contact windows, and metal etch) used to fabricate the diode of Problem 2.2? Assume that the areas shown are dark on the mask (a "light-field mask").
- (b)** Draw a set of alignment marks suitable for use with the alignment sequence of Problem 2.2(b).
- 2.5** In Table 2.5, 193-nm DUV lithography is shown as an alternative for 180-nm technology generation. (a) What value of NA would be required based upon Eq. 2.1? (b) What is the depth of field for this system?
- 2.6 (a)** What wavelength illumination is required to achieve $F = 0.25 \mu\text{m}$ with $NA = 1$ and without the use of phase-shifting masks? What is the value of DF corresponding to your value of NA ?
- (b)** Repeat for $NA = 0.5$.
- 2.7** Based upon the discussion in Section 2.4, what is the smallest feature size F that can be reproduced with a 193-nm optical source?
- 2.8** An extreme ultra violet (EUV) lithography source uses a 13-nm exposure wavelength. Based upon the discussion in Section 2.4, what is the smallest feature size F that can be reproduced with this source?

CHAPTER 3

Thermal Oxidation of Silicon

Upon exposure to oxygen, the surface of a silicon wafer oxidizes to form silicon dioxide. This native silicon dioxide film is a high-quality electrical insulator and can be used as a barrier material during impurity deposition. These two properties of silicon dioxide were the primary process factors leading to silicon becoming the dominant material in use today for the fabrication of integrated circuits. This chapter discusses the theory of oxide growth, the oxide growth processes, factors affecting oxide growth rate, impurity redistribution during oxidation, and techniques for selective oxidation of silicon. Methods for determining the thickness of the oxide film are also presented, and the SUPREM process simulation software is introduced.

3.1 THE OXIDATION PROCESS

Thermal oxidation of silicon is easily achieved by heating the wafer to a high temperature, typically 900 to 1200 °C, in an atmosphere containing either pure oxygen or water vapor. Both water vapor and oxygen move (diffuse) easily through silicon dioxide at these high temperatures. (See Fig. 3.1.) Oxygen arriving at the silicon surface can then combine with silicon to form silicon dioxide. The chemical reaction occurring at the silicon surface is



for dry oxygen and



for water vapor. Silicon is consumed as the oxide grows, and the resulting oxide expands during growth, as shown in Fig. 3.2. The final oxide layer is approximately 54% above the original surface of the silicon and 46% below the original surface.

FIGURE 3.1

Diffusivities of hydrogen, oxygen, sodium, and water vapor in silicon glass. Copyright John Wiley & Sons, Inc. Reprinted with permission from Ref. [4].

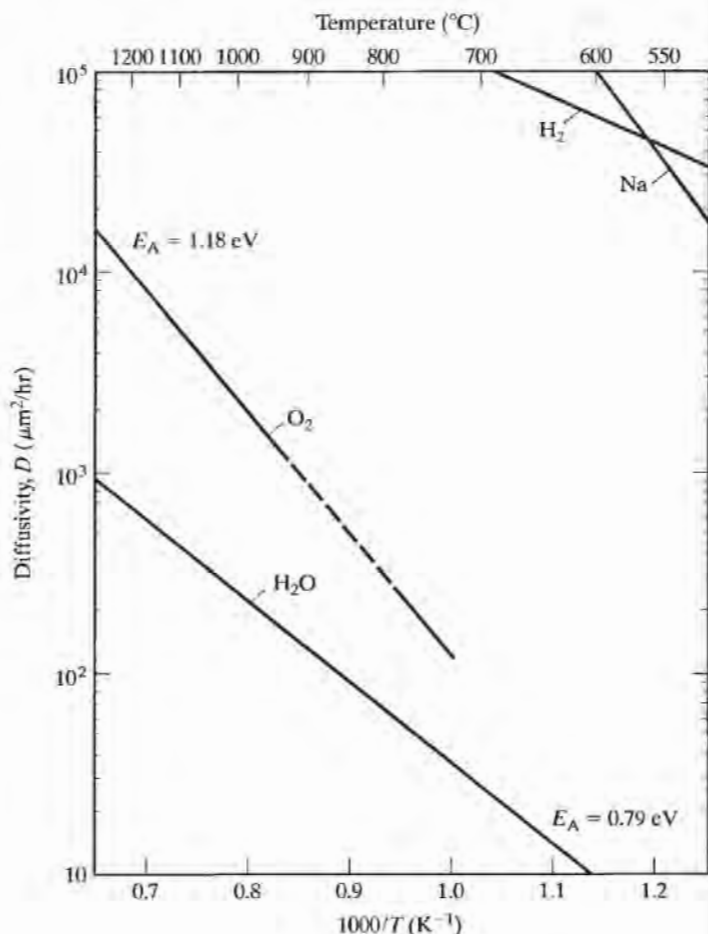
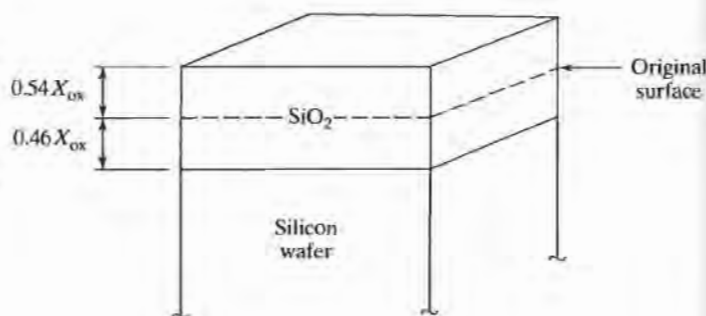


FIGURE 3.2

Formation of a silicon dioxide layer on the surface of a silicon wafer consumes silicon during growth of the layer. The oxide expands to fill a region approximately 54% above and 46% below the original surface of the wafer. The exact percentages depend on the density of the oxide.



3.2 MODELING OXIDATION

In order for oxidation to occur, oxygen must reach the silicon interface. As the oxide grows, oxygen must pass through more and more oxide, and the growth rate decreases as time goes on. A simple model for oxidation can be developed by assuming that

oxygen diffuses through the existing oxide layer. Fick's first law of diffusion states that the particle flow per unit area, J (called particle flux), is directly proportional to the concentration gradient of the particle:

$$J = -D \partial N(x, t) / \partial x, \quad (3.3)$$

where D is the diffusion coefficient and N is the particle concentration. The negative sign indicates that particles tend to move from a region of high concentration to a region of low concentration.

For our case of silicon oxidation, we will make the approximation that the oxygen flux passing through the oxide in Fig. 3.3 is constant everywhere in the oxide. (Oxygen does not accumulate in the oxide.) The oxygen flux J is then given by

$$J = -D(N_i - N_0) / X_o \quad (\text{number of particles/cm}^2 - \text{sec}), \quad (3.4)$$

where X_o is the thickness of the oxide at a given time, and N_0 and N_i are the concentrations of the oxidizing species in the oxide at the oxide surface and silicon dioxide-silicon interface, respectively. At the silicon dioxide-silicon interface, we assume that the oxidation rate is proportional to the concentration of the oxidizing species so that the flux at the interface is

$$J = k_s N_i, \quad (3.5)$$

where k_s is called the rate constant for the reaction at the Si-SiO₂ interface. Eliminating N_i using Eqs. (3.4) and (3.5), we find that the flux J becomes

$$J = DN_0 / (X_o + D/k_s). \quad (3.6)$$

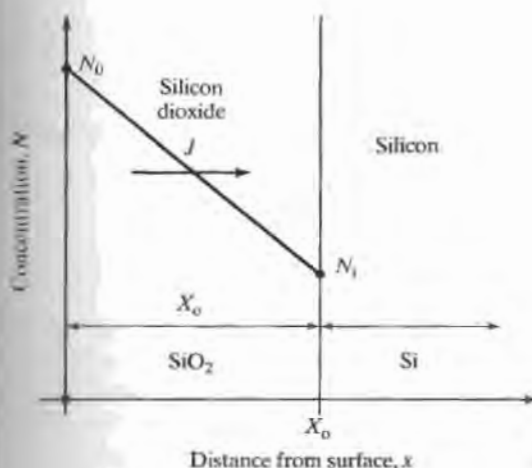


FIGURE 3.3

Model for thermal oxidation of silicon. X_o is the thickness of the silicon dioxide layer at any time t . J is the constant flux of oxygen diffusing through the layer, and N_0 and N_i represent the oxygen concentration at the oxide surface and silicon dioxide-silicon interface, respectively. Note that the oxide growth occurs at the silicon interface.

The rate of change of thickness of the oxide layer with time is then given by the oxidizing flux divided by the number of molecules M of the oxidizing species that are incorporated into a unit volume of the resulting oxide:

$$dX_o/dt = J/M = (DN_o/M)/(X_o + D/k_s). \quad (3.7)$$

This differential equation is easily solved using the boundary condition $X_o(t=0) = X_i$, which yields

$$t = X_o^2/B + X_o/(B/A) - \tau, \quad (3.8)$$

where $A = 2D/k_s$, $B = 2DN_o/M$, and $\tau = X_i^2/B + X_i/(B/A)$. X_i is the initial thickness of oxide on the wafer, and τ represents the time which would have been required to grow the initial oxide. A thin native oxide layer (10 to 20 Å) is always present on silicon due to atmospheric oxidation, or X_i may represent a thicker oxide grown during previous oxidation steps. Solving Eq. (3.8) for $X_o(t)$ yields

$$X_o(t) = 0.5A \left[\left\{ 1 + \frac{4B}{A^2}(t + \tau) \right\}^{1/2} - 1 \right]. \quad (3.9)$$

For short times with $(t + \tau) \ll A^2/4B$,

$$X_o(t) = (B/A)(t + \tau). \quad (3.10)$$

Oxide growth is proportional to time, and the ratio B/A is called the linear (growth) rate constant. In this region, growth rate is limited by the reaction at the silicon interface.

For long times with $(t + \tau) \gg A^2/4B$, $t \gg \tau$

$$X_o = \sqrt{Bt}. \quad (3.11)$$

Oxide growth is proportional to the square root of time, and B is called the *parabolic rate constant*. The oxidation rate is diffusion limited in this region.

3.3 FACTORS INFLUENCING OXIDATION RATE

There is good experimental agreement with this simple theory. Figures 3.4 (page 47) and 3.5 (page 48) show experimental data for the parabolic and linear rate constants. The rate-constant data follow straight lines when plotted on a semilogarithmic scale versus reciprocal temperature. This type of behavior occurs in many natural systems and is referred to as an Arrhenius relationship. A mathematical model for this behavior is as follows:

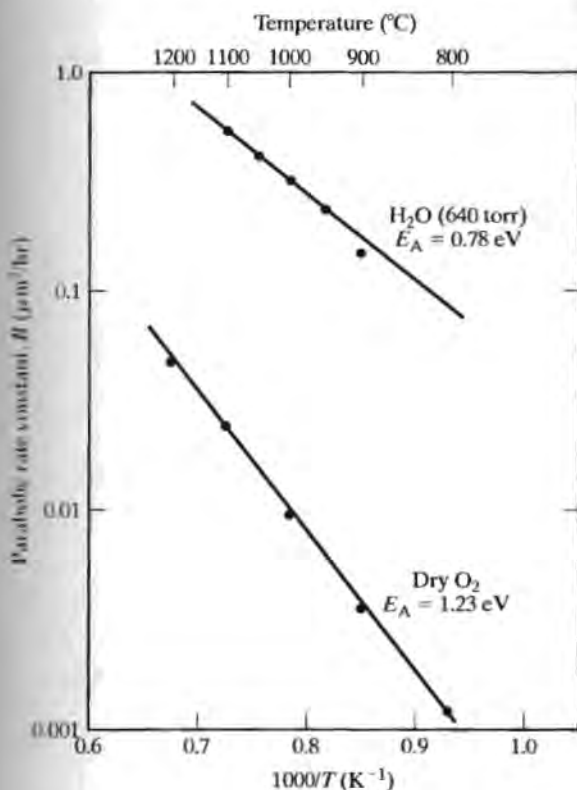


FIGURE 3.4

Dependence of the parabolic rate constant B on temperature for the thermal oxidation of silicon in pyrogenic H_2O (640 torr) or dry O_2 . Reprinted by permission of the publisher, The Electrochemical Society, Inc., from Ref. [10].

TABLE 3.1 Values for Coefficient D_0 and Activation Energy E_A for Wet and Dry Oxygen*

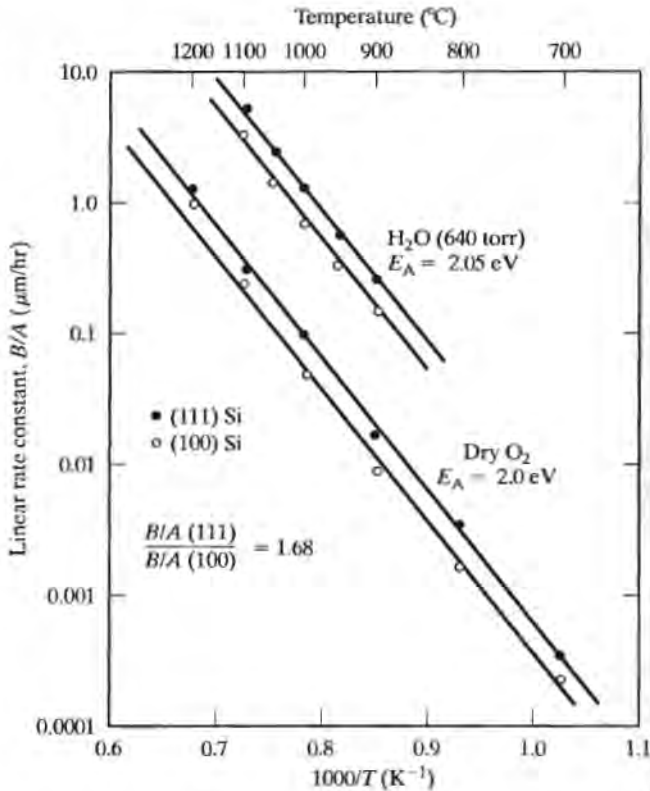
	Wet O_2 ($X_i = 0$ nm)		Dry O_2 ($X_i = 25$ nm)	
	D_0	E_A	D_0	E_A
<100> Silicon				
Linear (B/A)	$9.70 \times 10^7 \mu\text{m/hr}$	2.05 eV	$3.71 \times 10^6 \mu\text{m/hr}$	2.00 eV
Parabolic (B)	$386 \mu\text{m}^2/\text{hr}$	0.78 eV	$772 \mu\text{m}^2/\text{hr}$	1.23 eV
<111> Silicon				
Linear (B/A)	$1.63 \times 10^8 \mu\text{m/hr}$	2.05 eV	$6.23 \times 10^6 \mu\text{m/hr}$	2.00 eV
Parabolic (B)	$386 \mu\text{m}^2/\text{hr}$	0.78 eV	$772 \mu\text{m}^2/\text{hr}$	1.23 eV

*Data from Ref.[9]

Values for the coefficient D_0 and activation energy E_A for wet and dry oxygen are given in Table 3.1. For wet oxidation, a plot of the experimental data of oxide thickness versus oxidation time is consistent with an initial oxide thickness of approximately zero at $t = 0$. However, a similar plot for dry oxidation yields an initial oxide thickness of 250 Å for temperatures ranging from 800 to 1200 °C. Thus, a nonzero value for τ must be used in Eq. (3.8) for dry oxidation calculations. This large value of X_i indicates that our simple oxidation theory is not quite correct, and the reason for this value of X_i

FIGURE 3.5

Dependence of the linear rate constant B/A on temperature for the thermal oxidation of silicon in pyrogenic H_2O (640 torr) or dry O_2 . Reprinted by permission of the publisher, The Electrochemical Society, Inc., from Ref. [10].



is not well understood. Graphs of oxide growth versus time, calculated using the values from Table 3.1, are given in Figs. 3.6 and 3.7 on page 49.

Equation (3.12) indicates the strong dependence of oxide growth on temperature. A number of other factors affect the oxidation rate, including wet and dry oxidation, pressure, crystal orientation, and impurity doping. Water vapor has a much higher solubility than oxygen in silicon dioxide, which accounts for the much higher oxide growth rate in a wet atmosphere. Slower growth results in a denser, higher quality oxide and is usually used for MOS gate oxides. More rapid growth in wet oxygen is used for thicker masking layers.

Equation (3.8) shows that both the linear and parabolic rate constants are proportional to N_0 . N_0 is proportional to the partial pressure of the oxidizing species, so pressure can be used to control oxide growth rate. There is great interest in developing low-temperature processes for VLSI fabrication. High pressure is being used to increase oxidation rates at low temperatures. (See Fig. 3.8 on page 51.) In addition, very thin oxides (50 to 200 Å) are required for VLSI, and low-pressure oxidation is being investigated as a means of achieving controlled growth of very thin oxides.

Figures 3.4 through 3.7 also show the dependence of oxidation rate on substrate crystal orientation for the $\langle 111 \rangle$ and $\langle 100 \rangle$ materials most commonly used in bipolar and MOS processes, respectively. The crystal orientation changes the number of silicon bonds available at the silicon surface, which influences the oxide growth rate and quality of the silicon-silicon dioxide interface.

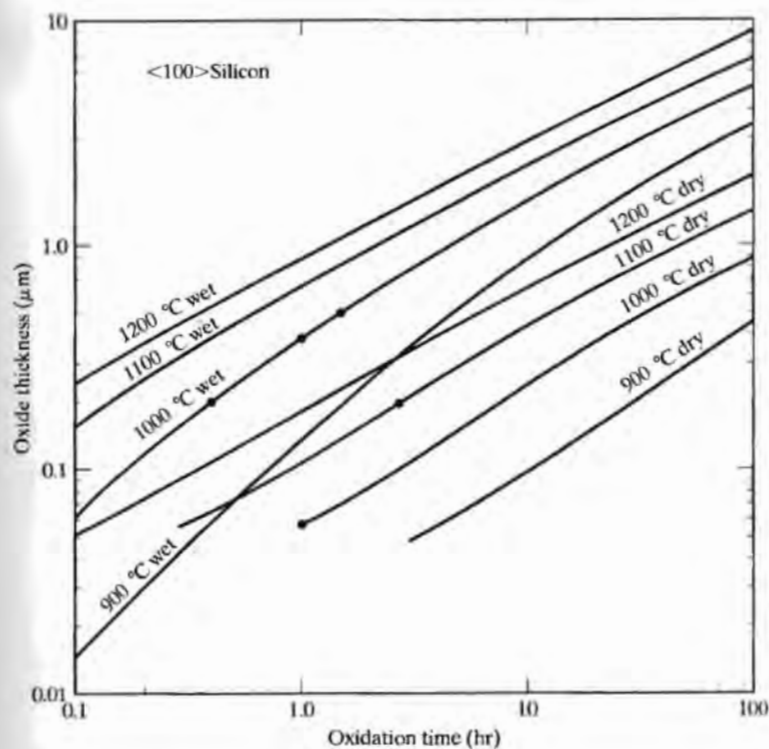


FIGURE 3.6

Wet and dry silicon dioxide growth for <100> silicon calculated using the data from Table 3.1. (The dots represent data used in examples.)

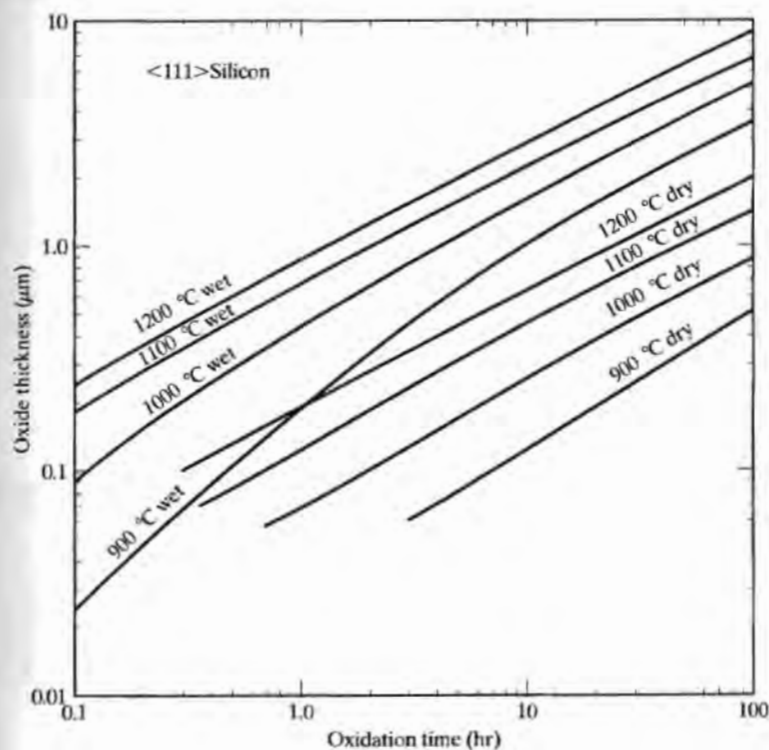
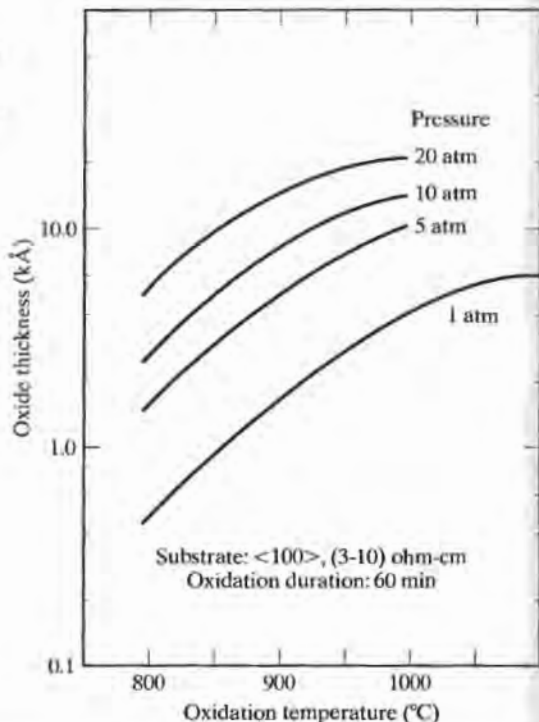


FIGURE 3.7

Wet and dry silicon dioxide growth for <111> silicon calculated using the data from Table 3.1.

FIGURE 3.8

Wet oxide growth at increased pressures. Reprinted with permission of Solid State Technology, published by Technical Publishing, a company of Dun and Bradstreet, from Ref. [12].



Example 3.1

According to Fig. 3.6, a 1-hr oxidation of <100> silicon at 1000 °C in dry oxygen will produce a silicon dioxide film approximately 580 Å (0.058 μm) thick. The same oxidation in wet oxygen will yield a film 3900 Å (0.39 μm) thick.

Example 3.2

A <100> wafer has a 2000-Å oxide on its surface.

- How long did it take to grow this oxide at 1100 °C in dry oxygen?
- The wafer is put back in the furnace in wet oxygen at 1000 °C. How long will it take to grow an additional 3000 Å of oxide? Solve this problem graphically using Figs. 3.6 and 3.7 as appropriate.
- Repeat part (b) using the oxidation theory presented in Eqs. (3.3) through (3.12).

Solution: (a) According to Fig. 3.6, it would take 2.8 hr to grow a 0.2-μm oxide in dry oxygen at 1100 °C.

(b) We can solve part (b) graphically using Fig. 3.6. The total oxide at the end of the oxidation would be 0.5 μm. If there were no oxide on the surface, it would take 1.5 hr to grow 0.5 μm. However, there is already a 0.2 μm oxide on the surface, and the wafer "thinks" that it has already been in the furnace for 0.4 hr. The time required to grow the additional 0.3 μm of oxide is the difference in these two times: $\Delta t = (1.5 - 0.4) \text{ hr} = 1.1 \text{ hr}$.

(c) From Table 3.1, $B = 3.86 \times 10^2 \exp(-0.78/kT) \text{ } \mu\text{m}^2/\text{hr}$ and $(B/A) = 0.97 \times 10^8 \exp(-2.05/kT) \text{ } \mu\text{m}/\text{hr}$. Using $T = 1273 \text{ K}$ and $k = 8.617 \times 10^{-5} \text{ eV/Kg}$, $B = 0.314 \text{ } \mu\text{m}^2/\text{hr}$

and $(B/A) = 0.738 \mu\text{m/hr}$. Using these values and an initial oxide thickness of $0.2 \mu\text{m}$ yields a value of 0.398 hr for the effective initial oxidation time τ . Using τ and a final oxide thickness of $0.5 \mu\text{m}$ yields an oxidation time of 1.08 hr . Note that both the values of t and τ are close to those found in part (b). Of course, the graphical results depend on our ability to interpolate logarithmic scales!

Heavy doping of silicon also changes its oxidation characteristics. Phosphorus doping increases the linear rate constant without altering the parabolic rate constant. Boron doping, on the other hand, increases the parabolic rate constant but has little effect on the linear rate constant. These effects are related to impurity redistribution during oxidation, which is discussed in the next section.

3.4 DOPANT REDISTRIBUTION DURING OXIDATION

During oxidation, the impurity concentration changes in the silicon near the silicon-silicon dioxide interface. Boron and gallium tend to be depleted from the surface, whereas phosphorus, arsenic, and antimony pile up at the surface.

Impurity depletion and pileup depend on both the diffusion coefficient and the segregation coefficient of the impurity in the oxide. The segregation coefficient m is equal to the ratio of the equilibrium concentration of the impurity in silicon to that of the impurity in the oxide. Various possibilities are depicted in Fig. 3.9 on page 52. The value of m for boron is temperature dependent and is less than 0.3 at normal diffusion temperatures. Boron also diffuses slowly through SiO_2 . Thus, boron is depleted from the silicon surface and remains in the oxide. (See Fig. 3.9(a).) The presence of hydrogen during oxide growth or impurity diffusion greatly enhances the diffusion of boron through oxide, resulting in enhanced depletion of boron at the silicon surface. (See Fig. 3.9(b).)

The value of m is approximately 10 for phosphorus, antimony, and arsenic. These elements are rejected by the oxide, and they diffuse slowly in the oxide, resulting in pileup at the silicon surface. (See Fig. 3.9(c).) In contrast, gallium has a segregation coefficient of 20, but it diffuses very rapidly through silicon dioxide. This combination causes gallium to deplete at the surface, as shown in Fig. 3.9(d).

The effects of boron depletion and phosphorus and arsenic pileup are particularly important in both bipolar and MOS processing. Process design must take both problems into account, and it is often necessary to add or change processing steps to overcome the effects of these phenomena.

3.5 MASKING PROPERTIES OF SILICON DIOXIDE

One of the most important properties of silicon dioxide is its ability to mask impurities during high-temperature diffusion. The diffusivities of antimony, arsenic, boron, and phosphorus in silicon dioxide are all orders of magnitude smaller than their corresponding values in silicon. Thus, SiO_2 films can be used effectively as a barrier layer to these elements. Relatively deep diffusion can take place in unprotected regions of silicon, whereas no significant impurity penetration will occur in regions covered by silicon dioxide.

FIGURE 3.9

The effects of oxidation on impurity profiles. (a) Slow diffusion in oxide (boron); (b) fast diffusion in oxide (boron with hydrogen ambient); (c) slow diffusion in oxide (phosphorus); (d) fast diffusion in oxide (gallium). C_B is the bulk concentration in the silicon. Copyright John Wiley & Sons, Inc. Reprinted with permission from Ref. [5].

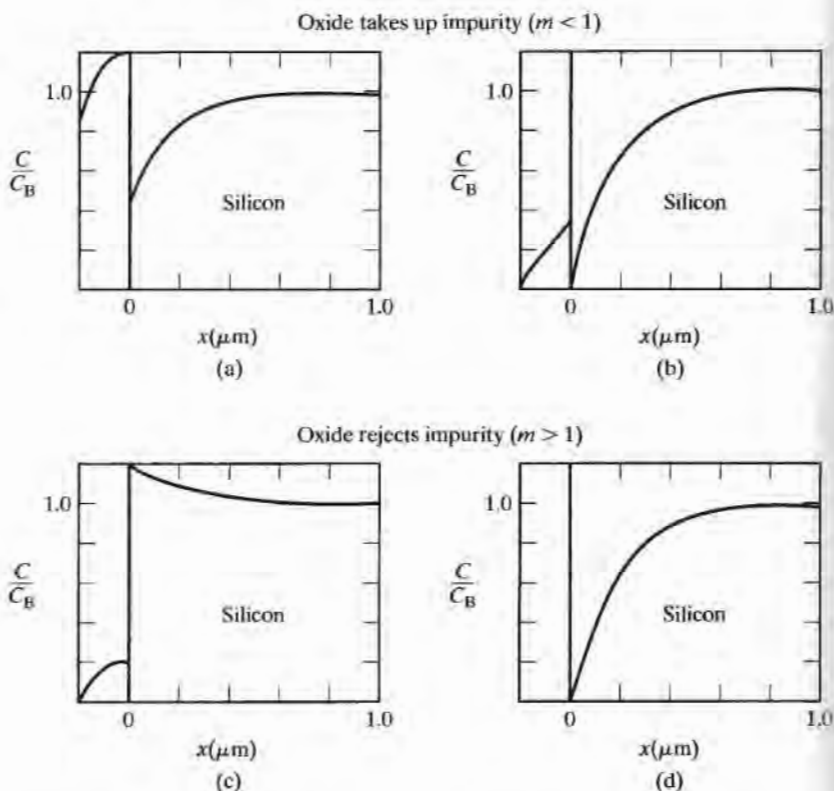


Figure 3.10 gives the SiO_2 thickness required to mask boron and phosphorus diffusions as a function of diffusion time and temperature. Note that silicon dioxide is much more effective in masking boron than in masking phosphorus. Arsenic and antimony diffuse more slowly than phosphorus, so an oxide thick enough to mask phosphorus is also sufficient to mask arsenic and antimony. Masking oxide thicknesses of 0.5 to 1.0 μm are typical in IC processes. The masking oxide would be considered to have failed if the impurity level under the mask were to reach a significant fraction (10%) of the background concentration in the silicon.

The graph for boron is valid for an environment that contains no hydrogen! As mentioned earlier, the presence of hydrogen greatly enhances the boron diffusivity. Wet oxidation releases hydrogen, and care must be taken to avoid boron diffusion in the presence of water vapor.

As mentioned in Section 3.4, gallium diffuses rapidly through SiO_2 , as does aluminum, and silicon dioxide cannot be used as a barrier for these elements. However, silicon nitride can be used effectively to prevent diffusion of these impurities.

3.6 TECHNOLOGY OF OXIDATION

Thermal oxidation of silicon is typically carried out in a high-temperature furnace. The furnace walls may be made of quartz, polycrystalline silicon, or silicon carbide and are specially fabricated to prevent sodium contamination during oxidation. The furnaces

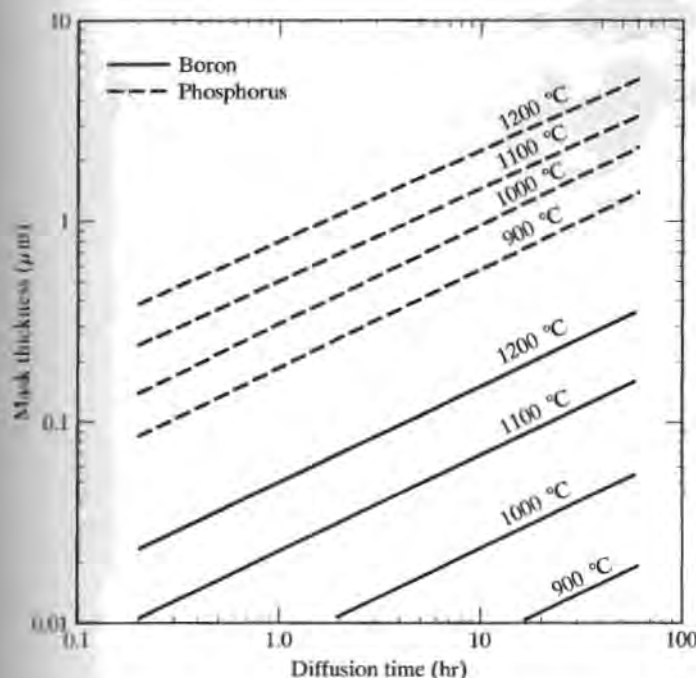


FIGURE 3.10

Thickness of silicon dioxide needed to mask boron and phosphorus diffusions as a function of diffusion time and temperature.

used for many years were manufactured with horizontal oxidation tubes. The wafers are placed upright on edge in a quartz boat and pushed slowly into the furnace. The furnace is maintained at a temperature between 800 and 1200 °C. Three-zone resistance-heated furnaces maintain the temperature within a fraction of a degree over a distance of 0.5 m in the center zone. Vertical furnaces are now commonly used for processing wafers with diameters of 150 mm and larger. A photograph of horizontal and vertical furnaces used for oxidation and diffusion appears in Fig. 3.11.

The furnace is continually purged with an inert gas such as nitrogen prior to oxidation. Oxidation begins by introducing the oxidizing species into the furnace in gaseous form. Extremely high-purity oxygen is available and is used for dry oxidation. Water vapor may be introduced by passing oxygen through a bubbler containing deionized water heated to 95 °C. The oxygen serves as a transport gas to carry the water vapor into the furnace. High-purity water vapor can also be obtained by burning hydrogen and oxygen in the furnace tube. Steam is not often used because it tends to pit the silicon surface.

3.7 OXIDE QUALITY

Wet oxidation is used to grow relatively thick oxides used for masking. An oxidation growth cycle usually consists of a sequence of dry-wet-dry oxidations. Most of the oxide is grown during the wet oxidation phase, since the growth rate is much higher in the presence of water. Dry oxidation results in a higher density oxide than that achieved with wet oxidation. Higher density in turn results in a higher breakdown voltage (5 to 10 MV/cm). To maintain good process control, the thin gate oxides (<1000 Å) of MOS devices are usually formed using dry oxidation.

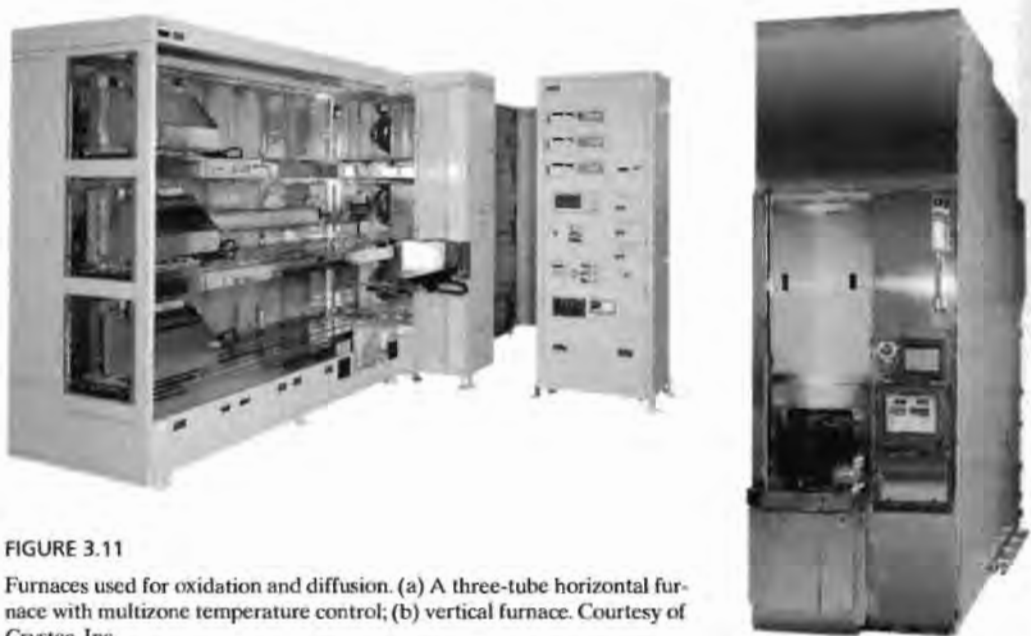


FIGURE 3.11

Furnaces used for oxidation and diffusion. (a) A three-tube horizontal furnace with multizone temperature control; (b) vertical furnace. Courtesy of Crystec, Inc.

MOS devices are usually fabricated on wafers having a $\langle 100 \rangle$ surface orientation. The $\langle 100 \rangle$ orientation results in the smallest number of unsatisfied silicon bonds at the Si-SiO₂ interface, and the choice of the $\langle 100 \rangle$ orientation yields the lowest number of interface traps.*

Sodium ions are highly mobile in SiO₂ films (see Fig. 3.1), and contamination of MOS gate oxides was a difficult problem to overcome in the early days of the integrated-circuit industry.* Bipolar devices are much more tolerant of oxide contamination than MOS devices, and this was a primary factor in the early dominance of bipolar integrated circuits.

Sodium-ion contamination results in mobile positive charge in the oxide. In addition, a substantial level of positive fixed oxide exists at the Si-SiO₂ interface.* These charge centers attract electrons to the surface of MOS transistors, resulting in a negative shift in the threshold voltage of the MOS devices. NMOS devices tend to become depletion-mode devices. PMOS devices remain enhancement-mode devices, but have more negative threshold voltages. The first successful MOS processes were therefore PMOS processes. As the industry was able to improve overall oxide quality, NMOS processes became dominant because of the mobility advantage of electrons over holes.

It was discovered that the effects of sodium contamination can be greatly reduced by adding chlorine during oxidation. Chlorine is incorporated into the oxide and immobilizes the sodium ions. A small amount (6% or less) of anhydrous HCl can be added to the oxidizing gas. Gaseous chlorine, oxygen, or nitrogen can also be bubbled through trichloroethylene (C₂HCl₃). It should also be noted that the presence of chlorine during dry oxidation results in an increase in both the linear and parabolic rate constants.

*See Volume IV in the Modular Series on Solid State Devices, *Field Effect Devices*, Chapter 4, for an excellent discussion of oxide quality.

SELECTIVE OXIDATION AND SHALLOW TRENCH FORMATION

The oxidation processes described above generally form an oxide film over the complete surface of the silicon wafer. The ability to selectively oxidize the silicon surface has become very important in high-density bipolar and MOS processes. Selective oxidation processes result in improved device packing density and more planar final structures. The techniques utilized for localized oxidation of silicon are generally referred to as LOCOS processes.

Oxygen and water vapor do not diffuse well through silicon nitride. Figure 3.12 shows a MOS process using selective oxidation in which silicon nitride is used as an oxidation barrier. A thin layer (10 to 20 nm) of silicon dioxide is first grown on the wafer to protect the silicon surface. Next, a layer of silicon nitride is deposited over the surface and patterned using photolithography. The wafer then goes through a thermal oxidation step. Oxide grows wherever the wafer is not protected by silicon nitride. This process results in the so-called *semirecessed oxide structure*.

Some oxide growth occurs under the edges of the nitride and causes the nitride to bend up at the edges of the masked area. The penetration of the oxide underneath the nitride results in a "bird's beak" structure. Formation of the bird's beak in Fig. 3.12 leads to loss of geometry control in VLSI structures, so minimization of the bird's beak phenomenon is an important goal in VLSI process design.

A *fully recessed oxide* can be formed by etching the silicon prior to oxidation. This process can yield a very planar surface after the removal of the nitride mask. However, subsequent processing reduces the advantage of this process over the semi-recessed version, and most processes today use some form of semirecessed oxidation.

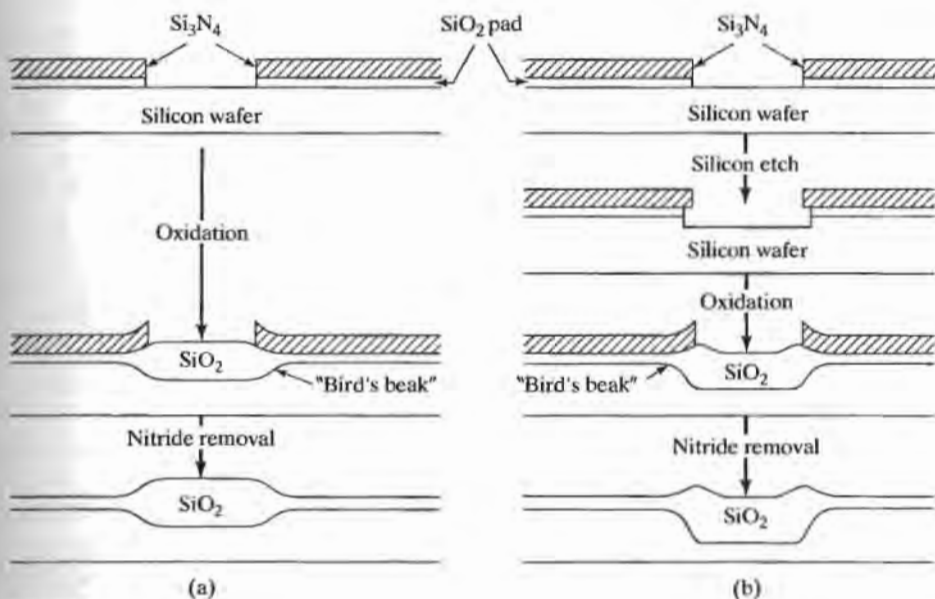


FIGURE 3.12

Cross section depicting process sequence for local oxidation of silicon (LOCOS): (a) semirecessed and (b) fully recessed structures.

3.8.1 Trench Isolation

Some form of shallow or deep refilled trench isolation is utilized in most of today's advanced MOS and bipolar processes. Figure 3.13(a) [18] on page 57 depicts formation of deep trenches filled with polysilicon in combination with a LOCOS field oxidation. A thin oxide pad is grown on silicon followed by deposition of a silicon nitride layer. Lithography is used to define openings in the nitride where trenches will be formed. The trenches are etched using reactive-ion etching and can be quite deep with very high aspect ratios. The surface of the trench is passivated with a thin layer of thermally grown oxide, and then the trench is refilled with deposited polysilicon. The final structure is produced by etching back any excess polysilicon, using a lithography step to remove the nitride layer where oxidation is desired, and growing the semirecessed oxide layer. The polysilicon may be doped, and similar structures are used to form trench capacitors for use as storage elements in some DRAM technologies [19]. The capacitor is formed between the polysilicon and the substrate, and the trenches may be as much as 5–10 μm in depth.*

Shallow trench isolation is used to provide isolation between transistors in most MOS and bipolar technologies (see Chapters 9 and 10) with feature sizes below 0.5 μm . Figure 3.13(b) depicts one possible process flow for forming shallow trenches [20]. A shallow trench with tapered sidewalls is etched in the silicon following patterning of the nitride layer. The pad oxide may be etched away slightly to round the corners of the final structure. A thin oxide layer is grown as a liner on the trench walls, and the trench is then filled with an oxide deposited using decomposition of TEOS or via a high-density plasma deposition.* A process called chemical mechanical polishing (CMP), discussed in more detail in the next section, is used to remove the excess oxide and create a highly planar surface. In the final step, the nitride may be removed, leaving the shallow trench isolation between two silicon regions. Figure 3.14 on page 58 shows the use of both deep and shallow trench isolation in an advanced bipolar process.

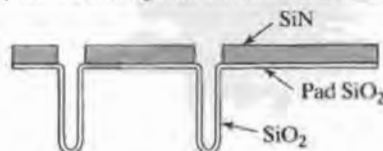
3.8.2 Chemical Mechanical Polishing (CMP)

CMP [21–24] was introduced into fabrication processing during the early 1990s and is now widely used in both bipolar and MOS processes to achieve the highly planar topologies required in deep submicron lithography. A conceptual diagram of a CMP apparatus is given in Fig. 3.15 on page 58. The wafer is mounted on a carrier and is brought into contact with a polishing pad mounted on a rotating platen. A liquid slurry is continuously dispensed onto the surface of the polishing pad. A combination of the vertical force between the wafer and the abrasive pad as well as the chemical action of the slurry is used to polish the surface to a highly planar state. In the case of formation of the shallow trenches, the nitride layer serves as a polishing stop. Polishing terminates when the nitride layer is fully exposed.

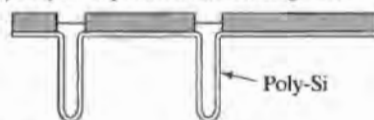
Actual structures differ slightly from the idealized planar surface, as depicted in Fig. 3.16 on page 59. A slight amount of “dishing” of the oxide and corner rounding can both occur due to polishing rate differences between the various materials. Some erosion of the nitride layer may also occur prior to process endpoint detection.

* See Chapter 6.

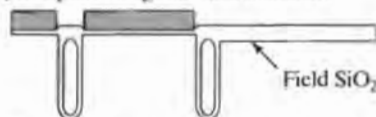
(1) Trench etching with SiN mask and oxidation



(2) Poly-Si deposition and etching back

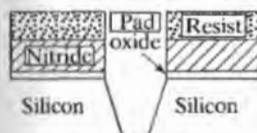


(3) SiN patterning and field oxidation

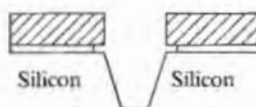


Fabrication procedure of trench isolation and field oxide.

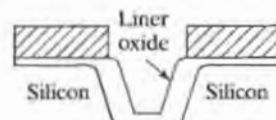
(a) Deep-trench process



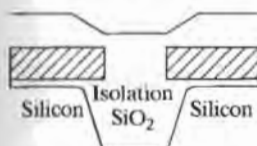
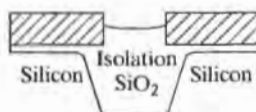
(1) Stack & trench etch



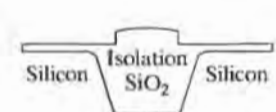
(2) Pad oxide undercut



(3) Liner oxidation

(4) CVD Oxide gap fill
(e.g. HDP, TEOS-O₂)

(5) CMP & HF dip

(6) H₂ PO₄ Nitride strip

(b) Steps in a typical STI process flow

FIGURE 3.13

Trench isolation structures. (a) Deep trench isolation. Copyright 1996 IEEE. Reprinted with permission from Ref. [18]. (b) Shallow trench isolation. Copyright 1998 IEEE. Reprinted with permission from Ref. [20].

3.9 OXIDE THICKNESS CHARACTERIZATION

One of the simplest methods for determining the thickness of an oxide is to compare the color of the wafer with the reference color chart in Table 3.2 on page 60. When a wafer is illuminated with white light perpendicular to the surface, the light penetrates the oxide film and is reflected by the underlying silicon wafer. Constructive interference causes enhancement of a certain wavelength in the reflected light, and the color

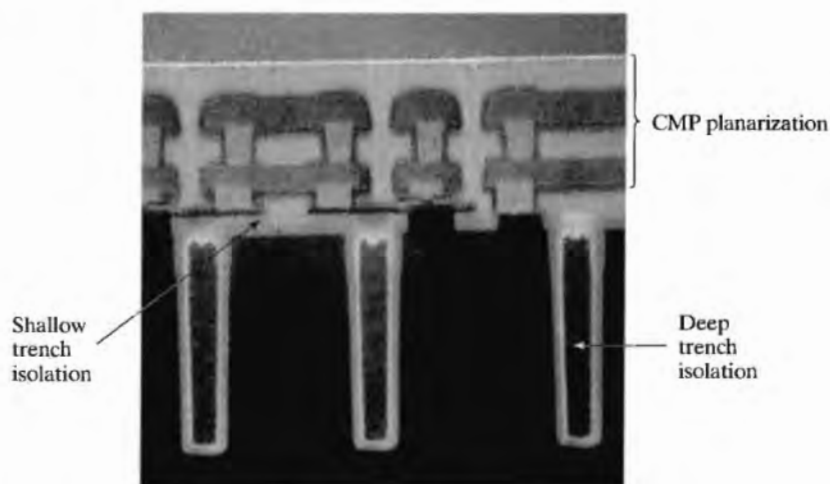


FIGURE 3.14

Microphotograph of actual deep and shallow trench applied to SiGe HBT technology. Copyright 1998 IEEE. Reprinted with permission from Ref. [31].

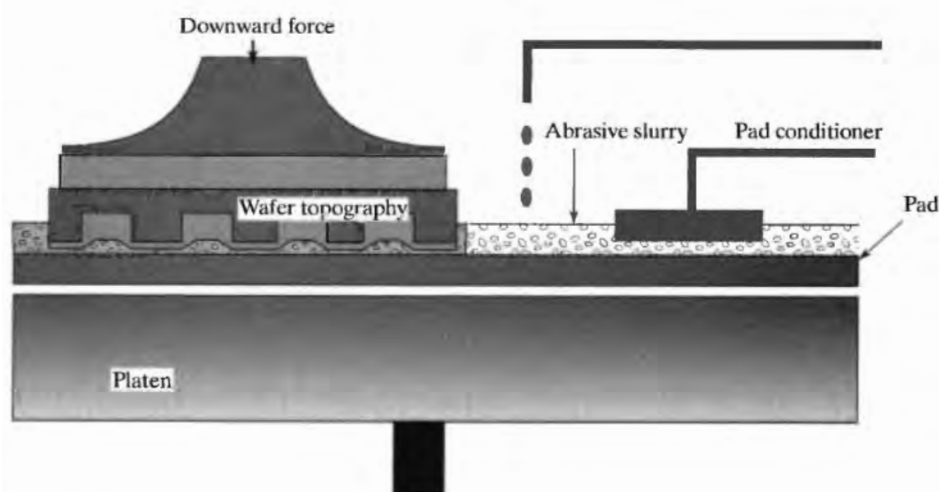


FIGURE 3.15

Chemical mechanical polishing technique.



(a)



(b)

FIGURE 3.16

Multilevel metallization fabricated with chemical mechanical polishing.

(a) SEM of 6-level thin-wire copper. First-level copper is connected by Tungsten studs to Tungsten local interconnect. (b) SEM of 6-level copper with low RC metallization on levels 5 and 6. Copyright 1997 IEEE. Reprinted with permission from Ref. [24].

of the wafer corresponds to the enhanced wavelength. Constructive interference occurs when the path length in the oxide ($2X_o$) is equal to an integer multiple of one wavelength of light in the oxide.

$$2X_o = k\lambda/n, \quad (3.13)$$

where the number k is any integer greater than zero and n is the refractive index of the oxide ($n = 1.46$ for SiO_2).

As an example, a wafer with a $5000\text{-}\text{\AA}$ silicon dioxide layer will appear blue green. Color-chart comparisons are quite subjective, and the colors vary periodically with thickness. In addition, care must be exercised to determine the color from a position perpendicular to the wafer. The color chart (Table 3.2) is only valid for vertical illumination with fluorescent light.

Accurate thickness measurement can be achieved with an instrument called an *ellipsometer*, and this instrument is often used to make an accurate reference color chart. Polarized monochromatic light is used to illuminate the wafer at an angle to the surface. Light is reflected from both the oxide and silicon surfaces. The differences in polarization are measured, and the oxide thickness can then be calculated [17].

A mechanical surface profiler can also be used to measure film thickness. The oxide is partially etched from the surface of a test wafer to expose a step between the wafer and oxide surfaces. A stylus is mechanically scanned over the surface of the wafer, and thickness variations are recorded by a computer. Films ranging from less than $0.01\text{ }\mu\text{m}$ to more than $5\text{ }\mu\text{m}$ can be measured with this instrument.

TABLE 3.2 Color Chart for Thermally Grown SiO_2 Films Observed Perpendicularly Under Daylight Fluorescent Lighting. Copyright 1964 by International Business Machines Corporation; reprinted with permission from Ref. [11]

Film Thickness (μm)	Color and Comments	Film Thickness (μm)	Color and Comments
0.05	Tan	0.63	Violet red
0.07	Brown	0.68	"Bluish" (not blue but borderline between violet and blue green; appears more like a mixture between violet red and blue green and looks grayish)
0.10	Dark violet to red violet	0.72	Blue green to green (quite broad)
0.12	Royal blue	0.77	"Yellowish"
0.15	Light blue to metallic blue	0.80	Orange (rather broad for orange)
0.17	Metallic to very light yellow green	0.82	Salmon
0.20	Light gold or yellow; slightly metallic	0.85	Dull, light red violet
0.22	Gold with slight yellow orange	0.86	Violet
0.25	Orange to melon	0.87	Blue violet
0.27	Red violet	0.89	Blue
0.30	Blue to violet blue	0.92	Blue green
0.31	Blue	0.95	Dull yellow green
0.32	Blue to blue green	0.97	Yellow to "yellowish"
0.34	Light green	0.99	Orange
0.35	Green to yellow green	1.00	Carnation pink
0.36	Yellow green	1.02	Violet red
0.37	Green yellow	1.05	Red violet
0.39	Yellow	1.06	Violet
0.41	Light orange	1.07	Blue violet
0.42	Carnation pink	1.10	Green
0.44	Violet red	1.11	Yellow green
0.46	Red violet	1.12	Green
0.47	Violet	1.18	Violet
0.48	Blue violet	1.19	Red violet
0.49	Blue	1.21	Violet red
0.50	Blue green	1.24	Carnation pink to salmon
0.52	Green (broad)	1.25	Orange
0.54	Yellow green	1.28	"Yellowish"
0.56	Green yellow	1.32	Sky blue to green blue
0.57	Yellow to "yellowish" (not yellow but is in the position where yellow is to be expected; at times appears to be light creamy gray or metallic)	1.40	Orange
0.58	Light orange or yellow to pink	1.45	Violet
0.60	Carnation pink	1.46	Blue violet
		1.50	Blue
		1.54	Dull yellow green

Accurate film thickness measurements can also be achieved using light-interference effects in microscopy, and automated interference-based equipment is commercially available for thin-film characterization.

3.10 PROCESS SIMULATION

As the scale of integrated circuits is reduced, accurate knowledge of one-, two-, and three-dimensional structural details of the process becomes more and more important. At the same time, experimental examination and characterization of the structures is a difficult and time-consuming task in deep submicron fabrication. For these reasons, computer simulation continues to grow in importance throughout the VLSI fabrication process.

Sophisticated computer programs that can predict the results of various fabrication steps have been available for many years [25–30]. These programs solve the generalized differential equations that model various fabrication processes and include the ability to simulate multidimensional oxide growth with its attendant moving Si–SiO₂ boundary, impurity segregation during oxide growth, and dopant evaporation from the surface as described in this chapter, as well as the various deposition, diffusion, epitaxial growth, and ion implantation processes studied in upcoming chapters. The detailed structures resulting from etching and recessed oxidation can also be simulated. Other programs have been developed to model lithography processes such as photoresist exposure and development.

One of the most widely used of these programs is the Stanford University Process Engineering Modeling program (SUPREM) [25–27] and its commercial implementations. The use of SUPREM requires specification of the process steps including times, temperature profiles, and other ambient conditions for oxidation, diffusion, ion implantation, film deposition, and etching. The program can model the one- and two-dimensional structures resulting from oxidation, as well as predicting the impurity profiles in the substrate, oxide, and polysilicon layers.

An example of a one-dimensional oxidation simulation is given in the SUPREM IV listing in Table 3.3 on page 63. The input listing describes a complex dry-wet-dry oxidation cycle, including the ramp-up of the furnace from one temperature to another and various ambient gas conditions at different steps in the oxidation cycle. The output data includes the oxide thickness and impurity dose in the oxide. Graphical output of the data is shown in Fig. 3.17. Incorporation of boron in the oxide and its depletion from the substrate are both clearly evident in the plotted results. Note that the oxide extends both above and below the original silicon surface defined as $x = 0$ on the horizontal axis. Most of the oxide growth occurs during the steam oxidation. The width of the oxide agrees well with a simple estimate from Fig. 3.6, based upon a single 5-hr wet oxidation at 1100 C: $X_o = 1.5 \mu\text{m}$.

SUMMARY

Silicon dioxide provides a high-quality insulating barrier on the surface of the silicon wafer. In addition, this layer can serve as a barrier layer during subsequent impurity-diffusion process steps. These two factors have allowed silicon to become the dominant semiconductor material in use today.

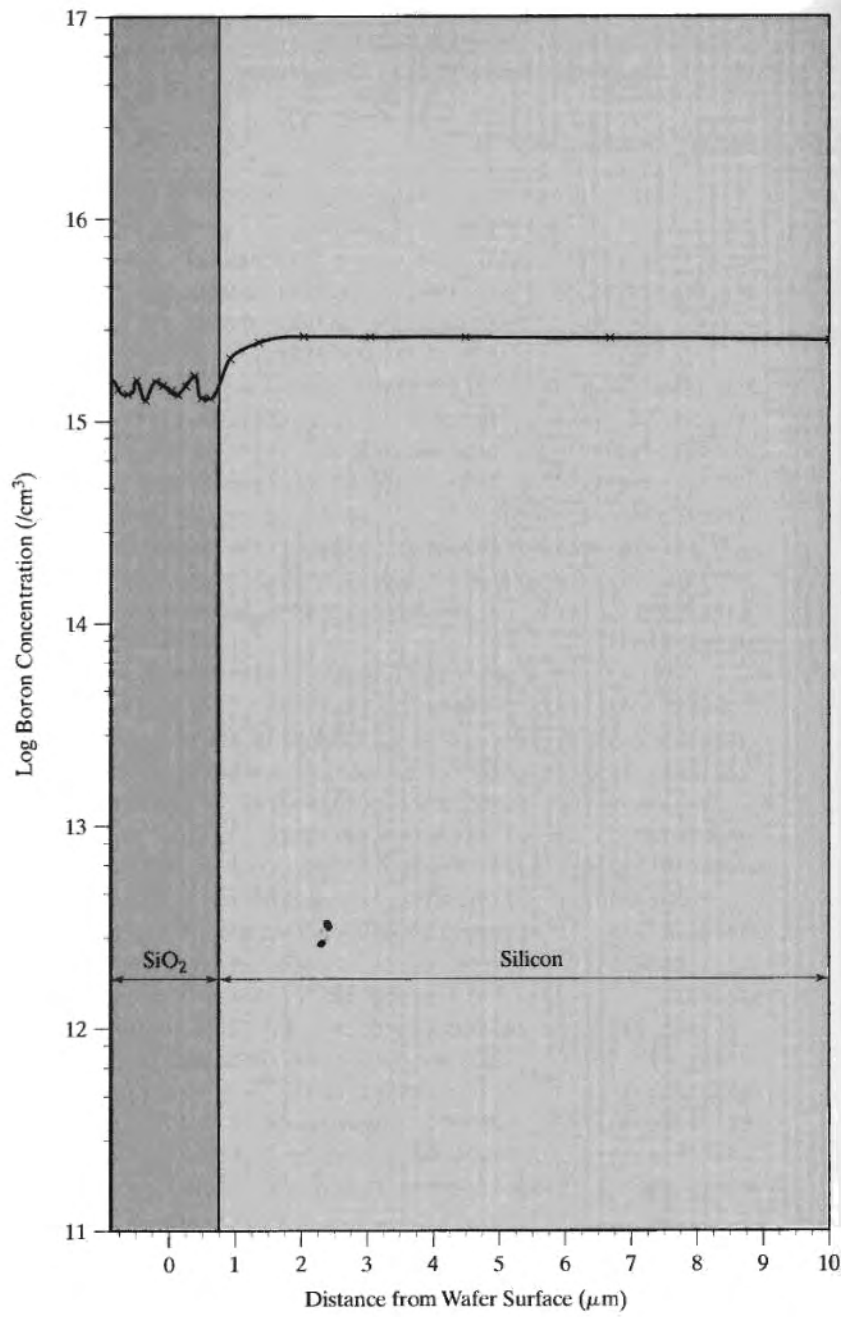


FIGURE 3.17
Results of SUPREM simulation of oxide growth on boron doped silicon wafers.

TABLE 3.3 SUPREM-IV Simulation Example

```

$ Multistep Oxidation

$ Use Automatic Grid Generation and Adaptive Grid

INITIALIZE <100> BORON=5 RESISTIV
DIFFUSION TEMP=950 TIME=30 F.N2=5
DIFFUSION TEMP=950 TIME=30 T.FINAL=1100 F.O2=5
DIFFUSION TEMP=1100 TIME=300 STEAM
DIFFUSION TEMP=1100 TIME=60 F.O2=5
DIFFUSION TEMP=1100 TIME=60 T.FINAL=800 F.N2=5
$ Print layer information
----
----
$ Plot results
----
----

```

A native oxide layer several tens of angstroms thick forms on the surface of silicon immediately upon exposure to oxygen even at room temperature. The thickness of this oxide layer may be readily measured from the accumulation-region capacitance of a MOS test capacitor. Thicker layers of silicon dioxide are conveniently grown in high-temperature oxidation furnaces using both wet and dry oxygen. Oxidation occurs much more rapidly in wet oxygen than in dry oxygen. However, the dry-oxygen environment produces a higher quality oxide and is usually used for the growth of MOS gate oxides. Thin oxides grow in direct proportion to time. However, as the oxide becomes thicker, growth rate slows and becomes proportional to the square root of time. These two growth regions are characterized by the linear and parabolic growth-rate constants.

Oxide cleanness is extremely important for MOS processes, and great care is exercised to prevent sodium contamination of the oxide. The addition of chlorine during oxidation improves oxide quality. Finally, oxidation alters the impurity distribution at the surface of the silicon wafer. Boron tends to be depleted from the silicon surface, whereas phosphorus tends to pile up at the silicon surface.

Oxidation thickness can be accurately measured using ellipsometers, interference microscopes, and mechanical surface profilers or can be estimated from the apparent color of the oxide under vertical illumination with white light.

REFERENCES

- [1] W. R. Runyon and K. E. Bean, *Semiconductor Integrated Circuit Processing Technology*, Addison Wesley, Reading, MA, 1990.
- [2] S. A. Campbell, *The Science and Engineering of Microelectronic Fabrication*, Oxford University Press, 1996.
- [3] S. M. Sze, *VLSI Technology*, McGraw-Hill, New York, 1983.
- [4] S. K. Ghandhi, *VLSI Fabrication Principles*, John Wiley & Sons, New York, 1983.
- [5] R. A. Colclaser, *Microelectronics—Processing and Device Design*, John Wiley & Sons, New York, 1980.

- [6] W. R. Runyan, *Semiconductor Measurements and Instrumentation*, McGraw-Hill, New York, 1975.
- [7] S. K. Ghandhi, *The Theory and Practice of Microelectronics*, John Wiley & Sons, New York, 1968.
- [8] B. E. Deal and A. S. Grove, "General Relationship for the Thermal Oxidation of Silicon," *Journal of Applied Physics*, 36, 3770-3778 (December, 1965).
- [9] B. E. Deal, "Thermal Oxidation Kinetics of Silicon in Pyrogenic H_2O and 5% HCl/H_2O Mixtures," *Journal of the Electrochemical Society*, 125, 576-579 (April, 1978).
- [10] B. E. Deal, "The Oxidation of Silicon in Dry Oxygen, Wet Oxygen and Steam," *Journal of the Electrochemical Society*, 110, 527-533 (June, 1963).
- [11] W. A. Pliskin and E. E. Conrad, "Nondestructive Determination of Thickness and Refractive Index of Transparent Films," *IBM Journal of Research and Development*, 8, 43-51 (January, 1964).
- [12] S. C. Su, "Low Temperature Silicon Processing Techniques for VLSIC Fabrication," *Solid-State Technology*, 24, 72-82 (March, 1981).
- [13] A. S. Grove, O. Leistiko, and C. T. Sah, "Redistribution of Acceptor and Donor Impurities During Thermal Oxidation of Silicon," *Journal of Applied Physics*, 35, 2695-2701 (September, 1964).
- [14] A. S. Grove, *Physics and Technology of Semiconductor Devices*, John Wiley & Sons, New York, 1967.
- [15] E. H. Nicollian and J. R. Brews, *MOS Physics and Technology*, John Wiley & Sons, New York, 1982.
- [16] M. Ghezzi and D. M. Brown, "Diffusivity Summary of B, Ga, P, As and Sb in SiO_2 ," *Journal of the Electrochemical Society*, 120, 146-148 (January, 1973).
- [17] E. Passaglia, R. R. Stromberg, and J. Kruger, Eds., *Ellipsometry in the Measurement of Surfaces and Thin Films*, National Bureau of Standards, Miscellaneous Publication #256, 1964.
- [18] H. Kotaki et al., "Novel Bulk Dynamic Threshold Voltage MOSFET (B-DTMO) with Advanced Isolation (SITOS) and Gate to Shallow-Well Contact (SSS-C) Processes for Ultra Low Power Dual Gate CMOS," *IEEE IEDM Technical Digest*, pp. 459-462, December 1996.
- [19] K. P. Muller, B. Flietner, C. L. Hwang, R. L. Kleinhenz, T. Nakao, R. Ranade, Y. Tsunashima, and T. Mii, "Trench Storage Node Technology for Gigabit DRAM Generations," *IEEE IEDM Technical Digest*, pp. 507-510, December 1996.
- [20] M. Nandakumar, A. Chatterjee, S. Sridhar, K. Joyner, M. Rodder, and I-C. Chen, "Shallow Trench Isolation for Advanced ULSI CMOS Technologies," *IEEE IEDM Technical Digest*, pp. 133-136, December 1998.
- [21] J. M. Steigerwald et al., "Pattern Geometry in Chemical-Mechanical Polishing of Inlaid Copper Structures," *Journal of the Electrochemical Society*, 141, no. 10: 2842-2848, Oct. 1994.
- [22] B. Stine et al., "Rapid Characterization and Modeling of Pattern-Dependent Variation in Chemical-Mechanical Polishing," *IEEE Trans. Semiconductor Manufacturing*, 11, no. 1, February 1998.
- [23] H. Nojo, M. Kodera, and R. Nakata, "Slurry Engineering for Self-Stopping Dishing Free SiO_2 -CMP," *IEEE IEDM Technical Digest*, pp. 349-352, December 1996.
- [24] D. Edelstein et al., "Full Copper Wiring in a Sub-0.25 μm CMOS ULSI Technology," *IEEE IEDM Technical Digest*, pp. 773-776, December 1997.
- [25] D. A. Antoniadis and R. W. Dutton, "Models for Computer Simulation of Complete IC Fabrication Processes," *IEEE Journal of Solid-State Circuits*, SC-14, pp. 412-422, April 1979.
- [26] C. P. Ho, J. D. Plummer, S. E. Hansen, and R. W. Dutton, "VLSI Process Modeling—SUPREM III," *IEEE Trans. Electron Devices*, ED-30, pp. 1438-1453, November 1983.

- [27] D. Chin, M. Kump, H.G. Lee, and R. W. Dutton, "Process Design Using Coupled 2D Process and Device Simulators," *IEEE IEDM Technical Digest*, pp. 223–226, December 1986.
- [28] C. D. Maldonado, F. Z. Custode, S. A. Louie, and R. K. Pancholy, "Two-Dimensional Simulation of a 2 μ m CMOS Process Using ROMANS II," *IEEE Trans. Electron Devices*, ED-30, pp. 1462–1469, November 1983.
- [29] M. E. Law, C. S. Rafferty, and R. W. Dutton, "New n-well Fabrication Techniques Based upon 2D Process Simulation," *IEEE IEDM Technical Digest*, pp. 518–521, December 1986.
- [30] R. W. Dutton, "Modeling and Simulation for VLSI," *IEEE IEDM Technical Digest*, pp. 2–7, December 1986.
- [31] John D. Cressler, "Re-engineering Silicon: Si-Ge Heterojunction Bipolar Technology," *IEEE Spectrum*, pp. 49–55, March 1995.

PROBLEMS

- 3.1 How long does it take to grow 100 nm of oxide in wet oxygen at 1000 °C (assume 100 silicon)? In dry oxygen? Which process would be preferred?
- 3.2 A 1.2- μ m silicon dioxide film is grown on a <100> silicon wafer in wet oxygen at 1100 °C. How long does it take to grow the first 0.4 μ m? The second 0.4 μ m? The final 0.4 μ m?
- 3.3 Derive Eq. (3.8) by solving differential Eq. (3.7).
- 3.4 A 3- μ m silicon dioxide film is grown on a <100> silicon wafer in wet oxygen at 1150 °C. How long does it take to grow the initial 1 μ m oxide? The second micron? The final micron?
- 3.5 The gate oxide for a CMOS process on <100> silicon is to have a thickness of 10 nm (100 Å). Calculate the time required to grow this oxide at 850 °C in wet oxygen using Eq. 3.8. Repeat for 1000 °C. (Be careful!) Do either of these possible processes seem to be controllable? If so, which one.
- 3.6 A 2- μ m SiO₂ film is needed as the initial oxide on a <100> silicon wafer. (a) Find the growth time in wet oxygen at 1150 °C using Fig. 3.6. (b) Use Eq. (3.8) to calculate the growth time.
- 3.7 A 1- μ m SiO₂ film is needed as the initial oxide on a <100> silicon wafer. (a) Find the growth time in wet oxygen at 1050 °C using Fig. 3.6. (b) Repeat the calculation for dry oxygen. (c) Use Eq. (3.8) to calculate the growth times.
- 3.8 A dry-wet-dry oxidation cycle of 30 min./120 min./30 min. is performed at 1100 °C. (a) What is the final oxide thickness for a <100> silicon wafer? Use Eq. 3.8. (b) What is the final oxide thickness for a <111> silicon wafer?
- 3.9 An one-hour dry oxidation at 1000 °C is followed by a 5-hour wet oxidation at 1100 °C. (a) Calculate the oxide thickness after each step for a <100> wafer. (b) Find the final oxide thickness graphically.
- 3.10 An one-hour dry oxidation at 1100 °C is followed by a 5-hour wet oxidation at 1100 °C. (a) Calculate the oxide thickness after each step for a <111> wafer. (b) Find the final oxide thickness graphically.
- 3.11 A square window is etched through 200 nm of oxide prior to a second oxidation, as in Example 3.2. The second oxidation grows 300 nm of oxide in the thick oxide region. Make a scale drawing of the cross section of the wafer after the second oxidation. What are the colors of the various regions under vertical illumination by white light?
- 3.12 A 10- μ m square window is etched through a 1- μ m thick oxide on a silicon wafer. The wafer is reoxidized to grow a new 1- μ m thick oxide film in the window. (a) Draw a cross section of the wafer following the second oxidation. (b) What are the colors of the oxide under vertical illumination by white light?

- 3.13 A $\langle 100 \rangle$ silicon wafer has 400 nm of oxide on its surface. How long will it take to grow an additional $1\mu\text{m}$ of oxide in wet oxygen at 1100°C ? Compare graphical and mathematical results. What is the color of the final oxide under vertical illumination by white light?
- 3.14 How much oxide is needed to mask a 4-hr boron diffusion at 1150°C ? A 1-hr phosphorus diffusion at 1050°C ?
- 3.15 The isolation diffusion in a bipolar process uses a 15-hour boron diffusion at 1150°C into a $\langle 111 \rangle$ silicon wafer. How much oxide is required as a barrier layer for this diffusion?
- 3.16 The n -well in a $\langle 100 \rangle$ CMOS process is formed with a 20-hour phosphorus diffusion at 1200°C . How much oxide is required as a barrier layer for this diffusion?
- 3.17 What is the color of the oxide in Problem 3.7? (b) How about in Problem 3.6?
- 3.18 Yellow light has a wavelength of approximately $0.57\mu\text{m}$. Calculate the thicknesses of silicon dioxide that will appear yellow under vertical illumination by white light. Consider oxide thicknesses less than $1.5\mu\text{m}$. Compare with the color chart (Table 3.2).
- 3.19 Write a computer program to calculate the linear and parabolic rate constants for wet and dry oxidation for temperatures of 950, 1000, 1050, 1100, 1150, and 1200 C. Assume $\langle 100 \rangle$ silicon.
- 3.20 Write a computer program to calculate the time required to grow a given thickness of oxide, based on the theory of Section 3.2. The user should be able to specify desired oxide thickness, wet or dry oxidation conditions, temperature, and orientation of the silicon wafer.
- 3.21 (a) Use SUPREM to simulate the oxide growth in Problem 3.8 on a $\langle 100 \rangle$ wafer doped with 10^{15} boron atoms/ cm^3 . Plot the final doping profiles in the silicon and oxide. (b) Repeat for a wafer doped with 10^{15} arsenic atoms/ cm^3 . (c) Compare the oxide thickness to hand calculations.
- 3.22 (a) Use SUPREM to simulate the oxide growth in Problem 3.8 on a $\langle 111 \rangle$ wafer doped with 3×10^{15} boron atoms/ cm^3 . (b) Compare the oxide thickness to hand calculations.
- 3.23 (a) Use SUPREM to simulate the oxide growth in Problem 3.6 on a $5\Omega\text{-cm}$ $\langle 100 \rangle$ boron-doped wafer. Plot the concentration of boron in the oxide and substrate. (b) Repeat for a $5\Omega\text{-cm}$ wafer doped with phosphorus atoms.
- 3.24 (a) Use SUPREM to simulate the oxide growth in Problem 3.7 on a $\langle 100 \rangle$ wafer doped with 5×10^{15} boron atoms/ cm^3 . Plot the concentration of boron in the oxide and substrate. (b) Repeat for a wafer with a doping of 5×10^{15} phosphorus atoms/ cm^3 .
- 3.25 Use SUPREM to calculate the thicknesses of the oxides in the two regions in Problem 3.12.

CHAPTER 4

Diffusion

High-temperature diffusion has historically been one of the most important processing steps used in the fabrication of monolithic integrated circuits. For many years, diffusion was the primary method of introducing impurities such as boron, phosphorus, and antimony into silicon to control the majority-carrier type and resistivity of layers formed in the wafer. Today, diffusion is used in the formation of “deep” layers exceeding a few tenths of a micron in depth. However, most deposition steps utilize the ion-implantation and rapid thermal annealing processes that will be explored in Chapter 5. We must still study the diffusion process in order to understand its limitations and the various problems associated with redistribution of impurities as they are added to silicon. In this chapter, we explore the theoretical and practical aspects of the diffusion process, the characterization of diffused layer sheet resistance, and the determination of junction depth. Physical diffusion systems and solid, liquid, and gaseous impurity sources are all discussed.

4.1 THE DIFFUSION PROCESS

The diffusion process begins with the deposition of a shallow high-concentration layer of the desired impurity in the silicon surface through windows etched in the protective barrier layer. At high temperatures (900 to 1200 °C), the impurity atoms move from the surface into the silicon crystal via the *substitutional* or *interstitial* diffusion mechanisms illustrated in Fig. 4.1 on page 68.

In the case of substitutional diffusion, the impurity atom hops from one crystal lattice site to another. The impurity atom thereby “substitutes” for a silicon atom in the lattice. Vacancies must be present in the silicon lattice in order for the substitutional process to occur. Statistically, a certain number of vacancies will always exist in the lattice. At high temperatures, vacancies may also be created by displacing silicon atoms from their normal lattice positions into the vacant *interstitial* space between lattice sites. The substitutional diffusion process in which silicon atoms are displaced into interstitial sites is called *interstitialcy* diffusion.

Considerable space exists between atoms in the silicon lattice, and certain impurity atoms diffuse through the crystal by jumping from one interstitial site to another. Since this mechanism does not require the presence of vacancies, interstitial diffusion

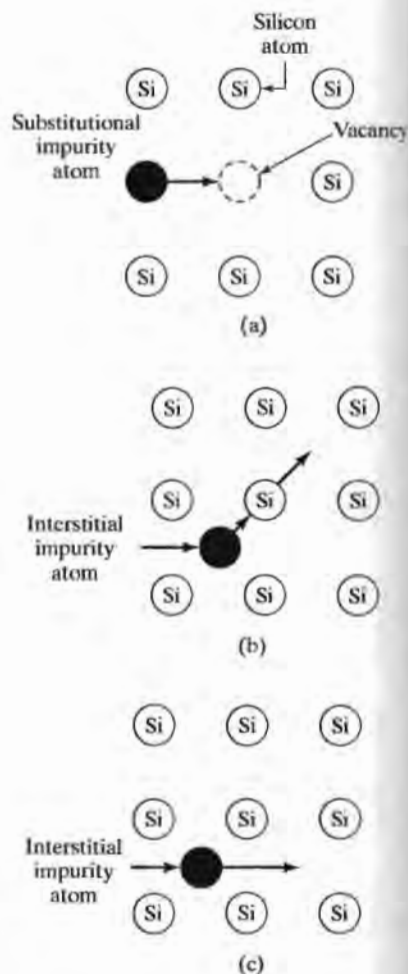


FIGURE 4.1

Atomic diffusion in a two-dimensional lattice. (a) Substitutional diffusion, in which the impurity moves among vacancies in the lattice; (b) interstitialcy mechanism, in which the impurity atom replaces a silicon atom in the lattice, and the silicon atom is displaced to an interstitial site; (c) interstitial diffusion, in which impurity atoms do not replace atoms in the crystal lattice.

proceeds much more rapidly than substitutional diffusion. The rapid diffusion rate makes interstitial diffusion difficult to control.

Impurity atoms need to occupy substitutional sites in the lattice in order to provide electrons or holes for conduction, as described in Volume I of this series [1]. Substitutional diffusion proceeds at a relatively low rate, because the supply of vacancies is limited, but this slow diffusion rate is actually an advantage, because it permits good control of the diffusion process.

4.2 MATHEMATICAL MODEL FOR DIFFUSION

The basic one-dimensional diffusion process follows *Fick's first law* of diffusion, presented in Chapter 3,

$$J = -D \frac{\partial N}{\partial x} \quad (4.1)$$

where J is the particle flux of the donor or acceptor impurity species, N is the concentration of the impurity, and D is the diffusion coefficient.

Fick's second law of diffusion may be derived using the continuity equation for the particle flux:

$$\partial N / \partial t = - \partial J / \partial x \quad (4.2)$$

Equation (4.2) states that the rate of increase of concentration with time is equal to the negative of the divergence of the particle flux. For the one-dimensional case, the divergence is equal to the gradient. Combining Eqs. (4.1) and (4.2) yields Fick's second law of diffusion:

$$\partial N / \partial t = D \partial^2 N / \partial x^2 \quad (4.3)$$

Here the diffusion coefficient D is assumed to be independent of position. This assumption is violated at high impurity concentrations. (See Section 4.6.3)

The partial differential equation in Eq. (4.3) can be solved by variable separation or Laplace transform techniques. Two specific types of boundary conditions are important in modeling impurity diffusion in silicon. The first is the *constant-source diffusion*, in which the surface concentration is held constant throughout the diffusion. The second is called a *limited-source diffusion*, in which a fixed quantity of the impurity species is deposited in a thin layer in the surface of the silicon.

4.2.1 Constant-Source Diffusion

During a constant-source diffusion, the impurity concentration is held constant at the surface of the wafer. Under this boundary condition, the solution to Eq. (4.3) is given by

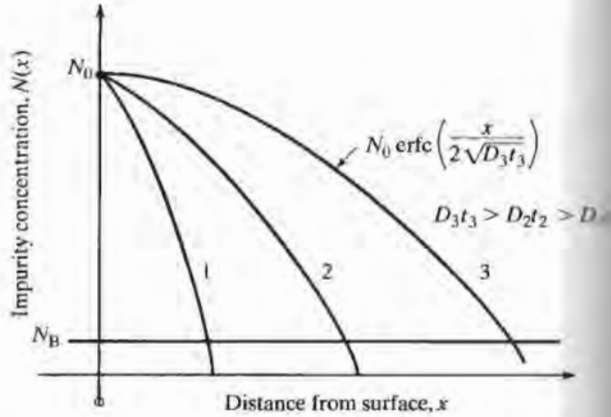
$$N(x, t) = N_0 \operatorname{erfc}(x/2\sqrt{Dt}) \quad (4.4)$$

for a semiinfinite wafer in which N_0 is the impurity concentration at the wafer surface ($x = 0$). Such a diffusion is called a *complementary error function* (erfc) diffusion and is shown graphically in Fig. 4.2 on page 70. As time progresses, the diffusion front proceeds further and further into the wafer with the surface concentration remaining constant. The total number of impurity atoms per unit area in the silicon is called the *dose*, Q , with units of atoms/cm². Q increases with time, and an external impurity source must supply a continual flow of impurity atoms to the surface of the wafer. The dose is found by integrating the diffused impurity concentration throughout the silicon wafer.

$$Q = \int_0^{\infty} N(x, t) dx = 2N_0\sqrt{Dt/\pi} \quad (4.5)$$

FIGURE 4.2

A constant-source diffusion results in a complementary error function impurity distribution. The surface concentration N_0 remains constant, and the diffusion moves deeper into the silicon wafer as the Dt product increases. Dt can change as a result of increasing diffusion time, increasing diffusion temperature, or a combination of both.



4.2.2 Limited-Source Diffusion

A limited-source diffusion is modeled mathematically using an impulse function at the silicon surface as the initial boundary condition. The magnitude of the impulse is equal to the dose Q . For this boundary condition in a semi-infinite wafer, the solution to Eq. (4.3) is given by the Gaussian distribution,

$$N(x, t) = (Q/\sqrt{\pi Dt})\exp-(x/2\sqrt{Dt})^2, \tag{4.4}$$

which is displayed graphically in Fig. 4.3. The dose remains constant throughout the limited-source diffusion process. As the diffusion front moves into the wafer, the surface concentration must decrease, so that the area under the curve can remain constant with time.

On a normalized logarithmic plot, the shapes of the Gaussian and complementary error function curves appear similar, as illustrated in Fig. 4.4. The erfc curve, however, falls off more rapidly than the Gaussian curve.

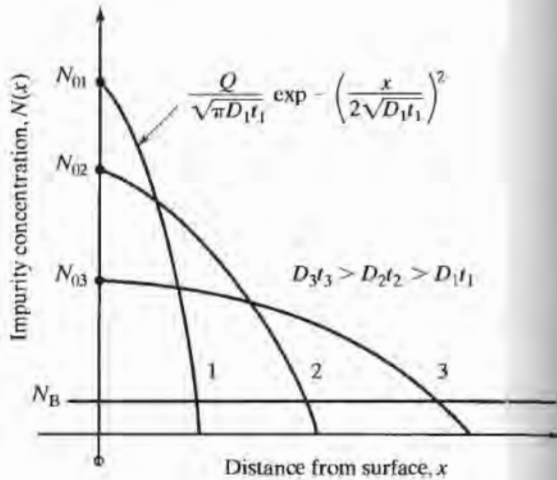


FIGURE 4.3

A Gaussian distribution results from a limited-source diffusion. As the Dt product increases, the diffusion front moves more deeply into the wafer, and the surface concentration decreases. The area (impurity dose) under each of the three curves is the same.

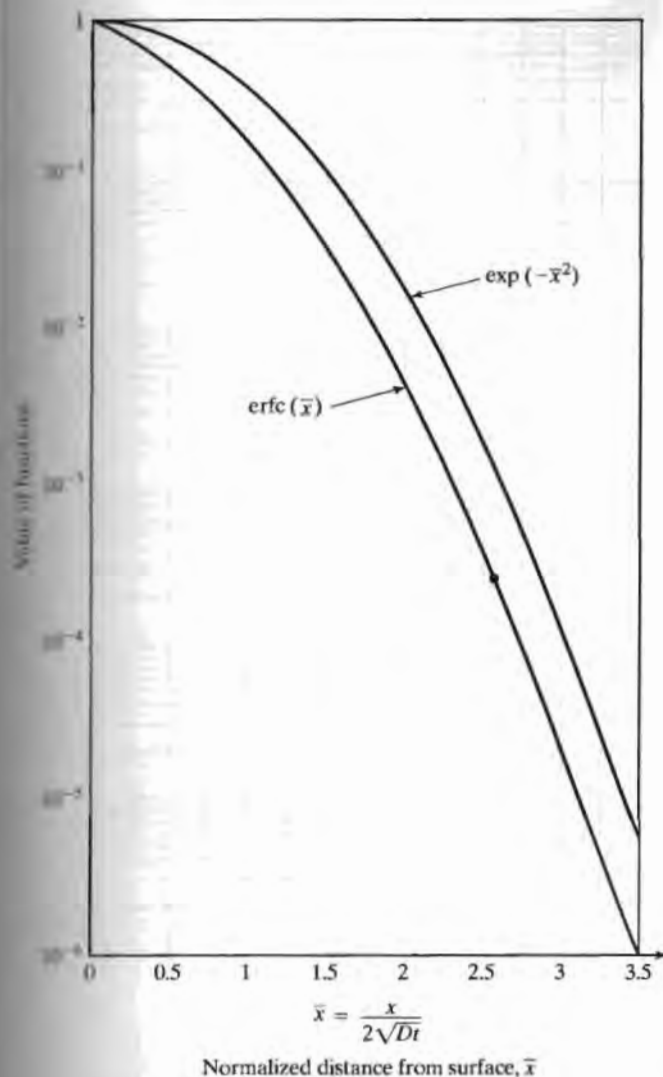


FIGURE 4.4

A graph comparing the Gaussian and complementary error function (erfc) profiles. We use this curve to evaluate the erfc and its inverse.

4.2.3 Two-Step Diffusion

A short constant-source diffusion is often followed by a limited-source diffusion, resulting in a two-step diffusion process. The constant-source diffusion step is used to establish a known dose in a shallow layer on the surface of the silicon and is called the *predeposition* step. The fixed dose approximates an impulse and serves as the impurity source for the second diffusion step.

The second diffusion is called the *drive-in* step and is used to move the diffusion front to the desired depth. If the Dt product for the drive-in step is much greater than the Dt product for the predeposition step, the resulting impurity profile is closely approximated by a Gaussian distribution, Eq. (4.6). If the Dt product for the drive-in step is much less than the Dt product for the predeposition step, the resulting impurity

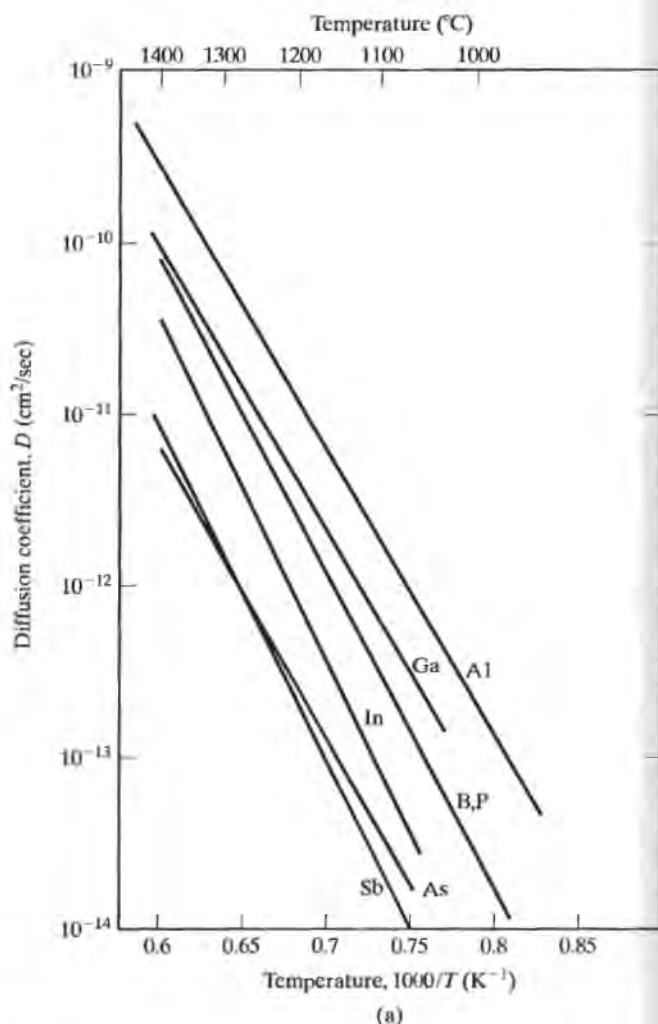


FIGURE 4.5

Diffusion constants in silicon for (a) substitutional diffusers (above) and (b) interstitial diffusers (next page). Copyright John Wiley & Sons, Inc.; reprinted with permission from Ref. [28].

profile is closely approximated by a complementary error function distribution, Eq. (4.4). An integral equation solution to the diffusion equation also exists for diffusion conditions that do not satisfy either inequality. (See Ref. [25].)

4.3 THE DIFFUSION COEFFICIENT

Figure 4.5 shows the temperature dependence of the diffusion coefficient D for (a) substitutional and (b) interstitial diffusers in silicon. The large difference between these coefficients is readily apparent. To achieve reasonable diffusion times with substitutional diffusers, temperatures in the range of 900 to 1200 $^{\circ}\text{C}$ are typically used. Interstitial diffusers are difficult to control, because of their large diffusion coefficients. (See Problem 4.19.)

Diffusion coefficients depend exponentially on temperature and follow the Arrhenius behavior introduced in Chapter 3:

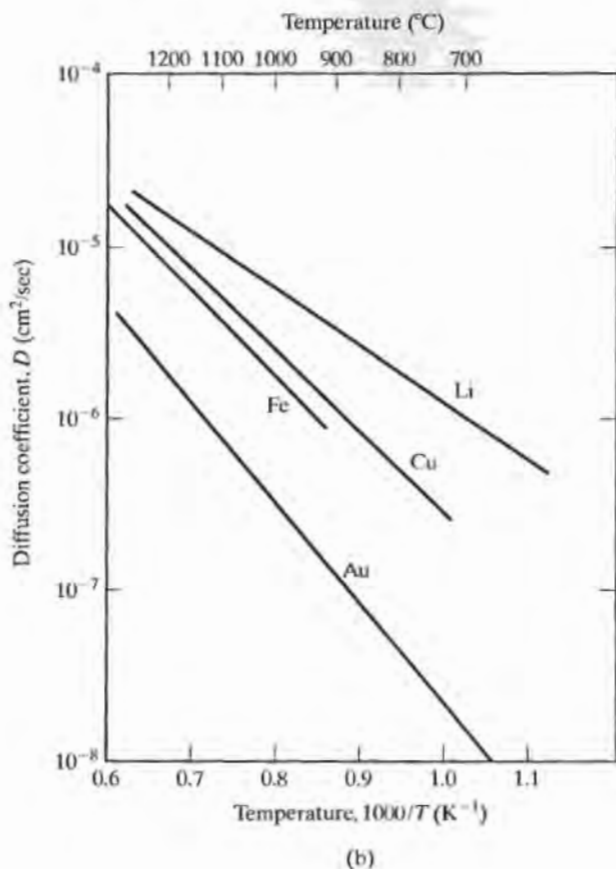


FIGURE 4.5
Continued.

$$D = D_0 \exp(-E_A/kT). \quad (4.7)$$

Values for D_0 and E_A can be determined from Fig. 4.5. Typical values for a number of impurities are given in Table 4.1.

Wide variability exists in diffusion coefficient data reported in the literature. We will use Eq. (4.7) and Table 4.1 in the examples and problems throughout the rest of this book. In general, calculations based on Eq. (4.7) and Table 4.1 can be used as guides. Most processes are then experimentally calibrated under the specific diffusion conditions in each laboratory.

Example 4.1

Calculate the diffusion coefficient for boron at 1100°C .

Solution: From Table 4.1, $D_0 = 10.5 \text{ cm}^2/\text{sec}$ and $E_A = 3.69 \text{ eV}$. $T = 1373 \text{ K}$.

$$D = 10.5 \exp \left[-\frac{3.69}{(8.614 \times 10^{-5})(1373)} \right] = 2.96 \times 10^{-13} \text{ cm}^2/\text{sec}.$$

TABLE 4.1 Typical Diffusion Coefficient Values for a Number of Impurities.

Element	$D_0(\text{cm}^2/\text{sec})$	$E_A(\text{eV})$
B	10.5	3.69
Al	8.00	3.47
Ga	3.60	3.51
In	16.5	3.90
P	10.5	3.69
As	0.32	3.56
Sb	5.60	3.95

4.4 SUCCESSIVE DIFFUSIONS

We are ultimately interested in the final impurity distribution after all processing is complete. A wafer typically goes through many time-temperature cycles during predeposition, drive-in, oxide growth, CVD, etc. For example, the base diffusion in a bipolar transistor will be followed by several high-temperature oxidations, as well as the emitter predeposition and drive-in cycles. These steps take place at different temperatures for different lengths of time. The effect of these steps is determined by calculating the total Dt product, $(Dt)_{\text{tot}}$, for the diffusion. $(Dt)_{\text{tot}}$ is equal to the sum of the Dt products for all high-temperature cycles affecting the diffusion:

$$(Dt)_{\text{tot}} = \sum_i D_i t_i \quad (4.8)$$

D_i and t_i are the diffusion coefficient and time associated with the i th processing step. $(Dt)_{\text{tot}}$ is then used in Eq. (4.4) or Eq. (4.6) to determine the final impurity distribution.

Example 4.2

Calculate $(Dt)_{\text{tot}}$ for a boron diffusion of 2 hours at 1100 °C followed by 5 hours at 1150 °C.

Solution: From Ex. 4.1 at $T = 1100$ °C, $D = 2.96 \times 10^{-13} \text{ cm}^2/\text{sec}$. For $T = 1150$ °C,

$$D = 10.5 \exp - \frac{3.69}{(8.614 \times 10^{-5})(1423)} = 8.86 \times 10^{-13} \text{ cm}^2/\text{sec},$$

$$(Dt)_{\text{tot}} = (2.96 \times 10^{-13} \text{ cm}^2/\text{sec})(7200 \text{ sec}) + (8.86 \times 10^{-13} \text{ cm}^2/\text{sec})(18000 \text{ sec})$$

$$(Dt)_{\text{tot}} = 1.81 \times 10^{-8} \text{ cm}^2.$$

4.5 SOLID-SOLUBILITY LIMITS

At a given temperature, there is an upper limit to the amount of an impurity that can be absorbed by silicon. This quantity is called the *solid-solubility limit* for the impurity and is indicated by the solid lines in Fig. 4.6 for boron, phosphorus, antimony, and

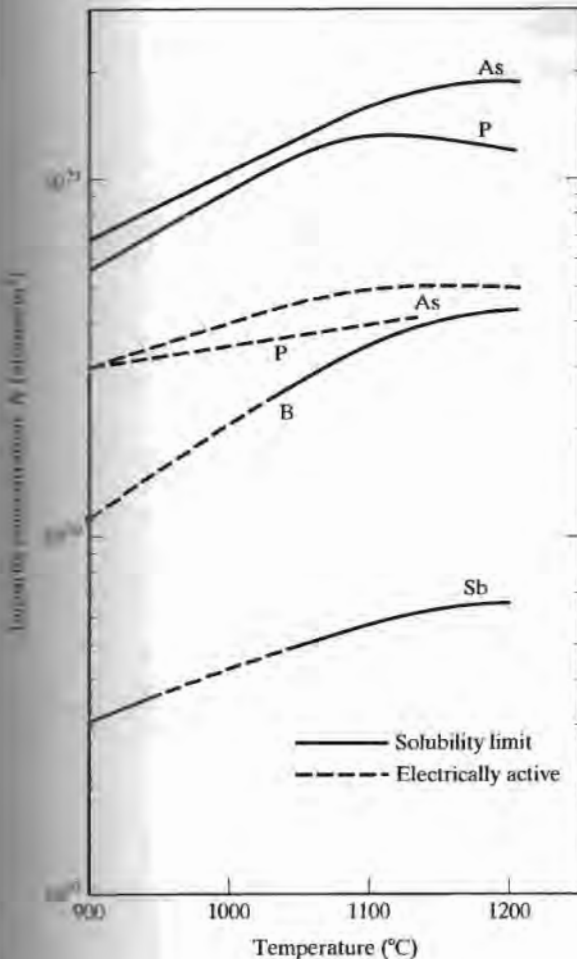


FIGURE 4.6

The solid-solubility and electrically active impurity-concentration limits in silicon for antimony, arsenic, boron, and phosphorus. Reprinted with permission from Ref. [29]. This paper was originally presented at the 1977 Spring Meeting of The Electrochemical Society, Inc., held in Philadelphia, Pennsylvania.

arsenic at normal diffusion temperatures. As can be seen in the figure, surface concentrations achieved through solid-solubility-limited diffusions will be quite high. For example, the solid-solubility limit of boron is approximately $3.3 \times 10^{20}/\text{cm}^3$ at 1100°C , and $1.2 \times 10^{21}/\text{cm}^3$ for phosphorus at the same temperature. High concentrations are desired for the emitter and subcollector diffusions in bipolar transistors and the source and drain diffusions in MOSFETs. However, solid-solubility-limited concentrations are too heavy for the base regions of bipolar transistors and for many resistors. The two-step diffusion process described in Section 4.2.3 overcomes this problem.

At high concentrations, only a fraction of the impurities actually contribute holes or electrons for conduction. The dotted lines in Fig. 4.6 show the “electrically active” portion of the impurity concentration. These curves will be referred to again in Section 4.7.2.

4.6 JUNCTION FORMATION AND CHARACTERIZATION

4.6.1 Vertical Diffusion and Junction Formation

The goal of most diffusions is to form pn junctions by converting p -type material to n -type material or vice versa. In Fig. 4.7, for example, the wafer is uniformly doped n -type material with a concentration indicated by N_B , and the diffusing impurity is boron. The point at which the diffused impurity profile intersects the background concentration is the *metallurgical junction depth*, x_j . The net impurity concentration at x_j is zero. Setting $N(x)$ equal to the background concentration N_B at $x = x_j$ yields

$$x_j = 2\sqrt{Dt \ln(N_0/N_B)} \quad (4.9)$$

and

$$x_j = 2\sqrt{Dt} \operatorname{erfc}^{-1}(N_B/N_0) \quad (4.10)$$

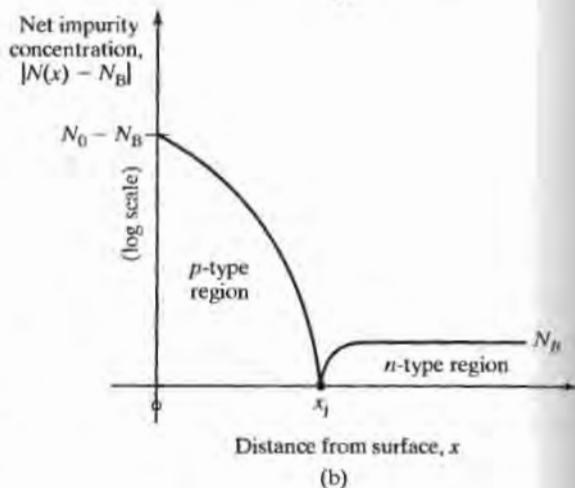
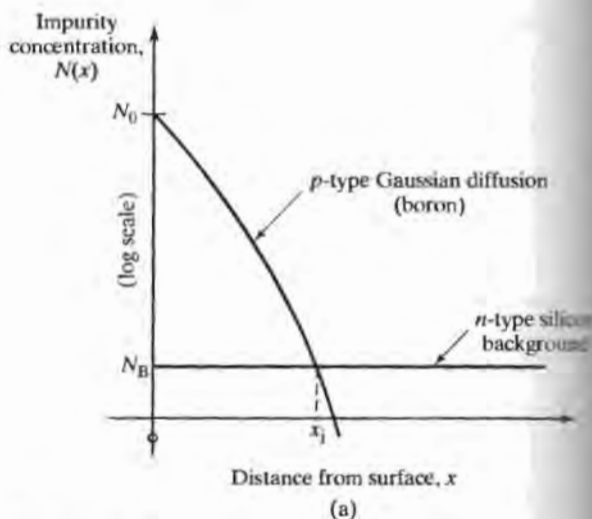


FIGURE 4.7

Formation of a pn junction by diffusion. (a) An example of a p -type Gaussian diffusion into a uniformly doped n -type wafer; (b) net impurity concentration in the wafer. The metallurgical junction occurs at the point $x = x_j$ where the net concentration is zero. The material is converted to p -type to the left of x_j and remains n -type to the right of x_j .

for the Gaussian and complementary error function distributions, respectively. In Fig. 4.7, the boron concentration N exceeds N_B to the left of the junction, and this region is p -type. To the right of x_j , N is less than N_B , and this region remains n -type.

We can use our scientific calculators to evaluate Eq. (4.9), and we will learn to evaluate the complementary error function expression using Fig. 4.4. To calculate the junction depth, we must know the background concentration N_B of the original wafer. Figure 4.8 gives the resistivity of n - and p -type silicon as a function of doping concentration. The background concentration can be determined using this figure when uniform concentrations of either donor or acceptor impurities are present in the silicon wafer.

Example 4.3

A boron diffusion is used to form the base region of an npn transistor in a 0.18-ohm-cm n -type silicon wafer. A solid-solubility-limited boron predeposition is performed at 900 °C for 15 min followed by a 5-hr drive-in at 1100 °C. Find the surface concentration and junction depth (a) following the predeposition step and (b) following the drive-in step.

Solution: The predeposition step is a solid-solubility-limited constant-source diffusion. Using Fig. 4.6, we see that the boron surface concentration is approximately $1.1 \times 10^{20}/\text{cm}^3$. The temperature of 900 °C equals 1173 K, which yields a diffusion coefficient $D_1 = 1.45 \times 10^{-15} \text{ cm}^2/\text{sec}$, and $t_1 = 900 \text{ sec}$ (15 min). The constant-source diffusion results in an erfc profile, and the impurity profile following predeposition is given by

$$N(x) = 1.1 \times 10^{20} \operatorname{erfc}(x/2\sqrt{D_1 t_1}) \text{ boron atoms}/\text{cm}^3.$$

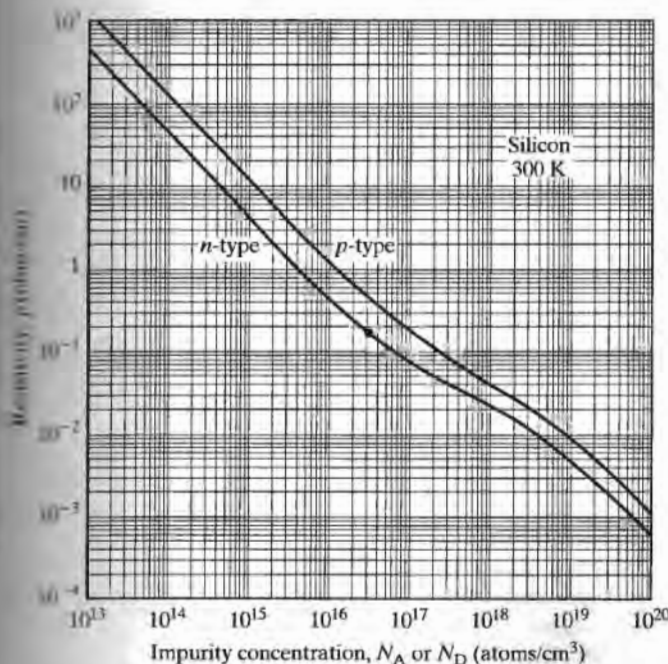


FIGURE 4.8

Room-temperature resistivity in n - and p -type silicon as a function of impurity concentration. (Note that these curves are valid for either donor or acceptor impurities but not for compensated material containing both types of impurities.)

Copyright 1987 Addison-Wesley Publishing Company. Reprinted with permission from Ref. [3].

To find the junction depth x_j , we must find the point at which the concentration $N(x)$ is equal to the background concentration N_B . Using Fig. 4.8, we find that a 0.18-ohm-cm n -type wafer corresponds to a doping concentration of $3 \times 10^{16}/\text{cm}^3$. Thus,

$$1.1 \times 10^{20} \operatorname{erfc}(x_j/2\sqrt{D_1 t_1}) = 3 \times 10^{16}.$$

Solving for x_j yields

$$x_j = 2\sqrt{D_1 t_1} \operatorname{erfc}^{-1}(0.000273) = 2(\sqrt{1.31 \times 10^{-12}})(2.57)\text{cm} = 0.0587 \mu\text{m}.$$

The value of $\operatorname{erfc}^{-1}(0.000273)$ is found with the aid of Fig. 4.4. The value 2.73×10^{-4} corresponds to the y -axis value of the complementary error function, and the corresponding value for the normalized distance is $x = 2.57$.

The dose in silicon is needed for the drive-in step and is equal to

$$Q = 2N_0\sqrt{D_1 t_1/\pi} = 2(1.1 \times 10^{20})\sqrt{(1.45 \times 10^{-15})(900)/\pi} \text{ boron atoms/cm}^2$$

$$Q = 1.42 \times 10^{14} \text{ boron atoms/cm}^2$$

At the drive-in temperature of 1100°C (1373 K), $D_2 = 2.96 \times 10^{-13} \text{ cm}^2/\text{sec}$, and the drive-in time of 5 hr = 18000 sec. Assuming that a Gaussian profile results from the drive-in step, the final profile is given by

$$N(x) = 1.1 \times 10^{18} \exp - (x/2\sqrt{D_2 t_2})^2 \text{ boron atoms/cm}^3. \quad (4.3.1)$$

Setting Eq. (4.3.1) equal to the background concentration yields the final junction depth of 2.77 μm . Figure 4.9 on page 79 shows the concentrations at various points in the diffusion process.

We must check our assumption that the drive-in step results in a Gaussian profile. The Dt product for the predeposition step is $1.31 \times 10^{-12} \text{ cm}^2$, and the Dt product for the drive-in step is $5.33 \times 10^{-9} \text{ cm}^2$. Thus, $D_2 t_2 \gg D_1 t_1$, and our assumption is justified.

4.6.2 Lateral Diffusion

During diffusion, impurities not only diffuse vertically, but also move laterally under the edge of any diffusion barrier. Figure 4.10 on page 80 presents the results of computer simulation of the two-dimensional diffusion process. The normalized impurity concentrations can be used to find the ratio of lateral to vertical diffusion. Lateral diffusion is an important effect-coupling device and process design and was an important factor driving the development of self-aligned polysilicon-gate MOS processes. The interaction of lateral diffusion and device layout will be discussed in greater detail in Chapters 9 and 10.

Example 4.4

An erfc diffusion results in a junction depth of 2 μm and a surface concentration of $1 \times 10^{20}/\text{cm}^3$. The background concentration of the wafer is $1 \times 10^{16}/\text{cm}^3$. What is the lateral diffusion underneath the edge of the mask?

Solution: The junction occurs at $N = N_B$, and $N_0/N_B = 10^4$. Using Fig. 4.10(a), we find that the ratio of lateral diffusion to vertical diffusion is 2.4/2.75, or 0.87. The lateral junction depth is therefore 1.74 μm .

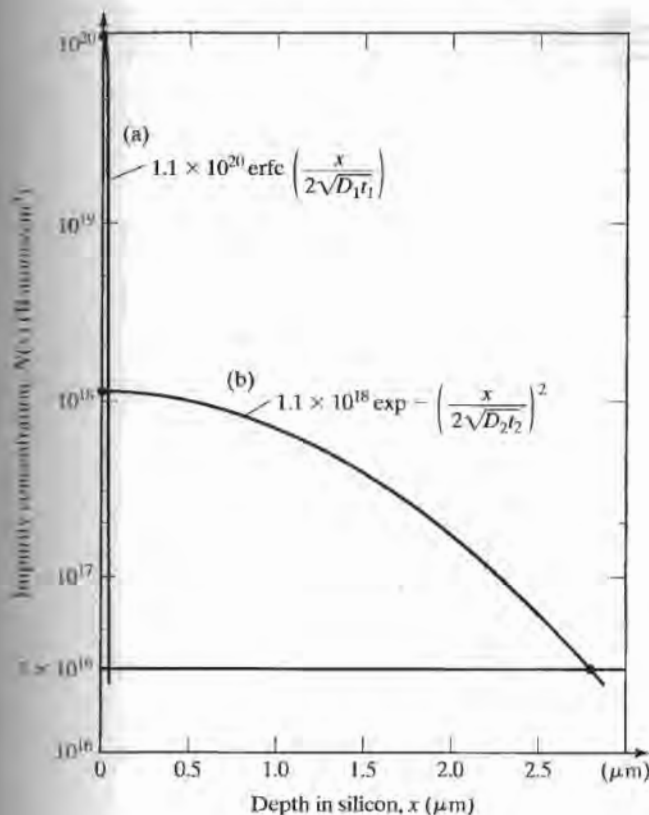


FIGURE 4.9

Calculated boron impurity profiles for Example 4.3. (a) Following the predeposition step at 900°C for 15 min; (b) following a subsequent 5-hr drive-in step at 1100°C. The final junction depth is 2.77 μm with a surface concentration of $1.1 \times 10^{18}/\text{cm}^3$. The initial profile approximates an impulse.

4.6.3 Concentration-Dependent Diffusion

Diffusion follows the theory of Section 4.3 as long as the impurity concentration remains below the value of the intrinsic-carrier concentration n_i at the diffusion temperature. Above this concentration, the diffusion coefficient becomes concentration dependent, and each of the common impurities exhibits a different behavior.

The diffusion equation can be solved analytically for linear, parabolic, and cubic dependencies of the diffusion coefficient on concentration. The results are presented in Fig. 4.11 on page 81, in which D_{sur} represents the diffusion coefficient at the surface. In general, concentration-dependent diffusion results in a much more abrupt profile than for the case of a constant-diffusion coefficient. These highly abrupt profiles are actually of value in formation of the shallow junctions desired in scaled VLSI devices.

Boron and arsenic can be modeled by the first-order dependence in Fig. 4.11, resulting in the analytical relations between junction depth, sheet resistance, total dose, and surface concentration given in Table 4.2 [7–9] on page 81.

High-concentration phosphorus diffusion results in a more complicated profile than that of boron or arsenic. Figure 4.12 on page 82 depicts typical shallow phosphorus diffusion profiles. As phosphorus diffuses into the wafer, the diffusion coefficient becomes enhanced at concentrations below approximately $10^{19}/\text{cm}^3$, resulting in a distinct “kink” in the profile. The kink effect represents a practical limitation to the use of phosphorus for the source–drain diffusions and emitter diffusions of shallow MOS and

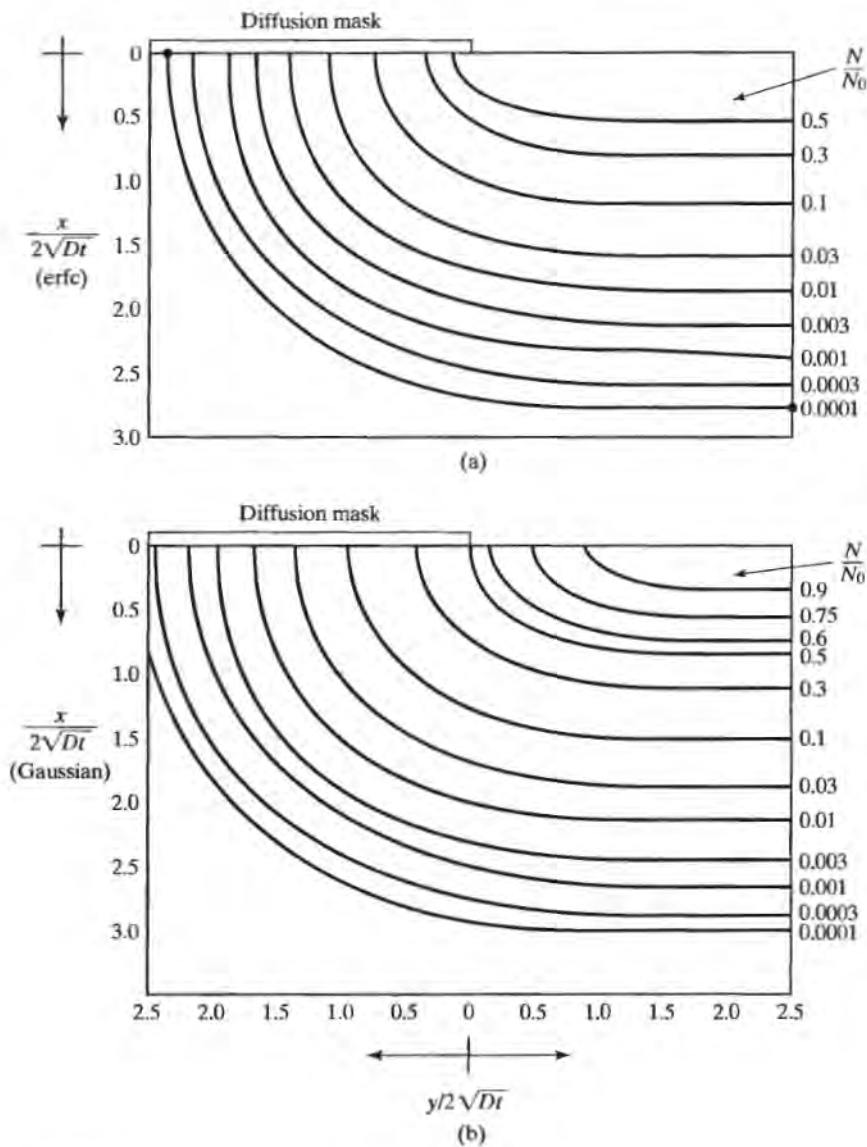


FIGURE 4.10

Normalized two-dimensional complementary error function and Gaussian diffusions near the edge of a window in the barrier layer. Copyright 1965 by International Business Machines Corporation; reprinted with permission from Ref. [4].

bipolar devices. Most MOS and bipolar VLSI processes now use arsenic to avoid this problem. Complex mathematical models describing the diffusion of phosphorus may be found in Refs. [10] and [11].

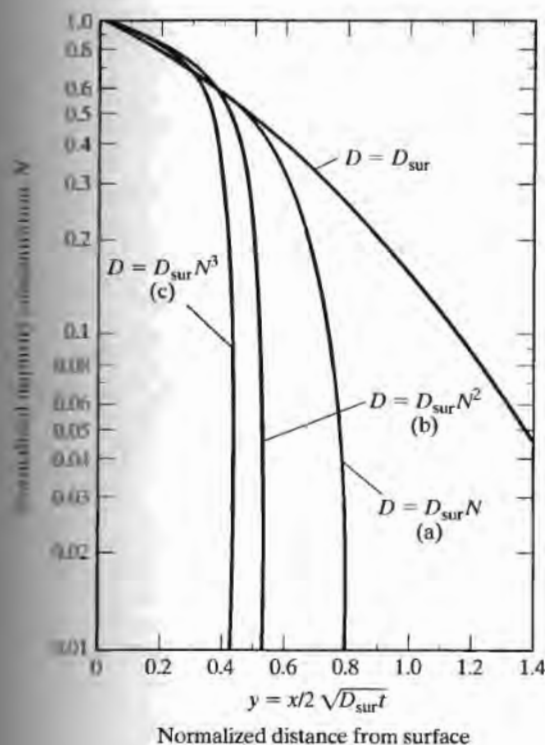


FIGURE 4.11

Diffusion profiles for concentration-dependent diffusion. Copyright 1963 by the American Physical Society. Reprinted with permission from Ref. [6].

Sec. 4.2 Properties of High-Concentration Arsenic and Boron Diffusions

Parameter	x_j (cm)	D (cm ² /sec)	N_0 (cm ⁻³)	Q (cm ⁻²)
As	$2.29 \sqrt{N_0 D t / n_i^*}$	$22.9 \exp(-4.1/kT)$	$1.56 \times 10^{17} (R_s x_j)^{-1}$	$0.55 N_0 x_j$
B	$2.45 \sqrt{N_0 D t / n_i^*}$	$3.17 \exp(-3.59/kT)$	$2.78 \times 10^{17} (R_s x_j)^{-1}$	$0.67 N_0 x_j$

The value of n_i must be calculated at the diffusion temperature.

SHEET RESISTANCE

In diffused layers, resistivity is a strong function of depth. For circuit and device design, it is convenient to work with a new parameter, R_s , called *sheet resistance*, which eliminates the need to know the details of the diffused-layer profile.

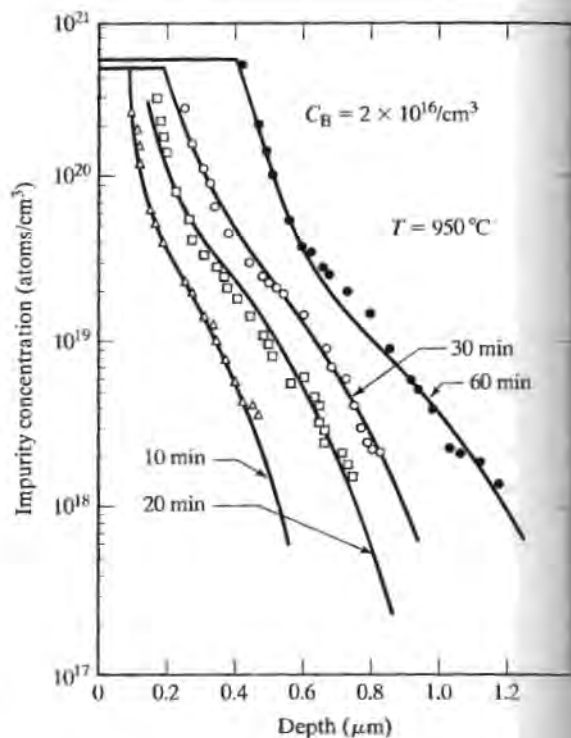


FIGURE 4.12

Shallow phosphorus diffusion profiles for constant-source diffusions at 950 °C. Copyright 1969 IEEE. Reprinted with permission from Ref. [10].

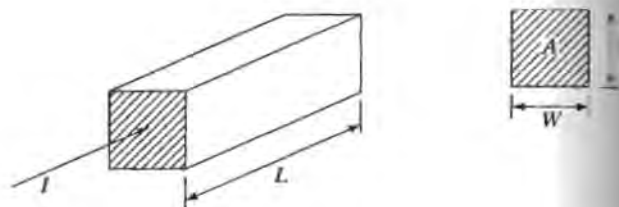
4.7.1 Sheet-Resistance Definition

Let us first consider the resistance R of the rectangular block of uniformly doped material in Fig. 4.13. R is given by

$$R = \rho L / A, \quad (4.11)$$

FIGURE 4.13

Resistance of a block of material having uniform resistivity. A uniform current distribution is entering the material perpendicular to the end of the block. The ratio of resistivity to thickness is called the *sheet resistance* of the material.



$$R = \rho \frac{L}{A} \quad \rho = \frac{1}{\sigma} \quad \sigma = q(\mu_n n + \mu_p p)$$

where ρ is the material's resistivity, and L and A represent the length and cross-sectional area of the block, respectively. Resistance is proportional to the material resistivity. If the length of the block is made longer, the resistance increases, and the resistance is inversely proportional to the cross-sectional area.

Using W as the width of the sample and t as the thickness of the sample, the resistance may be rewritten as

$$R = (\rho/t)(L/W) = R_s(L/W), \quad (4.12)$$

where $R_s = (\rho/t)$ is called the *sheet resistance* of the layer of material. Given the sheet resistance R_s , a circuit designer need specify only the length and width of the resistor to define its value. Strictly speaking, the unit for sheet resistance is the ohm, since the ratio L/W is unitless. To avoid confusion between R and R_s , sheet resistance is given the special descriptive unit of ohms per square. The ratio L/W represents the number of unit squares of material in the resistor.

Figure 4.14 shows top and side views of two typical dumbbell-shaped resistors with top contacts at the ends. The body of each resistor is seven "squares" long. If the sheet resistance of the diffusion were 50 ohms per square, each resistor would have a resistance of 350 ohms. The portion of the resistor surrounding the contacts also contributes to the total resistance of the structure. Figure 4.15 on page 84 presents the effective number of squares contributed by various end and corner configurations. For the resistor of Fig. 4.14, each end adds approximately 0.65 squares to the resistor and the total resistance would be 415 Ω . Note that lateral diffusion under the edges of the mask may change both the geometry of the contacts, as well as the number of squares in the body of the resistor. (See Problems 4.8–4.10.)

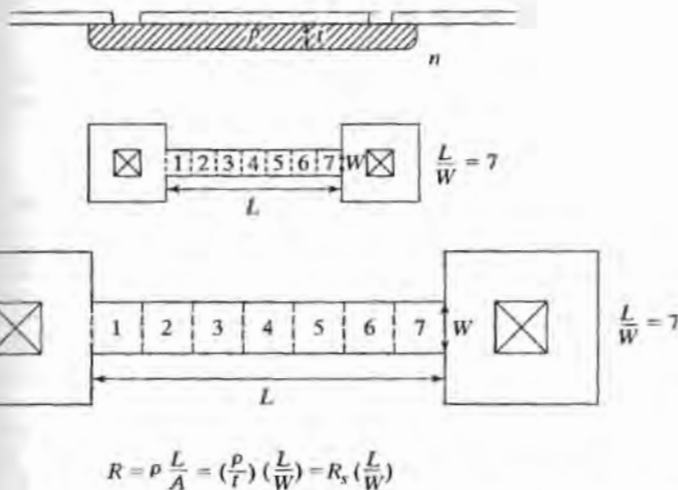


FIGURE 4.14

Top and side views of two diffused resistors of different physical size having equal values of resistance. Each resistor has a ratio L/W equal to 7 squares. Each end of the resistor contributes approximately 0.65 additional squares.

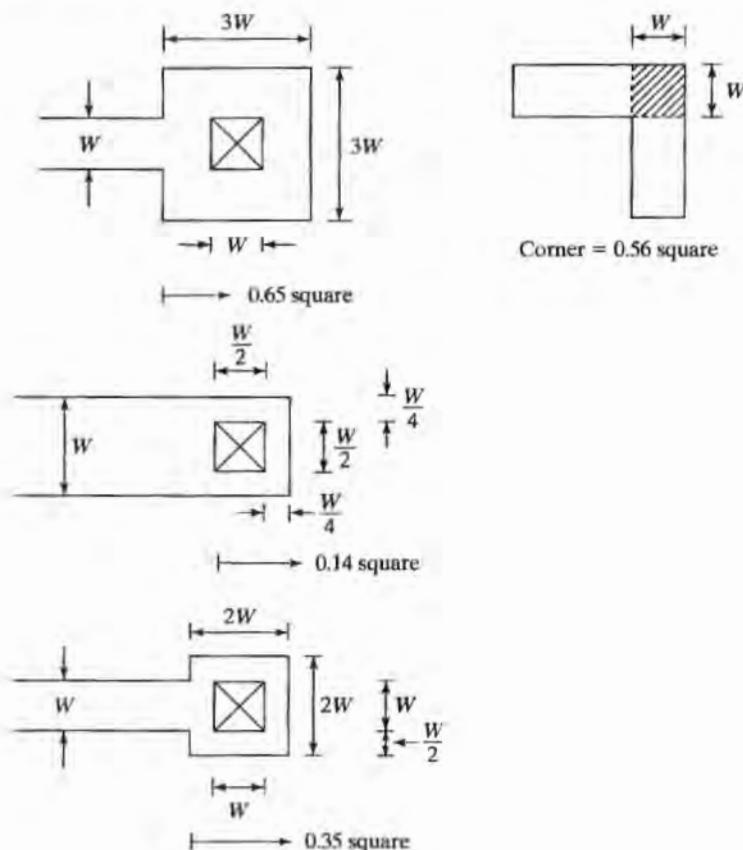


FIGURE 4.15

Effective square contributions of various resistor end and corner configurations.

Example 4.5

A diffusion with a sheet resistance of $200 \Omega/\square$ is used to fabricate a $5 \text{ k}\Omega$ resistor with the dumbbell shape similar to that in Fig. 4.14. How many squares are required in the body of the resistor?

Solution: The total number of squares required is $N_{\text{tot}} = (5000 \Omega) / (200 \Omega/\square) = 25 \square$. Subtracting the contribution of the two contacts gives $N = 25 - 2(0.65) = 23.7 \square$ for the body of the resistor. Note that the resistor body need not be an integer number of squares.

Irvin's Curves

From Section 4.2, we know that the impurity concentration resulting from a diffusion varies rapidly between the surface and the junction. Thus, ρ is a function of depth for diffused resistors. For diffused layers, we define the sheet resistance R_s in terms of the average resistivity of the layer

$$\bar{\rho} = \frac{1}{\bar{\sigma}} = \frac{1}{\frac{1}{x_j} \int_0^{x_j} \sigma(x) dx}$$

$$R_s = \frac{\bar{\rho}}{x_j} = \left[\int_0^{x_j} \sigma(x) dx \right]^{-1}.$$

In extrinsic material, this expression can be approximated by

$$R_s = \left[\int_0^{x_j} q\mu N(x) dx \right]^{-1}, \quad (4.13)$$

where x_j is the junction depth, μ is the majority-carrier mobility, and $N(x)$ is the net impurity concentration. We neglect the depletion of charge carriers near the junction x_j .

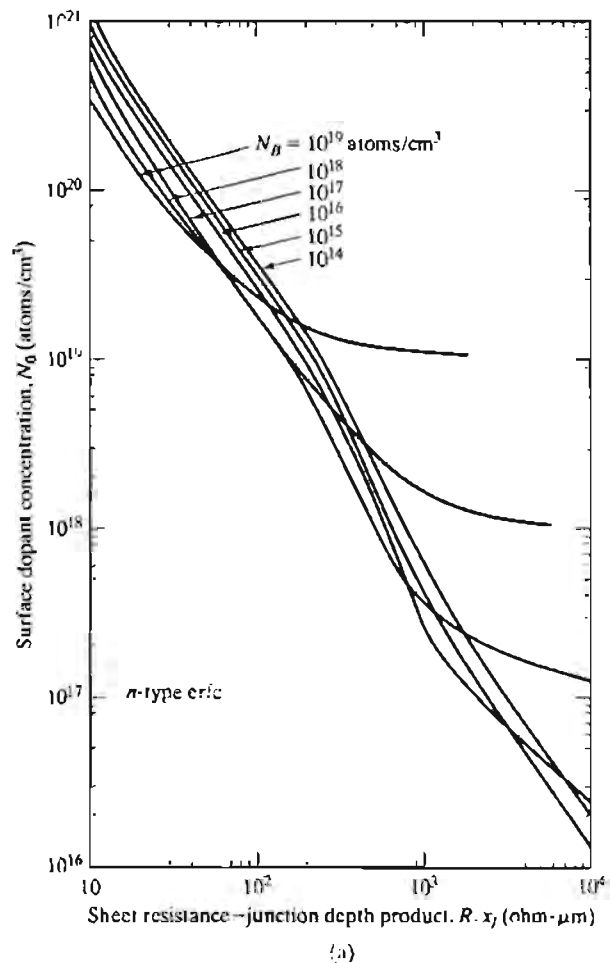
For a given diffusion profile, sheet resistance is uniquely related to the surface concentration of the diffused layer and the background concentration of the wafer. Equation (4.13) was evaluated numerically by Irvin [5], and a number of Irvin's results have been combined into Figs. 4.16(a)–(d) [2] on page 86–87. These figures plot surface concentration versus the $R_s x_j$ product and are used to find the sheet resistance and surface concentration of diffused layers.

Example 4.6

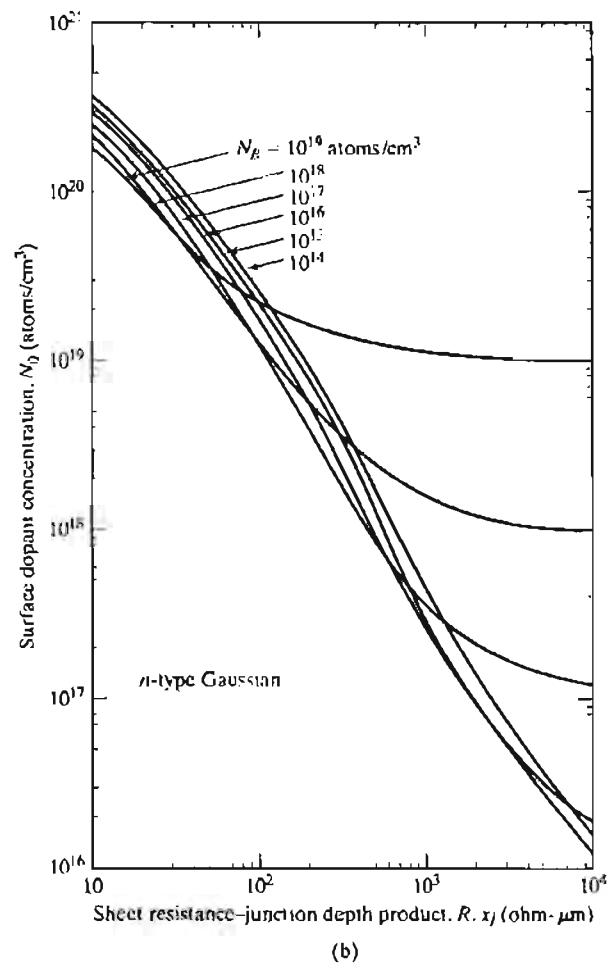
Find the sheet resistance of the boron diffusion from Example 4.2.

Solution: From Example 4.3, the background concentration of the wafer is $3 \times 10^{16}/\text{cm}^3$, the surface concentration is $1.1 \times 10^{18}/\text{cm}^3$, and the junction depth is $2.77 \mu\text{m}$. The diffusion resulted in a p -type Gaussian layer. Using Fig. 4.16(d), we find that the $R_s x_j$ product is found to be approximately $800 \text{ ohm}\cdot\mu\text{m}$. Dividing by a junction depth of $2.77 \mu\text{m}$ yields a sheet resistance of 289 ohms/square .

Sheet resistance is an electrical quantity that depends on the majority-carrier concentration. As shown in Fig. 4.6, the electrically active impurity concentration for phosphorus and arsenic is considerably less than the total impurity concentration at high doping levels. In order to use Irvin's curves at high doping levels, the vertical axis, which is labeled "surface dopant density," should be interpreted to be the electrically active dopant concentration at the surface.



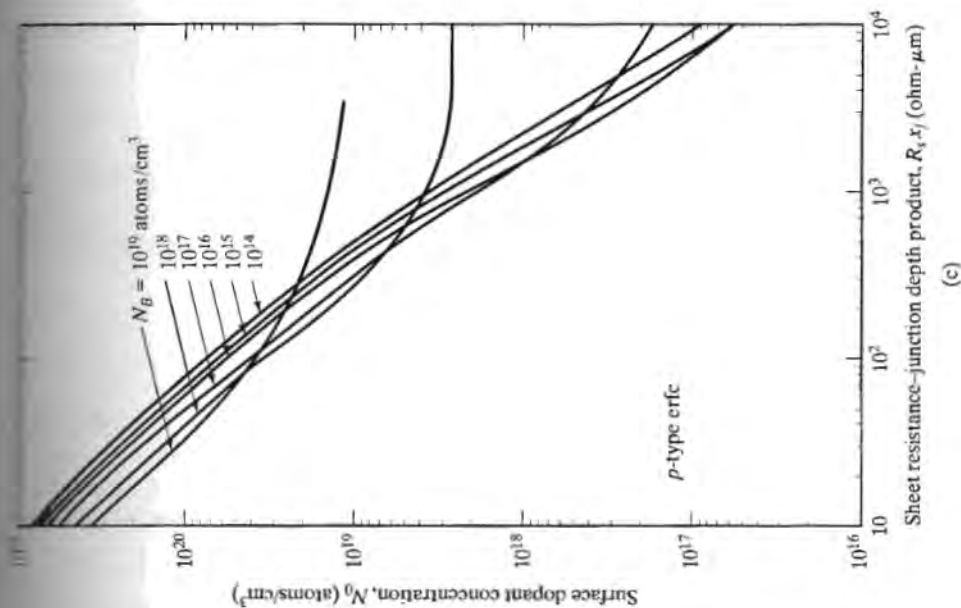
(a)



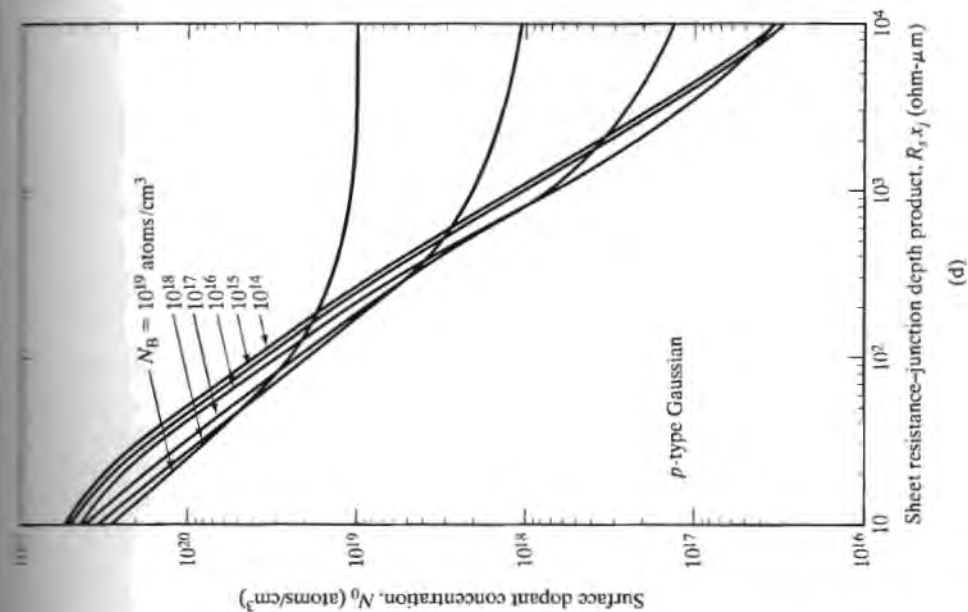
(b)

FIGURE 4.16

Surface impurity concentration versus the sheet resistance-junction depth product for different silicon background concentrations at 300 K. (a) *n*-type erfc distribution; (b) *n*-type Gaussian distribution; (c) *p*-type erfc distribution; (d) *p*-type Gaussian distribution. After Ref. [2]. Reprinted from Ref. [5] with permission from the *AT&T Technical Journal*. Copyright 1962 AT&T.



(c)

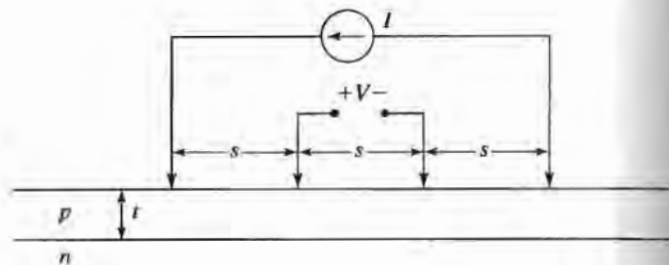


(d)

FIGURE 4.16
Continued.

FIGURE 4.17

Four-point probe with probe spacing s used for direct measurement of bulk wafer resistivity and the sheet resistance of thin diffused layers. A known current is forced through the outer probes, and the voltage developed is measured across the inner probes. (See Eqs. (4.14) through (4.16).)



4.7.3 The Four-Point Probe

A special instrument called a *four-point probe* may be used to measure the bulk resistivity of starting wafers and the sheet resistance of shallow diffused layers. As shown schematically in Fig. 4.17, a fixed current is injected into the wafer through the two outer probes, and the resulting voltage is measured between the two inner probes. If probes with a uniform spacing s are placed on an infinite slab of material, then the resistivity is given by

$$\rho = 2\pi s V / I \text{ ohm-meters for } t \gg s \quad (4.14)$$

and

$$\rho = (\pi t / \ln 2) V / I \text{ ohm-meters for } s \gg t. \quad (4.15)$$

For shallow layers, Eq. (4.15) gives the sheet resistance as

$$R_s = \rho / t = (\pi / \ln 2) V / I = 4.53 V / I \text{ ohm-meters for } s \gg t. \quad (4.16)$$

The approximation used in Eqs. (4.15) and (4.16) is easily met for shallow diffused layers in silicon. Unfortunately, silicon wafers are often thinner than the probe spacing s , and the approximation in Eq. (4.14) is not valid. Correction factors are given in Fig. 4.18 for thin wafers and for small-diameter wafers [12].

4.7.4 Van der Pauw's Method

The sheet resistance of an arbitrarily shaped sample of material may be measured by placing four contacts on the periphery of the sample. A current is injected through one pair of the contacts, and the voltage is measured across another pair of contacts. Van der Pauw [13,14] demonstrated that two of these measurements can be related by Eq. (4.17):

$$\exp(-\pi t R_{AB,CD} / \rho) + \exp(-\pi t R_{BC,DA} / \rho) = 1. \quad (4.17)$$

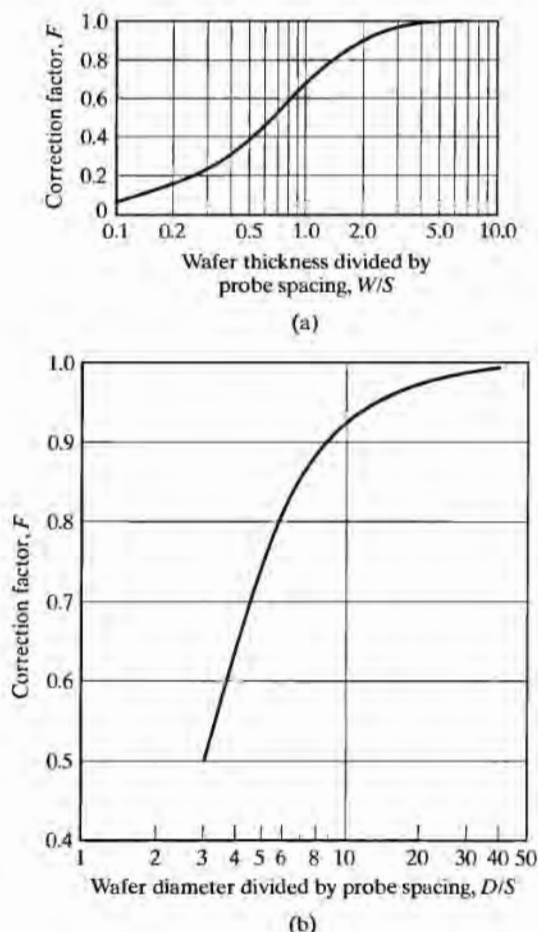


FIGURE 4.18

Four-point-probe correction factors, F , used to correct for (a) wafers which are relatively thick compared to the probe spacing s and (b) wafers of finite diameter. In each case $\rho = F\rho_{\text{measured}}$. (a) Copyright 1975 by McGraw-Hill Book Company. Reprinted with permission from Ref. [12]. (b) Reprinted from Ref. [30] with permission from the AT&T Technical Journal. Copyright 1958 AT&T.

Here $R_{AB,CD} = V_{CD}/I_{AB}$ and $R_{BC,DA} = V_{DA}/I_{BC}$. For a symmetrical structure such as a square or a circle,

$$R_{AB,CD} = R_{BC,DA}$$

and

$$R_s = \rho/t = (\pi/\ln 2)V_{CD}/I_{AB}. \quad (4.18)$$

Note the similarity to Eq. (4.16).

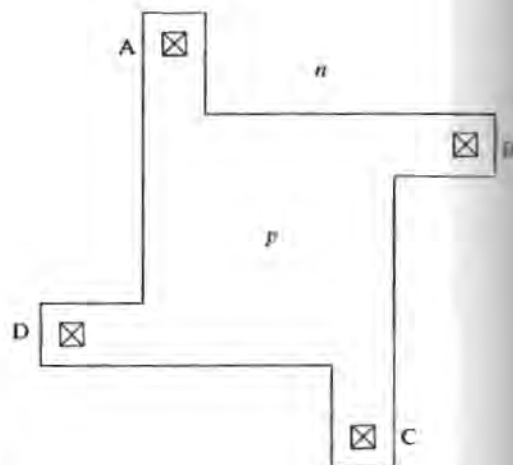


FIGURE 4.19

A simple van der Pauw test structure used to measure the sheet resistance of a diffused layer. Sheet resistance is calculated using Eq. (4.18).

Specially designed sheet-resistance test structures are often included on wafers so that the sheet resistances of *n*-type and *p*-type diffusions can be measured after final processing of the wafer. A sample structure is shown in Fig. 4.19.

4.8 JUNCTION-DEPTH AND IMPURITY PROFILE MEASUREMENT

Test wafers are normally processed in parallel with the actual IC wafers. No masking is done on the test wafer so that diffusion may take place across its full surface. The test wafer provides a large area for experimental characterization of junction depth. Alternatively, special test dice replace a few of the normal die sites on each wafer. These test dice provide an array of test structures for monitoring process and device characteristics during the various phases of the process.

4.8.1 Grove-and-Stain and Angle-Lap Methods

Two relatively simple methods can be used to measure the junction depth of diffused layers. In the first, known as the *groove-and-stain* method, a cylindrical groove is mechanically ground into the surface of the wafer, as shown in Fig. 4.20. If the radius R of the grinding tool is known, the junction depth x is easily found to be

$$x_j = \sqrt{(R^2 - b^2)} - \sqrt{(R^2 - a^2)}. \quad (4.19)$$

If the radius R is much larger than both distances a and b , then the junction depth is given approximately by

$$x_j = (a^2 - b^2)/2R = (a + b)(a - b)/2R. \quad (4.20)$$

After the grooving operation, the junction is delineated using a chemical etchant that stains the *pn* junction. Concentrated hydrofluoric acid with 0.1 to 0.5% nitric acid

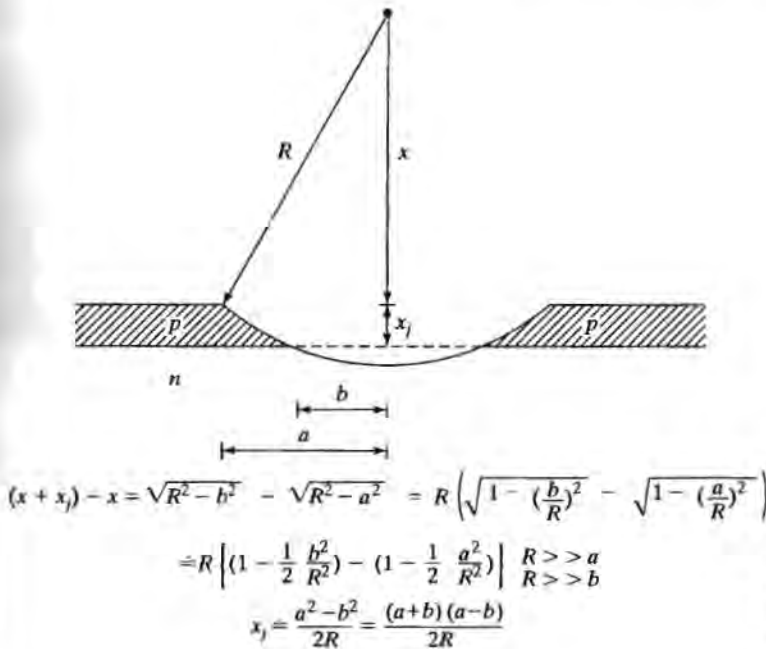


FIGURE 4.20

Junction-depth measurement by the groove-and-stain technique. The distances a and b are measured through a microscope, and the junction depth is calculated using Eq. (4.20).

can be used as a stain that is enhanced through exposure to high-intensity light [12]. The distances a and b are measured through a microscope, and the junction depth is calculated using Eq. (4.20).

The second technique is the *angle-lap* method. A piece of the wafer is mounted on a special fixture that permits the edge of the wafer to be lapped at an angle between 1 and 5°, as depicted in Fig. 4.21 on page 92. The junction depth is magnified so that the distance on the lapped surface is given by

$$x_j = d \tan \theta = N\lambda/2, \quad (4.21)$$

where θ is the angle of the fixture. An optically flat piece of glass is placed over the lapped region, and the test structure is illuminated with a collimated monochromatic beam of light with wavelength λ , typically from a sodium vapor lamp. The resulting interference pattern has fringe lines that are approximately 0.29 μm apart. The number N of fringes is counted through a microscope, and the junction depth may be found using Eq. (4.21). The usefulness of this method becomes limited for very shallow junctions. The analytic techniques discussed in the next section provide more general characterization capability for shallow structures.

4.8.2 Impurity-Profile Measurement

Spreading-resistance measurements and *Secondary Ion Mass Spectroscopy* (SIMS) are two techniques that are widely used for measurement of impurity profiles in semiconductors. Note that both methods are destructive; that is, they modify or destroy the region being characterized.

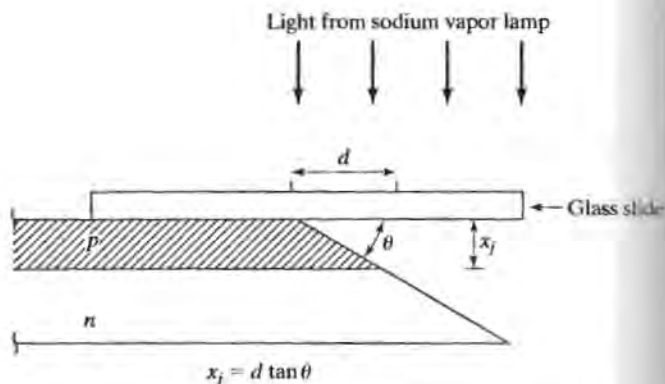
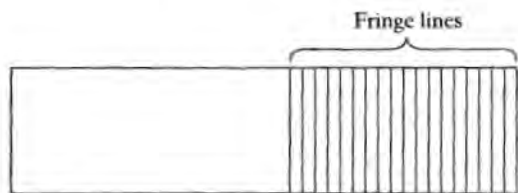


FIGURE 4.21
Junction depth measurement by the angle-lap and stain method. Interference fringe lines are used to measure the distance d , which is related to the junction depth using Eq. (4.21).



In the spreading-resistance method, a region of the semiconductor is angle lapped in a manner similar to that used for junction measurement discussed in Section 4.8.1. The resistivity of the layer is measured as a function of depth on the angle-lapped surface using a two-point probe. From this information, the impurity concentration and impurity type can be calculated. Figure 4.22 on page 93 presents an example of the impurity profile and junction depths determined from spreading resistance measurements.

In the SIMS method [15–16] depicted in Fig. 4.23(a) on page 94, a low-energy (1–20 keV) ion beam of say, cesium or oxygen, is used to remove (sputter) atoms from the surface, one or two atomic layers at a time. A small percentage of the atoms that are removed from the surface are ionized, and these ions are collected and analyzed by a mass spectrometer, which identifies the atomic species. The analysis is performed continuously during the sputtering process, and a profile of atomic distribution versus depth is produced as shown in Fig. 4.23(b). Mass removal proceeds at a rate of 2–5 Å/sec and can be used to a depth of a few microns. SIMS is the only surface-analysis tool with the sensitivity needed to characterize impurity profiles in silicon. Examples of typical sensitivities achievable with SIMS are provided in Table 4.3.

TABLE 4.3 SIMS Analysis in Silicon.

Element	Ion Beam	Sensitivity
Arsenic	Cesium	$5 \times 10^{14}/\text{cm}^3$
Boron	Oxygen	$1 \times 10^{13}/\text{cm}^3$
Phosphorus	Cesium	$5 \times 10^{15}/\text{cm}^3$
Oxygen	Cesium	$1 \times 10^{17}/\text{cm}^3$

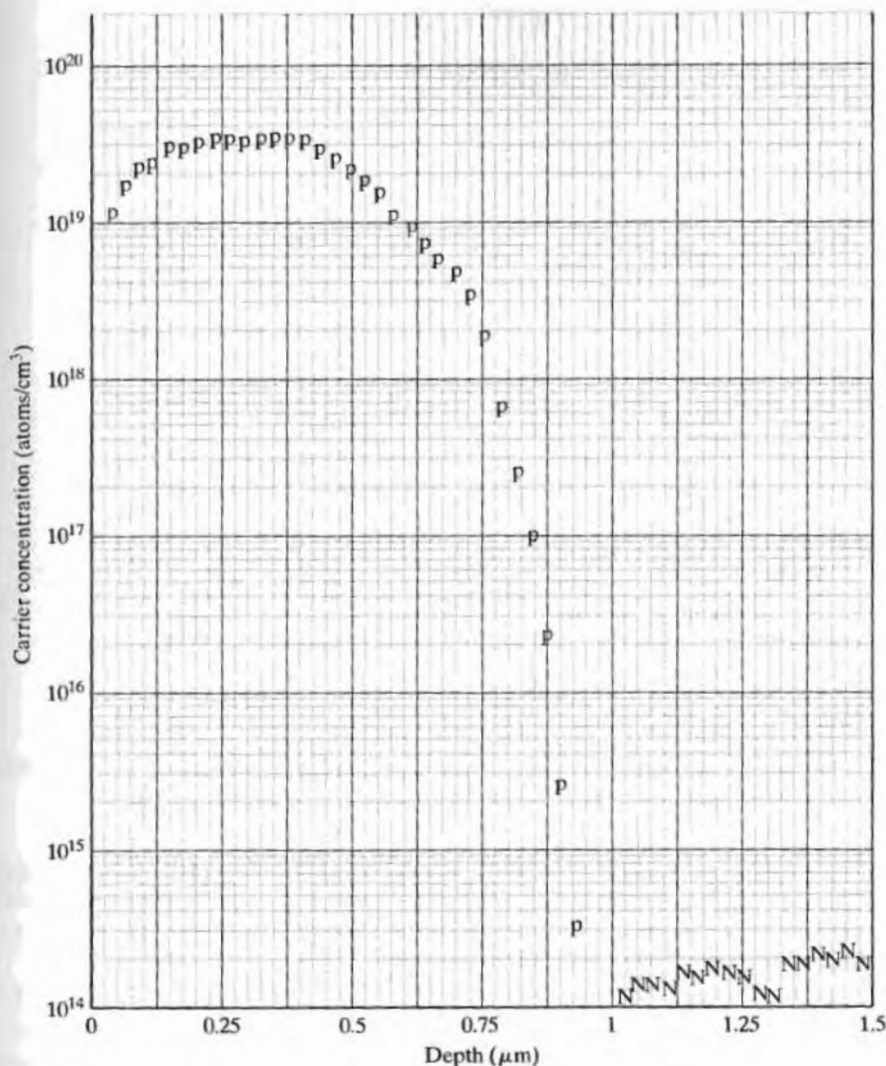


FIGURE 4.22

Example of an impurity profile measured using the spreading resistance method.

DIFFUSION SIMULATION

The SUPREM program introduced in Chapter 3 includes complete models for diffusion. SUPREM can simulate simple one-dimensional diffusion, as well as highly complex two-dimensional diffusion through a mask window, as depicted in Fig. 4.10. (See Problem 4.20.)

As a simple example, a portion of the input description of the two-step diffusion from Ex. 4.3 is given next, along with a plot of the corresponding output data in Fig.

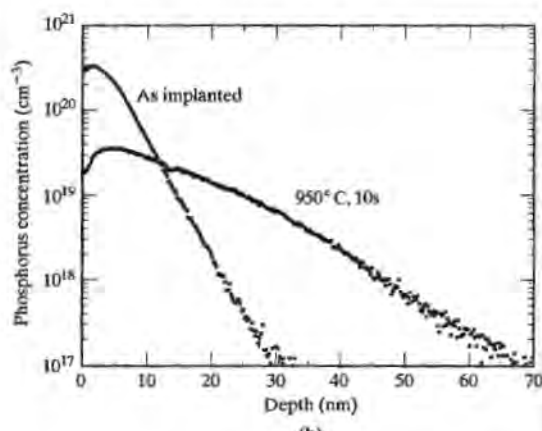
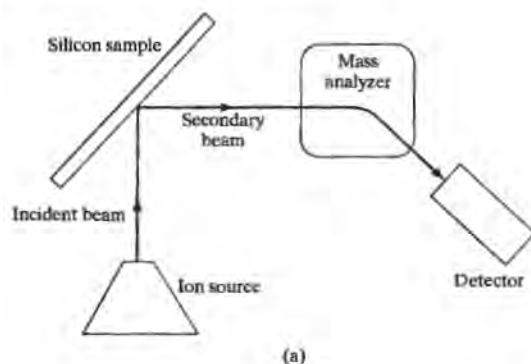


FIGURE 4.23

(a) Concept of a SIMS analysis system. (b) Example of an impurity profile measured using the SIMS analysis. Copyright 1997 IEEE. Reprinted with permission from Ref. [17].

4.24 on page 95. The input file defines the starting material to be a phosphorus doped <100> wafer with a resistivity of 0.18 $\Omega\text{-cm}$. The predeposition takes place at 900 $^{\circ}\text{C}$ for 15 minutes and the drive-in occurs at 1100 $^{\circ}\text{C}$ for 300 minutes. (Output control statements are not included in the listings.)

Following the predeposition step, the SUPREM simulation results predict that the surface concentration will be $N_0 = 3 \times 10^{20}/\text{cm}^3$ and that the junction depth will be $x_j = 0.1 \mu\text{m}$. After the drive-in step, the final values of N_0 and x_j are predicted to be $10^{17}/\text{cm}^3$ and 2.0 μm , respectively. The sheet resistance of the diffused layer is estimated to be approximately 500 Ω/\square . The simulation results show depletion of both boron and phosphorus near the wafer surface due to out-diffusion and indicate that the peak of the boron profile is actually below the silicon surface. These are features that we do not attempt to include in our basic hand analyses, and they show the power and importance of using the sophisticated computer simulation tools.

```

TWO STEP DIFFUSION
INITIALIZE <100> PHOS=0.18 RESISTIVITY
DIFFUSE TEMP=900 TIME=15 BORON=1E21
...
...
DIFFUSION TEMP=1100 TIME=300
...
...

```

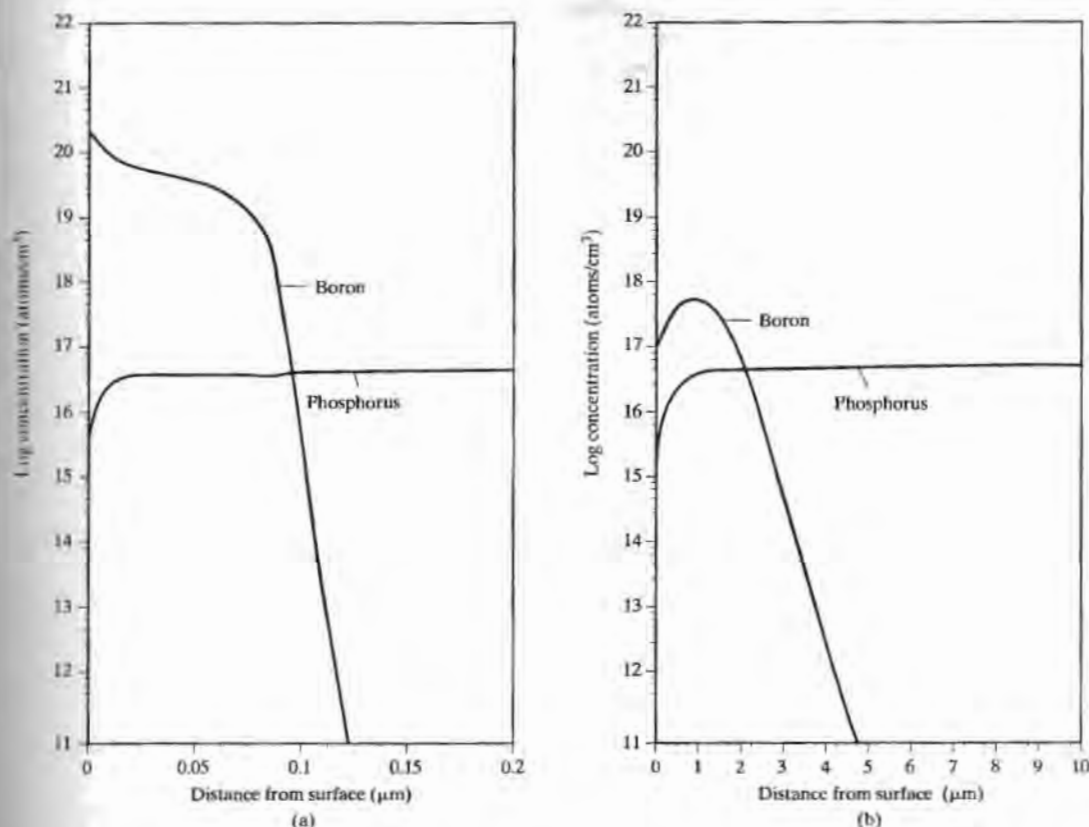


FIGURE 4.24

SUPREM simulation results for two-step boron diffusion into the phosphorus doped wafer from Ex. 4.3.

4.10 DIFFUSION SYSTEMS

The open-furnace-tube system using solid, liquid, or gaseous sources, as depicted in Fig. 4.25, yields good reproducibility and is a common diffusion technology used in IC fabrication. Three-zone horizontal furnaces can be used for diffusion. Wafers are placed in a quartz boat and positioned in the center zone of the furnace, where they are heated to a high temperature. Impurities are transported to the silicon surface, and then diffuse into the wafer.

Most common silicon dopants can be applied using liquid spin-on sources. These spin-on dopants are versatile, safe, and easy to apply, but the uniformity is often poorer than with other impurity sources. To achieve good quality control, most production systems use other solid, liquid, or gaseous impurity sources.

In one type of solid-source system, carrier gases (usually N_2 or O_2) flow at a controlled rate over a source boat placed in the furnace tube. The carrier gas picks up the vapor from the source and transports the dopant species to the wafer, where it is deposited on the surface of the wafer. The temperature of the source is controlled to maintain the desired vapor pressure. The source can be placed in a low-temperature section of the furnace or may be external to the furnace. Solid boron and phosphorus

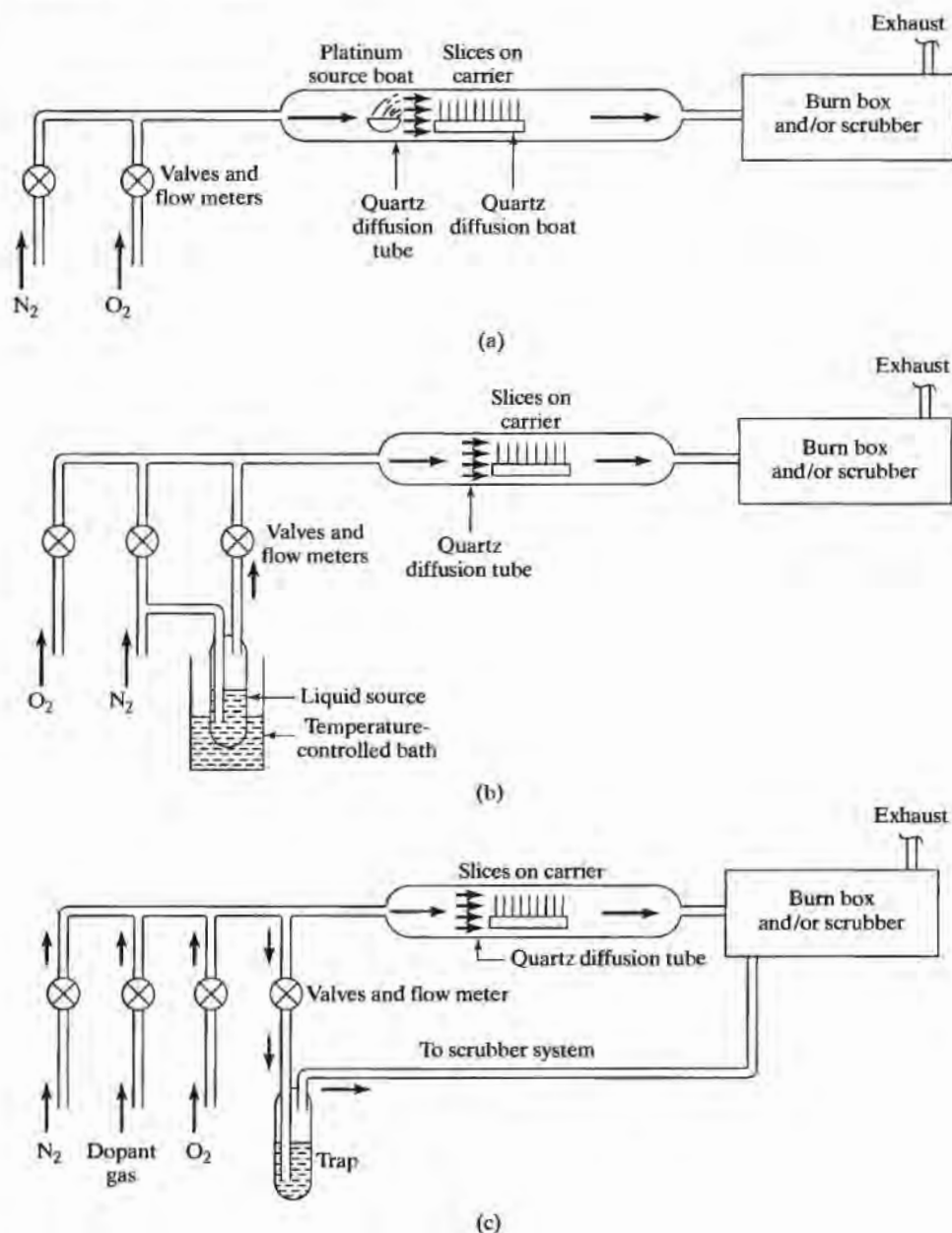


FIGURE 4.25

Open-furnace-tube diffusion systems. (a) Solid source in a platinum source boat in the rear of diffusion tube; (b) liquid-source system with carrier gas passing through a bubbler; (c) diffusion system using gaseous impurity sources. Copyright John Wiley and Sons. Reprinted with permission from Ref. [26].

impurity sources are also available in wafer form and are placed in the boat between adjacent pairs of silicon wafers.

In liquid-source systems, a carrier gas passes through a bubbler, where it picks up the vapor of the liquid source. The gas carries the vapor into the furnace tube, where it reacts with the surface of the silicon wafer.

Gas-source systems supply the dopant species directly to the furnace tube in the gaseous state. The common gas sources are extremely toxic, and additional input purging and trapping systems are required to ensure that all the source gas is removed from the system before wafer entry or removal. In addition, most diffusion processes either do not use all of the source gas or produce undesirable reaction by-products. Therefore, the output of diffusion systems must be processed by burning or by chemical or water scrubbing before being exhausted into the atmosphere.

Boron is the only commonly used *p*-type dopant. The diffusion coefficients of aluminum and gallium are quite high in silicon dioxide, and these elements cannot be masked effectively by SiO_2 . Indium is not used, because it is a relatively deep-level acceptor ($E_A - E_V = 0.14 \text{ eV}$).

In contrast, antimony, phosphorus, and arsenic can all be masked by silicon dioxide and are all routinely used as *n*-type dopants in silicon processing.

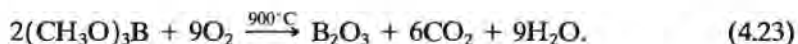
4.10.1 Boron Diffusion

Boron has a high solubility in silicon and can achieve active surface concentrations up to $4 \times 10^{20}/\text{cm}^3$. (See Fig. 4.6.) Elemental boron is inert up to temperatures exceeding the melting point of silicon. A surface reaction with boron trioxide (B_2O_3) is used to introduce boron to the silicon surface:



An excess amount of boron trioxide can cause a brown boron skin to form that is very difficult to remove with most acids. Boron skin formation can be minimized by performing the diffusions in an oxidizing atmosphere containing 3 to 10% oxygen. In a two-step diffusion, the boron predeposition step is commonly followed by a short wet-oxidation step to assist in removal of the boron skin prior to drive-in.

Common solid sources of boron include trimethylborate (TMB) and boron nitride wafers. TMB is a solid with high vapor pressure at room temperature. The TMB source is normally placed outside the diffusion furnace and cooled below room temperature during use. TMB vapor reacts in the furnace tube with oxygen to form boron trioxide, water, and carbon dioxide:



Unreacted TMB must be scrubbed from the exhaust stream.

Boron nitride is a solid source available in wafer form. Activated wafers are placed in every third slot in the same quartz boat used to hold the silicon wafers. A silicon wafer faces each side of the oxidized boron nitride wafer, and boron trioxide is transferred directly to the surface of the silicon wafer during high-temperature diffusion. A small flow of inert gas such as nitrogen is used to keep contaminants out of the tube during diffusion.

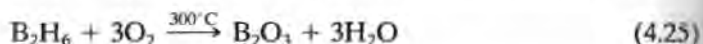
The most common liquid source for boron is boron tribromide (BBr_3). The reaction is



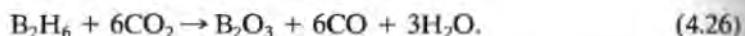
Free bromine combines easily with metallic impurities and is useful in removing (gettering) metallic impurities during diffusion. Bromine, as well as unused boron tribromide, is in the exhaust stream, so the outlet gases must be carefully cleaned.

The primary gaseous source of boron is diborane (B_2H_6). Diborane is a highly poisonous and explosive gas. Table 4.4 summarizes the ACGIH recommendations for the maximum permissible exposure to the common gases used as diffusion sources. Extreme care must be taken in using these gases. To reduce the risk of handling, diborane is usually diluted with 99.9% argon or nitrogen by volume.

Diborane oxidizes in either oxygen or carbon dioxide to form boron trioxide:



and



Both systems must provide a means for purging diborane from the input to the diffusion tube, and the output must be scrubbed to eliminate residual diborane and carbon monoxide.

4.10.2 Phosphorus Diffusion

Phosphorus has a higher solubility in silicon than does boron, and surface concentrations in the low $10^{21}/\text{cm}^3$ range can be achieved during high-temperature diffusion. Phosphorus is introduced into silicon through the reaction of phosphorus pentoxide at the wafer surface:



Solid P_2O_5 wafers can be used as a solid source for phosphorus, as can ammonium monophosphate ($\text{NH}_4\text{H}_2\text{PO}_4$) and ammonium diphosphate [$(\text{NH}_4)_2\text{H}_2\text{P}_2\text{O}_7$] in wafer form. However, the most popular diffusion systems use either liquid or gaseous sources. Phosphorus oxychloride (POCl_3) is a liquid at room temperature. A carrier gas is passed through a bubbler and brings the vapor into the diffusion furnace. The gas stream also contains oxygen, and P_2O_5 is deposited on the surface of the wafers:



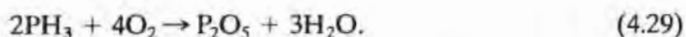
Liberated chlorine gas serves as a gettering agent, and Cl_2 and POCl_3 must be removed from the exhaust stream.

TABLE 4.4 Threshold Limit Recommendations for Common Gaseous Sources [24] *

Source	8-h exposure level (ppm)	Life-threatening exposure	Comments
Borane (B_2H_6)	0.10	160 ppm for 15 min	Colorless, sickly sweet, extremely toxic, flammable.
Phosphine (PH_3)	0.30	400 ppm for 30 min	Colorless, decaying fish odor, extremely toxic, flammable. A few minutes' exposure to 2000 ppm can be lethal.
Arsine (AsH_3)	0.05	6–15 ppm for 30 min	Colorless, garlic odor, extremely toxic. A few minutes' exposure to 500 ppm can be lethal.
Disilane (SiH_4)	0.50	Unknown	Repulsive odor, burns in air, explosive, poorly understood.
Dichlorosilane (SiH_2Cl_2)	5.00	...	Colorless, flammable, toxic. Irritating odor provides adequate warning for voluntary withdrawal from contaminated areas.

*Data from the 1979 American Conference of Governmental Hygienists (ACGIH).

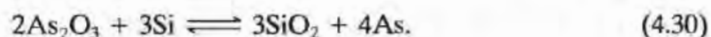
Phosphine, PH_3 , is a highly toxic, explosive gas used as the gaseous source for phosphorus. It is also supplied in dilute form with 99.9% argon or nitrogen. Phosphine is oxidized with oxygen in the furnace:



Unreacted phosphine must be cleaned from the exhaust gases, and the gas delivery system must be able to purge phosphine from the input to the tube.

4.10.3 Arsenic Diffusion

Arsenic has the highest solubility of any of the common dopants in silicon, with surface concentrations reaching $2 \times 10^{21}/\text{cm}^3$. The surface reaction involves arsenic trioxide:

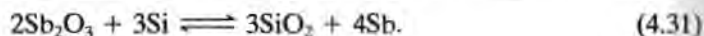


Oxide vapors can be carried into the furnace tube from a solid diffusion source by a nitrogen carrier gas. However, evaporation of arsenic from the surface limits surface concentrations to below $3 \times 10^{19}/\text{cm}^3$. The exhaust must be carefully cleaned, because of the presence of arsenic.

Arsine gas may be used as a source, but it is extremely toxic and also produces relatively low surface concentrations. The problems with arsenic deposition and safety delayed its widespread use in silicon processing until ion implantation was developed in the early 1970s. Ion implantation is now the preferred technique for introducing arsenic into silicon. (Chapter 5 is devoted to the subject of ion implantation.)

4.10.4 Antimony Diffusion

Antimony, like arsenic, has a low diffusion coefficient and has been used for a long time for buried layers in bipolar processes. Antimony trioxide is a solid source that is placed in a two-zone furnace in which the source is maintained at a temperature of 600 to 650 °C. Antimony is introduced at the silicon surface as in the other cases:



A liquid source, antimony pentachloride (Sb_2Cl_5), has been successfully used with oxygen as a carrier gas passing through a bubbler. The gas stibine (SbH_3) is unstable and is not used for antimony diffusion.

4.11 GETTERING

A process called gettering is often used to improve the quality of the silicon wafer as one of the first steps in the fabrication process. Gettering is used to remove unwanted impurities, typically heavy metals such as copper, gold, iron, and nickel, from the surface where the devices will be fabricated. These unwanted impurities can reduce both lifetime and mobility in silicon. The heavy metals tend to be fast diffusers in silicon and have high solubility in heavily doped *n*-type silicon. In general, the gettering techniques attempt to provide locations away from the active device regions where the undesired impurities can precipitate and be immobilized.

A number of different backside gettering techniques have been used successfully. In one of the earliest techniques, a heavily doped phosphorus diffusion is applied to the back of the wafer. In older bipolar processes, this gettering step often occurred naturally during the emitter diffusion step. Damage to the backside of the wafer is also effective and can be introduced by sandblasting the back of the wafer. Internal wafer stress can assist the gettering process and can be introduced into the wafer by depositing thin layers (0.1–0.5 μm) of either silicon nitride or low-temperature polysilicon on the backside of the wafer. Implantation of argon atoms is another method used to introduce damage into the wafer.

Oxygen, at levels as high as $10^{18}/\text{cm}^3$, is incorporated into silicon wafers during the crystal growth process. This oxygen can combine with other undesired impurities to form precipitates, and this technique is commonly referred to as intrinsic gettering. To be effective, the oxygen level must be controlled reasonably well during crystal growth, and a specific heat treatment cycle must be utilized. Defects in the original crystal structure or those introduced by epitaxial growth can also provide gettering sites. An excellent introduction to gettering can be found in reference [31].

SUMMARY

In Chapter 4, we have discussed the formation of *pn* junctions using high-temperature diffusion. Mathematical models for diffusion have been presented, and the behavior of common *n*- and *p*-type dopants in silicon has been discussed. A key parameter governing the diffusion process is the diffusion coefficient, which is highly temperature dependent, following an Arrhenius relationship. Boron, phosphorus, antimony, and arsenic all have reasonable diffusion coefficients in silicon at temperatures between 900 and 1200° C, and they can be conveniently masked by a barrier layer of silicon dioxide. Gallium and aluminum are not easily masked by SiO_2 and are seldom used, and indium is not used, because of its large activation energy. At high concentrations, diffusion coefficients become concentration-dependent, causing diffused profiles to differ substantially from predictions of simple theories.

Two types of diffusions are most often used. If the surface concentration is maintained constant throughout the diffusion process, then a complementary error function (erfc) distribution is obtained. In the erfc case, the surface concentration is usually set by the solid-solubility limit of the impurity in silicon. If a fixed dose of impurity is diffused into silicon, a Gaussian diffusion profile is achieved. These two cases are often combined in a two-step process to obtain lower surface concentrations than those achievable with a solid-solubility-limited diffusion. As the complexity of fabrication processes grows, simulation with process modeling programs such as SUPREM is becoming ever more important. The SUPREM program includes comprehensive models for single- and multi-dimensional diffusion, and it represents an extremely useful tool for modeling advanced device structures.

The concept of sheet resistance has been introduced, and Irvin's curves have been used to relate the sheet resistance, junction depth, and surface concentration of diffused layers. Techniques for calculating and measuring junction depth have also been presented. Resistor fabrication has been discussed, including end and corner effects, as well as the effects of lateral diffusion under the edges of diffusion barriers.

High-temperature open-furnace diffusion systems are routinely used for diffusion with solid, liquid, and gaseous impurity sources. Boron, phosphorus, and antimony are all easily introduced into silicon using high-temperature diffusion. However, arsenic deposition by diffusion is much more difficult, and today it is usually accomplished using ion implantation. (See Chapter 5.) As with many chemicals used in IC fabrication, some of the sources used for diffusion are extremely toxic and must be handled with great care.

REFERENCES

- [1] R. F. Pierret, *Semiconductor Fundamentals*, Volume I in the Modular Series on Solid State Devices, Addison-Wesley, Reading, MA, 1983.
- [2] R. A. Colclaser, *Microelectronics: Processing and Device Design*, John Wiley & Sons, New York, 1980.
- [3] R. F. Pierret, *Advanced Semiconductor Fundamentals*, Volume VI in the Modular Series on Solid State Devices, Addison-Wesley, Reading, MA, 1987.
- [4] D. P. Kennedy and R. R. O'Brien, "Analysis of the Impurity Atom Distribution Near the Diffusion Mask for a Planar p - n Junction," *IBM Journal of Research and Development*, 9, 179-186 (May, 1965).
- [5] J. C. Irvin, "Resistivity of Bulk Silicon and of Diffused Layers in Silicon," *Bell System Technical Journal*, 41, 387-410 (March, 1962).
- [6] L. R. Weisberg and J. Blanc, "Diffusion with Interstitial-Substitutional Equilibrium: Zinc in GaAs," *Physical Review*, 131, 1548-1552 (August 15, 1963).
- [7] R. B. Fair, "Boron Diffusion in Silicon—Concentration and Orientation Dependence, Background Effects, and Profile Estimation," *Journal of the Electrochemical Society*, 122, 800-805 (June, 1975).
- [8] R. B. Fair, "Profile Estimation of High-Concentration Arsenic Diffusions in Silicon," *Journal of Applied Physics*, 43, 1278-1280 (March, 1972).
- [9] R. B. Fair and J. C. C. Tsai, "Profile Parameters of Implanted-Diffused Arsenic Layers in Silicon," *Journal of the Electrochemical Society*, 123, 583-586 (1976).
- [10] J. C. C. Tsai, "Shallow Phosphorus Diffusion Profiles in Silicon," *Proceedings of the IEEE*, 57, 1499-1506 (September, 1969).
- [11] R. B. Fair and J. C. C. Tsai, "A Quantitative Model for the Diffusion of Phosphorus in Silicon and the Emitter Dip Effect," *Journal of the Electrochemical Society*, 124, 1107-1118 (July, 1977).
- [12] W. R. Runyan, *Semiconductor Measurements and Instrumentation*, McGraw-Hill, New York, 1975.
- [13] L. J. van der Pauw, "A Method of Measuring Specific Resistivity and Hall Effect of Discs of Arbitrary Shape," *Philips Research Reports*, 13, 1-9 (February, 1958).
- [14] R. Chwang, B. J. Smith, and C. R. Crowell, "Contact Size Effects on the van der Pauw Method for Resistivity and Hall Coefficient Measurements," *Solid-State Electronics*, 17, 1217-1227 (December, 1974).
- [15] P. F. Kane, and G. B. Larrabee, *Characterization of Semiconductor Materials*, McGraw-Hill Book Company, New York, 1970.
- [16] (a) C. W. White and W. H. Cristie, "The Use of RBS and SIMS to Measure Dopant Profile Changes in Silicon by Pulsed Laser Annealing," *Solid-State Technology*, pp. 109-116, September 1980. (b) J. M. Anthony et al., "Super SIMS for Ultra Sensitive Impurity Analysis," *Proceedings of the Materials Research Society Symposium* 69, pp. 311-316, 1986.
- [17] A. Agarwal et al., "Boron-enhanced-Diffusion of Boron: The Limiting Factor for Ultra-shallow Junctions," *IEEE IEDM Technical Digest*, pp. 467-470, December 1997.
- [18] D. A. Antoniadis and R. W. Dutton, "Models for Computer Simulation of Complete IC Fabrication Processes," *IEEE Journal of Solid State Circuits*, SC-14, 412-422 (April, 1979).
- [19] C. P. Ho, J. D. Plummer, S. E. Hansen, and R. W. Dutton, "VLSI Process Modeling—SUPREM III," *IEEE Trans. Electron Devices*, ED-30, 1438-1453 (November, 1983).

- [20] D. Chin, M. Kump, H. G. Lee, and R. W. Dutton, "Process Design Using Coupled 2D Process and Device Simulators," *IEEE IEDM Digest*, pp. 223–226 (December, 1980).
- [21] C. D. Maldonado, F. Z. Custode, S. A. Louie, and R. K. Pancholy, "Two Dimensional Simulation of a 2 μm CMOS Process Using ROMANS II," *IEEE Trans. Electron Devices*, 30, 1462–1469 (November, 1983).
- [22] M. E. Law, C. S. Rafferty, and R. W. Dutton, "New n -well Fabrication Techniques Based on 2D Process Simulation," *IEEE IEDM Digest*, pp. 518–521 (December, 1986).
- [23] R. W. Dutton, "Modeling and Simulation for VLSI," *IEEE IEDM Digest*, pp. 2–7 (December, 1986).
- [24] *Matheson Gas Data Book*, Matheson Gas Products, 1980.
- [25] A. B. Glaser and G. E. Subak-Sharpe, *Integrated Circuit Engineering*, Addison-Wesley, Reading, MA, 1979.
- [26] S. K. Ghandhi, *VLSI Fabrication Principles*, John Wiley & Sons, New York, 1983.
- [27] S. M. Sze, Ed., *VLSI Technology*, McGraw-Hill, New York, 1983.
- [28] S. K. Ghandhi, *The Theory and Practice of Microelectronics*, John Wiley & Sons, New York, 1968.
- [29] R. B. Fair, "Recent Advances in Implantation and Diffusion Modeling for the Design and Process Control of Bipolar ICs," *Semiconductor Silicon 1977*, PV 77-2, pp. 968–985.
- [30] F. M. Smits, "Measurement of Sheet Resistivities with the Four Point Probe," *Bell System Technical Journal*, 37, 711–718 (May, 1958).
- [31] W. R. Runyan and K. E. Bean, *Semiconductor Integrated Circuit Processing Technology*, Addison Wesley Publishing Company, Reading, MA, 1990.

PROBLEMS

- 4.1 A phosphorus diffusion has a surface concentration of $5 \times 10^{18}/\text{cm}^3$, and the background concentration of the p -type wafer is $1 \times 10^{15}/\text{cm}^3$. The Dt product for the diffusion is 10^{-8} cm^2 .
 - (a) Find the junction depth for a Gaussian distribution.
 - (b) Find the junction depth for an erfc profile.
 - (c) What is the sheet resistance of the two diffusions?
 - (d) Draw a graph of the two profiles.
- 4.2 A 5-hr boron diffusion is to be performed at 1100°C .
 - (a) What thickness of silicon dioxide is required to mask this diffusion?
 - (b) Repeat part (a) for phosphorus.
- 4.3 A boron diffusion into a 1-ohm-cm n -type wafer results in a Gaussian profile with a surface concentration of $5 \times 10^{18}/\text{cm}^3$ and a junction depth of $4 \mu\text{m}$.
 - (a) How long did the diffusion take if the diffusion temperature was 1100°C ?
 - (b) What was the sheet resistance of the layer?
 - (c) What is the dose in the layer?
 - (d) The boron dose was deposited by a solid-solubility-limited diffusion. Design a diffusion schedule (temperature and time) for this predeposition step.

- 4.4 The boron diffusion in Problem 4.3 is followed by a solid-solubility-limited phosphorus diffusion for 30 min at 950 °C. Assume that the boron profile does not change during the phosphorus diffusion.
- Find the junction depth of the new phosphorus layer. Assume an erfc profile.
 - Find the junction depth based on the concentration-dependent diffusion data presented in Fig. 4.21.
 - Calculate the total Dt product for Prob. 4.3 and compare the result to the Dt product for this problem. Is the assumption in the problem statement justified?
- 4.5 The p -well in a CMOS process is to be formed by a two-step boron diffusion into a 5-ohm-cm n -type substrate. The sheet resistance of the well is 1000 ohms per square and the junction depth is 7.5 μm .
- Design a reasonable diffusion schedule for the drive-in step that produces the p -well.
 - What is the final surface concentration in the p -well?
 - What is the dose required to form the well?
 - Can this dose be achieved using a solid-solubility-limited diffusion with diffusion temperatures of 900 °C or above? Explain.
- 4.6 (a) Calculate the Dt product required to form a 0.2- μm -deep source-drain diffusion for an NMOS transistor using a solid-solubility-limited arsenic deposition at 1000° C into a wafer with a background concentration of $3 \times 10^{16}/\text{cm}^3$.
- What is the diffusion time? Does this time seem like a reasonable process?
 - Recalculate Dt based upon the model in Fig. 4.20(e) and Table 4.2.
- 4.7 The channel length of a silicon-gate NMOS transistor is the spacing between the source and drain diffusions as shown in Fig. P4.7a. The spacing between the source and drain diffusion openings is 3 μm on the masking oxide used to make the transistor. The source/drain junctions are diffused to a depth of 0.5 μm using a constant-source diffusion. The surface concentration is $1 \times 10^{20}/\text{cm}^3$ and the wafer has a concentration of $1 \times 10^{16}/\text{cm}^3$. What is the channel length in the actual device after the diffusion is completed?

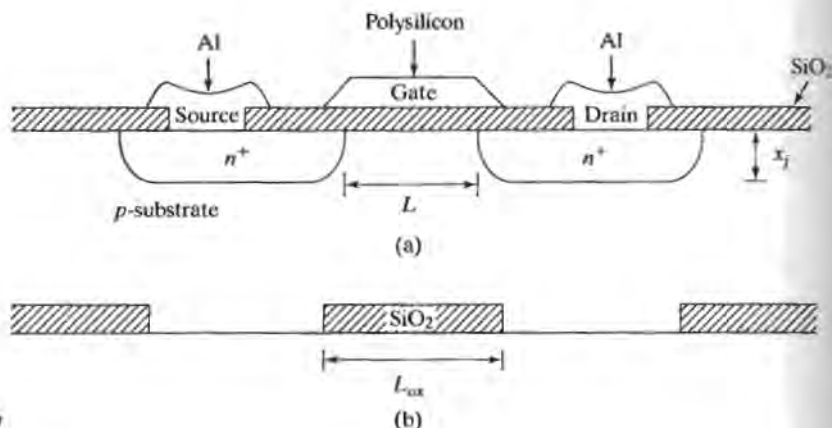


FIGURE P4.7

- 4.8 (a)** What is the total number of squares in the resistor shown in Fig. P4.8, assuming that its geometry is specified precisely by the mask dimensions?
- (b)** The resistor is actually formed from a p -type base diffusion with a $6\text{-}\mu\text{m}$ junction depth. What is the actual number of squares in this resistor, assuming that the lateral diffusion under the edge of the mask is $5\text{ }\mu\text{m}$.
- (c)** What would be the resistance of the resistors in parts (a) and (b) if the surface concentration of the base diffusion was 5×10^{18} boron atoms/cm³, the bulk concentration 10^{15} /cm³, and the junction depth $6\text{ }\mu\text{m}$.

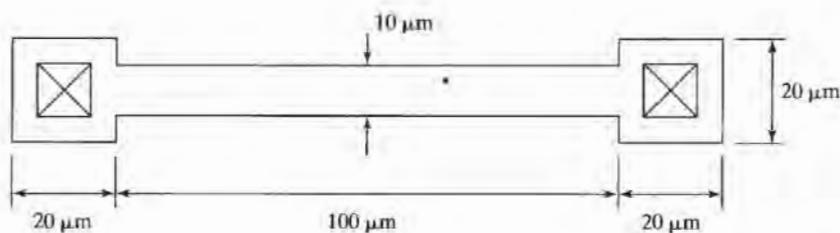


FIGURE P4.8

- 4.9 (a)** What is the total number of squares in the resistor drawn in Fig. P4.9, assuming that its horizontal geometry is identical to that on the mask in the figure.
- (b)** The resistor is actually formed from a $3\text{-}\mu\text{m}$ -deep diffusion that has a surface concentration of 3×10^{18} /cm³ in a wafer with background concentration of 10^{16} /cm³. What is actual number of squares in the resistor? (Use Fig. 4.10.)
- (c)** What is the resistance of the resistor if the diffusion is an n -type Gaussian layer?

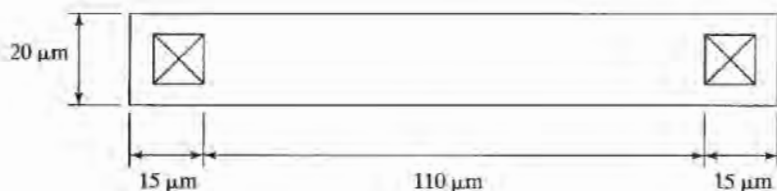


FIGURE P4.9

- 4.10** (a) Draw cross sections of the resistor in Fig. P4.10 through A'-A and B'-B if the lateral diffusion = 0.5λ and vertical diffusion = λ .
- (b) What is the total number of squares in the resistor as drawn on the masks?
- (c) What is the actual number of squares in the resistor for the diffusions given in part (a)?
- (d) Suppose the resistor is formed from a $2\text{-}\mu\text{m}$ -deep diffusion that has a surface concentration of $10^{19}/\text{cm}^3$ into a background concentration of $10^{16}/\text{cm}^3$. What is the actual number of squares in the resistor if $\lambda = 2\text{ }\mu\text{m}$? (Use Fig. 4.10.)
- (e) A processing error occurred and the actual junction depth in part (b) is discovered to be $3\text{ }\mu\text{m}$. Estimate the new number of squares in the resistor.

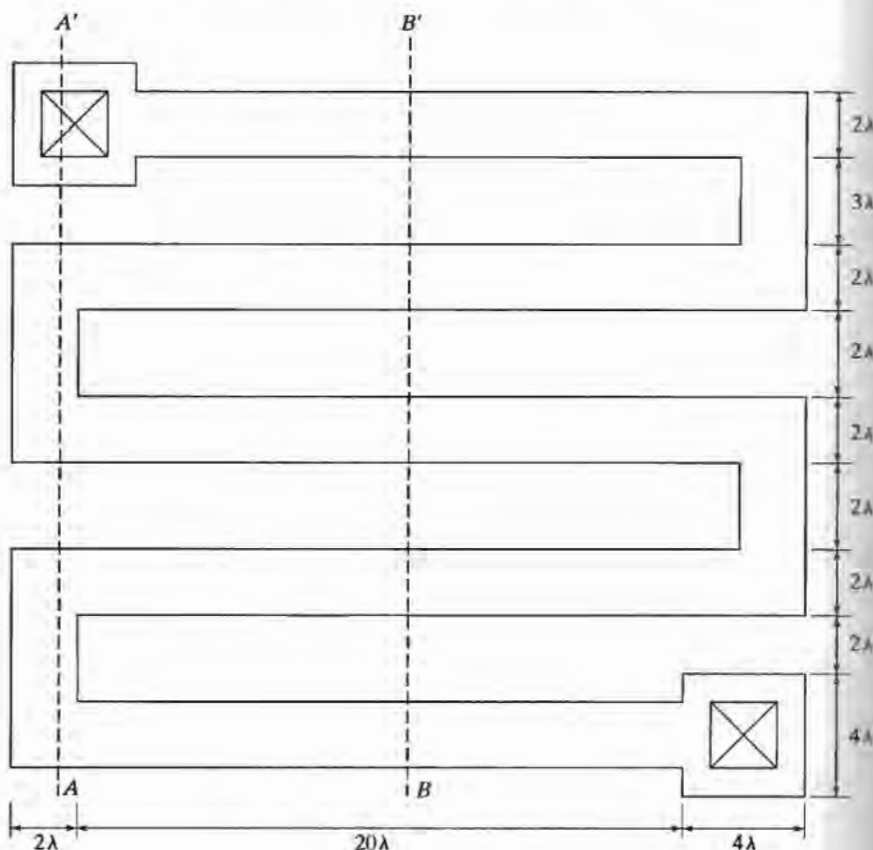


FIGURE P4.10

- 4.11** In practice, wafers are slowly pushed into and pulled out of the furnace, or the furnace temperature may be changed with time. Assume that the furnace temperature is being ramped down with time: $T = T_0 - Rt$, where T_0 is the initial temperature and R is the temperature change per second. Show that the effective Dt product defined by

$$(Dt)_{\text{eff}} = \int_0^{t_0} D(t) dt$$

where t_0 is the ramp-down time is given by

$$(Dt)_{\text{eff}} = D(T_0)(kT_0^2/RE_A)$$

where

$$D(T_0) = D_0 \exp(-E_A/kT_0)$$

- 4.12 Determine the sensitivity of junction depth to changes in furnace temperature by calculating $(dx_j/x_j)/(dT/T)$ for a Gaussian diffusion profile. What fractional change in junction depth will occur at 1100° C if the furnace temperature is in error by 10 °C?
- 4.13 Design (choose times and temperatures) a two-step diffusion to form a 5- μm -deep n -type layer with a surface concentration of $5 \times 10^{16}/\text{cm}^3$ in a 10 $\Omega\text{-cm}$ p -type substrate. (This layer could be the n -well for a CMOS process.)
- 4.14 Use Eq. (4.13) to calculate the sheet resistance of a p -type Gaussian diffusion having a surface concentration of $2 \times 10^{18}/\text{cm}^3$ and a junction depth of 2 μm . Assume that the hole mobility has a constant value of 300 $\text{cm}^2/\text{V}\cdot\text{sec}$.
- 4.15 An n -type Gaussian diffusion has a surface concentration of $10^{20}/\text{cm}^3$ and a junction depth of 2 μm . (a) Calculate the sheet resistance contribution of the first 0.5 μm of the layer, the second 0.5 μm , the third 0.5 μm , and the last 0.5 μm . (b) What is the sheet resistance of the total diffusion? Assume that the electron mobility has a constant value of 100 $\text{cm}^2/\text{V}\cdot\text{sec}$. (c) Compare the calculation to the prediction from Irvin's curves.
- 4.16 Rework Example 4.3 using the concentration-dependent boron diffusion expressions for the predeposition calculations. Find the new surface concentration and junction depth following the drive-in step, and compare the results with those presented in the example. At 900 °C, $n_i = 4 \times 10^{16}/\text{cm}^3$.
- 4.17 Derive the expressions for the Gaussian and complementary-error-function solutions to the diffusion equation.
- 4.18 What is the minimum sheet resistance to be expected from shallow arsenic- and boron-doped regions if the regions are 1 μm deep? 0.25 μm deep? Make use of Figs. 4.6 and 4.16. Assume that the region is uniformly doped. Compare your results to the equations presented in Section 4.6.
- 4.19 Gold is diffused into a silicon wafer using a constant-source diffusion with a surface concentration of $10^{18}/\text{cm}^3$. How long does it take the gold to diffuse completely through a silicon wafer 400 μm thick with a background concentration of $10^{16}/\text{cm}^3$ at a temperature of 1000° C?
- 4.20 Use SUPREM to simulate the (two-dimensional) results of the two-step diffusion from Ex. 4.3 through a 5- μm -wide opening in a silicon dioxide barrier layer. Plot the results.
- 4.21 (a) Use SUPREM to simulate the diffusion profile of Example 4.3. Compare the simulation results with those given in the example.
 (b) Follow the boron diffusion by the growth of a 500-nm layer of oxide in wet oxygen at 1100° C. Discuss what has happened to the boron concentration at the Si-SiO₂ interface.
 (c) Add a 30-min solid-solubility-limited phosphorus diffusion at 1000° C.
 (d) The phosphorus diffusion created a new pn junction. Update the hand calculations for the boron impurity profile of Ex. 4.3 and estimate the location of both pn junctions with the aid of Fig. 4.21. Compare your results with those of SUPREM in part (c).
- 4.22 A current of 10 μA is injected into a van der Pauw structure having a sheet resistance of 300 Ω/\square . What is the voltage that should be measured at the second set of terminals?

4.23 A gas cylinder contains 100 ft^3 of a mixture of diborane and argon. The diborane represents 0.1% by volume. An accident occurs and the complete cylinder is released into a room measuring $10 \times 12 \times 8 \text{ ft}$.

- (a) What will be the equilibrium concentration of diborane in the room in ppm?
- (b) Compare this level with the toxic level based on Table 4.3.
- (c) Would your answer to part (b) change if the gas cylinder contained arsine?

4.24 (a) Numerically calculate the sheet resistance of the diffusion in Problem 4.15 if the electron mobility can be described by

$$\mu_n = \left[92 + \frac{1270}{1 + \left(\frac{N}{1.3 \times 10^{17}} \right)^{.091}} \right] \frac{\text{cm}^2}{\text{V-sec}}$$

- (b) Repeat for a similar p-type Gaussian layer in an n-type substrate if the hole mobility is given by

$$\mu_p = \left[48 + \frac{447}{1 + \left(\frac{N}{6.3 \times 10^{16}} \right)^{.076}} \right] \frac{\text{cm}^2}{\text{V-sec}}$$

CHAPTER 5

Ion Implantation

Ion implantation offers many advantages over diffusion for the introduction of impurity atoms into the silicon wafer and has become a workhorse technology in modern IC fabrication. In this chapter, we will first discuss ion implantation technology and mathematical modeling of the impurity distributions obtained with ion implantation. We will subsequently explore deviations from the model caused by nonideal behavior and will discuss annealing techniques used to remove crystal damage caused by the implantation process.

5.1 IMPLANTATION TECHNOLOGY

An ion implanter is a high-voltage particle accelerator producing a high-velocity beam of impurity ions that can penetrate the surface of silicon target wafers. The following list shows the basic parts of the system, shown schematically in Fig. 5.1, beginning with the impurity-source end of the system.

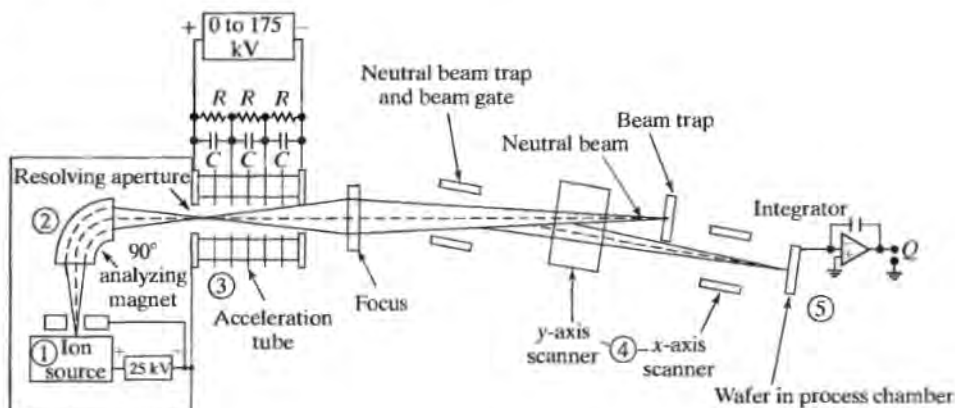


FIGURE 5.1

Schematic drawing of a typical ion implanter showing (1) ion source, (2) mass spectrometer, (3) high-voltage acceleration column, (4) x- and y-axis deflection system, and (5) target chamber.

1. **Ion Source.** The ion source operates at a high voltage (25 kV) and produces a plasma containing the desired impurity, as well as other undesired species. Arsine, phosphine, and diborane, as well as other gases, can be used in the source. Solids can be sputtered in special ion sources, and this technique offers a wide degree of flexibility in the choice of impurity.
2. **Mass Spectrometer.** An analyzer magnet bends the ion beam through a right angle to select the desired impurity ion. The selected ion passes through an aperture slit into the main accelerator column.
3. **High-Voltage Accelerator.** The accelerator column adds energy to the beam (up to 5 MeV) and accelerates the ions to their final velocity. Both the accelerator column and the ion source are operated at a high voltage relative to the target. For protection from high voltage and possible X-ray emission, the ion source and accelerator are mounted within a protective shield.
4. **Scanning System.** x - and y -axis deflection plates are used to scan the beam across the wafer to produce uniform implantation and to build up the desired dose. The beam is bent slightly to prevent neutral particles from hitting the target.
5. **Target Chamber.** Silicon wafers serve as targets for the ion beam. For safety, the target area is maintained near ground potential. The complete implanter system is operated under vacuum conditions.

The analyzer magnet is used to select the desired impurity ions from the output of the source. A charged particle moving with velocity \mathbf{v} through a magnetic field \mathbf{B} will experience a force \mathbf{F} , given by

$$\mathbf{F} = q(\mathbf{v} \times \mathbf{B}) \quad (5.1)$$

The force will tend to move the particle in a circle and the centrifugal force will balance \mathbf{F} . For the case where \mathbf{B} is perpendicular to \mathbf{v} , $q\mathbf{v} \times \mathbf{B} = m|\mathbf{v}|^2/r$, where $m|\mathbf{v}|^2/2 = qV$ and V is the accelerator voltage. Thus, the magnitude of the magnetic field \mathbf{B} may be adjusted to select an ion species with a given mass:

$$|\mathbf{B}| = \sqrt{(2mV/qr^2)} \quad (5.2)$$

The ion source in the figure operates at a constant potential (25 keV) so that the voltage V is known, and an ion species is selected by changing the dc current supplying the analyzer magnet. The selected impurity is then accelerated to its final velocity in the high-voltage column.

The silicon wafer is maintained in good electrical contact with the target holder, so electrons can readily flow to or from the wafer to neutralize the implanted ions. This electron current is integrated over time to measure the total dose Q from the implanter given by

$$Q = \frac{1}{mqA} \int_0^T I dt \quad (5.3)$$

where I is the beam current in amperes, A is the wafer area, $n = 1$ for singly ionized ions and 2 for doubly ionized species, and T is the implantation time. The use of a doubly ionized species increases the energy capability of the machine by a factor of 2, since $E = nqV$.

The target wafers can be maintained at relatively low temperatures during the implantation. Low-temperature processing prevents undesired spreading of impurities by diffusion, which is very important in VLSI fabrication. Another advantage of ion implantation is the ability to use a much wider range of impurity species than possible with diffusion. In principle, any element that can be ionized can be introduced into the wafer using implantation.

A production-level ion implanter costs millions of dollars, and this high cost is its greatest disadvantage. However, the advantages of flexibility and tight process control have far outweighed the disadvantage of cost, and ion implantation is now used routinely throughout bipolar and MOS integrated-circuit fabrication.

5.2 MATHEMATICAL MODEL FOR ION IMPLANTATION

As an ion enters the surface of the wafer, it collides with atoms in the lattice and interacts with electrons in the crystal. Each nuclear or electronic interaction reduces the energy of the ion until it finally comes to rest within the target. Interaction with the crystal is a statistical process, and the implanted impurity profile can be approximated by the Gaussian distribution function illustrated in Fig. 5.2. The distribution is described mathematically by

$$N(x) = N_p \exp \left[-\frac{(x - R_p)^2}{2\Delta R_p^2} \right] \quad (5.4)$$

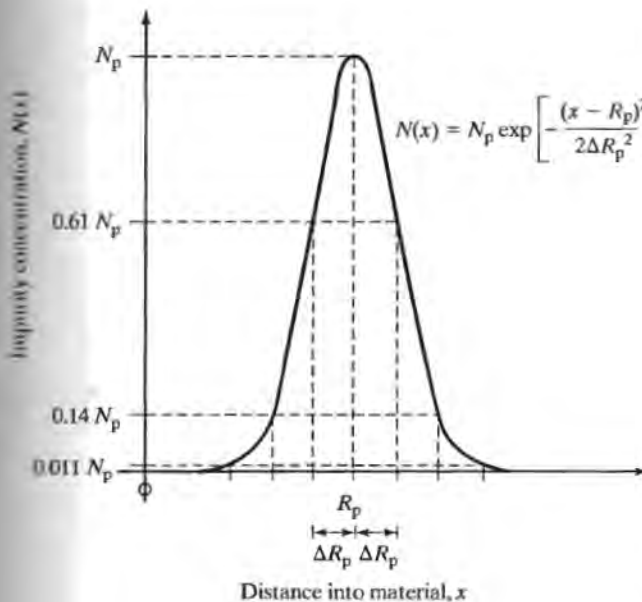


FIGURE 5.2

Gaussian distribution resulting from ion implantation. The impurity is shown implanted completely below the wafer surface ($x = 0$).

R_p is called the *projected range* and is equal to the average distance an ion travels before it stops. The peak concentration N_p occurs at $x = R_p$. Because of the statistical nature of the process, some ions will be "lucky" and will penetrate beyond the projected range R_p , and some will be "unlucky" and will not make it as far as R_p . The spread of the distribution is characterized by the standard deviation, ΔR_p , called the *straggle*.

The area under the impurity distribution curve is the implanted dose Q , defined by

$$Q = \int_0^{\infty} N(x) dx \quad (5.5)$$

For an implant completely contained within the silicon, the dose is equal to

$$Q = \sqrt{2\pi} N_p \Delta R_p \quad (5.6)$$

Implanted doses typically range from $10^{10}/\text{cm}^2$ to $10^{18}/\text{cm}^2$. For example, ion implantation is often used to replace the predeposition step in a two-step diffusion process. Doses in the range of 10^{10} to $10^{13}/\text{cm}^2$ are required for threshold adjustment in MOS technologies and are almost impossible to achieve using diffusion. CMOS well formation is another example where the dose control of the ion implanter is a distinct advantage. However, doses exceeding $10^{15}/\text{cm}^2$ are quite large and can be time-consuming to produce using ion implantation. As a reference for comparison, the silicon lattice atomic sheet density is approximately 7×10^{14} silicon atoms/ cm^2 on the $\langle 100 \rangle$ surface.

The implanted dose can be controlled within a few percent, and this tight control represents a major advantage of ion implantation. For example, resistors can be fabricated with absolute tolerances of a few percent in carefully controlled processes using ion implantation, whereas the same resistors would have an absolute tolerance exceeding 20% if they were formed using only diffusion.

The projected range of a given ion is a function of the energy of the ion, and of the mass and atomic number of both the ion and the target material. A theory for range and straggle was developed by Lindhard, Scharff, and Schiott and is called the *LSS theory* [1]. This theory assumes that the implantation goes into an amorphous material in which the atoms of the target material are randomly positioned. Figure 5.3 displays the results of LSS calculations for the projected range and straggle for antimony, boron, phosphorus, and arsenic in amorphous silicon and silicon dioxide. For the moment, we will assume that these results are also valid for crystalline silicon. Deviations from the LSS theory will be discussed in Section 5.5.

Range and straggle are roughly proportional to ion energy over a wide range, although nonlinear behavior is clearly evident in the figure. For a given energy, the lighter elements strike the silicon wafer with a higher velocity and penetrate more deeply into the wafer. The results indicate that the projected ranges in Si and SiO_2 are essentially the same, and we will assume that the stopping power of silicon dioxide is equal to that of silicon. Figure 5.3 also gives values for the transverse straggle ΔR_{\perp} , which will be discussed in the next section.

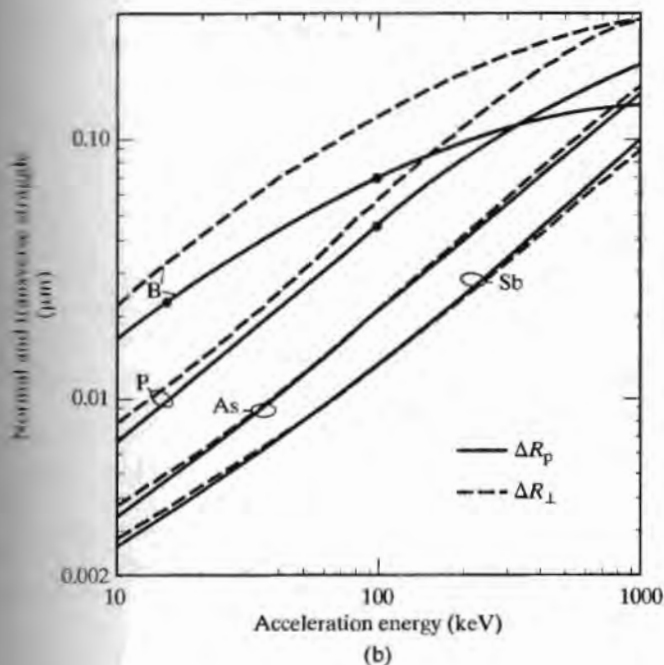
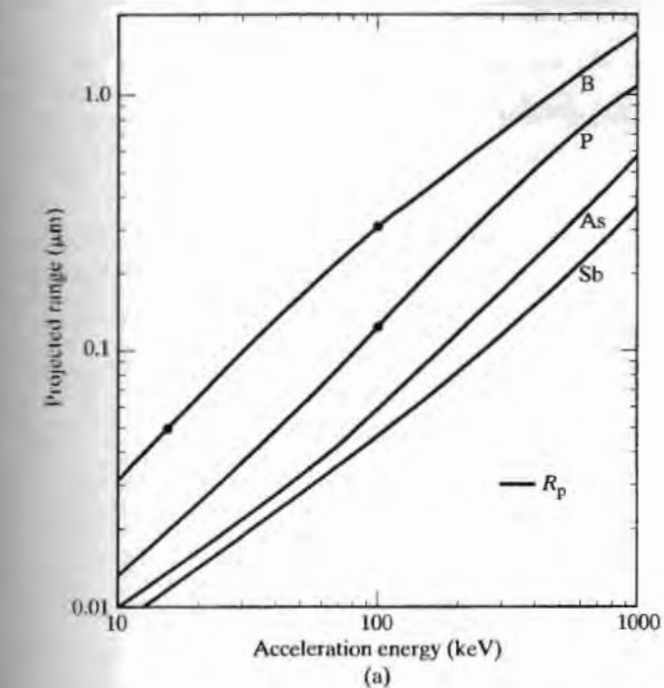


FIGURE 5.3

Projected range and straggle calculations based on LSS theory. (a) Projected range R_p for boron, phosphorus, arsenic, and antimony in amorphous silicon. Results for SiO_2 and for silicon are virtually identical. (b) Vertical ΔR_p and transverse ΔR_\perp straggle for boron, phosphorus, arsenic, and antimony. Reprinted with permission from Ref. [2]. (Copyright van Nostrand Reinhold Company, Inc.)

Example 5.1

Phosphorus with an energy of 100 keV is implanted into a silicon wafer. **(a)** What are the range and straggle associated with this implantation? **(b)** What should the implanted dose be if a peak concentration of $1 \times 10^{17}/\text{cm}^3$ is desired? **(c)** What length of time is required to implant this dose into a 200-mm wafer using a 2 μA beam current with singly ionized phosphorus?

Solution: Using Fig. 5.3, we find that the range and straggle are 0.12 μm and 0.045 μm , respectively. The dose and peak concentration are related by Eq. (5.6). Note that this is an approximation, since the peak is only a little over $2\Delta R_p$ below the silicon surface.

$$Q = \sqrt{2\pi} N_p \Delta R_p = \sqrt{2\pi} (1 \times 10^{17}/\text{cm}^3) (4.5 \times 10^{-6} \text{cm}) = 1.13 \times 10^{12}/\text{cm}^2$$

Rearranging Eq. (5.3) assuming a constant beam current gives:

$$T = \frac{nqAQ}{I} = \frac{(1)(1.6 \times 10^{-19} \text{coul})(\pi)(10 \text{cm})^2(1.13 \times 10^{12}/\text{cm}^2)}{2 \times 10^{-6} \text{coul/sec}} = 28.4 \text{sec}$$

5.3 SELECTIVE IMPLANTATION

In most cases, we desire to implant impurities only in selected areas of the wafer. Windows are opened in a barrier material wherever impurity penetration is desired. In the center of the window, the impurity distribution is described by Eq. (5.4), but near the edges the distribution decreases and actually extends under the edge of the window, as shown in Fig. 5.4. The overall distribution can be modeled by[3]

$$N(x, y) = N(x)F(y)$$

$$F(y) = 0.5[\text{erfc}\{(y-a)/\sqrt{2}\Delta R_\perp\} - \text{erfc}\{(y+a)/\sqrt{2}\Delta R_\perp\}] \quad (5.7)$$

where $N(x)$ is given by Eq. (5.4) and $2a$ is the width of the window. The parameter ΔR_\perp is called the *transverse straggle* and characterizes the behavior of the distribution near the edge of the window. Figure 5.4 shows normalized impurity distributions near the barrier edge calculated with Eq. (5.7). Figure 5.3(b) gives values of both normal straggle ΔR_p and transverse straggle ΔR_\perp .

In order to mask the ion implantation, it is necessary to prevent the implanted impurity from changing the doping level in the silicon beneath the barrier layer. Figure 5.5 shows a silicon wafer with a layer of silicon dioxide on the surface. An impurity has been implanted into the wafer with the peak of the distribution in the silicon dioxide. To prevent significantly altering the doping in the silicon, we require that the implanted concentration be less than 1/10th the background concentration at the interface between the silicon and silicon dioxide:

$$N(X_0) < N_B/10 \quad (5.8)$$

or

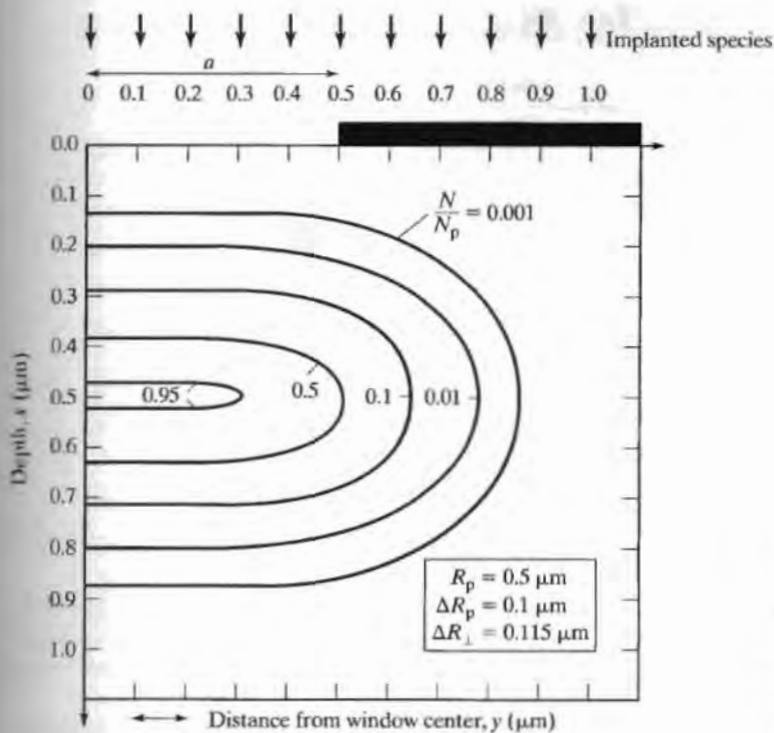


FIGURE 5.4

Contours of equal ion concentration for an implantation into silicon through a $1\text{ }\mu\text{m}$ window. The profiles are symmetrical about the x -axis and were calculated using Eq. (5.7), which is taken from Ref. [3].

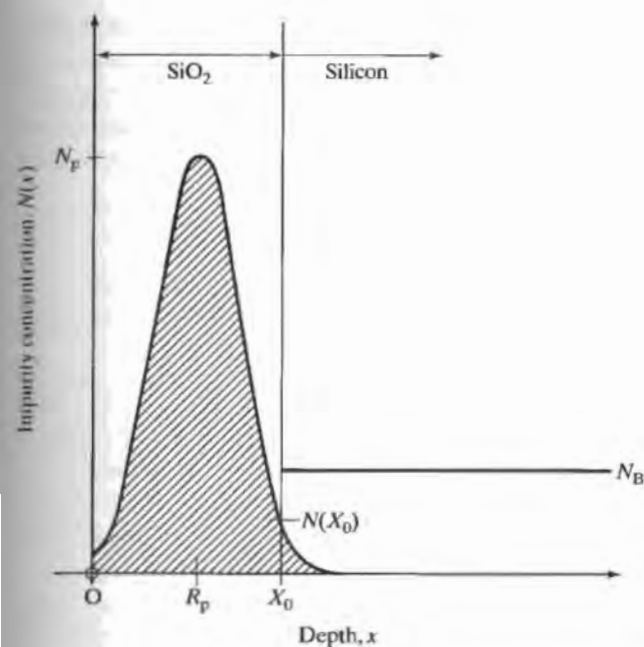


FIGURE 5.5

Implanted impurity profile with implant peak in the oxide. The barrier material must be thick enough to ensure that the concentration in the tail of the distribution is much less than N_B .

TABLE 5.1 Values of m for Various Values of N_p/N_B

N_p/N_B	m
10^1	3.0
10^2	3.7
10^3	4.3
10^4	4.8
10^5	5.3
10^6	5.7

$$N_p \exp\left[-\frac{(X_0 - R_p)^2}{2\Delta R_p^2}\right] < N_B/10 \quad (5.9)$$

Solving Eq. (5.9) for X_0 yields a minimum oxide thickness of

$$X_0 = R_p + \Delta R_p \sqrt{2 \ln(10N_p/N_B)} = R_p + m \Delta R_p \quad (5.10)$$

The oxide thickness must be at least equal to the projected range plus some multiple m times the straggle. Table 5.1 gives values of m for various ratios of peak concentration to background concentration. An oxide thickness equal to the projected range plus six times the straggle should mask most ion implantations.

Silicon dioxide and silicon nitride are routinely used as barrier materials during implantation. Since implantation is a low-temperature process, additional materials such as photoresist and aluminum, which cannot withstand high-temperature diffusion, may be used as barrier materials during the implantation.

Silicon nitride is more effective than silicon dioxide in stopping ions, and a silicon nitride barrier layer need only be 85% of the thickness of an SiO_2 barrier layer. On the other hand, photoresist is less effective in stopping ions, and a photoresist barrier layer must be 1.8 times the thickness of an SiO_2 layer under the same implantation conditions. Metals are of such a high density that even a very thin layer will mask most implantations.

Example 5.2

A boron implantation is to be performed through a 50-nm gate oxide so that the peak of the distribution is at the Si-SiO₂ interface. The dose of the implant in silicon is to be $1 \times 10^{13}/\text{cm}^2$. **(a)** What are the energy of the implant and the peak concentration at the interface? **(b)** How thick should the SiO₂ layer be in areas that are not to be implanted, if the background concentration is $1 \times 10^{16}/\text{cm}^3$? **(c)** Suppose the oxide is 50 nm thick everywhere. How much photoresist is required on top of the oxide to completely mask the ion implantation?

Solution: The projected range needs to be 0.05 μm in order to place the peak of the distribution at the Si-SiO₂ interface. Using Fig. 5.3(a), we find that the R_p of 0.05 μm requires an energy of 15 keV. Since the peak of the implant is at the interface, the total

dose will be twice the dose needed in silicon. The peak concentration is

$$N_p = Q/\Delta R_p \sqrt{2\pi} = 2 \times 10^{13}/(2.3 \times 10^{-6} \sqrt{2\pi}) = 3.5 \times 10^{18}/\text{cm}^3$$

where the straggle was found using Fig. 5.3(b). To completely mask the implantation, the tail of the distribution must be less than the background concentration at the interface. The minimum oxide thickness is found using Eq. (5.9):

$$X_0 = 0.05 + 0.023 \sqrt{2 \ln(10 \times 3.5 \times 10^{18}/10^{16})} \mu\text{m} = 0.14 \mu\text{m}$$

Since the oxide is 0.05 μm thick, the photoresist must provide a thickness equivalent to 0.09 μm . The resist thickness must be 1.8 times the needed thickness of SiO_2 to provide an equivalent barrier layer, so the photoresist should be at least 0.16 μm thick. This thickness requirement is easily met with most photoresist layers.

5.4 JUNCTION DEPTH AND SHEET RESISTANCE

Ion implantation is often used to form shallow pn junctions for various device applications. The implanted profile approximates a Gaussian distribution, and the junction depth may be found by equating the implanted distribution to the background concentration, as explained in Chapter 4:

$$N_p \exp\left[\frac{(X_j - R_p)^2}{2\Delta R_p^2}\right] = N_B$$

$$x_j = R_p \pm \Delta R_p \sqrt{2 \ln(N_p/N_B)} \quad (5.11)$$

Both roots may be meaningful, as indicated in Fig. 5.6, in which a deep subsurface implant has junctions occurring at two different depths, x_{j1} and x_{j2} .

Example 5.3

Boron is implanted into an n -type silicon wafer to a depth of 0.3 μm . Find the location of the junction if the peak concentration is $1 \times 10^{18}/\text{cm}^3$ and the doping of the wafer is $3 \times 10^{16}/\text{cm}^3$.

Solution: From Fig. 5.3(a), the implant energy is 100 keV. From Fig. 5.3(b), the straggle is 0.07 μm . Equating the Gaussian distribution to the background concentration yields

$$3 \times 10^{16} = 10^{18} \exp\left[-\frac{(X_j - R_p)^2}{2\Delta R_p^2}\right]$$

or

$$x_j = R_p \pm 2.65\Delta R_p$$

This yields junction depths of 0.12 μm and 0.49 μm .

The peak of an implantation is often positioned at the silicon surface. For this special case, we may use Irvin's curves for Gaussian distributions to find the sheet resistance of the implanted layer, as discussed in Chapter 4. These curves may also be

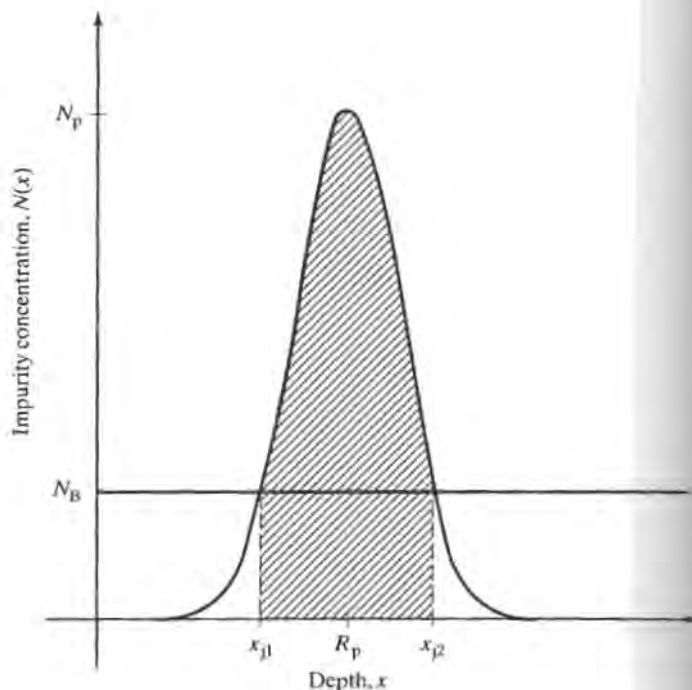


FIGURE 5.6

Junction formation by impurity implantation in silicon. Two pn junctions are formed at x_{j1} and x_{j2} .

used to find the sheet resistance of a layer which is completely below the surface. (See Problem 5.4.) Note that an implanted Gaussian impurity distribution will remain Gaussian through any subsequent high-temperature processing steps.

Diffused profiles generally have the maximum impurity concentration at the silicon surface. Ion-implantation techniques can be used to produce profiles with subsurface peaks or "retrograde" profiles that decrease toward the wafer surface. Multiple implant steps at different energies can also be used to build up more complicated impurity profiles.

5.5 CHANNELING, LATTICE DAMAGE, AND ANNEALING

5.5.1 Channeling

The LSS results of Section 5.2 are based on the assumption that the target material is amorphous, having a completely random order. This assumption is true of thermal SiO_2 , deposited Si_3N_4 and SiO_2 , and many thin metal films, but it is not valid for a crystalline substrate. The regular arrangement of atoms in the crystal lattice leaves a large amount of open space in the crystal. For example, Fig. 5.7 shows a view through the silicon lattice in the $\langle 110 \rangle$ direction. If the incoming ion flux is improperly oriented with respect to the crystal planes, the ions will tend to miss the silicon atoms in the lattice and will "channel" much more deeply into the material than the LSS theory predicts. However, electronic interactions will eventually stop the ions.

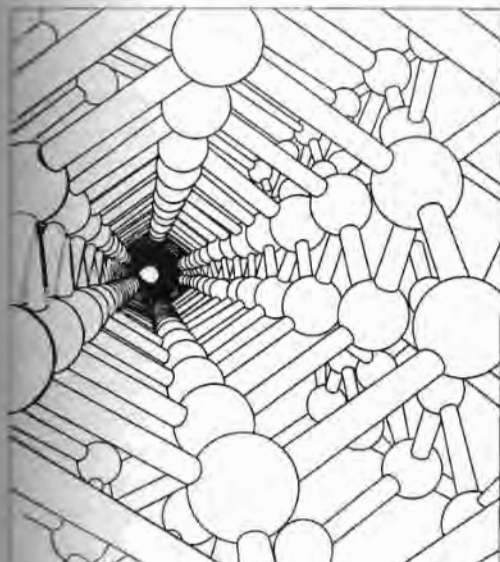


FIGURE 5.7

The silicon lattice viewed along the $\langle 110 \rangle$ axis. From *The Architecture of Molecules* by Linus Pauling and Roger Hayward. Copyright © 1964 W. H. Freeman and Company. Reprinted with permission from Refs. [4a] and [4b].

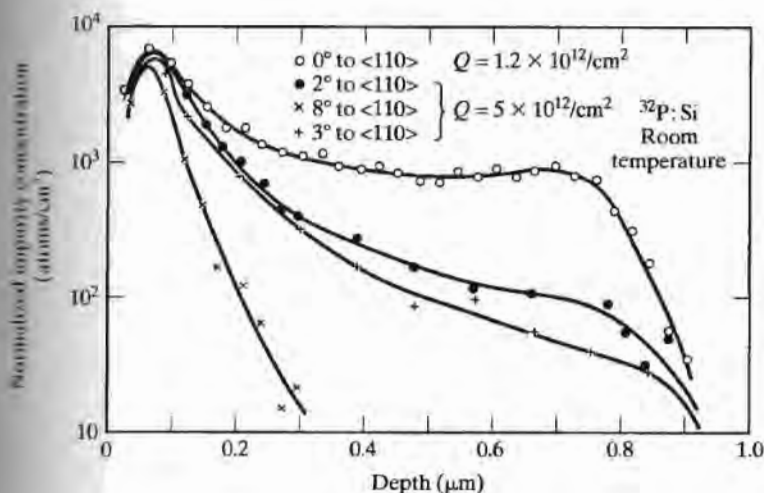


FIGURE 5.8

Phosphorus impurity profiles for 40-keV implantations at various angles from the $\langle 110 \rangle$ axis. Copyright 1968 by National Research Council of Canada. Reprinted with permission from Ref. [5].

The effects of channeling are demonstrated in Fig. 5.8. Phosphorus has been implanted at an energy of 40 keV into a silicon target with several orientations of the ion beam relative to the $\langle 100 \rangle$ silicon surface. The appearance of a random target can be achieved by tilting $\langle 100 \rangle$ silicon approximately 8° relative to the incoming beam. The results are represented by the x 's in Fig. 5.8. The range for this case compares well with the LSS calculations presented in Fig. 5.3. The open circles represent the boron profile implanted perpendicular to the $\langle 100 \rangle$ surface. Note that the range for the "channeled" case is almost twice that predicted by the LSS theory. Results for two other angles of incidence are given in Fig. 5.8, showing progressively less channeling as the angle is increased.

5.5.2 Lattice Damage and Annealing

During the implantation process, ion impact can knock atoms out of the silicon lattice, damaging the implanted region of the crystal. If the dose is high enough, the implanted layer will become amorphous. Figure 5.9 gives the dose required to produce an amorphous silicon layer for various impurities as a function of substrate temperature. The heavier the impurity, the lower the dose that is required to create an amorphous layer. At sufficiently high temperatures, an amorphous layer can no longer be formed. Note that damage from argon implantation was mentioned as a possible gettering technique at the end of Chapter 4.

Implantation damage can be removed by an "annealing" step. Following implantation, the wafer is heated to a temperature between 800 and 1000 °C for approximately 30 min. At these temperatures, silicon atoms can move back into lattice sites, and impurity atoms can enter substitutional sites in the lattice. After the annealing cycle, nearly all of the implanted dose becomes electrically active, except for impurity concentrations exceeding $10^{19}/\text{cm}^3$.

Unfortunately, annealing cycles of 30 min at temperatures approaching 1000 °C can cause considerable spreading of the implant by diffusion. It has been found that truly amorphous layers can actually be annealed at lower temperatures through the process of solid-phase epitaxy. The crystalline substrate seeds recrystallization of the amorphous layer, and epitaxial growth can proceed as rapidly as 500 Å/min at 600 °C. During solid-phase epitaxy, impurity atoms are incorporated into substitutional sites, and full activation is achieved at low temperatures.

Low-energy arsenic implantations produce shallow amorphous layers that can be annealed using solid-phase epitaxy to yield shallow, abrupt junctions that are ideal for

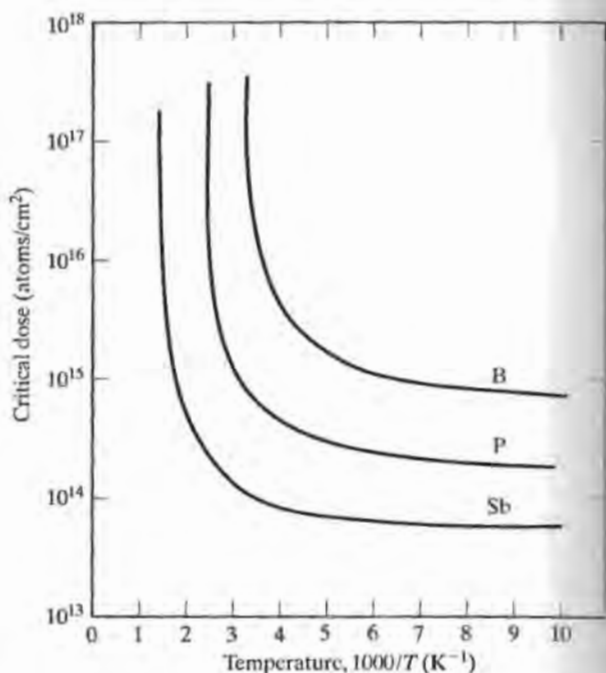


FIGURE 5.9

A plot of the dose required to form an amorphous layer on silicon versus reciprocal target temperature. Arsenic falls between phosphorus and antimony. Copyright 1970 by Plenum Publishing Corporation. Reprinted with permission from Ref. [6].

VLSI structures. Boron, however, is so light that it does not produce an amorphous layer even at relatively high doses, unless the substrate is deliberately cooled. (See Fig. 5.9.) Today, boron is often implanted using ions of the heavier BF_2 molecule. The lower-velocity implant results in shallow layers that can be annealed under solid-phase-epitaxy conditions.

5.5.3 Deviations from the Gaussian Theory

If we take a detailed look at the shape of the implanted impurity distribution, we find deviations from the ideal Gaussian profile. When light ions, such as boron, impact atoms of the silicon target, they experience a relatively large amount of backward scattering and fill in the distribution on the surface side of the peak, as in Fig. 5.10. Heavy atoms, such as arsenic, experience a larger amount of forward scattering and tend to fill in the profile on the substrate side of the peak. A number of methods have been proposed for mathematically modeling this behavior, such as the use of Pearson Type-IV distributions [8]. However, for common implant energies below 200 keV, the Gaussian theory provides a useful model of the impurity distribution. This is particularly true since the forward and backward scattering tend to alter the tails of the distribution where the concentration is well below the peak value.

5.6 SHALLOW IMPLANTATION

Deep submicron MOS devices require heavily doped source-drain regions with junction depths below 20 nm. In order to form these junctions, a new set of manufacturing tools was developed based upon low-energy implantation and rapid thermal annealing. A detailed understanding of a diffusion phenomenon termed Transient Enhanced Diffusion or TED was also required.

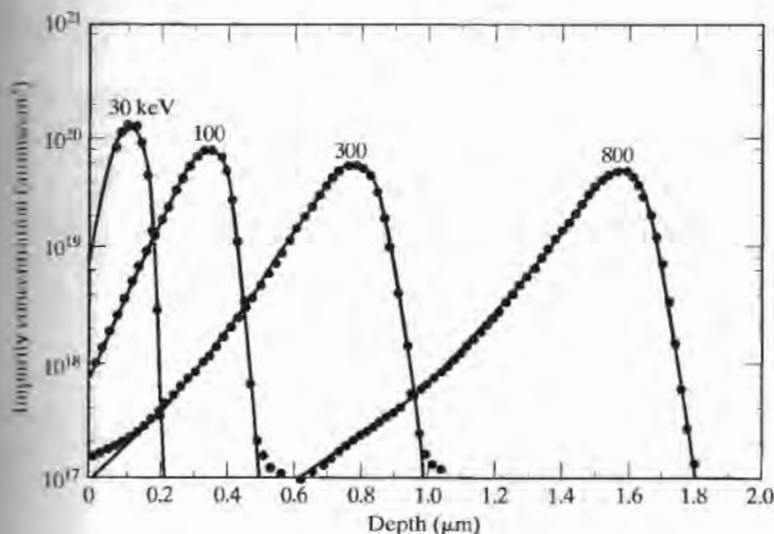


FIGURE 5.10

Measured boron impurity distributions compared with four-moment (Pearson IV) distribution functions. The boron was implanted into amorphous silicon without annealing. Reprinted with permission from Philips Journal of Research [8].

5.6.1 Low-Energy Implantation

To achieve shallow junctions, low-voltage ion implantation with energies in the range of 0.25–5 keV are utilized. Note that the classical ion-implanter described in Section 5.1 is not adequate for these implants, because of the high initial acceleration potential; specialized implanter systems have been developed specifically for low-energy ion implantation. In addition, new implantation species with high mass, such as decaborane $B_{10}H_{14}$ [12], and hence low resulting velocity, are being investigated. The “as-implanted” impurity profile distributions resulting from low-energy implantations can have peaks very near the surface with junction depths of less than 25 nm, as indicated by the SIMS data in Fig. 5.11. Because relatively high doses are used to achieve low sheet resistance for source–drain regions, even low-energy implants cause substantial damage to the crystal near the wafer surface, and this damage must be removed by annealing.

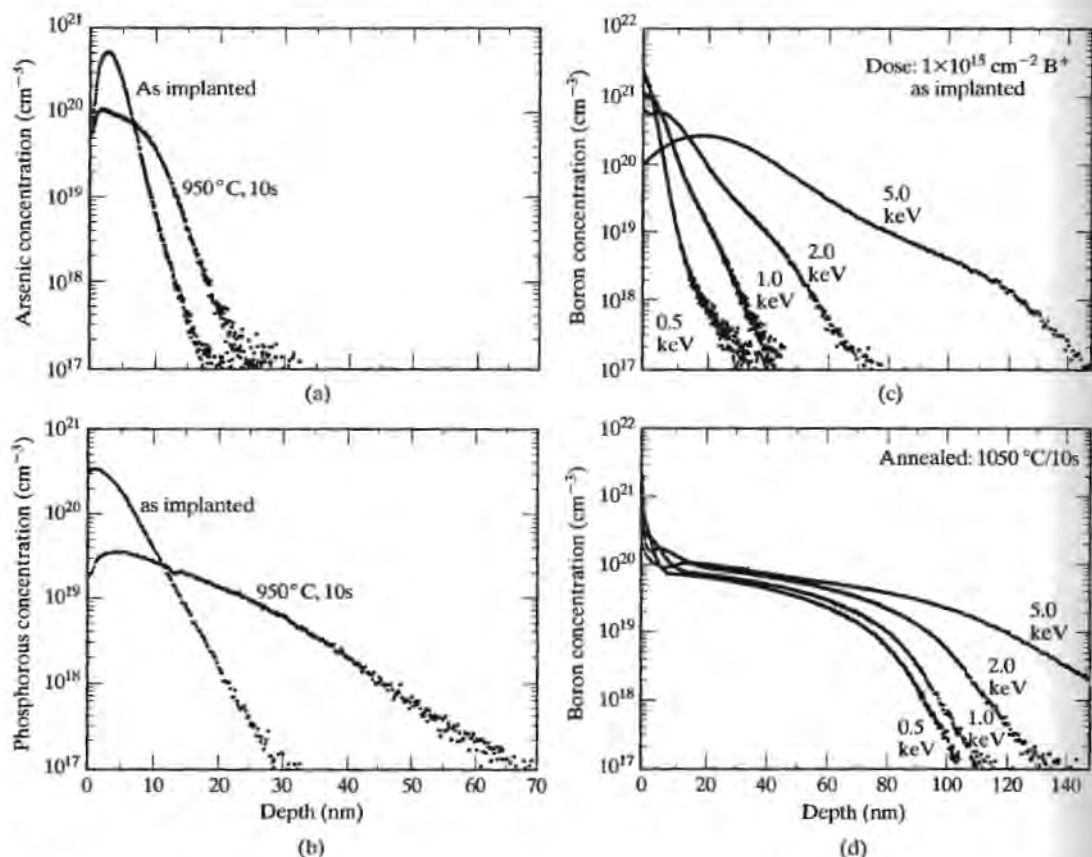


FIGURE 5.11

Examples of transient enhanced diffusion. SIMS data comparing as-implanted and annealed depth profiles from (a) $3 \times 10^{14}/cm^2$, 2 keV As^+ , and (b) $3 \times 10^{14}/cm^2$, 1 keV P^+ . Annealing conditions were 950 °C for 10 sec. SIMS depth profiles of $1 \times 10^{15}/cm^2$ B implanted at 0.5-, 1-, 2-, and 5-keV (c) as-implanted, and (d) after annealing at 1050 °C for 10 sec. Copyright 1997 IEEE. Reprinted with permission from Ref. [13].

5.6.2 Rapid Thermal Annealing

In addition to removing the damage caused by the implantation, the annealing step is required to electrically activate the implanted impurities. However, in order to minimize diffusion of the shallow implanted profiles, the Dt product associated with the annealing process must be kept as small as possible. Rapid Thermal Annealing (RTA) systems can achieve the desired results with annealing times that range from a few minutes down to only a few seconds. High-intensity lamps shown in Fig. 5.12(a) are used to rapidly heat the wafer to the desired annealing temperature (e.g., 950–1050 °C) in a very short time. Extremely rapid temperature ramp-up and ramp-down rates are achieved (e.g., 50 °C/sec or more). Using this rapid thermal-processing technique, the effective Dt products can be very small. However, even the short ramp time can be important if the dwell time at the upper temperature is short. (See Problems 5.19–5.22.)

Rapid thermal annealing represents only one form of rapid thermal processing. Similar systems are used to grow very thin oxide and nitride layers using processes termed rapid thermal oxidation (RTO) or rapid thermal nitridation (RTN). Silicide layers, to be discussed in Section 7.5, can also be formed using rapid thermal processing. An example of an RTP system appears in Fig. 5.12(b).

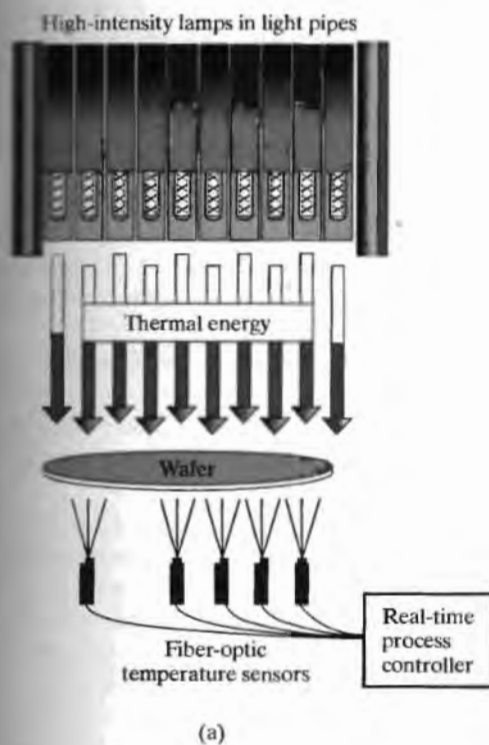


FIGURE 5.12

(a) Concept of a rapid thermal processing (RTP) system. (b) Applied Materials 300 mm RTP System. (Courtesy Applied Materials, Inc.)

5.6.3 Transient Enhanced Diffusion (TED)

During the investigation of shallow junction formation, it was discovered that the impurities diffuse considerably more than predicted by simple diffusion theory using the values of the Dt product calculated for the annealing process. It was found that the presence of damage to the silicon crystal from the implantation temporarily enhances the diffusion coefficient by a factor as large as 5 to 10 times [13–16]. This enhancement is a transient phenomenon which disappears as the damage is annealed out. However, the modification of the impurity profile can be substantial, particularly to the tails of the distribution, and the junction depth can change by a significant amount. Since the dose is constant, the peak concentration also falls, and the sheet resistance changes due to the impurity redistribution. An example of the profile redistribution due to TED is also presented in the SIMS data in Fig. 5.11 [13]. Because of its importance in small geometry devices, modeling of the TED phenomena has been added to many process simulation programs.

SUMMARY

Ion implantation uses a high-voltage accelerator to introduce impurity atoms into the surface of the silicon wafer, and it offers many advantages over deposition by high-temperature diffusion. Ion implantation is a low-temperature process minimizing impurity movement by diffusion, which has become very important to VLSI fabrication. Low-temperature processing also permits the use of a wide variety of materials as barrier layers to mask the implantation. Photoresist, oxide, nitride, aluminum, and other metal films can all be used, adding important increased flexibility to process design.

Ion implantation also permits the use of a much wider range of impurity species than diffusion. In principle, any element that can be ionized can be introduced into the wafer using implantation. Implantation offers much tighter control of the dose introduced into the wafer, and a much wider range of doses can be reproducibly achieved than possible with diffusion.

Diffused profiles almost always have the maximum impurity concentration at the surface. Ion-implantation techniques can be used to produce new profiles with subsurface peaks or retrograde profiles which decrease toward the wafer surface. Implantation can introduce impurities into very shallow layers near the surface, again a significant advantage for VLSI structures.

The main disadvantage of ion implantation is the cost of the equipment. Also, the ion implanter has trouble achieving high doses ($>10^{16}/\text{cm}^2$) in a time reasonable for high-volume production. High-current machines have been developed to overcome this latter problem. Overall, the flexibility and process control achievable with ion implantation have far outweighed the disadvantage of cost, and ion implantation is used routinely for state-of-the-art bipolar and MOS integrated-circuit fabrication.

Ion implantation results in profiles that can be modeled by a Gaussian distribution. The depth and width of the distribution depend on both the ion species and the energy of the implantation. To prevent channeling, implantation is normally performed at an angle of approximately 8° off the normal to the wafer surface.

The implantation process damages the surface, and an annealing step is required to remove the effects of the damage. Low doses may result in the need for annealing

800 to 1000 °C for 30 min. However, if the surface layer has become amorphous, annealing can be achieved through solid-phase epitaxy at temperatures of only 600 °C.

Very shallow implantations, that are required for deep submicron VLSI fabrication, can be achieved using low energy (< 5 keV) ion implantation. Rapid thermal annealing is then used to activate the implantation while maintaining a small Dt product. However, transient enhanced diffusion causes the implantations to spread more deeply than predicted by the Dt value and standard diffusion models.

REFERENCES

- [1] J. Lindhard, M. Scharff, and H. Schiott, "Range Concepts in Heavy Ion Ranges," *Mat.-Fys. Med. Dan. Vid. Selsk.*, 33, No. 14, 1963.
- [2] J. F. Gibbons, W. S. Johnson, and S. W. Mylroie, *Projected Range in Semiconductors*, 2nd ed., Dowden, Hutchinson, and Ross, New York, 1975.
- [3] S. Furukawa, H. Matsumura, and H. Ishiwara, "Lateral Distribution Theory of Implanted Ions," in S. Namba, Ed., *Ion Implantation in Semiconductors*, Japanese Society for the Promotion of Science, Kyoto, pp. 73, 1972.
- [4] (a) L. Pauling and R. Hayward, *The Architecture of Molecules*, W. H. Freeman, San Francisco, 1964. (b) S. M. Sze, Ed., *Semiconductor Devices Physics and Technology*, McGraw-Hill, New York, 1985.
- [5] G. Dearnaley, J. H. Freeman, G. A. Card, and M. A. Wilkins, "Implantation Profiles of ^{32}P Channeled into Silicon Crystals," *Canadian Journal of Physics*, 46, 587-595 (March 15, 1968).
- [6] F. F. Morehead and B. L. Crowder, "A Model for the Formation of Amorphous Si by Ion Implantation," pp. 25-30, in Eisen and Chadderton (see Source Listing 4).
- [7] B. L. Crowder and F. F. Morehead, Jr., "Annealing Characteristics of n-type Dopants in Ion-Implanted Silicon," *Applied Physics Letters*, 14, 313-315 (May 15, 1969).
- [8] W. K. Hofker, "Implantation of Boron in Silicon," *Philips Research Reports Supplements*, No. 8, 1975.
- [9] J. F. Gibbons, "Ion-Implantation in Semiconductors—Part I: Range Distribution Theory and Experiment," *Proceedings of the IEEE*, 56, 295-319, March 1968.
- [10] J. F. Gibbons, "Ion Implantation in Semiconductors—Part II: Damage Production and Annealing," *Proceedings of the IEEE*, 60, 1062-1096, September 1972.
- [11] T. Hirao, G. Fuse, K. Inoue, S. Takayanagi, Y. Yaegashi, S. Ichikawa, and T. Izumi, "Electrical Properties of Si Implanted with As Through SiO_2 Films," *Journal of Applied Physics*, 51, 262-268 (January 1980).
- [12] K. Goto, J. Matsuo, Y. Tada, T. Tanaka, Y. Momiyama, T. Sugii, and I. Yamada, "A High-Performance 50 nm PMOSFET Using Decaborane ($\text{B}_{10}\text{H}_{14}$) Ion Implantation and 2-step Activation Annealing Process," *IEEE IEDM Digest*, pp. 471-474, December 1997.
- [13] A. Agarwal, D. J. Eaglesham, H.-J. Gossman, L. Pelaz, S. B. Herner, D. C. Jacobson, T. E. Haynes, Y. Erokhin, and R. Simonton, "Boron-Enhanced-Diffusion of Boron: The Limiting Factor for Ultra Shallow Junctions," *IEEE IEDM Digest*, pp. 467-470, December 1997.
- [14] A. D. Lilak, S. K. Earles, K. S. Jones, M. E. Law, and M. D. Giles, "A Physics-Based Modeling Approach for the Simulation of Anomalous Boron Diffusion and Clustering Behavior," *IEEE IEDM Digest*, pp. 493-496, December 1997.
- [15] K. Suzuki, T. Miyashita, and Y. Tada, "Damage Calibration Concept and Novel B Cluster Reaction Model for B Transient Enhanced Diffusion Over Thermal Process Range from 600 °C (839 h) to 1100 °C (5 s) with Various Ion Implantation Doses and Energies," *IEEE IEDM Digest*, pp. 501-504, December 1997.

[16] S. S. Yu, H. W. Kennel, M. D. Giles, and P. A. Packan, "Simulation of Transient Enhanced Diffusion Using Computationally Efficient Models," *IEEE IEDM Digest*, pp. 509-512, December 1997.

SOURCE LISTING

- [1] J. W. Mayer, L. Eriksson, and J. A. Davies, *Ion-Implantation in Semiconductors*, Academic Press, New York, 1970.
- [2] G. Dearnaley, J. H. Freeman, R. S. Nelson, and J. Stephen, *Ion-Implantation*, North-Holland, New York, 1973.
- [3] G. Carter and W. A. Grant, *Ion-Implantation of Semiconductors*, John Wiley & Sons, New York, 1976.
- [4] F. Eisen and L. Chadderton, Eds., *Ion Implantation in Semiconductors*, First International Conference (Thousand Oaks, CA), Gordon and Breach, New York, 1970.
- [5] I. Ruge and J. Graul, Eds., *Ion Implantation in Semiconductors*, Second International Conference (Garmisch-Partenkirchen, Germany), Springer-Verlag, Berlin, 1972.
- [6] B. L. Crowder, Ed., *Ion Implantation in Semiconductors*, Third International Conference (Yorktown Heights, NY), Plenum, New York, 1973.
- [7] S. Namba, Ed., *Ion Implantation in Semiconductors*, Fourth International Conference (Osaka, Japan), Plenum, New York, 1975.
- [8] F. Chernow, J. Borders, and D. Bruce, Eds., *Ion Implantation in Semiconductors*, Fifth International Conference (Boulder, CO), Plenum, New York, 1976.

PROBLEMS

- 5.1 Boron is implanted with an energy of 60 keV through a 0.25- μm layer of silicon dioxide. The implanted dose is $1 \times 10^{14}/\text{cm}^2$.
 - (a) Find the boron concentration at the silicon-silicon dioxide interface.
 - (b) Find the dose in silicon.
 - (c) Determine the junction depth if the background concentration is $3 \times 10^{15}/\text{cm}^3$.
- 5.2 A measured boron dose of $2 \times 10^{15}/\text{cm}^2$ is implanted into the surface of a silicon wafer at an energy of 10 keV. What are the projected depth and straggle based upon Fig. 5.3? What is the junction depth if the implantation resulted in a Gaussian profile and the background concentration of the wafer is $10^{16}/\text{cm}^3$?
- 5.3 What energy is required to implant phosphorus through a 1- μm layer of silicon dioxide if the peak of the implant is to be at the Si-SiO₂ interface? What is the straggle?
- 5.4 An arsenic dose of $1 \times 10^{12}/\text{cm}^2$ is implanted through a 50-nm layer of silicon dioxide with the peak of the distribution at the Si-SiO₂ interface. A silicon nitride film on top of the silicon dioxide is to be used as a barrier material in the regions where arsenic is not desired. How thick should the nitride layer be if the background concentration is $1 \times 10^{15}/\text{cm}^3$?
- 5.5 An implantation will be used for the predeposition step for a boron base diffusion. The final layer is to be 5 μm deep with a sheet resistance of 125 ohms per square ($N_B = 10^{16}/\text{cm}^3$).

- (a) What is the dose required from the ion implanter if the boron will be implanted through a thin silicon dioxide layer so that the peak of the implanted distribution is at the silicon-silicon dioxide interface?
- (b) What drive-in time is required to produce the final base layer at a temperature of 1100°C ?
- 5.6 Repeat Problem 5.5 for a $2\text{-}\mu\text{m}$ -deep diffusion with a sheet resistance of 200 ohms/square.
- 5.7 Repeat Problem 5.5 for a $0.25\text{-}\mu\text{m}$ -deep diffusion with a sheet resistance of 250 ohms/square.
- 5.8 A phosphorus dose of $1 \times 10^{15}/\text{cm}^2$ is implanted into the surface of a silicon wafer at an energy of 20 keV. What are the projected depth and straggle based upon Fig. 5.3? What is the junction depth if the implantation resulted in a Gaussian profile and the background concentration of the wafer is $10^{16}/\text{cm}^3$?
- 5.9 (a) Use Irvin's curves to find the sheet resistance of a boron layer implanted completely below the surface in n -type silicon ($N_B = 10^{15}/\text{cm}^3$). Assume the layer has a peak concentration of $1 \times 10^{19}/\text{cm}^3$, and the range and straggle are $1.0\text{ }\mu\text{m}$ and $0.11\text{ }\mu\text{m}$, respectively. (Hint: Think about conductors in parallel.)
- (b) What is the dose of this implantation?
- (c) What was the energy used for this ion implantation?
- (d) At what depths are pn junctions located?
- 5.10 (a) Photoresist is used as an implantation barrier during a boron source/drain implantation at 50 keV. How thick a layer of photoresist should be utilized to block the implant?
- (b) Repeat the process for a phosphorus implant.
- (c) Repeat it for an arsenic implant.
- 5.11 What are the velocities of the following ions in an ion implanter if they are accelerated through a potential of 5 keV (a) B^+ ions? (b) $(\text{BF}_2)^{++}$ ions? (c) $(\text{B}_{10}\text{H}_{14})^+$ ions?
- 5.12 The source and drain regions of a self-aligned n -channel polysilicon-gate MOS transistor are to be formed using arsenic implantation. The dimensions of a cross section of the device are given in Fig. P5.12. Calculate the channel shrinkage caused by lateral straggle if the peak concentration of the implantation is $10^{20}/\text{cm}^3$ and the substrate doping is $10^{16}/\text{cm}^3$. Assume that the channel region is in the silicon immediately below the oxide. Use $R_p = 0.1\text{ }\mu\text{m}$, $\Delta R_p = 0.04\text{ }\mu\text{m}$, and $\Delta R_L = 0.022\text{ }\mu\text{m}$.

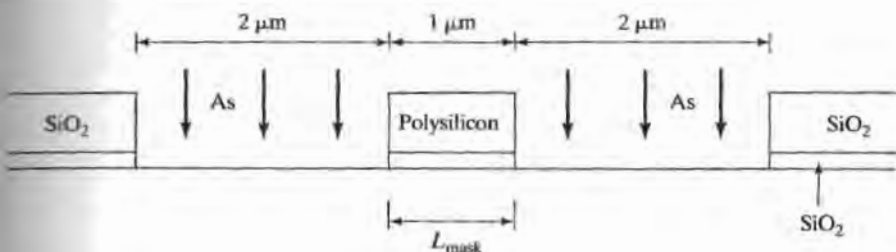


FIGURE 5.12

- 5.13** An implanted profile is formed by two boron implantations. The first uses an energy of 100 keV and the second an energy of 200 keV. The peak concentration of each distribution is $5 \times 10^{18}/\text{cm}^3$. Draw a graph of the composite profile and find the junction depth(s) if the phosphorus background concentration is $10^{16}/\text{cm}^3$. What are the doses of the two implant steps?
- 5.14** A high-energy (5 MeV) is used to implant oxygen deep below the silicon surface in order to form a buried SiO_2 layer. Assume that the desired SiO_2 layer is to be 0.2 μm wide.
- What is the oxygen dose required to be implanted in silicon?
 - What beam current is required if a 200-mm-diameter wafer is to be implanted in 15 min?
 - How much power is being supplied to the ion beam? Discuss what effects this implantation may have on the wafer.
- 5.15** The threshold voltage of an NMOS transistor may be increased by ion implantation of boron into the channel region. For shallow implantations, the voltage shift is given approximately by $\Delta V_T = qQ/C_{ox}$, where Q is the boron dose and $C_{ox} = \epsilon_0/X_0$. X_0 is the oxide thickness and ϵ_0 is the permittivity of silicon dioxide: $3.9 \times (8.854 \times 10^{-14} \text{ F/cm})$. What boron dose is required to shift the threshold by 0.75 V if the oxide thickness is 40 nm?
- 5.16** An ion implanter has a beam current of 10 μA . How long does it take to implant a boron dose of $10^{15}/\text{cm}^2$ into a wafer with a diameter of 200 mm?
- 5.17** Write a computer program to calculate the sheet resistance of an arbitrary Gaussian layer in silicon.
- 5.18** A boron dose of $1 \times 10^{15}/\text{cm}^2$ is implanted into a silicon wafer. (a) Use Fig. 5.9 to determine the minimum substrate temperature required to insure that an amorphous layer is formed so that solid-state epitaxy is possible. (b) What happens if the implanted species is changed to phosphorus?
- 5.19** An RTA system goes from 25 to 1050 $^\circ\text{C}$ at a rate of 50 $^\circ\text{C}/\text{sec}$, and remains at 1050 $^\circ\text{C}$ for 1 minute. It then returns to 25 $^\circ\text{C}$ at a rate of 50 $^\circ\text{C}/\text{sec}$. Numerically calculate the total Dt product for a boron diffusion using the data from Table 4.1 in Chapter 4. What is the Dt product considering only the time spent at 1050 $^\circ\text{C}$?
- 5.20** Repeat Prob. 5.19 if the time at 1050 $^\circ\text{C}$ is reduced to 5 seconds.
- 5.21** An RTA system goes from 25 to 980 $^\circ\text{C}$ at a rate of 40 $^\circ\text{C}/\text{sec}$ and stays at 980 $^\circ\text{C}$ for 2 minutes. It then returns to 25 $^\circ\text{C}$ at a rate of 40 $^\circ\text{C}/\text{sec}$. Numerically calculate the total Dt product for a phosphorus diffusion using the data from Table 4.1 in Chapter 4. What is the Dt product considering only the time spent at 980 $^\circ\text{C}$?
- 5.22** Repeat Prob. 5.21 reducing the time at 980 $^\circ\text{C}$ to 15 seconds.

CHAPTER 6

Film Deposition

Fabrication processes involve many steps in which thin films of various materials are deposited on the surface of the wafer. This chapter presents a survey of deposition processes, including evaporation, chemical vapor deposition, and sputtering, which are used to deposit metals, silicon and polysilicon, and dielectrics such as silicon dioxide and silicon nitride. Evaporation and sputtering require vacuum systems operating at low pressure, whereas chemical vapor deposition and epitaxy can be performed at either reduced or atmospheric pressure. An overview of vacuum systems and some results from the theory of ideal gases are also presented in this chapter.

EVAPORATION

Physical evaporation is one of the oldest methods of depositing metal films. Aluminum and gold are heated to the point of vaporization, and then evaporate to form a thin film covering the surface of the silicon wafer. To control the composition of the deposited material, evaporation is performed under vacuum conditions.

Figure 6.1 shows a basic vacuum deposition system consisting of a vacuum chamber, a mechanical roughing pump, a diffusion pump or turbomolecular pump, valves, vacuum gauges, and other instrumentation. In operation, the roughing pump is opened first, and the mechanical pump lowers the vacuum chamber pressure to an intermediate vacuum level of approximately 1 Pascal (Pa^1). If a higher vacuum level is needed, the roughing valve is closed, and the foreline and high-vacuum valves are opened. The roughing pump now maintains a vacuum on the output of the diffusion pump. A liquid-nitrogen (77 K) cold trap is used with the diffusion pump to reduce the pressure in the vacuum chamber to approximately 10^{-4} Pa. Ion and thermocouple gauges are used to monitor the pressure at a number of points in the vacuum system, and several other valves are used as vents to return the system to atmospheric pressure.

¹ 1 atm = 760 mm Hg = 760 torr = 1.013×10^5 Pa, 1 Pa = 1 N/m^2 = 0.0075 torr.

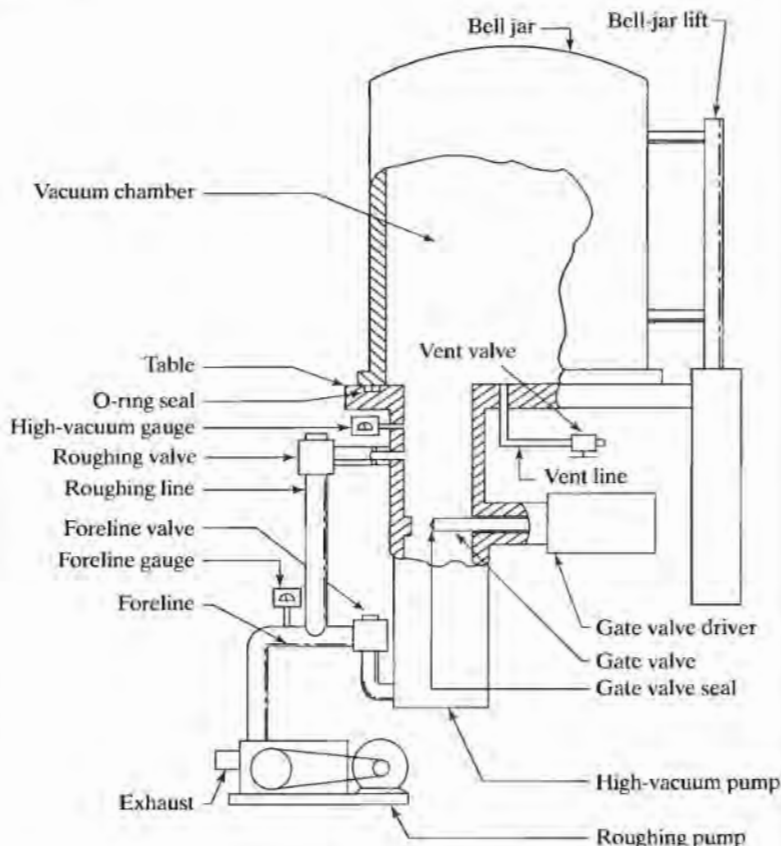


FIGURE 6.1

Typical vacuum system used for evaporation including vacuum chamber, roughing pump, high-vacuum pump, and various valves and vacuum gauges. Copyright 1987 McGraw-Hill Book Company. Reprinted with permission from Ref. [5].

6.1.1 Kinetic Gas Theory

Gases behave in an almost ideal manner at low pressure and are well described by the ideal gas law. Pressure P , volume V , and temperature T of one mole of a gas are related by

$$PV = N_{av}kT \quad (6.1)$$

where k is Boltzmann's constant² and N_{av} is Avogadro's number (6.02×10^{23} molecules/mole). The concentration of gas molecules is given by

$$n = \frac{N_{av}}{V} = \frac{P}{kT} \quad (6.2)$$

In some systems, the surface of the substrate must be kept extremely clean prior to deposition. The presence of even a small amount of oxygen or other elements will result in formation of a contamination layer on the surface of the substrate. The rate Φ

² k (Boltzmann's constant) = 1.38×10^{-23} J/K = 1.37×10^{-22} atm-cm³/K.

of formation of this layer is determined from the impingement rate of gas molecules hitting the substrate surface and is related to the pressure by

$$\Phi = \frac{P}{\sqrt{2\pi m k T}} \text{ (molecules/cm}^2\text{-sec)} \quad (6.3)$$

where m is the mass of the molecule. This can be reduced to

$$\Phi = \frac{2.63 \times 10^{20} P}{\sqrt{MT}} \text{ (molecules/cm}^2\text{-sec)} \quad (6.4)$$

where P is the pressure in Pa and M is the molecular weight (e.g., $M = 32$ for oxygen molecules). If we assume that each molecule sticks as it contacts the surface, then the time required to form a monolayer on the surface is given by

$$t = \frac{N_s}{\Phi} = \frac{N_s \sqrt{2\pi m k T}}{P} \quad (6.5)$$

where N_s is the number of molecules/cm² in the layer.

Example 6.2

Suppose the residual pressure of oxygen in the vacuum system is 1 Pa. How long does it take to deposit one atomic layer of oxygen on the surface of the wafer at 300 K?

Solution: The radius of an oxygen molecule is approximately 3.6 Å. If we assume close packing of the molecules on the surface, there will be approximately 2.2×10^{14} molecules/cm². At 300 K and 1 Pa, the impingement rate for oxygen is 2.7×10^{18} molecules/cm²-sec. One monolayer is deposited in 82 μsec (a very short period).

Pressure and temperature also determine another important film-deposition parameter called the *mean free path*, λ . The mean free path of a gas molecule is the average distance the molecule travels before it collides with another molecule. λ is given by

$$\lambda = \frac{kT}{\sqrt{2}\pi p d^2} \quad (6.6)$$

where d is the diameter of the gas molecule and is in the range of 2 to 5 Å. Evaporation is usually done at a background pressure near 10^{-4} Pa. At this pressure, a 4-Å molecule has a mean free path of approximately 60 m. Thus, during aluminum evaporation, for example, aluminum molecules do not interact with the background gases and tend to travel in a straight line from the evaporation source to the deposition target.

On the other hand, sputtering, which will be discussed in Section 6.4, uses argon gas at a pressure of approximately 100 Pa, using the same radius results in a mean free path of only 60 μm. Thus, the material being deposited tends to scatter often with the argon atoms and arrives at the target from random directions.

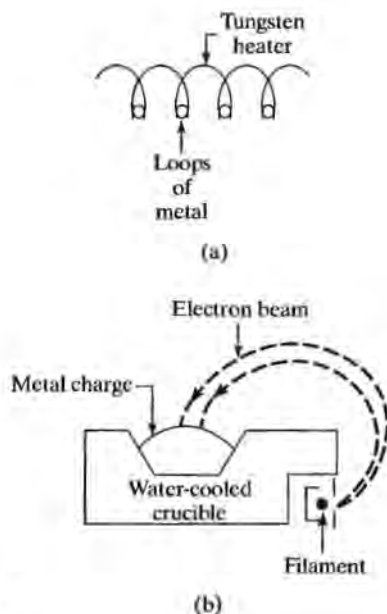


FIGURE 6.2

Two forms of evaporation sources. (a) Filament evaporation, in which loops of wire hang from a heated filament; (b) electron-beam source in which a beam of electrons is focused on a metal charge. The beam is bent in a magnetic field.

6.1.2 Filament Evaporation

The simplest evaporator consists of a vacuum system containing a filament that can be heated to high temperature. In Fig. 6.2(a), small loops of a metal such as aluminum are hung from a filament formed of a refractory (high-temperature) metal such as tungsten. Evaporation is accomplished by gradually increasing the temperature of the filament until the aluminum melts and wets the filament. Filament temperature is then raised to evaporate the aluminum from the filament. The wafers are mounted near the filament and are covered by a thin film of the evaporating material.

Although filament evaporation systems are easy to set up, contamination levels can be high, particularly from the filament material. In addition, evaporation of composite materials cannot be easily controlled using a filament evaporator. The material with the lowest melting point tends to evaporate first, and the deposited film will not have the same composition as the source material. Thick films are difficult to achieve, since a limited supply of material is contained in the metal loops.

6.1.3 Electron-Beam Evaporation

In electron-beam (E-beam) evaporation systems (see Fig. 6.2(b)), the high-temperature filament is replaced with an electron beam. A high-intensity beam of electrons, with an energy up to 15 keV, is focused on a source target containing the material to be evaporated. The energy from the electron beam melts a region of the target. Material evaporates from the source and covers the silicon wafers with a thin layer.

The growth rate using a small planar source is given by

$$G = \frac{m}{\pi pr^2} \cos \phi \cos \theta (\text{cm/sec}) \quad (6.7)$$

for the geometrical setup in Fig. 6.3. ϕ is the angle measured from the normal to the

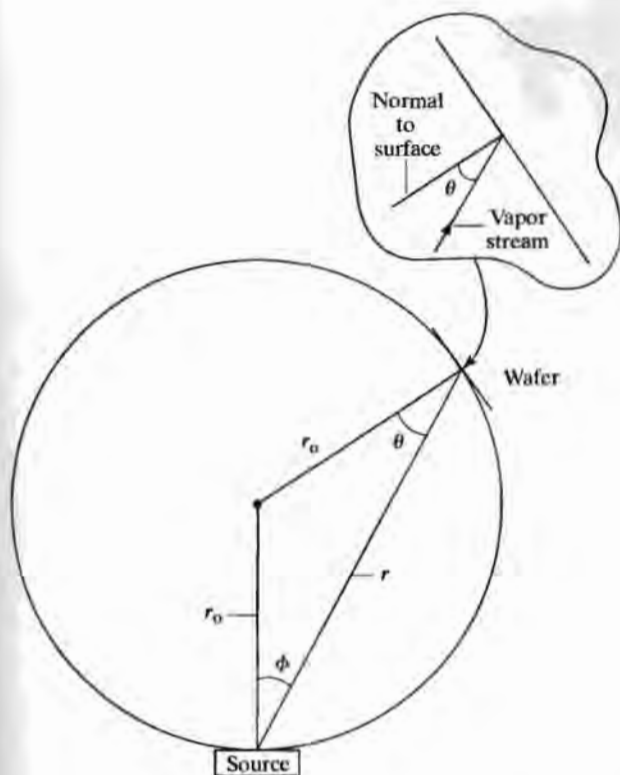


FIGURE 6.3

Geometry for evaporation in a system using a planetary substrate holder.

plane of the source, and θ is the angle of the substrate relative to the vapor stream. ρ and m are the density (g/cm^3) and mass evaporation rate (g/sec), respectively, of the material being deposited.

For batch deposition, a planetary substrate holder (Fig. 6.4) consisting of rotating sections of a sphere is used. Each substrate is positioned tangential to the surface of the sphere with radius r_0 , as in Fig. 6.3. Applying some geometry yields

$$\cos \theta = \cos \phi = \frac{r}{2r_0} \quad (6.8)$$

For the planetary substrate holder, G becomes independent of substrate position:

$$G = \frac{m}{4\pi\rho r_0^2} \quad (6.9)$$

The wafers are mounted above the source and are typically rotated around the source during deposition to ensure uniform coverage. The wafers are also often radiantly heated to improve adhesion and uniformity of the evaporated material. The source material sits in a water-cooled crucible, and its surface only comes in contact with the electron beam during the evaporation process. Purity is controlled by the purity of the original source material. The relatively large size of the source provides a virtually unlimited supply of material for evaporation, and the deposition rate is easily controlled by changing the current and energy of the electron beam.



FIGURE 6.4

Photograph of a laboratory E-beam evaporation system with a planetary substrate holder which rotates simultaneously around two axes.

One method of monitoring the deposition rate uses a quartz crystal, which is covered by the evaporating material during deposition. The resonant frequency of the crystal shifts in proportion to the thickness of the deposited film. By monitoring the resonant frequency of the crystal, the deposition rate may be measured with an accuracy of better than 1 Å/sec. Dual electron beams with dual targets may be used to coevaporate composite materials in E-beam evaporation systems.

X-ray radiation can be generated in an electron-beam system for acceleration voltages exceeding 5 to 10 keV, and substrates may suffer some radiation damage from both energetic electrons and X-rays. The damage can usually be annealed out during subsequent process steps. However, the radiation effects are of great concern to MOS process designers, and sputtering has replaced electron-beam evaporation in many steps in manufacturing processes.

6.1.4 Flash Evaporation

Flash evaporation uses a fine wire as the source material, and a high-temperature ceramic bar is used to evaporate the wire. The wire is fed continuously and evaporates on contact with the ceramic bar. Flash evaporation can produce relatively thick films, as in an E-beam system, without problems associated with radiation damage.

6.1.5 Shadowing and Step Coverage

Because of the large mean free paths of gas molecules at low pressure, evaporation techniques tend to be directional in nature, and shadowing of patterns and poor step coverage can occur during deposition. Figure 6.5 illustrates the shadowing phenomenon that can occur with closely spaced features on the surface of an integrated circuit. In the fully shadowed region, there will be little deposition. In the partially shadowed region, there will be variation in film thickness. To minimize these effects, the planetary substrate holder of the electron-beam system continuously rotates the wafers during the film deposition.

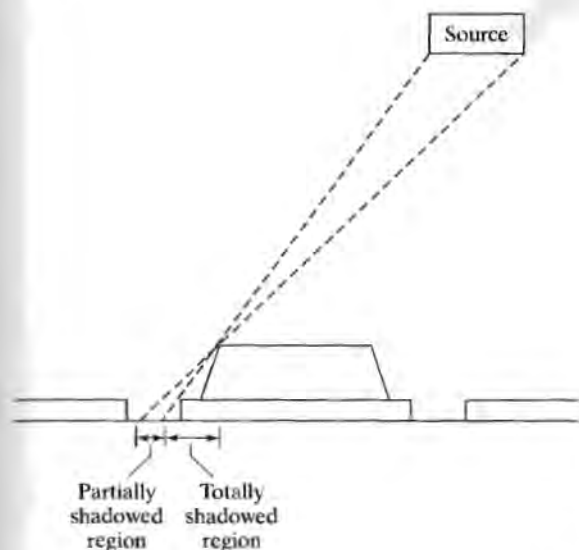


FIGURE 6.5

An example of the shadowing problem that can occur during low-pressure vacuum deposition in which the molecular mean free path is large.

6.2 SPUTTERING

Sputtering is achieved by bombarding a target with energetic ions, typically Ar^+ . Atoms at the surface of the target are knocked loose and transported to the substrate, where deposition occurs. Electrically conductive materials such as Al, W, and Ti can use a dc power source, in which the target acts as the cathode in a diode system. Sputtering of dielectrics such as silicon dioxide or aluminum oxide requires an RF power source to supply energy to the argon atoms. A diagram of a sputtering system is shown in Fig. 6.6.

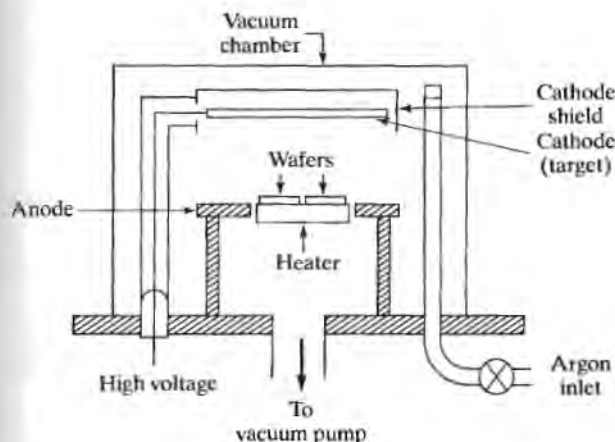


FIGURE 6.6

A dc sputtering system in which the target material acts as the cathode of a diode and the wafers are mounted on the system anode.

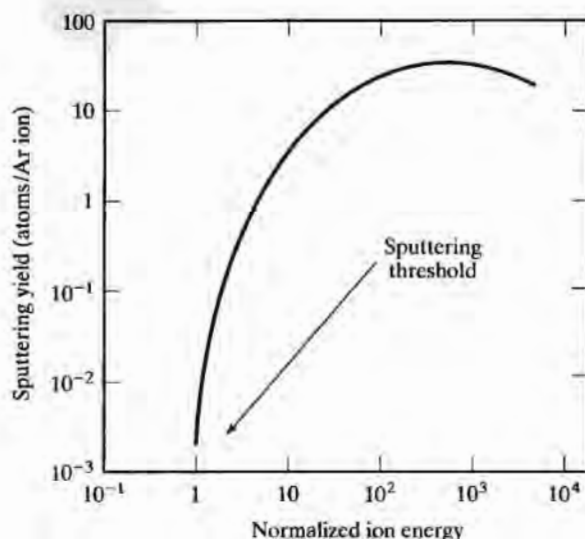


FIGURE 6.7

Sputtering yield versus ion energy for a dc sputtering system using argon.

In sputter deposition, there is a threshold energy that must be exceeded before sputtering occurs. The sputtering yield (Fig. 6.7) represents the number of atoms liberated from the target by each incident atom, and it increases rapidly with energy of the incident ions. Systems are usually operated with an energy large enough to ensure a sputtering yield of at least unity.

Sputtering can be used to deposit a broad range of materials. In addition, alloys may be deposited in which the deposited film has the same composition as the target. An example is the Al-Cu-Si alloy commonly used for metallization in integrated circuits. (We will discuss this alloy in Chapter 7.) As one might expect, sputtering results in the incorporation of some argon into the film, and heating of the substrate up to 350 °C can occur during the deposition process. Sputtering provides excellent coverage of the sharp topologies often encountered in integrated circuits.

Sputter etching (a reversal of the sputter deposition process) can be used to clean the substrate prior to film deposition, and the sputter etching process is often used to clean contact windows prior to metal deposition. Etching removes any residual oxide from the window and improves the contact between the metal and the underlying material.

6.3 CHEMICAL VAPOR DEPOSITION

Chemical vapor deposition (CVD) forms thin films on the surface of a substrate by thermal decomposition or reaction of gaseous compounds. The desired material is deposited directly from the gas phase onto the surface of the substrate. Polysilicon, silicon dioxide, and silicon nitride are routinely deposited using CVD techniques. In addition, refractory metals such as tungsten (W) can also be deposited using CVD.

Chemical vapor deposition can be performed at pressures for which the mean free path for gas molecules is quite small, and the use of relatively high temperatures can result in excellent conformal step coverage over a broad range of topological profiles.

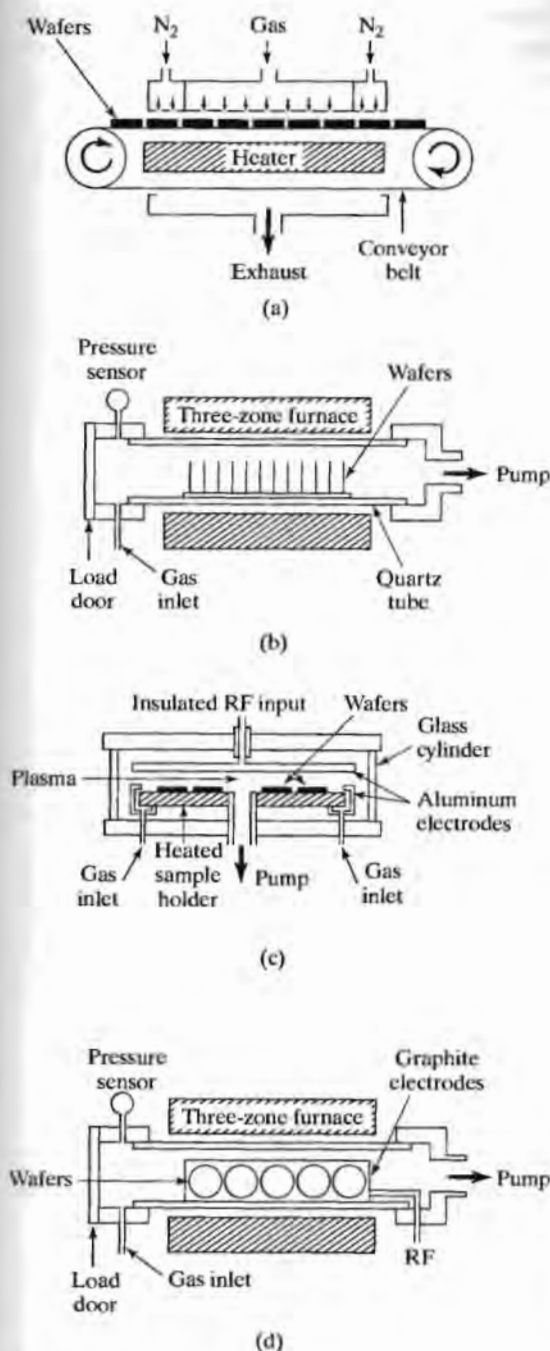


FIGURE 6.8

Four types of chemical vapor deposition (CVD) systems. (a) Atmospheric-pressure reactor; (b) hot-wall LPCVD system using a three-zone furnace tube; (c) parallel-plate plasma-enhanced CVD system; (d) PECVD system using a three-zone furnace tube. Copyright 1983 Bell Telephone Laboratories, Inc. Reprinted by permission from Ref. [2].

6.3.1 CVD Reactors

Several different types of CVD reactor systems are shown in Fig. 6.8. In Fig. 6.8(a), a continuous atmospheric-pressure (APCVD) reactor is shown. This type of reactor has been used for deposition of the silicon dioxide passivation layer as one of the last steps

in IC processing. The reactant gases flow through the center section of the reactor and are contained by nitrogen curtains at the ends. The substrates can be fed continuously through the system, and large-diameter wafers are easily handled. However, high gas-flow rates are required by the atmospheric-pressure reactor.

The hot-wall, low-pressure system of Fig. 6.8(b) is commonly used to deposit polysilicon, silicon dioxide, and silicon nitride and is referred to as a low-pressure CVD (LPCVD) system. The reactant gases are introduced into one end of a three-zone furnace tube and are pumped out the other end. Temperatures range from 300 to 1150 °C, and the pressure is typically 30 to 250 Pa. Excellent uniformity can be obtained with LPCVD systems, and several hundred wafers may be processed in a single run. Hot-wall systems have the disadvantage that the deposited film simultaneously coats the inside of the tube. The tube must be periodically cleaned or replaced to minimize problems with particulate matter. In spite of this problem, hot-wall LPCVD systems are in widespread use throughout the semiconductor industry. Vertical furnaces similar to that depicted in Fig. 3.11(b) are also utilized for chemical vapor deposition.

CVD reactions can also take place in a plasma reactor, as shown in Fig. 6.8(c). Formation of the plasma permits the reaction to take place at low temperatures, which is a primary advantage of plasma-enhanced CVD (PECVD) processes. In the parallel-plate system, the wafers lie on a grounded aluminum plate, which serves as the bottom electrode for establishing the plasma. The wafers can be heated up to 400 °C using high-intensity lamps or resistance heaters. The top electrode is a second aluminum plate placed in close proximity to the wafer surface. Gases are introduced along the outside of the system, flow radially across the wafers, and are pumped through an exhaust in the center. An RF signal is applied to the top plate to establish the plasma. The capacity of this type of system is limited, and wafers must be loaded manually. A major problem in VLSI fabrication is particulate matter that may fall from the upper plate onto the wafers.

The furnace-plasma system in Fig. 6.8(d) can handle a large number of wafers at one time. A special electrode assembly holds the wafers parallel to the gas flow. The plasma is established between alternating groups of electrodes supporting the wafers.

Optical excitation, usually with laser sources, can also be used to assist or replace the thermal energy required for CVD reactions, and this form of processing is referred to as photon-enhanced chemical vapor deposition.

6.3.2 Polysilicon Deposition

Silicon is deposited in an LPCVD system using thermal decomposition of silane:



Low-pressure systems (25 to 150 Pa) use either 100% silane or 20 to 30% silane diluted with nitrogen. A temperature between 600 and 650 °C results in deposition of polysilicon material at a rate of 100 to 200 Å/min. A less commonly used deposition occurs between 850 and 1050 °C in a hydrogen atmosphere. The higher temperature overcomes a reduction in deposition rate caused by the hydrogen carrier gas.

Polysilicon can be doped by diffusion or ion implantation or during deposition (in situ) by the addition of dopant gases such as phosphine, arsine, or diborane. The addition of diborane greatly increases the deposition rate, whereas the addition of phosphine or arsine substantially reduces the deposition rate.

Polysilicon is often deposited as undoped material and is then doped by diffusion. High-temperature diffusion occurs much more rapidly in polysilicon than in single-crystal silicon, and the polysilicon film is typically saturated with the dopant to achieve as low a resistivity as possible for interconnection purposes. Resistivities of 0.01 to 0.001 ohm-cm can be achieved in diffusion-doped polysilicon. Ion implantation typically yields a lower active-impurity density in the polysilicon film, and ion-implanted polysilicon exhibits a resistivity about 10 times higher than that achieved by high-temperature diffusion.

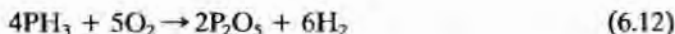
6.3.3 Silicon Dioxide Deposition

Silicon dioxide films can be deposited using a variety of reactions and temperature ranges, and the films can be doped or undoped. Phosphorus-doped oxide can be used as a passivation layer over a completed integrated circuit or as the insulating medium in multilevel metal processes (which will be discussed in the next chapter). Silicon dioxide containing 6 to 8% phosphorus by weight will soften and flow at temperatures between 1000 and 1100 °C. This "P-glass reflow" process can be used to improve step coverage and provide a smoother topography for later process steps. SiO₂ with lower concentrations of phosphorus will not reflow properly, and higher concentrations can corrode aluminum if moisture is present. Oxide doped with 5 to 15% by weight of various dopants can also be used as a diffusion source.

Deposition of silicon dioxide over aluminum must occur at a temperature below the silicon-aluminum eutectic point of 577 °C. (See Chapter 7.) A reaction between silane and oxygen between 300 and 500 °C is commonly used to deposit SiO₂:

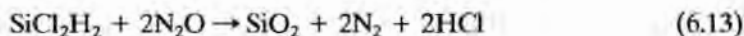


The oxide may be doped with phosphorus using phosphine:



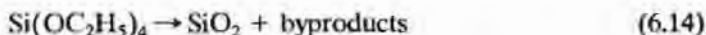
Oxide passivation layers can be deposited at atmospheric pressure using the continuous reactor of Fig. 6.8(a), or they can be deposited at reduced pressure in an LPCVD system, as in Fig. 6.8(b).

Deposition of SiO₂ films prior to metallization can be performed at higher temperatures, which gives a wider choice of reactions and results in better uniformity and step coverage. For example, a dichlorosilane reaction with nitrous oxide in an LPCVD system at approximately 900 °C,



can be used to deposit insulating layers of SiO₂ on wafer surfaces.

Decomposition of the vapor produced from a liquid source, tetraethylorthosilicate (TEOS), can also be used in an LPCVD system between 650 and 750 °C:



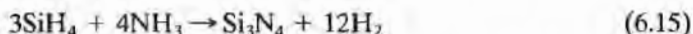
Deposition based on the decomposition of TEOS provides excellent uniformity and step coverage. Oxide doping may be accomplished in the LPCVD systems by adding phosphine, arsine, or diborane.

A comparison of some of the properties of various CVD oxides is given in Table 6.1.

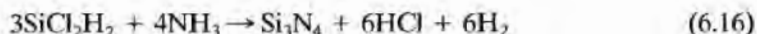
6.3.4 Silicon Nitride Deposition

As discussed in Chapter 3, silicon nitride is used as an oxidation mask in recessed oxide processes. Silicon nitride is also used as a final passivation layer, because it provides an excellent barrier to both moisture and sodium contamination. Composite films of oxide and nitride are being investigated for use as very thin gate insulators in scaled VLSI devices, and they are also used as the gate dielectric in electrically programmable memory devices.

Both silane and dichlorosilane will react with ammonia to produce silicon nitride. The silane reaction occurs between 700 and 900 °C at atmospheric pressure:

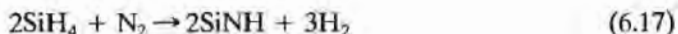


Dichlorosilane is used in an LPCVD system between 700 and 800 °C:

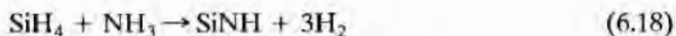


Thermal growth of silicon nitride is possible, but not very practical. Silicon nitride will form when silicon is exposed to ammonia at temperatures between 1000 and 1100 °C, but the growth rate is very low.

Plasma systems may be used for the deposition of silicon nitride. Silane will react with a nitrogen discharge to form plasma nitride (SiN):



Silane will also react with ammonia in an argon plasma:



LPCVD films are hydrogen-rich, containing up to 8% hydrogen. Plasma deposition does not produce stoichiometric silicon nitride films. Instead, the films contain as much as 20 to 25% hydrogen. LPCVD films have high internal tensile stresses, and

TABLE 6.1 Properties of Various Deposited Oxides (After Ref. [2])

Source	Deposition Temperature (°C)	Composition	Conformal Step Coverage	Dielectric Strength (MV/cm)	Etch Rate (Å/min) [100:1 H ₂ O:HF]
Silane	450	SiO ₂ (H)	No	8	60
Dichlorosilane	900	SiO ₂ (Cl)	Yes	10	30
TEOS	700	SiO ₂	Yes	10	30
Plasma	200	SiO _{1.9} (H)	No	5	400

films thicker than 2000 Å may crack because of this stress. On the other hand, plasma-deposited films have much lower tensile stresses.

The resistivity (10^{16} ohm-cm) and dielectric strength (10 MV/cm) of the LPCVD nitride films are better than those of most plasma films. Resistivity of plasma nitride can range from 10^6 to 10^{15} ohm-cm, depending on the amount of nitrogen in the film, while the dielectric strength ranges between 1 and 5 MV/cm.

6.3.5 CVD Metal Deposition

Many metals can be deposited by CVD processes. Molybdenum (Mo), tantalum (Ta), titanium (Ti), and tungsten (W) are all of interest in today's processes, because of their low resistivity and their ability to form silicides with silicon. (See Chapter 7.) Aluminum can be deposited from a metallorganic compound such as tri-isobutyl aluminum, but this technique has not been commonly used because many other excellent methods of aluminum deposition are available. Advanced metallization systems employ copper, and researchers are actively exploring CVD processes for copper deposition. However, at this writing, CVD copper processes suitable for use in manufacturing have not been developed, and copper is still deposited by standard electro- and electrolytic plating techniques similar to those utilized to produce printed circuit boards.

Tungsten can be deposited by thermal, plasma, or optically assisted decomposition of WF_6 :



or through reduction with hydrogen:



Mo, Ta, and Ti can be deposited in an LPCVD system through reaction with hydrogen. The reaction is the same for all three metals:



Here M stands for any one of the three metals previously mentioned.

6.4 EPITAXY

CVD processes can be used to deposit silicon onto the surface of a silicon wafer. Under appropriate conditions, the silicon wafer acts as a seed crystal, and a single-crystal silicon layer is grown on the surface of the wafer. The growth of a crystalline silicon layer from the vapor phase is called *vapor-phase epitaxy* (VPE), and it is the most common form of epitaxy used in silicon processing. In addition, *liquid-phase epitaxy* (LPE) and *molecular-beam epitaxy* (MBE) are being used widely in compound semiconductor technology.

Epitaxial growth was first used in IC processing to grow single-crystal *n*-type layers on *p*-type substrates for use in standard buried-collector bipolar processing. More recently, it has been introduced into CMOS VLSI processes where lightly doped layers are grown on heavily doped substrates of the same type (*n* on n^+ or *p* on p^+) to help suppress a circuit-failure mode called *latchup*.

6.4.1 VPE

Silicon epitaxial layers are commonly grown with silicon deposited from the gas phase. A basic model for the process is given in Fig. 6.9. At the silicon surface, the flux J_s of gas molecules is determined by

$$J_s = k_s N_s \quad (6.22)$$

where k_s is the surface-reaction rate constant and N_s is the surface concentration of the molecule involved in the reaction. In the steady state, this flux must equal the flux J_g of molecules diffusing in from the gas stream. The flux J_g may be approximated by

$$J_g = (\bar{D}_g/\delta)(N_g - N_s) = h_g(N_g - N_s) \quad (6.23)$$

where \bar{D} is an effective diffusion constant for the gas molecule and δ is the distance over which the diffusion is taking place. The ratio \bar{D}_g/δ is called the *vapor-phase mass-transfer coefficient*, h_g . Equating J_s and J_g yields the flux impinging on the surface of the wafer. The growth rate ν is equal to the flux divided by the number N of molecules incorporated per unit volume of film:

$$\nu = \frac{J_s}{N} = \frac{k_s h_g}{k_s + h_g} \frac{N_g}{N} \quad (6.24)$$

If $k_s \gg h_g$, then growth is said to be mass-transfer-limited, and

$$\nu = h_g \frac{N_g}{N} \quad (6.25)$$

If $h_g \gg k_s$, then growth is said to be surface-reaction-limited, and

$$\nu = k_s \frac{N_g}{N} \quad (6.26)$$

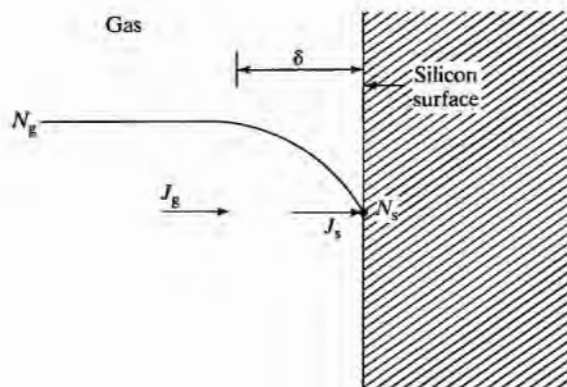


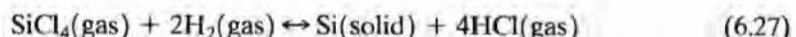
FIGURE 6.9

Model for the epitaxial growth process.

Figure 6.10 shows epitaxial growth rate as a function of temperature. Chemical reactions at the surface tend to follow an Arrhenius relationship characterized by an activation energy E_A , whereas the mass-transfer process tends to be independent of temperature. These two regions show up clearly in this figure. At low temperatures, the growth rate follows an Arrhenius relationship with an activation energy of approximately 1.5 eV. At higher temperatures, the growth rate becomes independent of temperature. To have good growth-rate control and to minimize sensitivity to variations in temperature, epitaxial growth conditions are usually chosen to yield a mass-transfer-limited growth rate.

Three common types of VPE reactors, the horizontal, pancake, and barrel systems, are shown in Fig. 6.11. The susceptor that supports the wafers is made of graphite and is heated by RF induction in the horizontal and vertical reactors and by radiant heating in the barrel reactor.

Silicon tetrachloride (SiCl_4), silane (SiH_4), dichlorosilane (SiH_2Cl_2), and trichlorosilane (SiHCl_3) have all been used for silicon VPE. Silicon tetrachloride has been widely used in industrial processing:



This reaction takes place at approximately 1200 °C and is reversible. If the carrier gas coming into the reactor contains hydrochloric acid, etching of the surface of the silicon wafer can occur. This in-situ etching process can be used to clean the wafer prior to the start of epitaxial deposition.

A second reaction competes with the epitaxial deposition process:



This second reaction also etches the silicon from the wafer surface. If the concentration of SiCl_4 is too high, etching of the wafer surface will take place rather than epitaxial

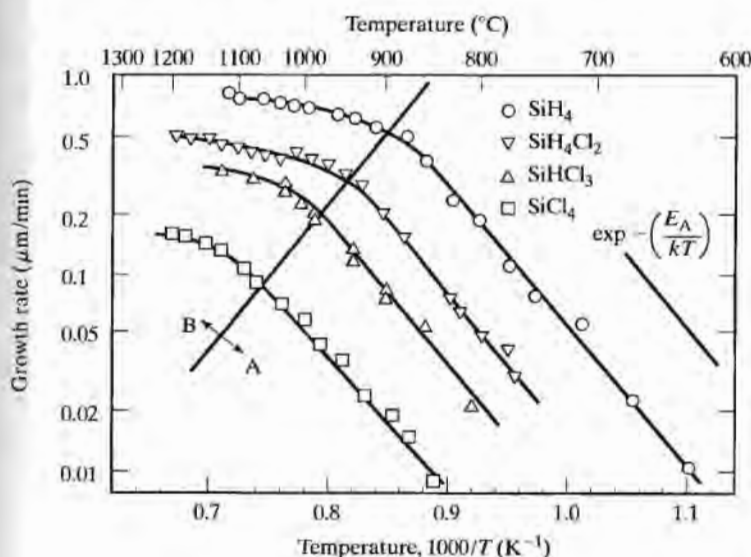


FIGURE 6.10

Temperature dependence of the silicon epitaxial growth process for four different sources. The growth rate is surface-reaction-limited in region A and is mass-transfer-limited in region B. Reprinted with permission from Philips Journal of Research from Ref. [3].

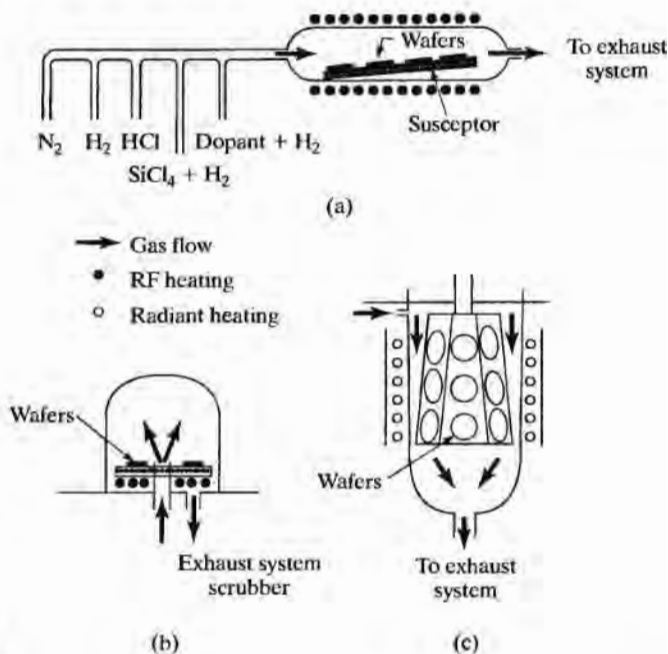


FIGURE 6.11

(a) Horizontal, (b) pancake, and (c) barrel susceptors commonly used for vapor-phase epitaxy. Copyright 1985 John Wiley & Sons, Inc. Reprinted with permission from Ref. [1].

deposition. Figure 6.12 shows the effect of $SiCl_4$ concentration on the growth of epitaxial silicon. The growth rate initially increases with increasing $SiCl_4$ concentration, peaks, and then decreases. Eventually, growth stops and the etching process becomes dominant. If the growth rate is too high, a polysilicon layer is deposited, rather than a layer of single-crystal silicon.

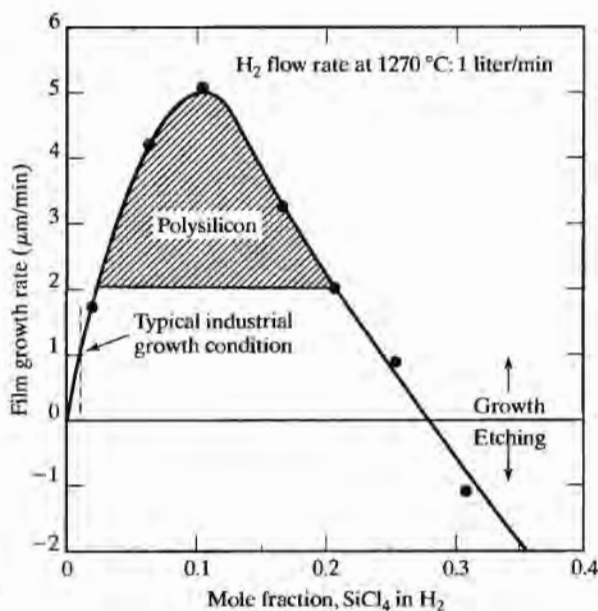


FIGURE 6.12

Silicon epitaxial growth rate as a function of $SiCl_4$ concentration. Polysilicon deposition occurs for growth rates exceeding $2 \mu m/min$. Etching of the surface will occur for mole fraction concentrations exceeding 28%. Copyright 1985 John Wiley & Sons, Inc. Reprinted with permission from Ref. [1].

Epitaxial growth can also be achieved by the pyrolytic decomposition of silane:



The reaction is not reversible and takes place at low temperatures. In addition, it avoids the formation of HCl gas as a reaction by-product. However, careful control of the reactor is needed to prevent formation of polysilicon rather than single-crystal silicon layers. The presence of any oxidizing species in the reactor can also lead to contamination of the epitaxial layer by silica dust.

6.4.2 Doping of Epitaxial Layers

Epitaxial layers may be doped during the growth process by adding impurities to the gas used for deposition. Arsine, diborane, and phosphine are the most convenient sources of the common impurities. The resistivity of the epitaxial layer is controlled by varying the partial pressure of the dopant species in the gas supplied to the reactor. The addition of arsine or phosphine tends to slow down the rate of epitaxial growth, while the addition of diborane tends to enhance the epitaxial growth rate.

Lightly doped epitaxial layers are often grown on more heavily doped substrates, and autodoping of the epitaxial layer can occur during growth. Impurities can evaporate from the wafer or may be liberated by chlorine etching of the surface during deposition. The impurities are incorporated into the gas stream, resulting in doping of the growing layer. As the epitaxial layer grows, less dopant is released from the wafer into the gas stream, and the impurity profile eventually reaches a constant level determined by the doping in the gas stream.

During deposition, the substrate also acts as a source of impurities which diffuse into the epitaxial layer. This out-diffusion will be discussed more fully in the next section. Both autodoping and out-diffusion cause the transition from the doping level of the substrate to that of the epitaxial layer to be less abrupt than desired. The effects of autodoping and out-diffusion are illustrated in Fig. 6.13.

6.4.3 Buried Layers

Out-diffusion is a common problem that occurs with the buried layer in bipolar transistors. In order to reduce the resistance in series with the collector of the bipolar transistor, heavily doped n -type regions are diffused into the substrate prior to the growth of an n -type epitaxial layer. During epitaxy, impurities diffuse upward from the heavily doped buried-layer regions.

Diffusion of impurities from the substrate during epitaxial growth is modeled by the diffusion equation with a moving boundary [4], as in Fig. 6.14,

$$D \frac{\partial^2 N}{\partial x^2} = \frac{\partial N}{\partial t} + v_x \frac{\partial N}{\partial x} \quad (6.30)$$

where v_x is the rate of growth of the epitaxial layer.

Two specific solutions of Eq. (6.30) are applicable to epitaxial layer growth. The first case is the growth of an undoped epitaxial layer on a uniformly doped substrate.

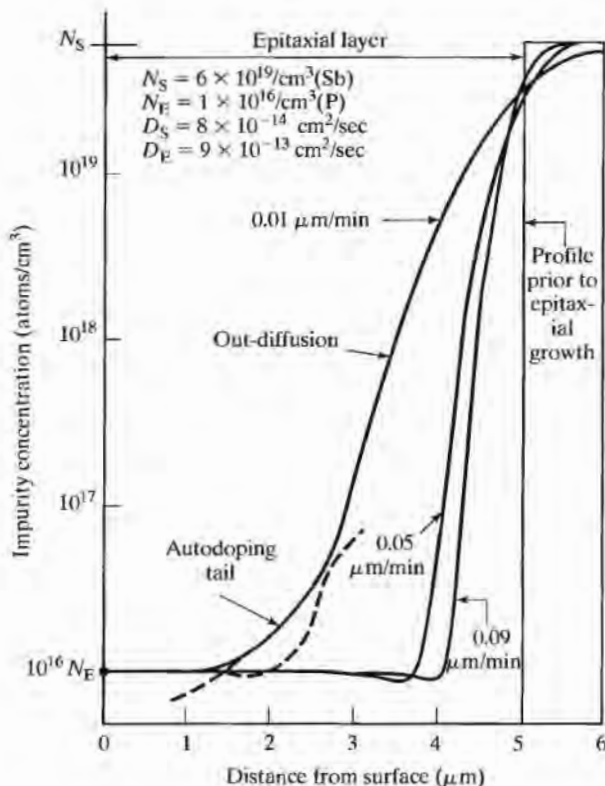


FIGURE 6.13

Redistribution of impurity atoms due to gas-phase autodoping and impurity out-diffusion during epitaxial layer growth. Out-diffusion is calculated using Eq. (6.33) for epitaxial growth of a phosphorus-doped layer at 1150 °C over an antimony-doped buried layer with a surface concentration of $6 \times 10^{19}/\text{cm}^3$. The three curves are for growth rates of 0.01, 0.05, and 0.09 $\mu\text{m}/\text{min}$. For clarity, the effects of autodoping are shown on only one curve.

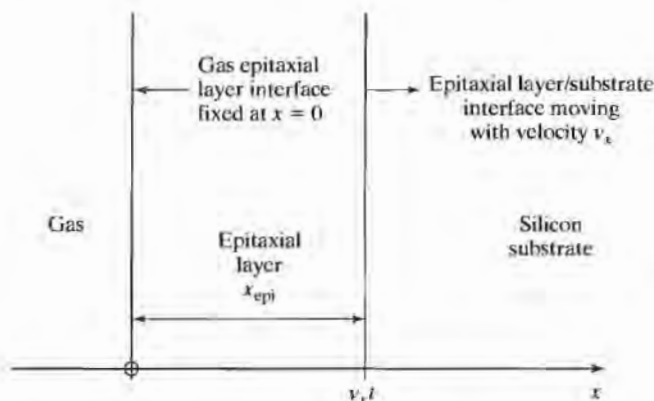


FIGURE 6.14

Geometrical model for the moving boundary value problem which describes the epitaxial growth process.

The boundary conditions are $N(x, 0) = N_s = N(\infty, t)$, and the flux $J_x = (h + v_x)N(0, t)$, where h is the mass-transfer coefficient, which characterizes the escape rate of dopant atoms from the silicon into the gas. Normally, $h \ll v_x$. A change of variables from x to $x' = x - v_x t$ simplifies Eq. (6.30) and gives an approximate solution for $N(x, t)$:

$$N_1(x, t) = \frac{N_s}{2} \left[1 + \operatorname{erf} \frac{x - x_{\text{epi}}}{2\sqrt{D_s t}} \right] \quad (6.31)$$

Equation (6.31) assumes that the epitaxial layer growth rate greatly exceeds the rate of movement of the diffusion front. Equation (6.31) is the exact solution for diffusion from one semi-infinite layer into a second semi-infinite layer.

The second case is the growth of a doped epitaxial layer on an undoped substrate. The boundary conditions for this case are $N(0, t) = N_E$ and $N(\infty, t) = 0 = N(x, 0)$. The solution of Eq. (6.30) for these boundary conditions is:

$$N_2(x, t) = \frac{N_E}{2} \left[\operatorname{erfc} \frac{x - x_{\text{epi}}}{2\sqrt{D_E t}} + \exp \frac{v_s x}{D_E} \operatorname{erfc} \frac{x + x_{\text{epi}}}{2\sqrt{D_E t}} \right] \quad (6.32)$$

where $x_{\text{epi}} = v_s t$ is the epitaxial layer thickness. Superposition of the solutions for these two cases gives a good approximation to diffusion which occurs during epitaxial growth:

$$N(x, t) = N_1(x, t) + N_2(x, t). \quad (6.33)$$

N_s represents the doping in the substrate, and N_E is the doping intentionally introduced into the epitaxial layer. D_s and D_E represent the diffusion coefficients of the impurity species in the substrate and epitaxial layer, respectively. Figure 6.13 shows diffusion profiles for a phosphorus-doped epitaxial layer grown at various rates on an antimony-doped substrate. The curves were produced using Eq. (6.33).

An additional problem occurs during epitaxial growth. The oxidation and lithographic processing steps used during formation of a buried layer result in a step of as much as 0.2 μm around the perimeter of the buried layer. Epitaxial growth on this non-planar surface causes the pattern to shift during growth, as illustrated in Fig. 6.15. Pattern shift is difficult to predict, may be as large as the epitaxial layer thickness, and must be accounted for during the design of subsequent mask levels.

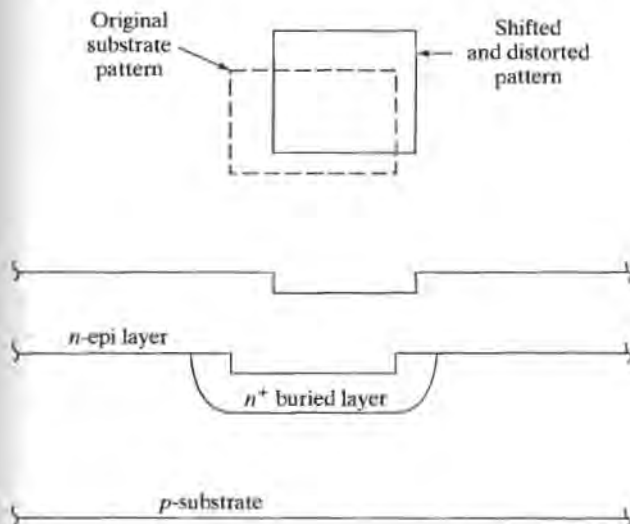


FIGURE 6.15

Pattern shift during epitaxial growth over an n^+ buried layer. The original pattern is shifted and distorted in shape.

6.4.4 Liquid-Phase and Molecular-Beam Epitaxy

In liquid-phase epitaxy, the substrate is brought into contact with a solution containing the material to be deposited in liquid form. The substrate acts as a seed for material crystallizing directly from the solute. Growth rates typically range between 0.1 and 1 $\mu\text{m}/\text{min}$.

In the molecular-beam epitaxy process, the crystalline layer is formed by deposition from a thermal beam of atoms or molecules. Deposition is performed in ultrahigh-vacuum conditions (10^{-8} Pa). Substrate temperatures during MBE range from 400 to 900 $^{\circ}\text{C}$, and the growth rate is relatively low (0.001 to 0.3 $\mu\text{m}/\text{min}$). The epitaxial layer is grown atomic layer by atomic layer, and many unique device structures can be fabricated by changing the material which is deposited between one layer and the next.

The throughput of MBE is relatively low. Plasma-assisted CVD processes, which promise to give many of the benefits of MBE with much higher throughput, continue to be investigated in research laboratories.

SUMMARY

Thin films of a very broad range of materials are used in IC fabrication. This chapter has presented an overview of film-deposition techniques, including physical evaporation, chemical vapor deposition (CVD), epitaxial growth, and sputtering. Most of these processes are performed at low pressure, and this chapter has presented an introduction to vacuum systems and a review of some important aspects of ideal gas theory.

Physical evaporation using filament or electron-beam evaporators can be used to deposit metals and other materials that can easily be melted. E-beam systems can operate at high power levels and melt high-temperature metals. However, E-beam evaporation may result in radiation damage to thin oxide layers at the surface of the wafer. In addition, it is difficult to deposit material compounds and alloys using evaporation. Finally, gas molecules at low pressures have large mean free paths, and evaporation has problems with shadowing and poor step coverage during film deposition.

Sputtering uses energetic ions such as argon to bombard a target material and dislodge atoms from the surface of the target. The dislodged atoms are deposited on the surface of the wafer. Direct-current sputtering systems can be used to deposit conductive materials, and RF sputtering can be used to deposit insulators. Sputtering can be used to deposit composite materials in which the deposited film maintains the same composition as the source material. Sputtering also uses higher pressures than evaporation. The much shorter mean free paths that result yield a deposition with freedom from shadowing and much better step coverage.

Low-pressure and atmospheric chemical vapor deposition (CVD) systems deposit films from chemical reactions taking place in a gas stream passing over the wafer. Polysilicon, silicon dioxide, silicon nitride, and metals can all be deposited using CVD techniques. A special type of CVD deposition called epitaxy results in the growth of single-crystal silicon films on the surface of silicon wafers. Out-diffusion and autodoping cause problems with impurity profile control during epitaxial layer growth.

In a modern bipolar or MOS fabrication process, one can expect to find evaporation, sputtering, and CVD techniques all used somewhere in the process flow.

REFERENCES

- [1] S. M. Sze, *Semiconductor Devices—Physics and Technology*, John Wiley & Sons, New York, 1985.
- [2] A. C. Adams, "Dielectric and Polysilicon Film Deposition," Chapter 3 in S. M. Sze, Ed., *VLSI Technology*, McGraw-Hill, New York, 1983.
- [3] F. C. Eversteyn, "Chemical-Reaction Engineering in Semiconductor Industry," *Philips Research Reports*, 29, 45–66 (February, 1974).
- [4] A. B. Glaser and G. E. Subak-Sharpe, *Integrated Circuit Engineering*, pp. 205–209, Addison-Wesley, Reading, MA, 1979.
- [5] W. S. Rusk, *Microelectronic Processing*, McGraw-Hill, New York, 1987.
- [6] W. R. Runyan and K. E. Bean, *Semiconductor Integrated Circuit Processing Technology*, Addison-Wesley Publishing Company, Reading, MA, 1990.
- [7] S. A. Campbell, *The Science and Engineering of Microelectronic Fabrication*, Oxford University Press, New York, 1996.

FURTHER READING

- [1] J. L. Vossen and W. Kern, Eds., *Thin Film Processes*, Academic Press, New York, 1978.
- [2] J. F. O'Hanlon, *A User's Guide to Vacuum Technology*, John Wiley & Sons, New York, 1980.
- [3] L. Holland, *Vacuum Deposition of Thin Films*, John Wiley & Sons, New York, 1961.
- [4] A. S. Grove, *Physics and Technology of Semiconductor Devices*, John Wiley & Sons, New York, 1967.
- [5] H. C. Theuerer, "Epitaxial Silicon Films by Hydrogen Reduction of SiCl_4 ," *Journal of the Electrochemical Society*, 108, 649–653, July 1961.
- [6] C. O. Thomas, D. Kahng, and R. C. Manz, "Impurity Distribution in Epitaxial Silicon Films," *Journal of the Electrochemical Society*, 109, 1055–1061, November 1962.
- [7] A. S. Grove, A. Roder, and C. T. Sah, "Impurity Distribution in Epitaxial Growth," *Journal of Applied Physics*, 36, 802–810, March 1965.
- [8] D. Kahng, C. O. Thomas, and R. C. Manz, "Epitaxial Silicon Junctions," *Journal of the Electrochemical Society*, 110, 394–400, May 1963.
- [9] W. H. Shepherd, "Vapor Phase Deposition and Etching of Silicon," *Journal of the Electrochemical Society*, 112, 988–994, October 1965.
- [10] G. R. Srinivasan, "Autodoping Effects in Silicon Epitaxy," *Journal of the Electrochemical Society*, 127, 1334–1342, June 1980.
- [11] J. C. Bean, "Silicon Molecular Beam Epitaxy as a VLSI Processing Technique," 1981 *IEEE IEDM Proceedings*, pp. 6–13, December 1981.

PROBLEMS

- 6.1 A silicon wafer sits on a bench in the laboratory at a temperature of 300 K and a pressure of 1 atm. Assume that the air consists of 100% oxygen. How long does it take to deposit one atomic layer of oxygen on the wafer surface, assuming 100% adhesion?
- 6.2 Repeat Problem 6.1, but this time the wafer is kept in a nitrogen-purged cabinet in which the oxygen content is less than 0.1% of the total gas content.
- 6.3 Calculate the impingement rate and mean free path for oxygen molecules ($M = 32$) at 300 K and a pressure of 10^{-4} Pa. What is this pressure in torr?
- 6.4 An ultrahigh-vacuum system operates at a pressure of 10^{-8} Pa. What is the concentration of residual air molecules in the chamber at 300 K?
- 6.5 A high-vacuum system has a residual nitrogen concentration of 1000 molecules/cm². What is the gas pressure at 300K?

- 6.6** The partial pressure of a material being deposited in a vacuum system must be well above the residual background gas pressure if reasonable deposition rates are to be achieved. What must the partial pressure of aluminum be to achieve a deposition rate of 100 nm/min? Assume close packing of spheres with a diameter of 5 Å, 100% adhesion of the impinging aluminum, and a temperature of 300 K.
- 6.7** A wafer 100 mm in diameter is mounted in an electron-beam evaporation system in which the spherical radius is 40 cm. Use Eq. (6.7) to estimate the worst-case variation in film thickness between the center and edges of the wafer for an evaporated aluminum film 1 μm thick.
- 6.8** Repeat Problem 6.7 for a 200-mm wafer mounted 50 cm from the e-beam source.
- 6.9** Electron-beam evaporation is going to be used to deposit a 0.6-μm-thick layer of aluminum on a 300-mm-diameter wafer. The thickness variation between the center and edges of the wafer is desired to be less than 0.05 μm. How far should the wafer be from the source?
- 6.10** An MBE system must operate under ultrahigh-vacuum conditions to prevent the formation of undesired atomic layers on the surface of the substrate. What pressure of oxygen can be permitted at 300 K if formation of a monolayer of contamination can be permitted after the sample has been in the chamber for no less than 4 hr?
- 6.11** (a) Calculate the growth rate of a silicon layer from an SiCl_4 source at 1200 °C. Use $h_g = 1$ cm/sec, $k_s = 2 \times 10^6 \exp(-1.9/kT)$ cm/sec, and $N_g = 3 \times 10^{16}$ atoms/cm³. (For silicon, $N = 5 \times 10^{22}$ /cm³.)
 (b) What is the change in growth rate if the temperature is increased by 25 °C?
 (c) At what temperature does $k_s = h_g$? What is the growth rate at this temperature?
 (d) What is the value of E_A in Fig. 6.10?
- 6.12** Use Eqs. (6.31) and (6.32) to model the case of a 10-μm *n*-type epitaxial layer ($N_E = 1 \times 10^{16}$ /cm³) grown on a *p*-type substrate ($N_S = 1 \times 10^{18}$ /cm³). Plot the impurity profile in the epitaxial layer and substrate assuming that the layer was grown at a rate of 0.2 μm/min at a temperature of 1200 °C. Assume that boron and phosphorus are the impurities. Find the location of the *pn* junction.
- 6.13** Compare and discuss the advantages and disadvantages of evaporation, sputtering, and chemical vapor deposition.
- 6.14** (a) A 1-kg source of aluminum is used in an E-beam evaporation system. How many 100-mm wafers can be coated with a 1-μm Al film before the source material is exhausted? Assume that 15% of the evaporated aluminum actually coats a wafer. (The rest is deposited on the inside of the electron-beam system.)
 (b) Repeat the process for a 300-mm wafer.
- 6.15** (a) A silicon wafer 100 mm in diameter is centered 200 mm above a small planar evaporation source. Calculate the ratio of thickness between the center and edges of the wafer using Eq. (6.7), following a 1-μm film deposition.
 (b) Repeat the process for 200- and 300-mm wafers centered 40 cm above the small planar source.
- 6.16** Advanced CMOS processes often use lightly doped epitaxial layers grown on heavily doped substrates. Use Eq. (6.31) to predict the dopant profile in the epilayer if an intrinsic silicon layer is grown on top of a substrate that has a uniform concentration of 10^{20} As atoms/cm³. Assume the layer thickness is 1 μm and that it is grown in SiCl_4 at 1,100 °C.
- 6.17** Repeat Problem 6.16 for the case where the intrinsic layer is grown on a substrate that has a uniform concentration of 2×10^{20} B atoms/cm³. Assume the layer thickness is 2 μm and that it is grown in SiH_4Cl_2 at 950 °C.

CHAPTER 7

Interconnections and Contacts

The previous six chapters focused on the various processes used to fabricate semiconductor devices in the silicon substrate. To complete the formation of an integrated circuit, one must interconnect the devices and finally get connections to the world outside the silicon chip. Until the 1970s, integrated circuits had two possible levels of interconnection: diffusions and metallization. The use of polysilicon as a gate material in MOS devices added a third level useful for interconnecting devices and circuits.

In this chapter, we discuss the various forms of interconnections and the problems associated with making good contacts between metal and silicon. Refractory metal silicides and multilevel metallization used in VLSI processes are discussed, and an additional method for depositing patterned films, called *liftoff*, is also introduced.

Copper has been introduced in IC processing because of its lower resistivity. The CMP process, introduced in Chapter 3, is combined with standard electroplating techniques to achieve highly planar, inlaid copper interconnections referred to as Damascene technology. Low dielectric constant interlevel films are used to reduce the capacitance of the interconnection levels.

7.1 INTERCONNECTIONS IN INTEGRATED CIRCUITS

As we found in previous chapters, aluminum, polysilicon, and diffused regions are all easily isolated from each other using an insulating layer of silicon dioxide. Thus, today's integrated circuits have three different materials that may cross over each other. To be useful as an interconnect, the materials must also provide as low a sheet resistance as possible in order to minimize voltage drops along the interconnect lines, as well as to minimize propagation delay caused by the resistance and capacitance of the line. Finally, low-resistance "ohmic" contacts must be made between the materials, and the interconnection lines must be reliable throughout long-term operation.

Figure 7.1 shows a simple MOS logic circuit illustrating how polysilicon, metal, and diffused interconnections may cross over or contact each other. Aluminum is used

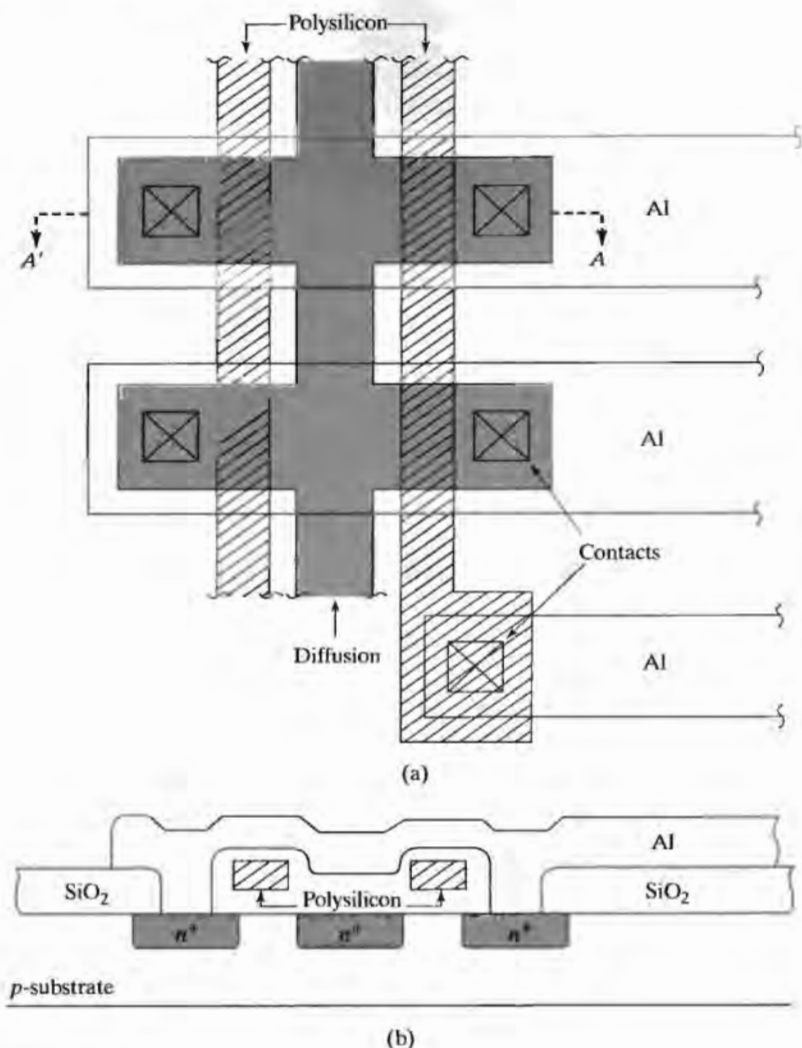


FIGURE 7.1

Portion of a MOS logic circuit showing the use of diffusion, polysilicon, and aluminum interconnections. (a) Top view; (b) cross section through $A'-A$.

to make contact to diffusions and polysilicon, and diffusions in various regions have been extended and merged together to form interconnections.

In this technology, polysilicon lines and diffused lines can only be connected together using the metal level. Improved circuit density can often be achieved by using "buted contacts" between polysilicon and diffusions or by changing the process to introduce "buried contacts" directly between the polysilicon and diffused layers. These two techniques will be examined later in this chapter.

7.2 METAL INTERCONNECTIONS AND CONTACT TECHNOLOGY

The requirement for low-resistivity materials leads one immediately to consider metals for use as interconnections, and the resistivities of common metals are compared in Table 7.1. Historically, aluminum and gold have been used with silicon IC processing. Gold requires the use of a multilayer sandwich involving other metals such as titanium or tungsten. Gold can be troublesome, because it is a rapid diffuser (see Fig. 4.5) in silicon and produces deep-level recombination centers in silicon that tend to significantly reduce the lifetime of free carriers. In addition, gold forms many problematic inter-metallic compounds. Because of these various problems, the use of gold is most often restricted to chip packaging technologies.

Aluminum is compatible with silicon IC processing and is the most common material in use today. It is relatively inexpensive, adheres well to silicon dioxide, and has a bulk resistivity of $2.7 \mu\Omega\text{-cm}$. However, care must be exercised to avoid a number of problems associated with the formation of good aluminum contacts to silicon.

Advanced multilevel metallization systems employ copper, because of its improved resistivity relative to aluminum. Copper has an even larger diffusion coefficient in silicon than gold and also causes lifetime reduction and leakage in silicon. Therefore, copper is generally not introduced into the fabrication sequence until one or two levels of aluminum metallization and interlevel dielectric levels have been formed above the semiconductor devices. These metallization and dielectric layers act as passivation layers to protect the active devices below.

7.2.1 Ohmic Contact Formation

We desire to form “ohmic” contacts between the metal and semiconductor. True ohmic contacts would exhibit a straight-line I - V characteristic with a low value of resistance (see Fig. 7.2(a)), as opposed to the I - V characteristic of a rectifying contact shown in Fig. 7.2(b). Figure 7.2(c) shows an I - V characteristic more representative of a practical ohmic contact to silicon. Although nonlinear near the origin, it develops only a small voltage across the contact at normal current levels.

Figure 7.3 shows a number of ways in which aluminum may contact semiconductor regions during device fabrication. Aluminum contact to p -type silicon normally results in

TABLE 7.1 Bulk Resistivity of Metals ($\mu\Omega\text{-cm}$)

Ag: Silver	1.6
Al: Aluminum	2.65
Au: Gold	2.2
Co: Cobalt	6
Cu: Copper	1.7
Mo: Molybdenum	5
Ni: Nickel	7
Pd: Palladium	10
Pt: Platinum	10.6
Ti: Titanium	50
W: Tungsten	5

Source: WebElements (<http://www.webelements.com>)

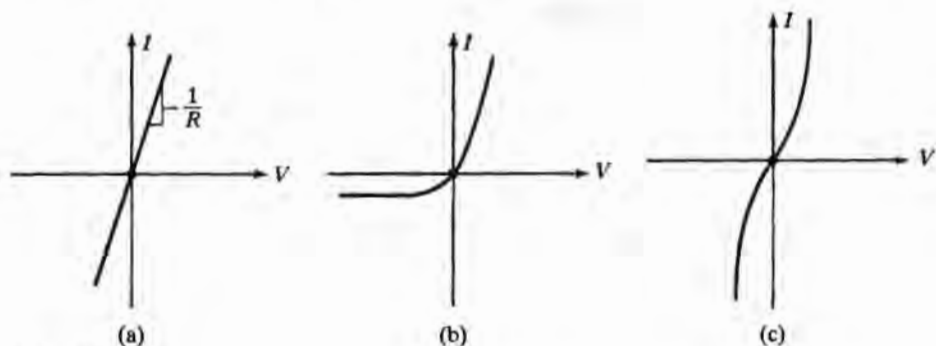


FIGURE 7.2

I-V characteristics of contacts between integrated-circuit materials. (a) Ideal ohmic contact; (b) rectifying contact; (c) practical nonlinear "ohmic" contact.

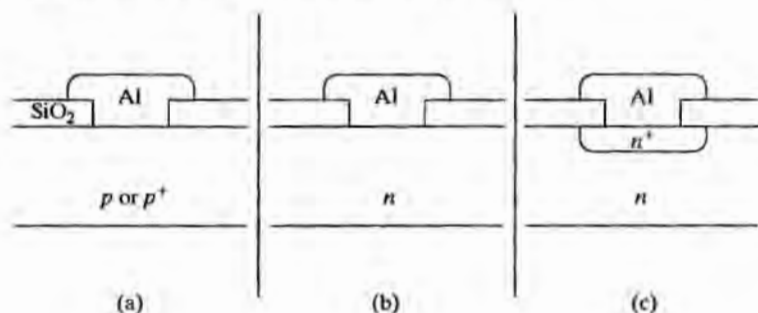


FIGURE 7.3

Three possible types of aluminum contacts to silicon. (a) Aluminum to *p*-type silicon forms an ohmic contact with an *I-V* characteristic approximating that in Fig. 7.2a; (b) aluminum to *n*-type silicon can form a rectifying contact (Schottky barrier diode) like that in Fig. 7.2b; (c) aluminum to *n*⁺ silicon yields a contact similar to that in Fig. 7.2c.

a good ohmic contact for doping levels exceeding $10^{16}/\text{cm}^3$. However, a problem arises in trying to contact *n*-type silicon, as shown in Fig. 7.3(b). For lightly doped *n*-type material, aluminum can form a metal-semiconductor "Schottky-barrier" diode rather than an ohmic contact. To prevent this rectifying contact from forming, an *n*⁺ diffusion is placed between the aluminum and any lightly doped *n*-type regions, as in Fig. 7.3(c). The resulting contact has an *I-V* characteristic similar to that in Fig. 7.2(c). This technique was used in forming the collector contact in the bipolar process shown in Fig. 1.6.

7.2.2 Aluminum-Silicon Eutectic Behavior

Silicon melts at a temperature of 1412 °C, and pure aluminum melts at 660 °C. However, aluminum and silicon together exhibit "eutectic" characteristics in which a mixture of the two materials lowers the melting point of the composite material to below that of either element. Figure 7.4 shows the phase diagram of the aluminum-silicon system at a pressure of 1 atm. The minimum melting temperature, or *eutectic temperature*, is 577 °C and corresponds to an 88.7% Al, 11.3% Si composition. Because of the relatively low eutectic

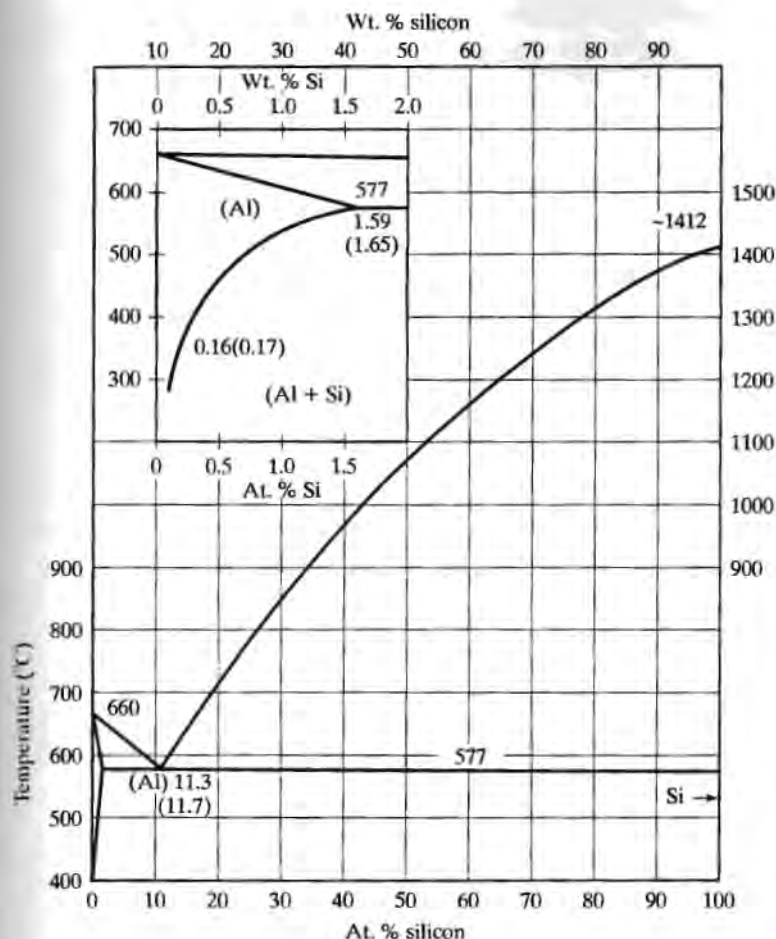


FIGURE 7.4

Phase diagram of the aluminum-silicon system. The silicon-aluminum eutectic point occurs at a temperature of 577 °C. At contact-alloying temperatures between 450 and 500 °C, aluminum will absorb from 0.5 to 1% silicon. Copyright 1958 McGraw-Hill Book Company, reprinted with permission from Ref. [1].

temperature of the Al-Si system, aluminum must be introduced into the IC process sequence after all high-temperature processing has been completed.

7.2.3 Aluminum Spiking and Junction Penetration

To ensure good contact formation, aluminum is normally annealed in an inert atmosphere at a temperature of 450 to 500 °C following deposition and patterning. Although this temperature is well below the eutectic temperature for silicon and aluminum, silicon still diffuses into the aluminum. The diffusion leads to a major problem associated with the formation of aluminum contacts to silicon, particularly for shallow junctions.

Anywhere a contact is made between aluminum and silicon, silicon will be absorbed by the aluminum during the annealing process. The amount of silicon absorbed will depend on the time and temperature involved in the annealing process, as well as the area of the contact. (See Problem 7.4.) To make matters worse, the silicon is not absorbed uniformly from the contact region. Instead, it tends to be supplied from a few points. As

the silicon is dissolved, spikes of aluminum form and penetrate the silicon contact region. If the contact is to a shallow junction, the spike may cause a junction short, as in Fig. 7.5.

The inset in Fig. 7.4 gives the solubility of silicon in aluminum. Between 400 °C and the eutectic temperature, the solubility of silicon in aluminum ranges from 0.25 to 1.5% by weight. To solve the spiking problem, silicon may be added to the aluminum film during deposition by coevaporation from two targets, or sputter deposition can be used with an aluminum target containing approximately 1% silicon. Both of these techniques deposit a layer in which the aluminum demand for silicon is satisfied, and the metallization does not absorb silicon from the substrate during subsequent annealing steps.

The junction penetration problem becomes particularly acute in high-density VLSI processes with extremely shallow junctions; another way to prevent spiking is to place a barrier material between the aluminum and silicon, as shown in Fig. 7.5. One possibility is to deposit a thin layer of polysilicon prior to aluminum deposition. The polysilicon will then supply the silicon needed to saturate the aluminum. Another alternative is to use a metal as a barrier. The metal must form a low-resistance contact with silicon, not react with aluminum, and be compatible with other process steps. Various semiconductor manufacturers have used a number of metals, including platinum, palladium, titanium, and tungsten.

7.2.4 Contact Resistance

There is a small resistance associated with an ohmic contact between two materials. To a first approximation, the *contact resistance* R_c is inversely proportional to the area of the contact:

$$R_c = \rho_c / A \quad (7.1)$$

where ρ_c is the specific contact resistivity in ohm-cm² and A is the area of the contact. For example, a $2 \times 2 \mu\text{m}$ contact with $\rho_c = 1 \mu\text{ohm-cm}^2$ yields a contact resistance of 25 ohms. Figure 7.6 shows the contact resistivity as a function of annealing temperature for several aluminum-silicon systems. It is evident why the 450 °C annealing process is used following aluminum deposition. Also note that the use of polysilicon under aluminum to prevent junction spiking yields a much poorer value of ρ_c .

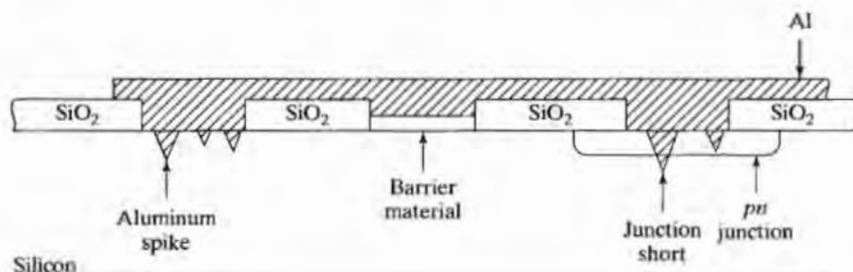


FIGURE 7.5

Aluminum spiking which occurs during aluminum-silicon alloying. Aluminum spikes can cause shorts in shallow junctions. Aluminum containing 1% silicon is often used to eliminate spiking. A barrier material of polysilicon or a metal such as titanium can also be used to prevent spiking.

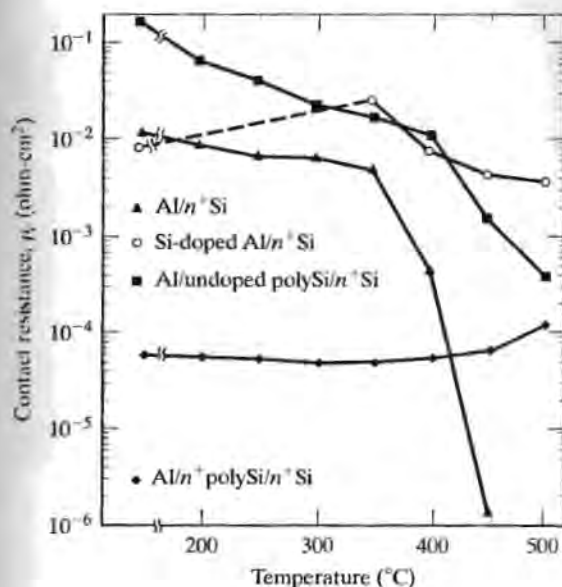


FIGURE 7.6

Contact resistivity of a variety of aluminum-silicon systems. An alloying temperature of 450 $^{\circ}\text{C}$ is typically used to obtain low-contact resistance for Al-Si contacts. Reprinted with permission from *Solid-State Electronics*, Vol. 23, pp. 255-262, M. Finetti et al., "Aluminum-Silicon Ohmic Contact on Shallow n^+ /p Junctions" [2]. Copyright 1980, Pergamon Press, Ltd.

7.2.5 Electromigration

Metal interconnections in integrated circuits are operated at relatively high current densities, and a very interesting failure mechanism develops in aluminum and other conductors. *Electromigration* is the movement of atoms in a metal film due to momentum transfer from the electrons carrying the current. Under high-current-density conditions, metal-atom movement causes voids in some regions and metal pileup, or *hillocks*, in other regions, as shown in Fig. 7.7. Voids can eventually result in open circuits, and pileup can cause short circuits between closely spaced conductors.

The mean time to failure (MTF) of a conductor due to electromigration has been experimentally related to current density, J , and temperature by

$$\text{MTF} \propto (J^{-2}) \exp(E_A/kT) \quad (7.2)$$

where E_A is an activation energy with a typical value of 0.4 to 0.5 eV for aluminum.

The most common method of improving aluminum resistance to electromigration is to add a small percentage of a heavier metal such as copper. Targets composed of 95% Al, 4% Cu, and 1% Si are routinely used in sputter deposition systems. The aluminum-copper-silicon alloy films simultaneously provide electromigration resistance and eliminate aluminum spiking.

Pure copper interconnections would be expected to exhibit a much higher electromigration resistance than aluminum, which is in fact the case, as shown in the results in Fig. 7.8, which compares electromigration performance of TiN clad copper lines with those formed of an AlCu alloy [6]. An order of magnitude in improvement is obtained.

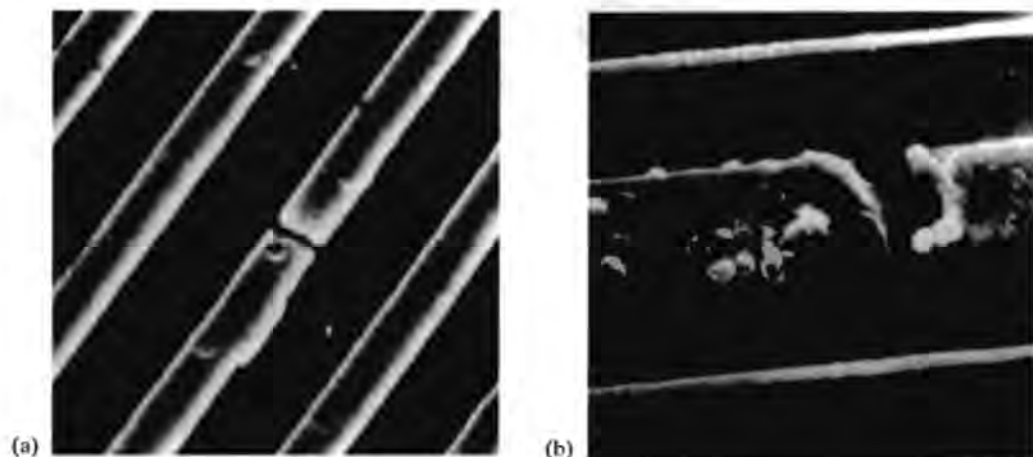


FIGURE 7.7

Scanning electron micrographs of aluminum interconnection failure caused by electromigration. (a) Sputtered aluminum with 0.5% copper; (b) evaporated aluminum with 0.5% copper. Copyright 1980, IEEE. Reprinted with permission from Ref. [3].

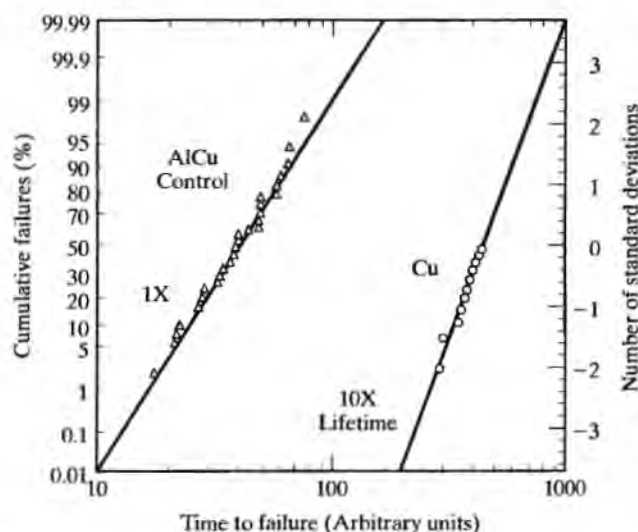


FIGURE 7.8

Electromigration performance improvement using copper metallization. Copyright 1997, IEEE. Reprinted with permission from Ref. [6].

7.3 DIFFUSED INTERCONNECTIONS

Diffused conductors with low sheet resistances represent the second available interconnect medium in basic IC technology. From Fig. 4.16, we can see that the minimum resistivity is approximately 1,000 $\mu\text{ohm-cm}$. For shallow structures measuring about 1 μm , the minimum obtainable sheet resistance is typically between 10 and 20 ohms per square. Such sheet resistances are obviously much higher than that of metal, and one must be selective in the use of diffusions for signal or power distribution.

The diffused line must really be modeled as a distributed RC structure, as illustrated in Fig. 7.9, when signal propagation is considered. The resistance, R , of diffused regions was discussed in detail in Chapter 4, and C represents capacitance of the

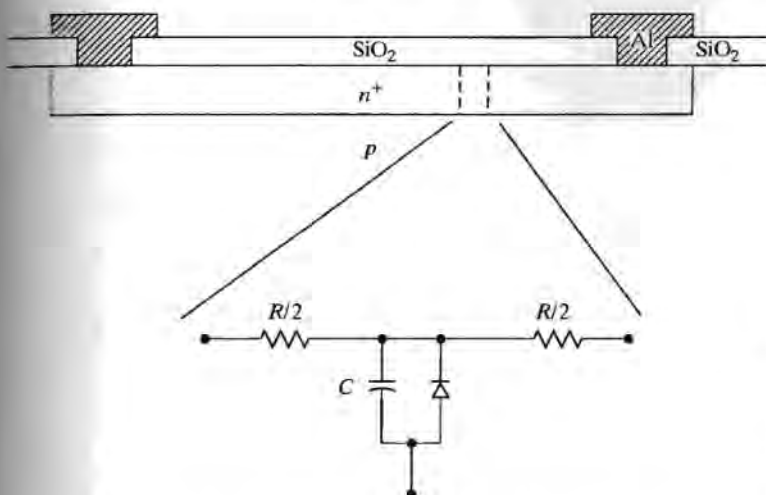


FIGURE 7.9

A lumped circuit model for a small section of an n^+ diffusion. The RC line delay limits the use of diffusions for high-speed signal distribution.

reverse-biased pn junction formed between the diffused region and the substrate. Heavily doped diffusions are normally used for interconnection purposes and can be approximated by a one-sided step junction in which the depletion layer extends predominantly into the substrate. The capacitance per unit area is given by

$$C = \sqrt{\frac{qN_sK_s\epsilon_0}{2(\phi_{bi} + V_R)}} \quad \phi_{bi} = (kT/q) \ln \left(\frac{N_s}{n_i} \right) + 0.56 \text{ V} \quad (7.3)$$

where N_s is the substrate doping, ϕ_{bi} is the built-in potential of the junction, and V_R is the reverse bias applied to the junction.

The relatively large RC product of long diffused lines results in substantial time delay for signals propagating down such a line. Hence, diffusions are more useful in interconnecting adjacent devices in integrated circuits. Figure 7.10 shows a three-input NMOS NOR-gate in which the source diffusions of the three input transistors are merged together as one diffusion. The three drains of the input devices, as well as the source of the depletion-mode load device, are also merged together as one diffusion. Figure 7.1 shows an example of the use of long diffused interconnection regions in a programmable-logic-array (PLA) structure.

7.4 POLYSILICON INTERCONNECTIONS AND BURIED CONTACTS

Heavily doped n -type polysilicon is the primary MOS transistor gate material in use today, and it provides an additional layer of interconnection that is easily insulated from other layers by thermal oxidation or insulator deposition. This extra level of interconnection greatly facilitates the layout of compact digital integrated circuits. Thin, heavily doped polysilicon layers have a minimum resistivity of approximately 300 $\mu\text{ohm-cm}$, and they suffer from the same sheet-resistance problems associated with shallow diffused interconnections (typically 20 to 30 ohms per square). Polysilicon lines have substantial capacitance to the substrate and exhibit RC delay problems similar to those of diffused interconnections.

7.4.1 Buried Contacts

In the polysilicon-gate processes presented thus far, the polysilicon acts as a barrier material during ion implantation or diffusion. Thus, a diffusion can never pass beneath a polysilicon line. In addition, contact windows to the diffusions are not opened until after polysilicon deposition. It is therefore necessary to use a metal link to connect between polysilicon and diffusion, as in Fig. 7.11(a). Interconnecting the diffusion to polysilicon in this manner requires two contact windows and an intervening space, both of which are wasteful of area.

In memory arrays, where density is extremely important, an extra mask step can be introduced into the process to permit direct contact between polysilicon and silicon, as shown in Fig. 7.11(b). Prior to polysilicon deposition, windows are opened in the thin gate oxide, permitting the polysilicon to contact the silicon surface. Diffusion of the n -type dopant from the heavily doped n^+ polysilicon merges with the adjacent ion-implanted n^+ regions, and the result is called a *buried contact*. The edge of the contact exhibits the lowest resistance since the impurity concentration is greatest in that region.

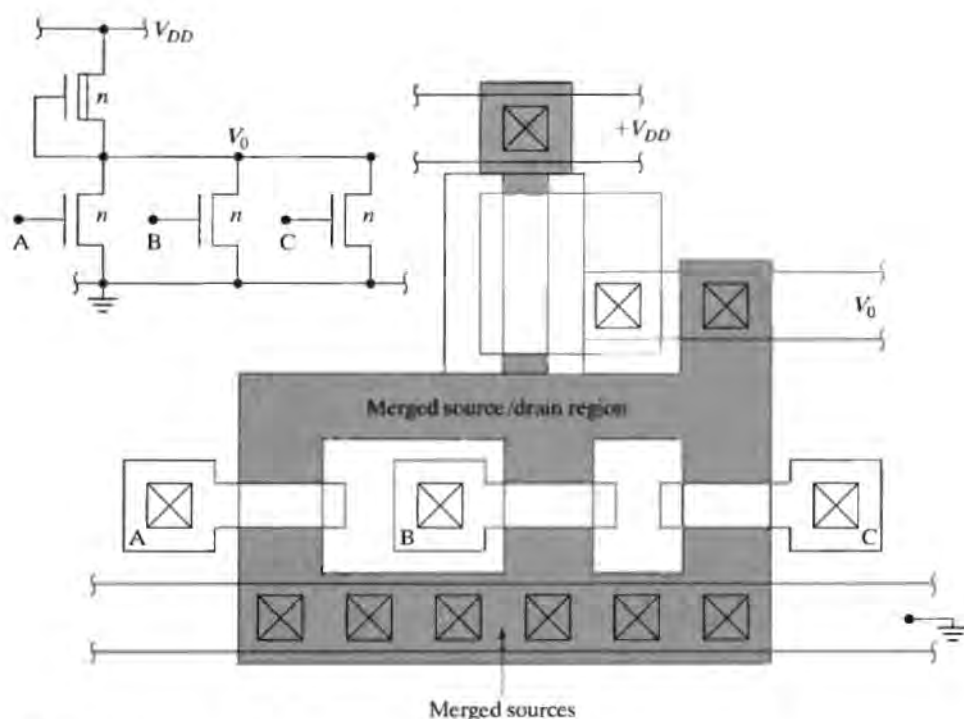


FIGURE 7.10

Layout of a three-input NMOS NOR-gate showing device interconnection through merging of adjacent source and drain diffusions.

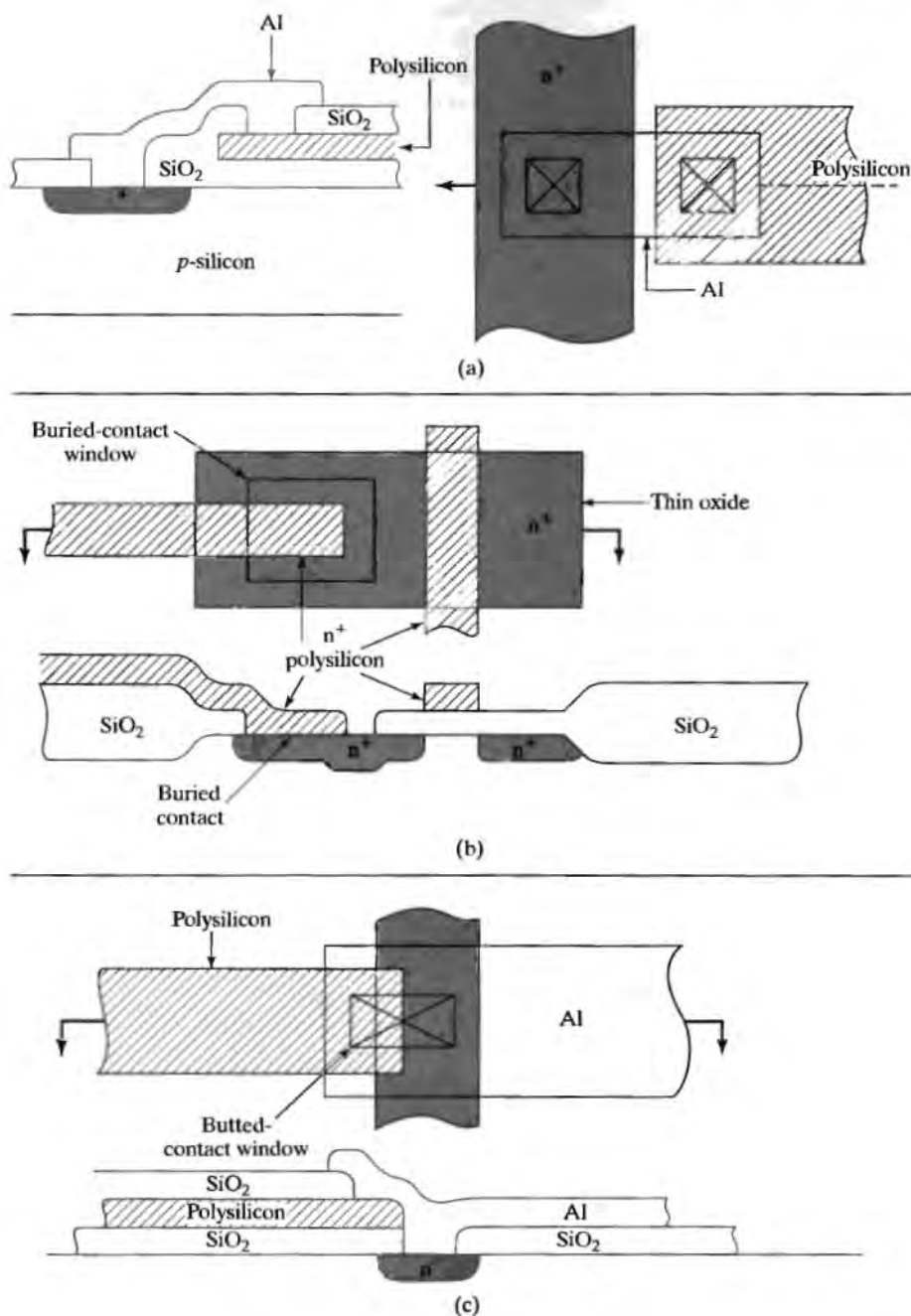


FIGURE 7.11

Three techniques for interconnecting polysilicon and n^+ diffusion. (a) Normal aluminum link requiring two contact regions and an intervening space; (b) buried-contact structure; (c) butted-contact structure.

7.4.2 Butted Contacts

Another method of conserving area is to form a “butted” contact as shown in Fig. 7.11(c). In this example, polysilicon is aligned with the edge of the diffusion contact window, and metal connects the diffusion and polysilicon together. The butted contact saves area by eliminating the space normally required between separate contact windows.

7.5 SILICIDES AND MULTILAYER-CONTACT TECHNOLOGY

The sheet resistance of both thin polysilicon and shallow diffusions cannot be reduced below 10 to 20 ohms per square, which greatly reduces their utility as an interconnection medium. Interconnect delays limit the speed of VLSI circuits, and as dice get larger and feature sizes get smaller, methods for improving these interconnections have had to be found.

7.5.1 Silicides, Polycides, and Salicides

A wide range of noble and refractory metals form compounds with silicon called *silicides*, and the sheet resistance of polysilicon and diffusion can be reduced by forming a low-resistivity, shunting silicide layer on their surfaces. A list of properties of possible silicides is given in Table 7.2. Several of the elements, including titanium, tungsten, platinum, and palladium, have been used in the formation of Schottky-barrier diodes in bipolar processes since the 1960s and are now used to form silicides for interconnection purposes.

TABLE 7.2 Properties of Some Silicides of Interest. Reprinted with permission of the American Institute of Physics from Ref [4].

Silicide	Starting Form	Sintering Temperature (°C)	Lowest Binary Eutectic Temperature (°C)	Specific Resistivity (μhm-cm)
CoSi ₂	Metal on polysilicon	900	1195	18–25
	Cosputtered alloy	900		
HfSi ₂	Metal on polysilicon	900	1300	45–50
MoSi ₂	Cosputtered alloy	1000	1410	100
NiSi ₂	Metal on polysilicon	900	966	50
	Cosputtered alloy	900		50–60
Pd ₂ Si	Metal on polysilicon	400	720	30–50
PtSi	Metal on polysilicon	600–800	830	28–35
TaSi ₂	Metal on polysilicon	1000	1385	35–45
	Cosputtered alloy	1000		50–55
TiSi ₂	Metal on polysilicon	900	1330	13–16
	Cosputtered alloy	900		25
WSi ₂	Cosputtered alloy	1000	1440	70
ZrSi ₂	Metal on polysilicon	900	1355	35–40

A structure with a silicide formed on top of the polysilicon gate, often called a *polycide*, is shown in Fig. 7.12. A layer of the desired metal is deposited using evaporation, sputtering, or CVD techniques. Upon heating of the structure to a temperature between 600 and 1000 °C, the metal reacts with the polysilicon to form the desired silicide. Coevaporation, cosputtering, or sputtering of a composite target may be used to simultaneously deposit both silicon and metal onto the polysilicon surface prior to the thermal treatment or "sintering" step. Silicides have resistivities in the range of 15 to 50 $\mu\text{ohm-cm}$.

Another feature of silicide layers is the ability to oxidize the surface following silicide formation. At high temperatures, silicon diffuses readily through the silicide layer and will combine with oxygen at the silicide surface to form an SiO_2 insulating layer.

The eutectic temperature of the silicide and silicon will limit the temperature of further processing steps, as in the case of aluminum. However, many silicides are stable at temperatures exceeding 1000 °C. Exceptions include the silicides of nickel (900 °C), platinum (800 °C), and palladium (700 °C).

Silicides are also used to reduce the effective sheet resistance of diffused interconnections. Figure 7.13 outlines a process for simultaneous formation of silicides on both the gate and source-drain regions of an MOS transistor. An oxide spacer is used to prevent silicide formation on the side of the gate, because such formation could cause a short between the gate and diffusions. The spacer is formed by first coating the surface with a CVD oxide, followed by a reactive-ion etching step. The oxide along the edge of the gate is thicker than over other regions, and some oxide is left on the side of the gate at the point when the oxide is completely removed from the source and drain regions and the top of the gate. Next, metal is deposited over the wafer. During sintering, silicide forms only in the regions where metal touches silicon or polysilicon. Unreacted metal may be removed with a selective etch that does not attack the silicide. The result is a silicide that is automatically self-aligned to the gate and source-drain regions. *Self-aligned silicides* are often called *salicides*.

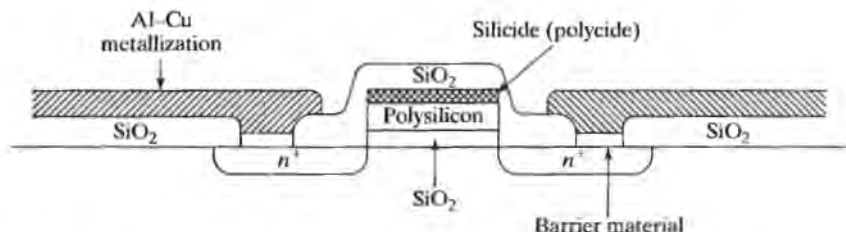


FIGURE 7.12

MOS structure showing the use of a "polycide" to reduce the sheet resistance of the polysilicon gate material and a barrier material to prevent aluminum spiking through shallow source-drain junctions.

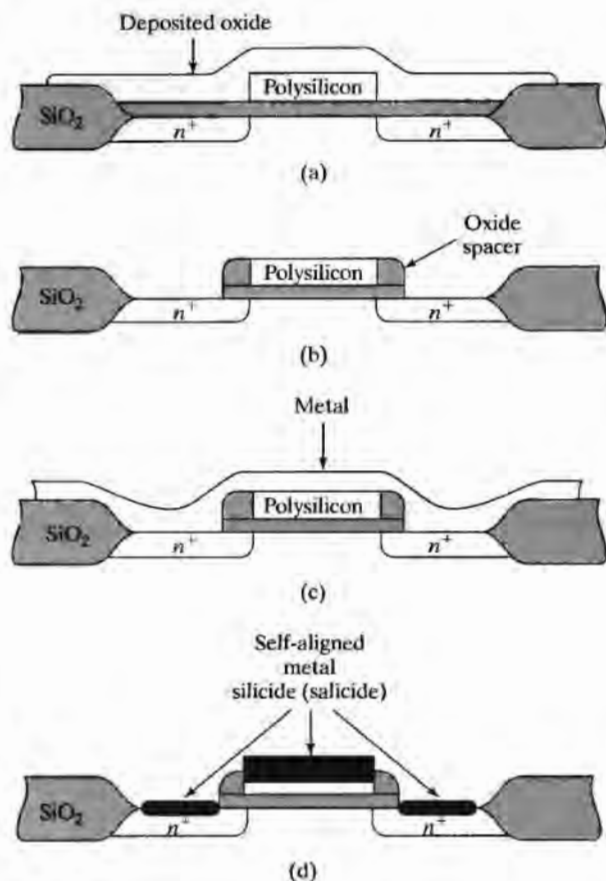


FIGURE 7.13

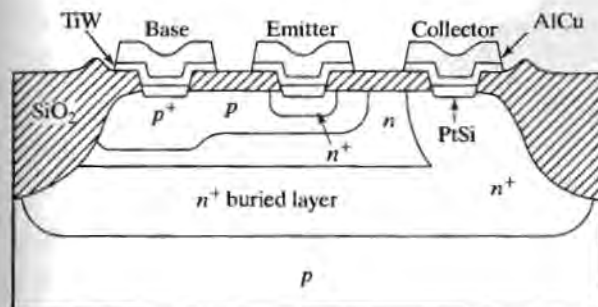
Use of self-aligned silicide ("salicide") in the formation of an MOS device. (a) Oxide is deposited over the normal MOS structure following polysilicon definition; (b) structure after reactive-ion etching leaving a sidewall oxide spacer; (c) metal is deposited over the structure and heated to form silicides; (d) unreacted metal is readily etched away, leaving silicide automatically aligned to gate and source-drain regions.

7.5.2 Barrier Metals and Multilayer Contacts

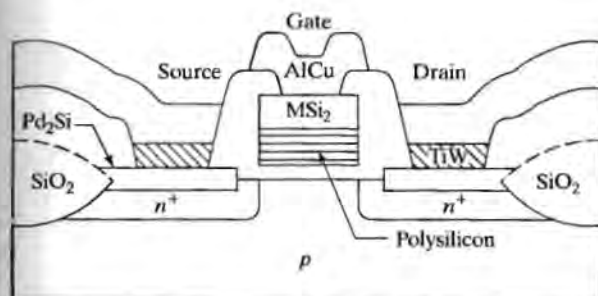
Aluminum contacts to silicides suffer from the same pitting and spiking problems associated with direct contact to silicon. To circumvent these problems, an intermediate layer of metal is used that prevents silicon diffusion. Figure 7.14 shows the application of titanium-tungsten (TiW) as a barrier metal over the silicides in the contact regions of both bipolar and MOS technologies. The final contact consists of a sandwich of a silicide over the diffusion, followed by the TiW diffusion barrier, and completed with aluminum-copper interconnection metallization. Multilayer contact structures are common in advanced, high-performance MOS, and bipolar technologies.

7.6 THE LIFTOFF PROCESS

The pattern definition processes which have been discussed previously have been "subtractive" processes, as illustrated in Fig. 7.15(a). The wafer is completely covered with a thin film layer, which is selectively protected with a masking layer such as photoresist. Wet or dry etching then removes the thin film material from the unprotected areas.



(a)



(b)

FIGURE 7.14

Device cross sections showing the use of silicide contacts in (a) bipolar and (b) MOS devices. Reprinted with permission from *Semiconductor International* magazine, August 1985[5]. Copyright 1985 by Cahners Publishing Co., Des Plaines, IL.

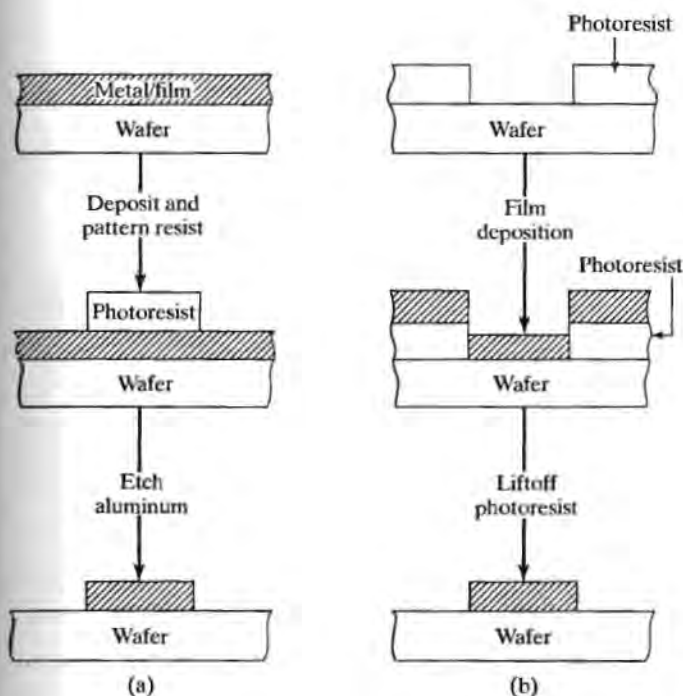


FIGURE 7.15

A comparison of interconnection formation by (a) subtractive etching and (b) additive metal liftoff.

The additive or *liftoff* process shown in Fig. 7.15(b) can also be used, in which the substrate is first covered with a photoresist layer patterned with openings where the final material is to appear. The thin film layer is deposited over the surface of the wafer. Any material deposited on top of the photoresist layer will be removed with the resist, leaving the patterned material on the substrate. For liftoff to work properly, there must be a very thin region or a gap between the upper and lower films. Otherwise, tearing and incomplete liftoff will occur.

The masking patterns for the liftoff and subtractive processes are the negatives of each other. This can be achieved by changing the mask from dark field to light field or by changing from negative to positive photoresist.

7.7 MULTILEVEL METALLIZATION

A single level of metal simply does not provide sufficient capability to fully interconnect complex VLSI chips. Many processes now use two or three levels of polysilicon, as well as several levels of metallization, in order to ensure wirability and provide adequate power distribution.

7.7.1 Basic Multilevel Metallization

A multilevel metal system is shown in Fig. 7.16. Standard processing is used through the deposition and patterning of the first level of metal. An interlevel dielectric, consisting

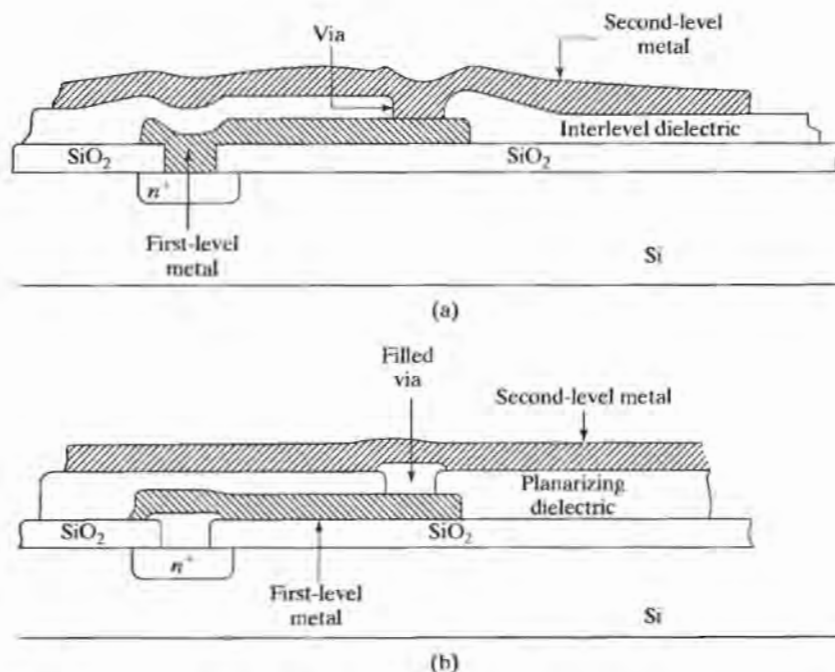


FIGURE 7.16

(a) Basic two-level metallization process may use polyimide, oxide, or nitride as an interlevel dielectric; (b) additional process steps may be added to fill the vias with metal prior to each metal deposition in order to achieve a more planar structure.

of CVD or sputtered SiO_2 , or a plastic-like material called *polyimide*, is then deposited over the first metal layer. The dielectric layer must provide good step coverage and should help smooth the topology. In addition, the layer must be free of pinholes and be a good insulator. Next, vias are opened in the dielectric layer, and the second level of metallization is deposited and patterned.

7.7.2 Planarized Metallization

The topology that results from the simple multilayer interconnect process of Fig. 7.16(a) simply cannot be utilized in submicron processes because of the depth-of-field limitations in the lithographic processes. The CMP process introduced in Chapter 3 is used to achieve highly planar layers. In the process flow in Fig. 7.16(b), a via filling technique is used to form the vias between metal layers. Tungsten is commonly used as the via metallization. The dielectric deposition, metallization, and CMP processes are repeated until the desired number of levels of interconnection is achieved. Integrated circuits with six levels of metal have been successfully fabricated using similar processes. An example of a planarized multilevel metal system employing aluminum metallization and tungsten “plugs” appears in Fig. 7.17 [7].

7.7.3 Low Dielectric Constant Interlevel Dielectrics

Propagation delay associated with interconnections is a critical issue in high-performance microprocessors, as well as other integrated circuits. The RC product associated with these interconnections can be decreased by reducing either the resistance or capacitance or both. Copper is being used to reduce the resistance term and is described in

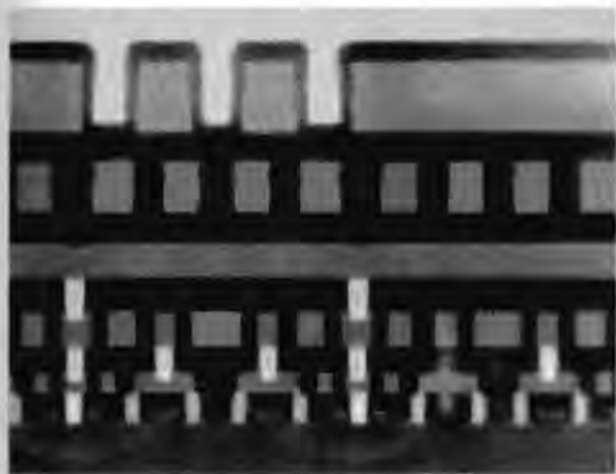


FIGURE 7.17

Multilevel aluminum metallization with tungsten plugs. Copyright 1998 IEEE. Reprinted with permission from Ref. [7].

detail in Section 7.8. Silicon dioxide, the most common interlevel dielectric, has a relative dielectric constant of 3.9. This value is fairly high, although much less than that of silicon itself ($\epsilon_r = 11.7$). Air isolated interconnects, with $\epsilon_r = 1$, have been utilized in GaAs circuits, but have not been successfully applied to silicon integrated circuits. Research teams are presently trying to identify dielectric materials that are compatible with silicon IC technology and have ϵ_r values in the 2.0–2.5 range. Fluorinated oxides, porous oxides and many polymer materials are under investigation [11–13].

7.8 COPPER INTERCONNECTS AND DAMASCENE PROCESSES

Because of its lower resistivity (see Table 7.1), copper is being used in place of other metals in multilevel metal systems. Unfortunately, copper is a deep-level impurity and a very rapid diffuser in silicon (see Fig. 4.5), and so great care must be exercised to prevent it from contaminating the silicon substrate and devices. The metallization techniques discussed so far have been subtractive processes in which the metal is deposited everywhere and then etched away where not desired. Manufacturable dry etching processes have not been developed for the removal of copper, so additive plating techniques are used. The Damascene processes use chemical mechanical polishing, as discussed in Chapter 3, to produce highly planar layers that may be used for multiple layers of interconnect.

7.8.1 Electroplated Copper Interconnect

Two methods of electroplating copper interconnect lines are shown in Fig. 7.18 [8]. The first involves plating through a mask. In Fig. 7.18(a), a conductive seed layer must first be deposited on the wafer that may already be planarized utilizing a CMP step. The seed layer provides an electrical path that is needed for current during the plating process. A masking layer such as photoresist is then deposited on the wafer and lithographically patterned. The wafer is immersed in the plating system, and a dc bias is applied between the solution and the seed layer. Copper plating occurs wherever the seed layer is exposed to the plating solution. Following the plating operation, the masking layer is removed and the seed layer etched away, leaving a copper interconnection line on top of the substrate. However, the lack of planarity of this structure limits its application in today's ICs.

7.8.2 Damascene Plating

The Damascene process of Fig. 7.18(b) is ideally suited to IC interconnect structures as pointed out in [8]. An insulating layer such as silicon dioxide is deposited on the surface of the substrate, and standard photolithography is used to define the desired interconnect pattern. Following seed-layer deposition, the entire surface is electroplated, filling the interconnect regions as well as covering the rest of the surface with copper. The excess copper is polished away using a CMP process step. The final structure is highly planar with metal lines inlaid in the insulator.¹

¹The metal interconnect structures produced by Damascene processing are embedded or inlaid in the insulating background material. The inlaid structures are reminiscent of inlaid metal from Damascus and hence the use of the name Damascene processes.

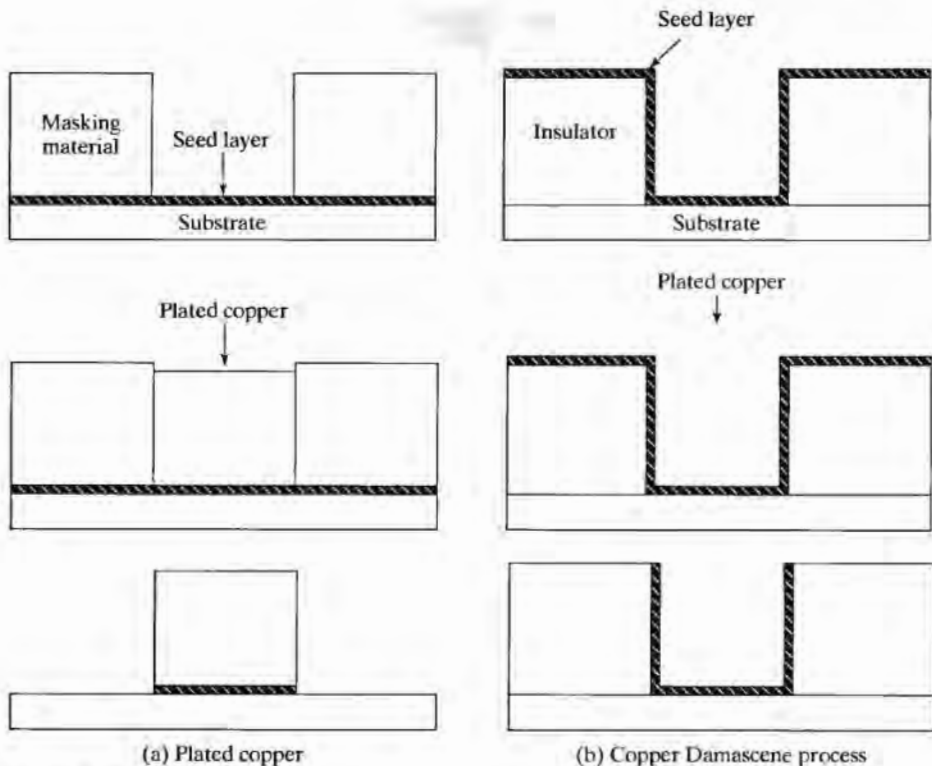


FIGURE 7.18

Plated copper and copper Damascene process steps. (a) Mask openings are defined over seed layer and copper is plated in the opening. A nonplanar structure results after plating mask removal. (b) A seed layer is deposited over the patterned insulator, and copper is plated over the entire structure. A planar surface results after the excess copper is lapped away.

7.8.3 Dual Damascene Structures

The full power of the Damascene technique is realized through the dual Damascene process, which forms interconnection lines and vias between interconnect levels at the same time, as illustrated in Fig. 7.19. The process begins with the substrate coated with a thin etch stop layer such as silicon nitride. Two layers of an insulator such as SiO_2 are deposited with a thin intervening etch stop layer. The insulator sandwich is capped with a final etch stop layer.

Windows are opened in the silicon nitride layer defining the via locations, and the insulator is etched away with the etch terminating on the silicon nitride layer. A new set of windows defining the interconnection lines are etched through the nitride. The nitride etch stop is also removed from the bottom of the via. The oxide is then etched simultaneously from the upper and lower levels of oxide. A barrier layer such as titanium nitride (TiN) may be deposited before the seed layer deposition to prevent interaction between the plated copper layer and the insulator. Electroplated copper then builds up the vias and the interconnect lines simultaneously. The structure is completed with CMP removal of the excess copper. The process can be

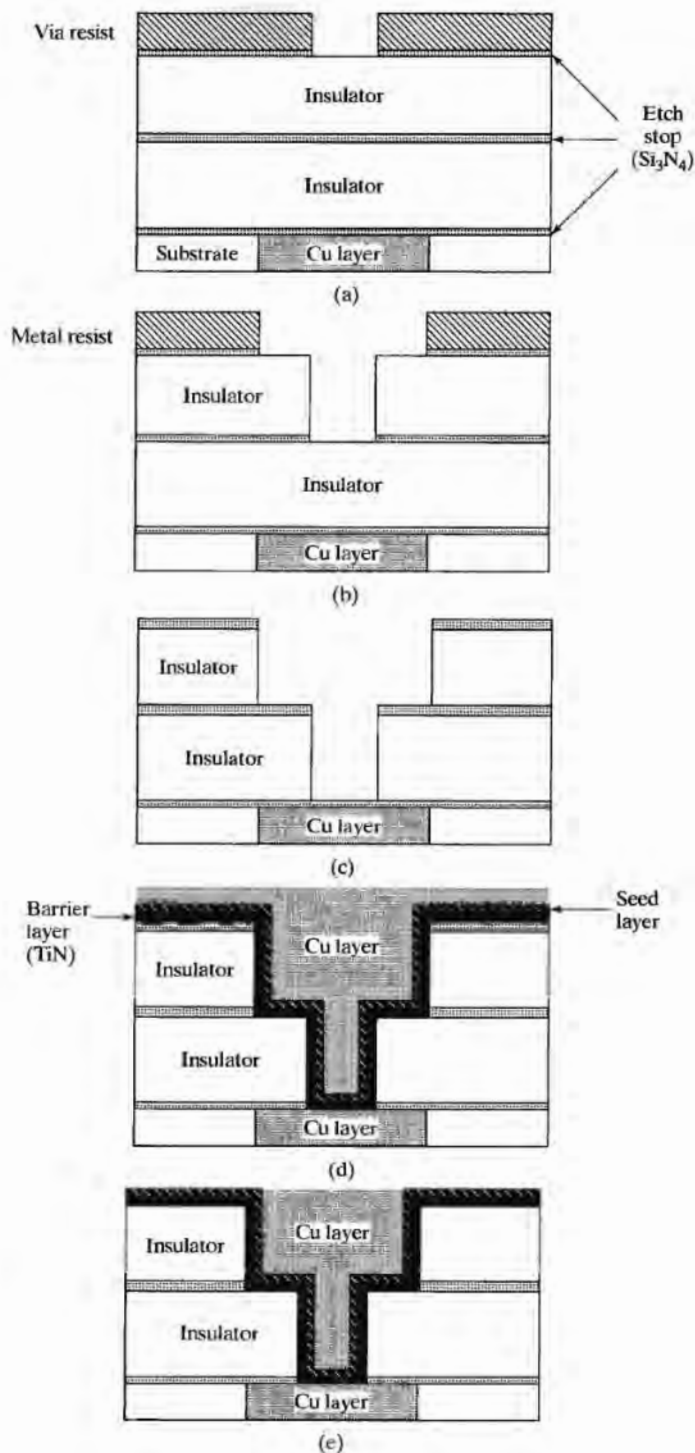
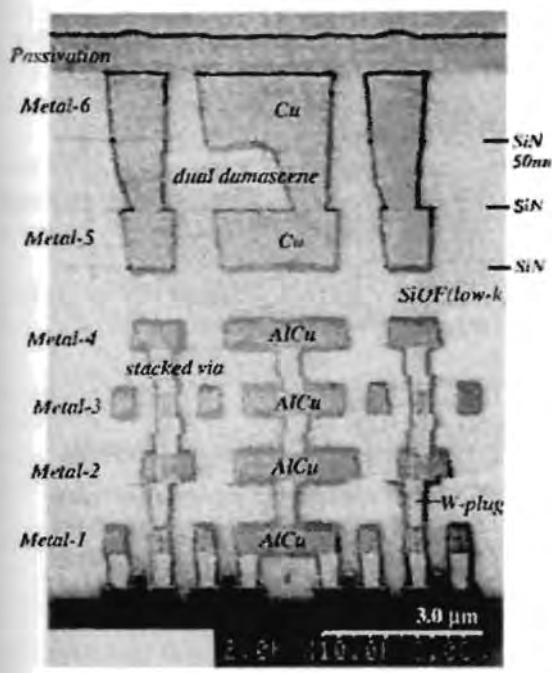


FIGURE 7.19

Dual Damascene process flow. (a) An insulator sandwich is first deposited and the upper nitride layer is patterned. The insulator layer is etched. The etch terminates on the silicon nitride etch stop. (b) The nitride layer is patterned and etched. (c) Following the next oxide etch step, two different width openings exist in the two oxide layers. (d) Barrier and seed layers are deposited and plated with copper. (e) Final structure following removal of excess copper.

repeated to build up additional layers of interconnect. The resulting copper interconnects are surrounded by a thin cladding layer of TiN.

Figure 7.20 shows two multilevel metal systems involving dual Damascene processing. The first [6], Fig. 7.20(a), uses four layers of aluminum copper interconnections with tungsten via plugs plus two levels of Damascene copper interconnections. The second [9], Fig. 7.20(b), shows the use of six levels of copper wiring.



(a)



(b)

FIGURE 7.20

Microphotographs of six-level metalization.

(a) Dual Damascene copper. Courtesy of Motorola, Inc. (b) Dual Damascene copper combined with aluminum-copper and tungsten plugs on the lower levels. Note planarity of both structures. Copyright 1997 IEEE. Reprinted with permission from Refs [6].

SUMMARY

In this chapter, we have explored the various types of interconnections used in modern integrated circuits, including diffusion, polysilicon, and metal. Diffusion and polysilicon have a relatively high sheet resistance, which often restricts their use to local interconnections. The formation of metal silicides on the surface of polysilicon lines and diffusions can substantially reduce the sheet resistance of these interconnections.

Problems relating to the formation of good ohmic contacts between aluminum and silicon have also been discussed. An n^+ layer is required between aluminum and n -type silicon to prevent formation of a Schottky-barrier diode instead of an ohmic contact. Aluminum penetration into silicon is a serious problem in forming contacts to shallow junctions. Metals such as tungsten and titanium are often used as silicon diffusion barriers to prevent aluminum penetration into contacts to silicon or silicides.

At high current densities, a failure mechanism called *electromigration* can cause open and short circuits to form in the metallization layers. Aluminum containing approximately 1% silicon and 4% copper is used to minimize aluminum spiking and electromigration, respectively.

Multilevel metal processes have been developed for integrated circuits which require more than one level of metallization. Some of today's processes contain up to three levels of polysilicon, and others use six or more levels of metallization.

The first successful approaches to multilevel metallization covered aluminum with layers of planarizing oxide or polyimide. However the topology achieved with this approach has proven to be too rugged for deep submicron processes. Damascene processes combine electroplated copper deposition with chemical mechanical polishing to achieve highly planar multilevel metallization.

REFERENCES

- [1] M. Hansen and A. Anderko, *Constitution of Binary Alloys*, McGraw-Hill, New York, 1958.
- [2] M. Finetti, P. Ostojia, S. Solmi, and G. Soncini, "Aluminum-Silicon Ohmic Contact on 'Shallow' n^+ /p Junctions," *Solid-State Electronics*, 23, 255-262, March 1980.
- [3] S. Vaidya, D. B. Fraser, and A. K. Sinha, "Electromigration Resistance of Fine Line Al for VLSI Applications," *Proceedings of 18th IEEE Reliability Physics Symposium*, pp. 165-167, 1980.
- [4] S. P. Murarka, "Refractory Silicides for Integrated Circuits," *Journal of Vacuum Science and Technology*, 17, 775-792, July/August 1980.
- [5] P. S. Ho, "VLSI Interconnect Metallization—Part 3," *Semiconductor International*, 128-133, August 1985.
- [6] S. Venkatesen et al., "A High-Performance 1.8-V, 0.20- μ m CMOS Technology with Copper Metallization," *IEEE IEDM Technical Digest*, pp. 769-772, December 1997.
- [7] S. Yang et al., "A High-Performance 180 nm Generation Logic Technology," *IEEE IEDM Technical Digest*, pp. 197-200, December 1998.
- [8] P. C. Andricacos, U. Uzoh, J. O. Dukovic, J. Horkans, and H. Deligianni, "Damascene Copper Electroplating for Chip Interconnections," *IBM Journal of Research and Development*, pp. 567-574, September 1998.
- [9] M. Igarashi et al., "The Best Combination of Aluminum and Copper Interconnects for a High-Performance 0.18 μ m CMOS Logic Device," *IEEE IEDM Technical Digest*, pp. 829-832, December 1998.

- [10] D. Edelstein et al., "Full Copper Wiring in a sub-0.25 μm CMOS ULSI Technology," *IEEE IEDM Technical Digest*, pp. 773–776, December 1997.
- [11] Y. Matsubara et al., "Low-k Fluorinated Amorphous Carbon Interlayer Technology for Quarter Micron Devices," *IEEE IEDM Technical Digest*, pp. 369–372, December 1996.
- [12] E. M. Zielinski et al., "Damascene Integration of Copper and Ultra-low-k Xerogel for High-Performance Interconnects," *IEEE IEDM Technical Digest*, pp. 936–938, December 1997.
- [13] Y. Matsubara et al., "RC Delay Reduction of 0.18 μm CMOS Technology Using Low Dielectric Constant Fluorinated Amorphous Carbon," *IEEE IEDM Technical Digest*, pp. 841–844, December 1998.

FURTHER READING

- [1] P. B. Ghate, J. C. Blair, and C. R. Fuller, "Metallization in Microelectronics," *Thin Solid Films*, 45, 69–84, August 15, 1977.
- [2] G. L. Schnable and R. S. Keen, "Aluminum Metallization—Advantages and Limitations for Integrated Circuit Applications," *Proceedings of the IEEE*, 57, 1570–1580, September 1969.
- [3] J. Black, "Physics of Electromigration," *Proceedings of the 12th IEEE Reliability Physics Symposium*, p. 142, 1974.
- [4] C. Y. Ting, "Silicide for Contacts and Interconnects," *IEEE IEDM Technical Digest*, pp. 110–113, December 1984.
- [5] S. P. Murarka, "Recent Advances in Silicide Technology," *Solid-State Technology*, 28, 181–185, September 1985.
- [6] S. Sachdev and R. Castellano, "CVD Tungsten and Tungsten Silicide for VLSI Applications," *Semiconductor International*, pp. 306–310, May 1985.
- [7] P. B. Ghate, J. C. Blair, C. R. Fuller, and G. E. McGuire, "Application of Ti:W Barrier Metallization for Integrated Circuits," *Thin Solid Films*, 53, 117–128, September 1, 1978.
- [8] R. A. M. Wolters and A. J. M. Nellissen, "Properties of Reactive Sputtered TiW," *Solid-State Technology*, 29, 131–136, February 1986.
- [9] T. Sakurai and T. Serikawa, "Lift-Off Metallization of Sputtered Al Alloy Films," *Journal of the Electrochemical Society*, 126, 1257–1260, July 1979.
- [10] T. Batchelder, "A Simple Metal Lift-Off Process," *Solid-State Technology*, 25, 111–114, February 1982.
- [11] S. A. Evans, S. A. Morris, L. A. Arledge, Jr., J. O. Engle, and C. R. Fuller, "A 1- μm Bipolar VLSI Technology," *IEEE Transactions on Electron Devices*, ED-27, 1373–1379, August 1980.
- [12] Y. Sasaki, O. Ozawa, and S. Kameyama, "Application of MoSi_2 to the Double-Level Interconnections of I^2L Circuits," *IEEE Transactions on Electron Devices*, ED-27, 1385–1389, August 1980.
- [13] R. A. Larsen, "A Silicon and Aluminum Dynamic Memory Technology," *IBM Journal of Research and Development*, 24, 268–282, May 1980.
- [14] J. M. Mikkelsen, L. A. Hall, A. K. Malhotra, S. D. Secombe, and M. S. Wilson, "An NMOS VLSI Process for Fabrication of a 32b CPU Chip," *IEEE ISSCC Digest*, 24, 106–107, February 1981.
- [15] P. B. Ghate, "Multilevel Interconnection Technology," *IEEE IEDM Digest*, 126–129, December 1984.

PROBLEMS

- 7.1** (a) What is the sheet resistance of a 1- μm -thick aluminum-copper-silicon line with a resistivity of 3.2 $\mu\text{ohm-cm}$?
- (b) What would be the resistance of a line 500 μm long and 10 μm wide?
- (c) What is the capacitance of this line to the substrate if it is on an oxide which is 1 μm thick? (Assume that you can use the parallel-plate capacitance formula.)
- (d) What is the RC product associated with this 500- μm line?
- 7.2** (a) Repeat Problem 7.1 for a polysilicon line with a resistivity of 500 $\mu\text{ohm-cm}$.
- (b) Repeat Problem 7.1 for a titanium silicide line with a resistivity of 25 $\mu\text{ohm-cm}$.
- (c) Repeat Problem 7.1 for a copper line with a resistivity of 1.7 $\mu\text{ohm-cm}$.
- 7.3** (a) Compute estimates of the sheet resistance of shallow arsenic and boron diffusions by assuming uniformly doped rectangular regions with the maximum achievable electrically active impurity concentrations. (See Fig. 4.6.) Use hole and electron mobilities of 75 and 100 $\text{cm}^2/\text{V-sec}$, respectively, and a depth of 0.25 μm .
- (b) Compare your answers with those for diffused lines obtained from Figs. 4.6 and 4.16. Use the maximum possible electrically active concentration for the boron and arsenic surface concentrations.
- 7.4** Suppose that a $500 \times 15 \mu\text{m}$ aluminum line makes contact with silicon through a $10 \times 10 \mu\text{m}$ contact window as shown in Fig. P7.4. The aluminum is 1 μm thick and is annealed at 450 $^\circ\text{C}$ for 30 min. Assume that the silicon will saturate the aluminum up to a distance \sqrt{Dt} from the contact. D is the diffusion coefficient of silicon in aluminum which follows an Arrhenius relationship with $D = 0.04 \text{ cm}^2/\text{sec}$ and $E_A = 0.92 \text{ eV}$. Assume that silicon is absorbed uniformly through the contact and that the density of aluminum and silicon is the same. How deep will the aluminum penetration into the silicon be?

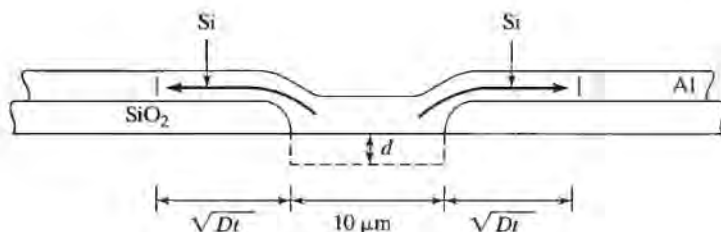


FIGURE P7.4

- 7.5 A certain process forms aluminum contacts to n^+ silicon through a $1 \times 1 \mu\text{m}$ contact window resulting in a contact resistance of 0.5 ohms.
- What is the specific contact resistivity for this contact?
 - What will the contact resistance be if the contact windows are reduced to $0.1 \times 0.1 \mu\text{m}$? Does this seem acceptable for a VLSI process?
- 7.6 Electromigration failures depend exponentially on temperature.
- What is the ratio of the MTFs of identical aluminum conductors operating at the same current density at 300 K and 400 K?
 - At 77 K, (liquid-nitrogen temperature) and 400 K? Use $E_A = 0.5 \text{ eV}$.
- 7.7 (a) What is the mean time to failure for the AlCu line in Fig. 7.8?
(b) How about for the copper line in the same figure?
- 7.8 An n^+ diffusion is used for interconnection. The surface concentration of the diffusion is $4 \times 10^{19}/\text{cm}^3$ and the junction depth is $4 \mu\text{m}$. The diffusion is formed in a p -type wafer with a background concentration of $1 \times 10^{15}/\text{cm}^3$.
- What is the sheet resistance of this diffusion?
 - Estimate the capacitance per unit length if the diffusion is $15 \mu\text{m}$ wide. Assume the rectangular geometry shown in Fig. P7.8 and use the step-junction-capacitance formula.
- 7.9 What is the maximum current that may be allowed to flow in an aluminum conductor $1 \mu\text{m}$ thick and $4 \mu\text{m}$ wide if the current density must not exceed $5 \times 10^5 \text{ A/cm}^2$?
- 7.10 What is the maximum current that may be permitted to flow in an aluminum conductor $0.25 \mu\text{m}$ thick and $0.5 \mu\text{m}$ wide if the current density cannot exceed 10^6 A/cm^2 ?
- 7.11 What is the resistance of a $0.25 \times 0.25 \mu\text{m}$ tungsten plug that is $1 \mu\text{m}$ high?

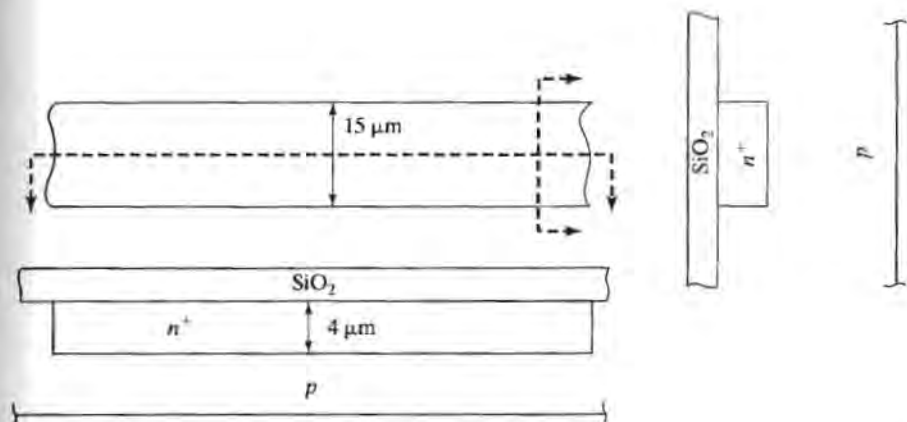


FIGURE P7.8

- 7.12** (a) What is the sheet resistance of a 0.5- μm -thick copper line with a resistivity of 1.7 $\mu\text{ohm-cm}$?
- (b) What is the resistance of a line that is 50 μm long and 0.5 μm wide?
- (c) What is the capacitance of this line to the substrate if a "low-K" dielectric is utilized with $\epsilon = 2\epsilon_0$?
- (d) What is the RC product associated with the 50- μm line?

CHAPTER 8

Packaging and Yield

The low cost normally associated with integrated circuits results from mass production in which many wafers, each containing a large number of IC dice, are all processed together. There may be tens to thousands of dice per wafer and 25 to 200 wafers per lot. After wafer processing is completed, however, the dice must be tested, separated, and assembled in packages that are easy to handle and to mount in electronic systems. The testing and assembly operations substantially increase the cost of the final product.

In this chapter, we first present an overview of testing and die separation. Then we discuss IC assembly, including die attachment, wire bonding, and a survey of the various types of packages used with integrated circuits.

The ultimate cost of the integrated circuit is related to the total yield of assembled and tested devices. In the early stages in the development of a new process or circuit, we are lucky if a few functional dice are found per wafer. Late in the life of a process with a mature circuit design, yields of 60 to 80% are not uncommon. A discussion of the dependence of yield on defect density and die size concludes this chapter.

8.1 TESTING

Following metallization and passivation-layer processing, each die on the wafer is tested for functionality. Special parametric test dice are placed at a number of sites on the wafer. At this stage, dc tests are used to verify that basic process parameters fall within acceptable limits. To perform the tests, a probe station lowers a ring of very fine, needle-sharp probes into contact with the pads on the test die. Test equipment is connected to the circuit through the probes and controlled by a computer system. If the wafer-screening operation shows that basic process and device parameters are within specification, functional testing of each die begins.

Under computer control, the probe station automatically steps across the wafer, performing functional testing at each die site. Defective dice are marked with a drop of ink. Later, when the dice are separated from the wafer, any die with an ink spot is discarded. It has become impossible to exhaustively test complex VLSI devices such as

microprocessors. Instead, a great deal of computer time is used to find a minimum sequence of tests that can be used to indicate die functionality. At the wafer-probe stage, functional testing is primarily static in nature. High-speed dynamic testing is difficult to do through the probes, so parametric speed tests are usually performed after die packaging is complete.

The ratio of functional dice to total dice on the wafer gives the *yield* for each wafer. Yield is directly related to the ultimate cost of the completed integrated circuit and will be discussed more fully in Section 8.7.

8.2 WAFER THINNING AND DIE SEPARATION

As mentioned in Chapter 1, wafers range from approximately 350–1250 microns in thickness. Large-diameter wafers must be thicker in order to maintain structural integrity and planarity during the wide range of processing steps encountered during IC fabrication. However, many applications require much thinner dice; a thickness of approximately 275 μm is commonly used by many manufacturers. Flash memory cards, credit cards, and medical electronics are just three examples of applications that use dice that have been thinned to only 125 μm . The thinning is done using back lapping and polishing processes similar to those for CMP described in Chapter 6.

Following initial functional testing, individual IC dice must be separated from the wafer. In one method used for many years, the wafer is mounted on a holder and automatically scribed in both the x - and y -directions using a diamond-tipped scribe. Scribing borders of 75 to 250 μm are formed around the periphery of the dice during fabrication. These borders are left free of oxide and metal and are aligned with crystal planes if possible. In $\langle 100 \rangle$ wafers, natural cleavage planes exist perpendicular to the surface of the wafer in directions both parallel and perpendicular to the wafer flat. For $\langle 111 \rangle$ wafers, a vertical cleavage plane runs parallel, but not perpendicular, to the wafer flat. This can lead to some separation and handling problems with $\langle 111 \rangle$ material.

Following scribing, the wafer is removed from the holder and placed upside down on a soft support. A roller applies pressure to the wafer, causing it to fracture along the scribe lines. Care is taken to ensure that the wafer cracks along the scribe lines to minimize die damage during separation, but there will always be some damage and loss of yield during the scribing and breaking steps.

Today, diamond saws are used for die separation. A wafer is placed in a holder on a sticky sheet of Mylar as shown in Fig. 8.1. The saw can be used either to scribe the wafer or to cut completely through the wafer. Following separation, the dice remain attached to the Mylar film.

8.3 DIE ATTACHMENT

Visual inspection is used to sort out dice that may have been damaged during die separation, and the inked dice are also discarded. The next step in the assembly process is to mount the good dice in packages.

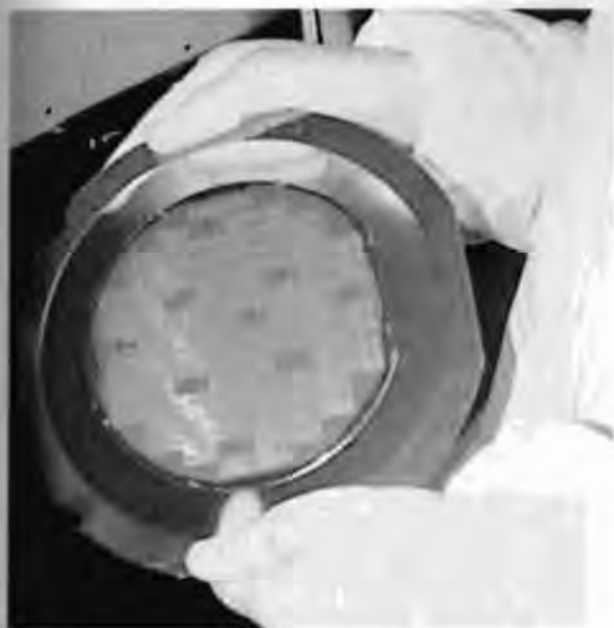


FIGURE 8.1

150-mm wafer mounted and ready for the dicing saw.

8.3.1 Epoxy Die Attachment

An epoxy cement may be used to attach the die to a package or “header.” However, epoxy is a poor thermal conductor and an electrical insulator. Alumina can be mixed with the epoxy to increase its thermal conductivity, and gold- or silver-filled epoxies are used to reduce the thermal resistance of the epoxy bonding material and to provide a low-resistance electrical connection between the die and the package.

8.3.2 Eutectic Die Attachment

The gold-silicon eutectic point occurs at a temperature of 370 °C for a mixture of approximately 3.6% Si and 96.4% Au. Gold can be deposited on the back of the wafer prior to die separation or can be in the form of a thin alloy “preform” placed between the die and package. To form a eutectic bond, the die and package are heated to 390 to 420 °C, and pressure is applied to the die in conjunction with an ultrasonic scrubbing motion. Eutectic bonding is possible with a number of other metal-alloy systems, including gold-tin and gold-germanium. A solder attachment technique is also used with semiconductor power devices.

8.4 WIRE BONDING

Wire bonding is the most widely used method for making electrical connections between the die and the package. The bonding areas on the die are large, square pads 100 to 125 μm on a side, located around the periphery of the die. Figure 8.2 shows a 2.5 \times 2.5 mm die with 100 μm wire bond pads on a 125 μm pitch. Many manufacturers do not put pads in the corners, because of reliability concerns. In particular, packaging-induced die stress and die cracking tend to be highest in the corners.

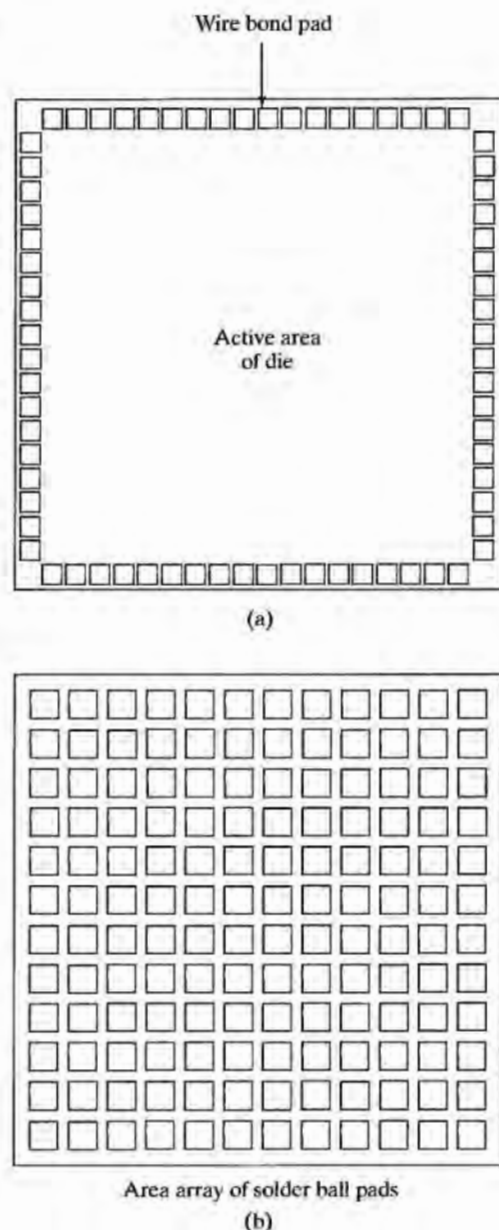
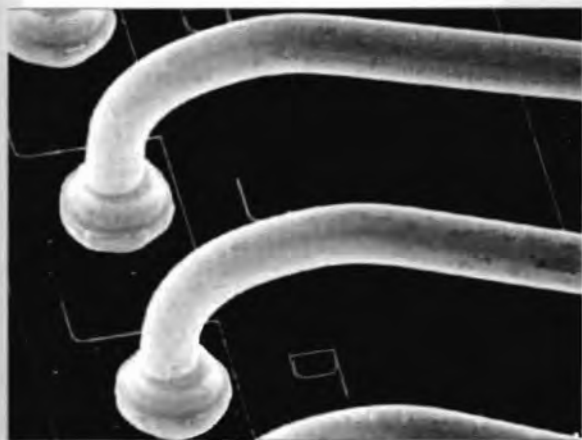


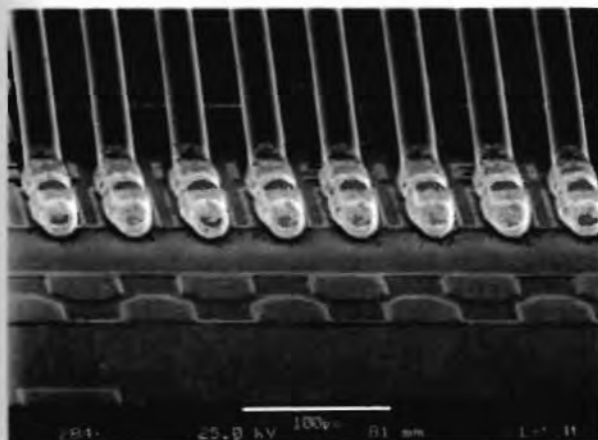
FIGURE 8.2

Die layout for (a) peripheral pads and (b) area array of pads.

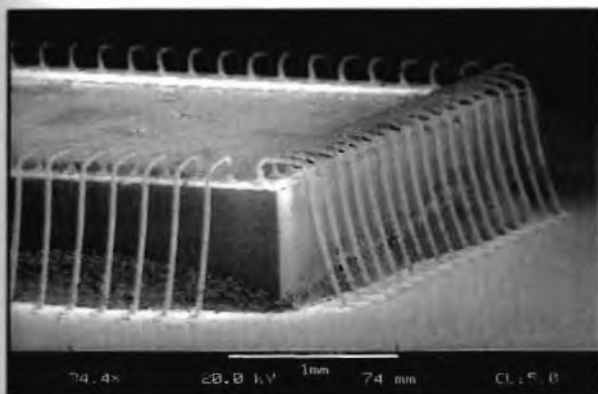
Thermocompression bonding was originally used with gold wire, and ultrasonic bonding is used with aluminum wire. A combination of the two is termed thermosonic bonding. Present high-density wire bonding can achieve a 50–70 μm pitch (see Fig. 8.3), and two or more staggered rows of wire bond pads are sometimes used to increase the number of possible pads/die. The ITRS goal is to reach a wire bond pitch of 40 μm by the year 2010. Fine wires interconnect the aluminum bonding pads on the IC die to the leads of the package.



(a)



(b)



(c)

FIGURE 8.3

(a) An SEM micrograph of gold ball bonding (b) SEM of high density gold ball bonding (c) SEM micrograph of wire bonded die. Courtesy of Kulicke and Soffa Industries, Inc. (K&S).

8.4.1 Thermocompression Bonding

The thermocompression bonding technique, also called *nail-head* or *ball bonding*, uses a combination of pressure and temperature to weld a fine gold wire to the aluminum bonding pads on the die and the gold-plated leads of the package. Figure 8.4 shows the steps used in forming either a thermocompression or thermosonic bond.

A fine gold wire, 15 to 75 μm in diameter, is fed from a spool through a heated capillary. A small hydrogen torch or electric spark melts the end of the wire, forming a gold ball two to three times the diameter of the wire. Under either manual or computer control, the ball is positioned over the bonding pad, the capillary is lowered, and the ball deforms into a "nail head" as a result of the pressure and heat from the capillary.

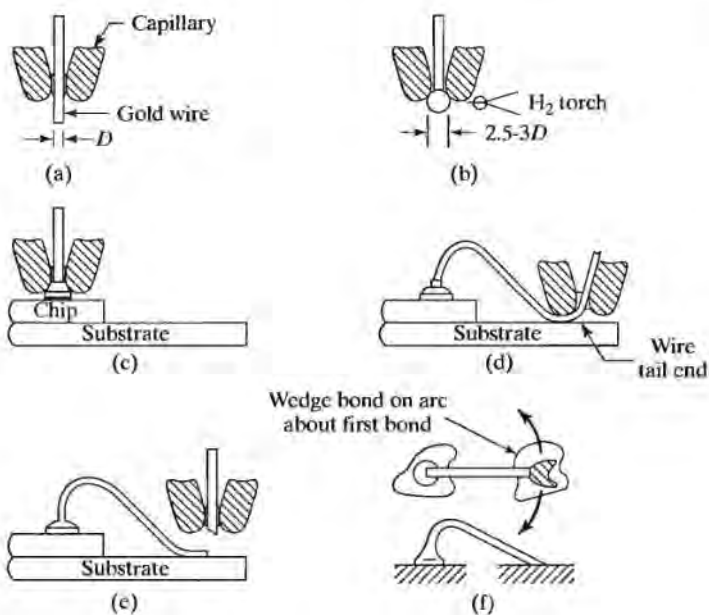
Next, the capillary is raised, and wire is fed from the spool as the tool is moved into position over the package. The second bond is a wedge bond produced by deforming the wire with the edge of the capillary. After the formation of the second bond, the capillary is raised and the wire is broken near the edge of the bond. Because of the symmetry of the bonding head, the bonder may move in any direction following formation of the nail-head bond. An SEM micrograph of a gold-ball bond is shown in Fig. 8.3(a).

A problem encountered in production of gold-aluminum bonds is formation of the "purple plague." Gold and aluminum react to form intermetallic compounds. One such compound, AuAl_2 , is purple in color, and its appearance has been associated with faulty bonds. However, this compound is highly conductive. The actual culprit is a low-conductivity, tan-colored compound, Au_2Al , which is also present.

During thermocompression bonding, the substrate is maintained at a temperature between 150 and 200 $^{\circ}\text{C}$. The temperature at the bonding interface ranges from 280 to 350 $^{\circ}\text{C}$, and significant formation of the Au-Al compounds can occur at these

FIGURE 8.4

Thermosonic ball-wedge bonding of a gold wire. (a) Gold wire in a capillary; (b) ball formation accomplished by passing a hydrogen torch over the end of a gold wire or by capacitance discharge; (c) bonding accomplished by simultaneously applying a vertical load on the ball while ultrasonically exciting the wire (the chip and substrate are heated to about 150 $^{\circ}\text{C}$); (d) a wire loop and a wedge bond ready to be formed; (e) the wire is broken at the wedge bond; (f) the geometry of the ball-wedge bond that allows high-speed bonding. Because the wedge can be on an arc from the ball, the bond head or package table does not have to rotate to form the wedge bond. Reprinted with permission from *Semiconductor International* magazine, May 1982 [1]. Copyright 1982 by Cahners Publishing Co., Des Plaines, IL.



temperatures. Limiting the die temperature during the bonding process helps to prevent the formation of the intermetallic compounds, and high-temperature processing and storage following bonding should be avoided.

In addition to the foregoing problem, many epoxy materials cannot withstand the temperatures encountered in thermocompression bonding, and thermocompression bonding has been replaced by the ultrasonic and thermosonic bonding techniques discussed in the next two sections.

8.4.2 Ultrasonic Bonding

Oxidation of aluminum wire at high temperatures makes it difficult to form a good ball at the end of the wire. An alternative process called *ultrasonic bonding*, which forms the bond through a combination of pressure and rapid mechanical vibration, is used with aluminum wire. Aluminum wire is fed from a spool through a hole in the ultrasonic bonding tool, as shown in Fig. 8.5. To form a bond, the bonding tool is lowered over the bonding position, and ultrasonic vibration at 20 to 60 kHz causes the metal to deform and flow easily under pressure at room temperature. Vibration also breaks through the oxide film that is always present on aluminum and results in formation of a clean, strong weld.

As the tool is raised after forming the second bond, a clamp is engaged and pulls and breaks the wire at a weak point just beyond the bond. To maintain proper wire alignment in the bonding tool, the ultrasonic bonder can move only in a front-to-back motion between the first and second bonds, and the package must be rotated 90° to permit complete bonding of rectangular dice.

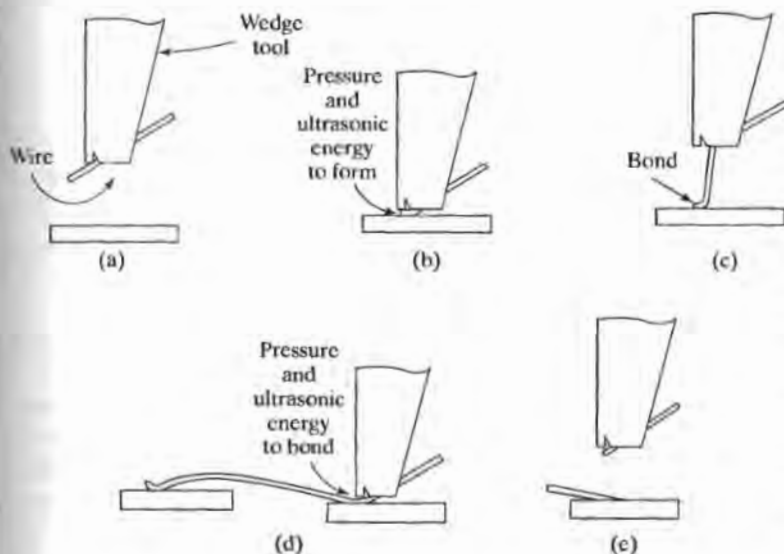


FIGURE 8.5

(a) In ultrasonic bonding, the tool guides wire to the package terminal; (b) pressure and ultrasonic energy form the bond; (c) and (d) the tool feeds out wire and repositions itself above the IC chip. The tool lowers and ultrasonically forms the second bond; (e) tool lifts, breaking the wire at the bond. Reprinted with permission from *Circuits Manufacturing*, January 1980. Copyright Morgan-Grampian 1980.

8.4.3 Thermosonic Bonding

Thermosonic bonding combines the best properties of ultrasonic and thermocompression bonding. The bonding procedure is the same as in thermocompression bonding, except that the substrate is maintained at a temperature of approximately 150 °C. Ultrasonic vibration causes the metal to flow under pressure and form a strong weld. The symmetrical bonding tool permits movement in any direction following the nail-head bond. Thermosonic bonding can be easily automated, and computer-controlled thermosonic bonders can produce 5 to 10 bonds per second.

8.5 PACKAGES

IC dice can be mounted in a wide array of packages. In this section, we discuss the round TO-style packages, dual-in-line packages (DIP), the pin-grid array (PGA), the leadless chip carrier (LCC), and packages used for surface mounting. Later in this chapter, we will look at flip-chip mounting, ball-grid arrays, and tape-automated bonding. The long-term trend is to attempt to achieve a package footprint that is similar in size to the semiconductor die itself. These packages are referred to as Chip Scale Packages (CSPs).

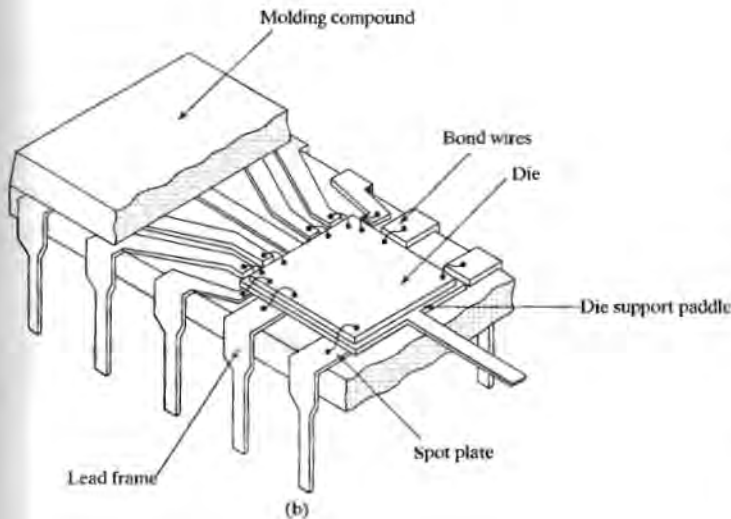
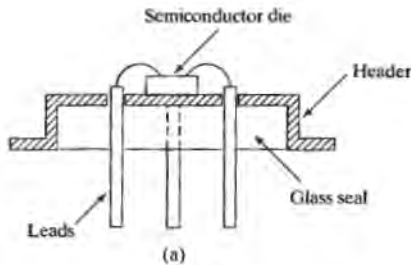
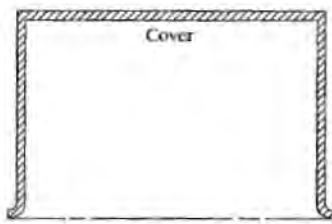
8.5.1 Circular TO-Style Packages

Figure 8.6(a) shows a round TO-type package that was one of the earliest IC packages. Different configurations of this package are available with 4 to 48 pins. The silicon die is attached to the center of the gold-plated header. Wire bonds connect the pads on the die to Kovar lead posts that protrude through the header and glass seal. Kovar is an iron–nickel alloy designed to have the same coefficient of thermal expansion as the glass seal. A metal cap is welded in place after die attachment and wire bonding.

8.5.2 DIPs

The DIP shown in Fig. 8.6(b) is extremely popular, because of its low cost and ease of use. Plastic and epoxy DIPs are the least expensive packages and are available with as few as 4 leads to more than 80 leads. In the postmolded DIP, a silicon die is first mounted on and wire-bonded to a metallic lead frame. Epoxy is then molded around both the die and the frame. As a result, the silicon die becomes an integral part of the package. Plastic DIPs have evolved into other forms of in-line packages, including single-in-line packages (SIPs) and zig-zag-in-line packages (ZIPs) shown in Fig. 8.6(c), as well as many surface-mount packages to be described in Section 8.5.5.

In a ceramic DIP, the die is mounted in a cavity on a gold-plated ceramic substrate and wire bonded to gold-plated Kovar leads. A ceramic or metal lid is used to seal the top of the cavity. Ceramic packages are considerably more expensive than plastic and are designed for use over a wider temperature range. In addition, ceramic packages may be hermetically sealed. A premolded plastic package similar to the ceramic package is also available.



Dual in-line package (DIP)



Single in-line package (SIP)



Zig-zag in-line package (ZIP)

(c)

FIGURE 8.6

(a) TO-style package; (b) ball- and wedge-bonded silicon die in a plastic DIP. The die support paddle may be connected to one of the external leads. For most commercial products, only the die paddle and the wedge-bond pads are selectively plated. The external leads are solder-plated or dipped after package molding. Copyright 1981, IEEE. Reprinted with permission from Ref. [3]. (c) DIP, SIP, and ZIP packages.

8.5.3 Pin-Grid Arrays (PGAs)

The DIP package is satisfactory for packaging IC dice with up to approximately 80 pins. The *pin-grid array* of Fig. 8.7 provides a package with a much higher pin density than that of the DIP package. The pins are placed in a regular x - y array, and the package can have hundreds of pins. Wire bonding is still used to connect the die to gold interconnection lines, which fan out to the array of pins. Other versions of this package use the flip-chip process (which will be discussed in Section 8.6).



FIGURE 8.7

An example of a PGA in which the chip faces upward in the cavity.

8.5.4 Leadless Chip Carriers (LCCs)

Figure 8.8 shows two types of leadless chip carriers. In each case, the die is mounted in a cavity in the middle of the package. Connections are made between the package and die using wire bonding, and the cavity is sealed with a cap of metal, ceramic, or epoxy. The package in Fig. 8.8(a) has contact pads only on the top surface of the chip carrier. The chip carrier is pressed tightly against contact fingers in a socket mounted on a printed-circuit board. Another type of LCC is shown in Fig. 8.8(b). Conductors are formed in grooves in the edges of the chip carrier and are again pressed tightly against contact pins in a socket on the next level of packaging.

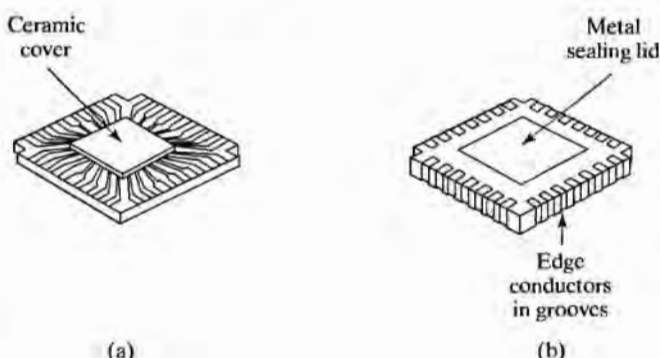
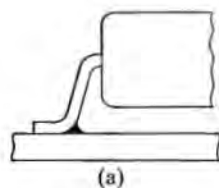


FIGURE 8.8

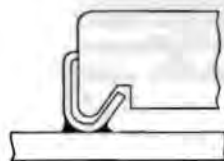
(a) Ceramic leadless chip carriers with top connections; (b) LCC with edge connections in grooves on the sides of the package.

8.5.5 Packages for Surface Mounting

The TO-, DIP-, and PGA-style packages are made for mounting in holes fabricated in printed-circuit boards. Surface-mount packages do not require holes. The gull-wing package shown in Fig. 8.9(a) has short-lead stubs bent away from the package, whereas the leads of the J-style package of Fig. 8.9(b) are bent back underneath the package. Both styles permit soldering of the package directly to the surface of a printed-circuit board or hybrid package. Several versions are shown in Fig. 8.9(c)–(h). The leadless chip carriers described in the previous section are also available with leads added for surface mounting.



(a)



(b)



(c)

Small outline transistor (SOT)



(d)

Small outline integrated circuit (SOIC)



(e)

Thin small outline package (TSOP)



(f)

Small outline J-leaded (SOJ)



(g)

Quad flat pack (QFP)



(h)

Plastic-leaded chip carrier (PLCC)/
Quad J-leaded pack (QJLP)

FIGURE 8.9

(a) Gull-wing and (b) J-lead surface-mount (c) Small outline transistor (d) Small outline IC (e) Thin small outline (f) Small outline J-leaded (g) Quad flat pack (h) Plastic-leaded chip carrier or Quad J-leaded packages.

8.6 FLIP-CHIP AND TAPE-AUTOMATED-BONDING PROCESSES

As can be envisioned from the preceding discussion, the die-mounting and wire-bonding processes involve a large number of manual operations and are therefore quite expensive. In fact, the cost of assembly and test may be many times the cost of a small die. The one-at-a-time nature of the wire-bonding process also leads to reduced reliability, and failure of wire bonds is one of the most common reliability problems in integrated circuits. The flip-chip and tape-automated-bonding processes were developed to permit batch fabrication of die-to-package interconnections.

8.6.1 Flip-Chip Technology

The *flip-chip* mounting process was developed at IBM during the 1960s [4]. It took almost three decades for the semiconductor industry to adopt flip-chip technology. For today and the foreseeable future, however, flip-chip technology is an integral part of the ITRS plan for high-density chip interconnection. The first step is to form a solder ball (see Fig. 8.10) on top of each bonding pad. A sandwich of Cr, Cu, and Au is sequentially evaporated through a mask to form a cap over each of the aluminum bonding pads. Chrome and copper provide a barrier and a good contact to the aluminum pad. Gold adheres well to chrome and acts as an oxidation barrier prior to solder deposition. Lead-tin solder is evaporated through a mask onto the Au-Cu-Cr cap, occupying an area slightly larger than the cap. The die is heated, causing the solder to recede from the oxide surface and form a solder ball on top of the Au-Cu-Cr bonding-pad cap.

After being tested and separated, the dice are placed face down on a substrate. Temperature is increased, causing the solder to reflow, and the die is bonded directly to the interconnections on the substrate. Solder balls provide functions of both electrical interconnection and die attachment. Hundreds of bonds can be formed simultaneously using this technique, and bonding pads may be placed anywhere on the surface of the die, rather than just around the edge. In addition, the bond between the die and the substrate is very short. The main disadvantages of this technique are the additional

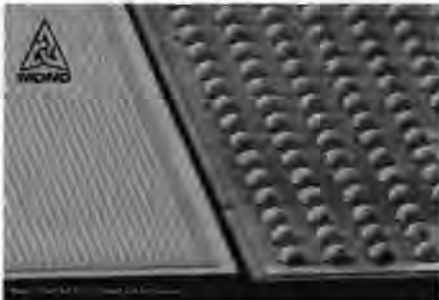
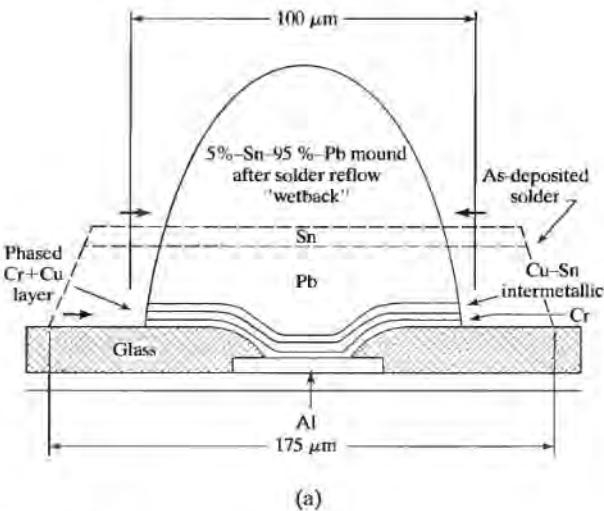


FIGURE 8.10
(a) Cross section through a solder ball before and after reflowing. Copyright 1969 by International Business Machines Corporation. Reprinted with permission from Ref. [4]. (b) Flip-chip Pb/Sn solder bumps in standard 250 μm pitch (right) and 50 μm . Courtesy of MCNC Optical and Electronic Packaging Group.

processing complexity, the higher thermal resistance between die and substrate, and the inability to visually inspect the completed solder joints. To enhance reliability, an under-fill material is used in the gap between the chip and substrate. The under-fill material is usually dispensed in liquid form and then cured. To be effective, the under-fill must achieve good adhesion at the interfaces with both the die and the substrate.

Two forms of solder-ball footprints are in common use at the die level. The first simply replaces wire bonds with solder balls around the periphery of the die. However, the real power of the flip-chip approach is achieved through the use of large arrays of solder balls placed over the full die area. For example, the die in Fig. 8.2(a) has only 72 pads, but an area array would achieve 144 interconnections using 125- μm pads on a 200- μm pitch as shown in Fig. 8.2(b).

The original IBM technique was a difficult and expensive process to implement, and these factors impeded its widespread adoption. More recently, a screen-printing approach has been used to deposit solder paste to form solder balls for low-volume applications. Gold stud bumping, shown in the photograph in Fig. 8.11, is yet another approach in which a gold wire is bonded to the chip pad, but the wire is removed after the first bond, leaving the gold bump on the pad. These chips may then be flipped over and mounted to form a flip-chip structure.



FIGURE 8.11

Bumps formed by modification of the wire-bond process. Courtesy of Kulicke and Soffa Industries, Inc. (K&S).

8.6.2 Ball-Grid Array (BGA)

Ball-grid arrays essentially apply the solder-ball approach to the package rather than to the chip as indicated in Fig. 8.12. Chips are mounted to the BGA substrate using either peripheral wire bonding or flip-chip technology, and then the chip is coated with some form of molding compound. Standard BGA array ball pitches are in the range of 1270 μm (50 mils). Fine-pitch ball-grid arrays (FBGA, or μBGA) are projected to decrease from a pitch of 500 μm in the year 2000 to 300 μm by 2010.

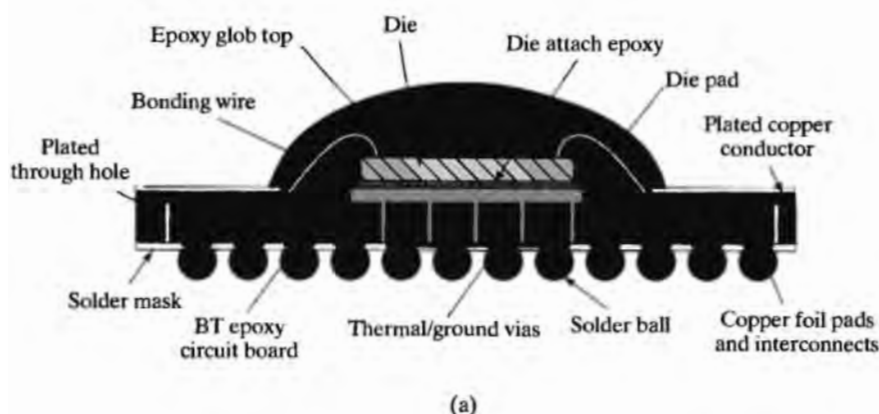


FIGURE 8.12

(a) Ball grid array cross section. (b) Intel microprocessor packaged using a BGA. Courtesy of Intel Corp.

8.6.3 The Tape-Automated-Bonding (TAB) Process

In *tape-automated bonding*, dice are attached to copper leads supported by a tape similar to 35-mm film. The film is initially coated with copper, and the leads are defined by lithography and etching. The lead pattern may contain hundreds of connections.

Die attachment requires a process similar in concept to the solder-ball technology discussed earlier. Gold bumps are formed on either the die or the tape and are used to bond the die to the leads on the tape. Figure 8.13 outlines the steps used to form a gold bump on a bonding pad [5]. A multilayer metal sandwich is deposited over the passivation oxide. Next, a relatively thick layer of photoresist is deposited, and windows are opened above the bonding pads. Electroplating is used to fill the openings with gold. The photoresist is removed, and the thin metal sandwich is etched away using wet or dry etching. The final result is a 25- μm -high gold bump standing above each pad. As in the flip-chip approach, bonding sites may be anywhere on the die.

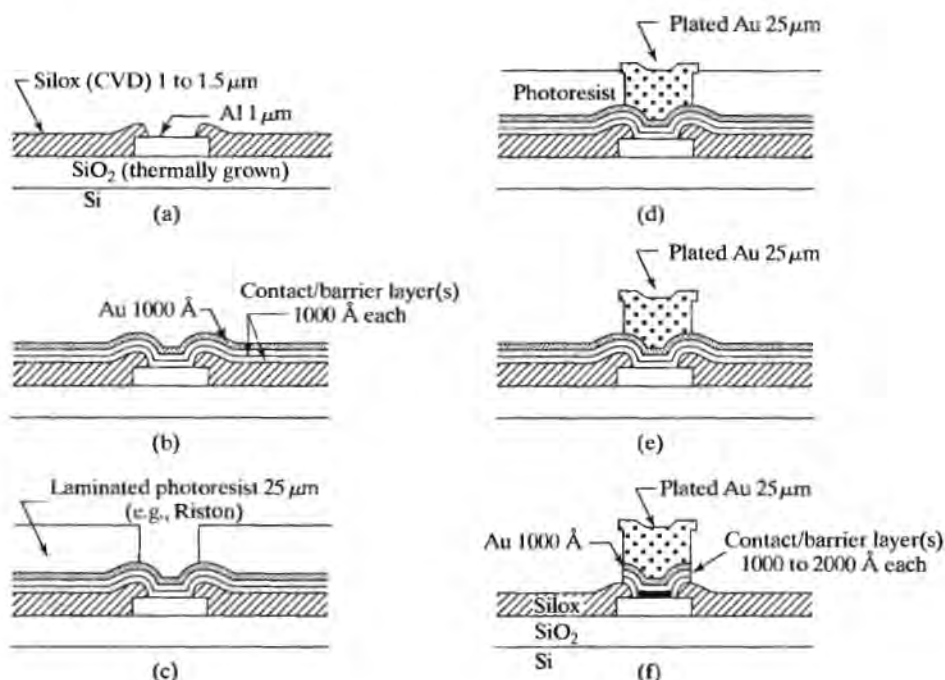


FIGURE 8.13

Process sequence for making gold bumps on aluminum metallurgy devices. (a) The wafer is cleaned and sputter-etched; (b) a contact/barrier layer (which also serves as a conductive film for electroplating) is sputter-deposited with a layer of gold for oxidation protection; (c) a thick-film photoresist (25 μm) is laminated and developed; (d) gold is electroplated to a height of approximately 25 μm to form the bumps; (e) the resist is stripped; (f) the sputter-deposited conductive film is removed chemically or by back-sputtering. Reprinted with permission of *Solid State Technology*, published by Technical Publishing, a company of Dun & Bradstreet, from Ref. [5].

The mounting process aligns the tape over the die, as in Fig. 8.14. A heated bonding head presses the tape against the die, forming thermocompression bonds. In a production process, a new die is brought under the bonding head and the tape indexes automatically to the next lead site.

TAB-mounted parts offer the advantage that they can be functionally tested and burned in once the dice are attached to the film. In addition, the IC passivation layer and gold bump completely seal the semiconductor surface.

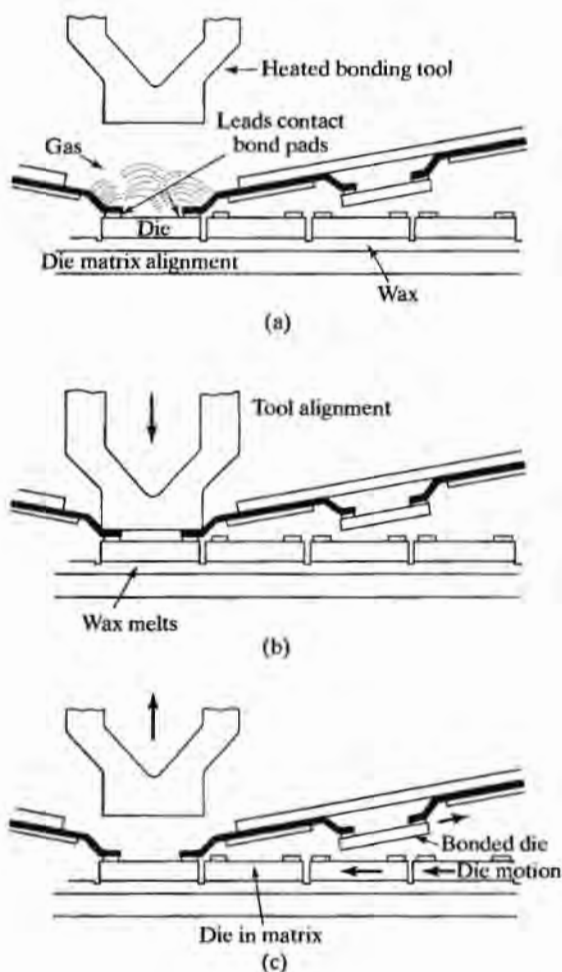


FIGURE 8.14

Tape-automated-bonding procedure. (a) Preformed leads of film are lowered into position and aligned above bonding pads on the die, which is held in place with a wax; (b) bonding tool descends and forms bond with pressure and heat; heat melts the wax, releasing the die; (c) tool and film are raised, lifting the bonded die clear so a new die can be moved into position and the process can be repeated. Reprinted with permission from *Small Precision Tools Bonding Handbook*. Copyright 1976.

8.6.4 Chip Scale Packages

Many of the technologies described, including wire bonding, TAB, and flip-chip mounting, are evolving to the Chip Scale Package (CSP) in which the goal is to achieve a package whose area is no larger than the die itself. Figures 8.15(a–b) present drawings of two possible approaches to the CSP, one based upon wire-bonded die and the second employing a form of TAB. Another approach referred to as chip-on-board attempts to eliminate the package altogether by mounting the bare die directly on the printed circuit board or flexible substrate as shown in Fig. 8.15(c). Either wire bonding or flip-chip mounting may be used, and the final structure is encapsulated with an epoxy glob-top coating.

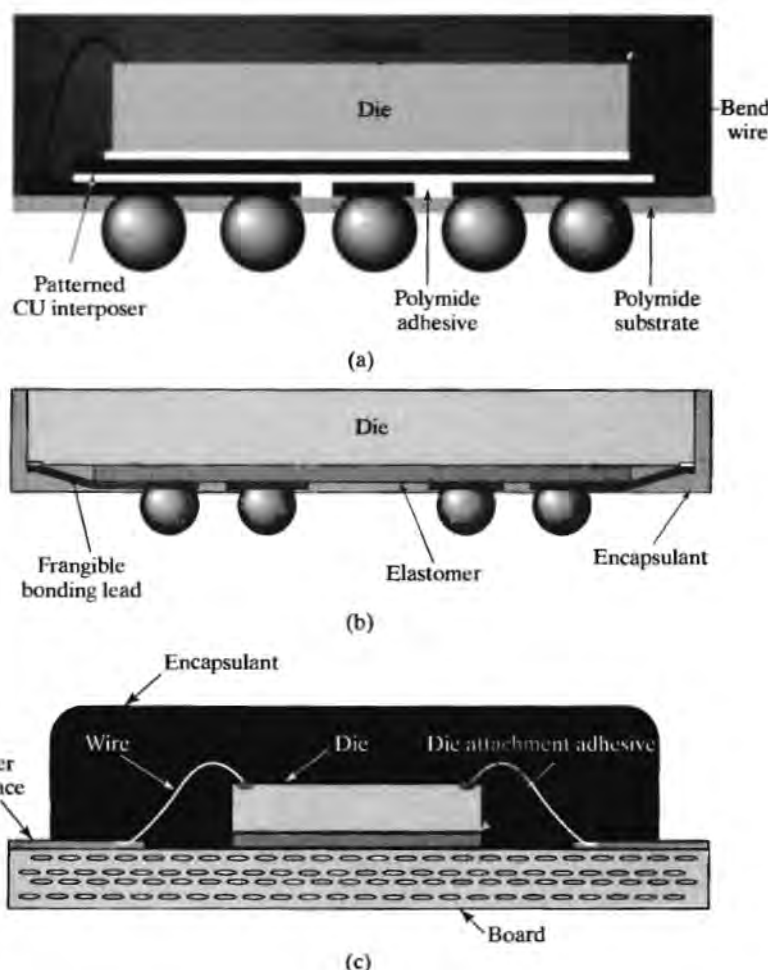


FIGURE 8.15

(a) Chip scale package using wire bonding. (b) Alternate form of a CSP. (c) Chip-on-board packaging.

8.7 YIELD

The manufacturer of integrated circuits is ultimately interested in how many finished chips will be available for sale. A substantial fraction of the dice on a given wafer will not be functional when they are tested at the wafer-probe step at the end of the process. Additional dice will be lost during the die separation and packaging operations, and a number of the packaged devices will fail final testing.

As mentioned earlier, the cost of packaging and testing is substantial and may be the dominant factor in the manufacturing cost of small die. For a large die with low yield, the manufacturing cost will be dominated by the wafer processing cost. A great deal of time has been spent attempting to model wafer yield associated with IC processes. Wafer yield is related to the complexity of the process and is strongly dependent on the area of the IC die.

8.7.1 Uniform Defect Densities

One can visualize how die area affects yield by looking at the wafer in Fig. 8.16, which has 120 die sites. The dots represent randomly distributed defects that have caused a die to fail testing at the wafer-probe step. In Fig. 8.16(a), there are 52 good dice out of the total of 120, giving a yield of 43%. If the die size were twice as large, as in Fig. 8.16(b), the yield would be reduced to 22% for this particular wafer.

An estimate of the yield of good dice can be found from a classical problem in probability theory in which n defects are randomly placed in N die sites. The probability P_k that a given die site contains exactly k defects is given by the binomial distribution:

$$P_k = \frac{n!}{k!(n-k)!} N^{-n} (N-1)^{n-k} \quad (8.1)$$

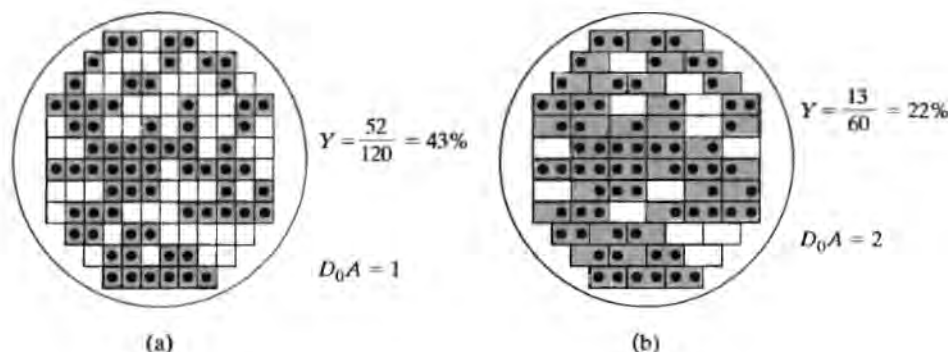


FIGURE 8.16

Illustrations of wafers, showing effect of die size on yield. Dots indicate the presence of a defective die location. (a) For a particular die size the yield is 43%; (b) if the die size were doubled, the yield would be only 22%.

For large n and N , Eq. (8.1) can be approximated by the Poisson distribution:

$$P_k = \frac{\lambda^k}{k!} \exp(-\lambda) \quad (8.2)$$

where $\lambda = n/N$. The yield is given by the probability that a die is found with no defects:

$$Y = P_0 = \exp(-\lambda) \quad (8.3)$$

The area of the wafer is equal to NA , where A is the area of one die. The density of defects, D_0 , is given by the total number of defects, n , divided by the area of all the chips, and the average number of defects per die, λ , is given by

$$\lambda = n/N = D_0 A, \quad \text{for } D_0 = n/NA \quad (8.4)$$

The yield based on the Poisson distribution then becomes

$$Y = \exp(-D_0 A) \quad (8.5)$$

This expression was used to predict early die yield, but was found to give too low an estimate for large dice with $D_0 A > 1$. Equation (8.5) implicitly assumes that the defect distribution is uniform across a given wafer and does not vary from wafer to wafer. However, it was quickly realized that these conditions are not realistic. Defect densities vary from wafer to wafer because of differences in handling and processing. On a given wafer, there are usually more defects around the edge of the wafer than in the center, and the defects tend to be found in clusters. These realizations led to investigation of nonuniform defect densities.

8.7.2 Nonuniform Defect Densities

Murphy [6] showed that the wafer yield for a nonuniform defect distribution can be calculated from

$$Y = \int_0^{\infty} \exp(-DA) f(D) dD \quad (8.6)$$

where $f(D)$ is the probability density for D . He considered several possible distributions, as shown in Fig. 8.17. The impulse function in Fig. 8.17(a) represents the case in which the defect density is the same everywhere, and substituting it for $f(D)$ in Eq. (8.6) yields Eq. (8.5). The triangular distribution in Fig. 8.17(b) is a simple approximation to a Gaussian distribution function and allows some wafers to have very few defects and others to have up to $2D_0$ defects. Application of Eq. (8.6) results in the following yield expression:

$$Y = \left[\frac{1 - \exp(-D_0 A)}{D_0 A} \right]^2 \quad (8.7)$$

A uniform distribution of defect densities is modeled by $f(D)$ in Fig. 8.17(c) and predicts a yield of

$$Y = \left[\frac{1 - \exp(-2D_0 A)}{2D_0 A} \right] \quad (8.8)$$

More complicated probability distributions have also been investigated, including the negative binomial and gamma distributions [7,8]. These result in the yield expression in Eq. (8.9)

$$Y = \left[1 + \frac{D_0 A}{\alpha} \right]^{-\alpha} \quad (8.9)$$

in which α represents a clustering parameter which ranges from 0.5 to 10. The 1999 ITRS is basing future defect density requirements on Eq. (8.9) with a clustering parameter of 5 [9]. To achieve a yield loss of less than 20% due to random defects, the ITRS projects the required defect density to be less than $0.1/\text{cm}^2$ by the year 2010 with a critical defect size of less than 30 nm.

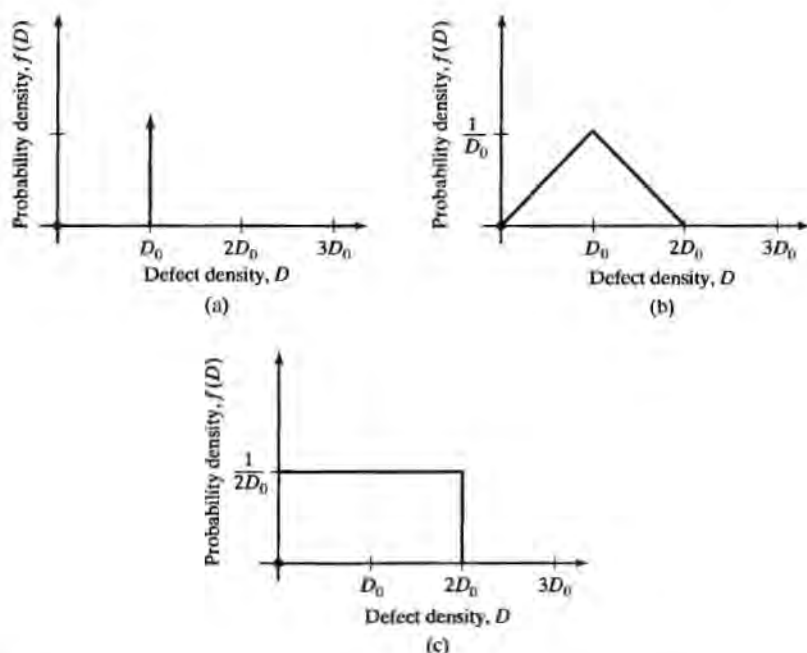


FIGURE 8.17

Possible defect probability density functions. (a) Impulse, where every wafer has exactly the same number of defects; (b) a triangular approximation to a Gaussian density; (c) a uniform density function.

Figure 8.18 plots the various yield functions versus D_0A , the average number of defects in a die of area A . Early yield estimates based on Poisson statistics are clearly much more pessimistic than those based on the other functions. However, the negative binomial yield model in Eq. (8.9) with $\alpha = 5$ is beginning to approach the exponential function. (See Problems 8.9 and 8.10.)

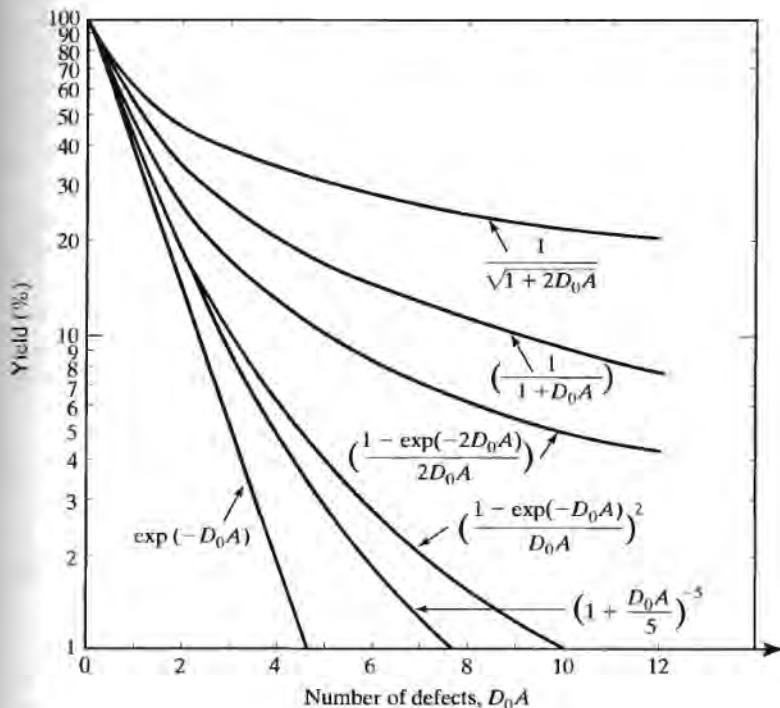


FIGURE 8.18

Theoretical yield curves for different defect densities. See Eqs. (8.5) through (8.9).

Example 8.1

A 150-mm wafer has a defect density of 10 defects/cm², and costs \$200 to process. The cost of assembly and testing is \$1.50 per die. (a) What is the total manufacturing cost for a 5 × 5 mm die in this process based on yield Eq. (8.7)? (The number of square dice per wafer is given approximately by $N = \pi(R - S)^2/S^2$, where R is the wafer radius and S is the length of the side of the die.) (b) The market price for this part is \$2.50. What must be the wafer yield needed for manufacturing cost to drop below the market price?

Solution: The area of the die is 0.25 cm, so the average number of defects per die is $D_0A = 2.5$. Equation (8.7) predicts a yield of 13.5%. The wafer has a radius of 75 mm and contains approximately 616 dice. So the average wafer will yield 83 good dice. The cost of the packaged dice will be $C = (\$200/83) + \$1.50 = \$3.91$. Getting the cost to the market price requires $\$2.50 = (\$200/ND) + \$1.50$. We must get $ND = 200$ good dice per wafer to break even, corresponding to a yield of $Y = 200/616 = 0.325$ or more.

SUMMARY

Following the completion of processing, wafers are screened by checking various processing and device parameters using special test sites on the wafer. If the parameters are within proper limits, each die on the wafer is tested for functionality, and bad dice are marked with a drop of ink.

Next, the dice are separated from the wafer using a diamond saw or a scribe-and-break process. Some die loss is caused by damage during the separation process. The remaining good dice are mounted in ceramic or plastic DIPs, LCCs, PGAs, surface-mount, or BGA packages using epoxy or eutectic die-attachment techniques.

Bonding pads on the die are connected to leads on the package using ultrasonic or thermosonic bonding of 15 to 75 μm aluminum or gold wire. Batch-fabricated flip-chip and TAB interconnection processes that permit simultaneous formation of hundreds or even thousands of bonds can also be used.

The final manufacturing cost of an integrated circuit is determined by the number of functional parts produced. The overall yield is the ratio of the number of working packaged dice to the original number of dice on the wafer. Yield loss is due to defects on the wafer, processing errors, damage during assembly, and lack of full functionality during final testing. The relationship between wafer yield and the size of an integrated-circuit die has been explored in detail. The larger the die size, the lower will be the number of good dice available from a wafer.

REFERENCES

- [1] J. W. Stafford, "The Implications of Destructive Wire Bond Pull and Ball Bond Shear Testing on Gold Ball-Wedge Wire Bond Reliability," *Semiconductor International*, p. 82, May 1982.
- [2] W. C. Till and J. T. Luxon, *Integrated Circuits: Materials, Devices and Fabrication*, Prentice-Hall, Englewood Cliffs, NJ, 1982.
- [3] J. R. Howell, "Reliability Study of Plastic Encapsulated Copper Lead Frame Epoxy Die Attach Packaging System," *Proceedings of the Reliability Physics Symposium, IEEE* pp. 104-110, 1981.
- [4] P. A. Totta and R. P. Sopher, "SLT Device Metallurgy and Its Monolithic Extension," *IBM Journal of Research and Development*, 13, 226-238, May 1969.
- [5] T. S. Liu, W. R. Rodrigues de Miranda, and P. R. Zipperlin, "A Review of Wafer Bumping for Tape Automated Bonding," *Solid State Technology*, 23, 71-76, March 1980.
- [6] B. T. Murphy, "Cost-Size Optima of Monolithic Integrated Circuits," *Proceedings of the IEEE*, 52, 1537-1545, December 1964.
- [7] R. B. Seeds, "Yield and Cost Analysis of Bipolar LSI," *IEEE IEDM Proceedings*, p. 12, October 1967.
- [8] C. H. Stapper, "On Yield, Fault Distributions, and Clustering of Particles," *IBM Journal of Research and Development*, 30, 326-338, May 1986.
- [9] *The International Technology Roadmap for Semiconductors*, The Semiconductor Industry Association (SIA), San Jose, CA: 1999. (<http://www.semichips.org>)

FURTHER READING

- [1] G. G. Harman and J. Albers, "The Ultrasonic Welding Mechanism as Applied to Aluminum- and Gold-Wire Bonding in Microelectronics," *IEEE Transactions on Parts, Hybrids and Packaging*, PHP-13, 406-412, December 1977.
- [2] K. I. Johnson, M. H. Scott, and D. A. Edson, "Ultrasonic Wire Welding—Part I: Wedge-Wedge Bonding of Aluminum Wires," *Solid State Technology*, 20, 50-56, March 1977.
- [3] K. I. Johnson, M. H. Scott, and D. A. Edson, "Ultrasonic Wire Welding—Part II: Ball-Wedge Wire Welding," *Solid State Technology*, 20, 91-95, April 1977.
- [4] C. Plough, D. Davis, and H. Lawler, "High Reliability Aluminum Wire Bonding," *Proceedings of the Electronic Components Conference, IEEE*, pp. 157-165, 1969.
- [5] N. Ahmed and J. J. Svitak, "Characterization of Gold-Gold Thermocompression Bonding," *Proceedings of the Electronic Components Conference, IEEE*, pp. 92-97, 1976.
- [6] L. S. Goldmann, "Geometric Optimization of Controlled Collapse Interconnections," *IBM Journal of Research and Development*, 13, 251-265, May 1969.
- [7] K. C. Norris and A. H. Landzberg, "Reliability of Controlled Collapse Interconnections," *IBM Journal of Research and Development*, 13, 266-271, May 1969.
- [8] J. E. Price, "A New Look at Yield of Integrated Circuits," *Proceedings of the IEEE*, 58, 1290-1291, August 1970.
- [9] C. H. Stapper, Jr., "On a Composite Model to the IC Yield Problem," *IEEE Journal of Solid-State Circuits*, SC-10, 537-539, December 1975.
- [10] C. H. Stapper, "The Effects of Wafer to Wafer Defect Density Variations on Integrated Circuit Defect and Fault Distributions," *IBM Journal of Research and Development*, 29, 87-97, January 1985.

PROBLEMS

- 8.1 Make a list of at least ten process or device parameters which could easily be monitored using a special test site on a wafer.
- 8.2 A simple microprocessor contains 115 flip-flops and hence 2^{115} possible states. If a tester can perform a new static test every 100 nsec, how many years will it take to test every state in the microprocessor chip? If the wafer has 100 dice, how long will it take to test the wafer?
- 8.3 (a) How many pads can be placed on a 10×15 mm die if a single row of pads is used? Assume the use of 100- μ m pads on a 125- μ m pitch with no pads in the corners.
(b) Repeat for 75- μ m pads on a 100- μ m pitch.
(c) How many solder balls would fit on the same die using an area array if the ball pads were 125- μ m pads on a 200- μ m pitch?
- 8.4 Compare the four yield formulas for a large VLSI die in which $D_0A = 10$ defects. Assume a clustering parameter of 1.0. How many good dice can we expect from 100- and 150-mm-diameter wafers using the different yield expressions? (The number of square dice per wafer can be estimated from $N = \pi (R - S)^2 / S^2$, where R is the radius of the wafer and S is the length of one side of the die.) Assume $S = 5$ mm.
- 8.5 What is the wafer yield for the defect map in Fig. 8.11 if the die is four times the size of that in Fig. 8.11(a)? What is the yield predicted by Poisson statistics? Assume the data from Fig. 8.11 is best represented by Eq. (8.9). What value of clustering parameter best fits the data?
- 8.6 A new circuit design is estimated to require a die which is 5×8 mm and will be fabricated on a wafer 125 mm in diameter. The process is achieving a defect density of 10 defects/cm², and the wafer processing cost is \$250.

- (a) What will be the cost of the final product if testing and packaging adds \$1.60 to the completed product?
- (b) The circuit design could be partitioned into two chips rather than one, but each die will increase in area by 15% in order to accommodate additional pads and I/O circuitry. If the testing and packaging cost remains the same, what is the cost of the two-chip set? Base your answers on Eq. (8.8). (See Problem 8.4 for the number of dice per wafer.)
- 8.7 (a) Repeat Problem 8.5 for a defect density of 5 defects/cm² and a wafer cost of \$150.
 (b) Repeat Problem 8.5 for a defect density of 5 defects/cm² and a wafer cost of \$300.
- 8.8 A die has an area of 25 mm² and is being manufactured on a 100-mm-diameter wafer using a process rated at 2 defects/cm². A new process is being developed which allows the die area to be reduced by a factor of 2. However, because of the smaller feature sizes, the new process costs 30% more and is presently achieving only 10 defects/cm².
- (a) Is it economical to switch to this new process?
- (b) At what defect density does the cost of the new die equal the cost of the old die?
- (c) Based on your judgment, would you recommend switching to the new process even if it is not now economical? Why?
- (d) At what die size is the cost the same in either process? Use Eq. (8.8) for this problem.
- 8.9 What is the limit of the yield distribution in Eq. (8.9) as the clustering parameter approaches infinity?
- 8.10 Compare the predictions of yield equations 8.5 and 8.9 for D_0A ranging from 1 to 10 with $\alpha = 5$ and $\alpha = 5,000$.
- 8.11 Suppose $D_0 = 0.1/\text{cm}^2$. What is the average number of defects on 150 mm, 200 mm, and 300 mm wafers?
- 8.12 Suppose that going from 100-mm wafers to 150-mm wafers changes the wafer processing cost from \$150/wafer to \$250/wafer, and the defect density remains constant at 10 defects/cm². What two die sizes give the same die cost? Use Eq. (8.9) with a cluster factor of 2. Use a calculator or computer to find the answer by iteration.
- 8.13 What would be the die yield in Fig. 8.11(b) if the defect positions were the same but the die pattern was rotated by 90°? How many good dice with four times the area of that in Fig. 8.11(a) would now exist?
- 8.14 A Gaussian probability density function for defect density is given by

$$f(D) = \frac{2}{D_0\sqrt{\pi}} \exp\left[-\frac{2(D-D_0)^2}{D_0}\right], \text{ for } 0 \leq D \leq 2D_0, \text{ and } 0 \text{ otherwise}$$

Calculate the yield Y for various values of D_0A and compare your results to those of the triangular distribution given in Eq. (8.7). (You may want to use a calculator or computer to perform the iteration.)

- 8.15 The wafers shown in Fig. 8.11 actually have 120 defects placed randomly on the wafer. Obviously, some chips must have several defects. Use Eq. (8.1) to predict how many dice will have exactly 1, 2, 3, 4, and 5 defects.
- 8.16 What defect density is required to achieve a yield of 70% for a 10×15 mm die if the process is characterized by a cluster parameter of 5? (b) Repeat for 80% yield. (c) Repeat for 90% yield.
- 8.17 What defect density is required to achieve a yield of 75% for a 20×20 mm die if the process is characterized by a cluster parameter of 6? (b) Repeat for 85% yield.

MOS Process Integration

In Chapter 9, we explore a number of relationships between process and device design and circuit layout. Processes are usually developed to provide devices with the highest possible performance in a specific circuit application, and one must understand the circuit environment and its relation to device parameters and device layout.

In this chapter, we look at a number of basic concerns in MOS process design, including channel-length control; layout ground rules and ground-rule design; source-drain breakdown and punch-through voltages; and threshold-voltage adjustment. Metal-gate technology is discussed, and the important advantages of self-aligned silicon-gate technologies are presented. Discussions of CMOS and silicon-on-insulator technologies complete the chapter.

3.1 BASIC MOS DEVICE CONSIDERATIONS

To explore the relationship between MOS process design and basic device behavior, we begin by discussing the static current-voltage relationship for the MOS transistor, as developed in Volume IV of this series [1]. The cross section of two NMOS transistors is shown in Fig. 9.1. In the linear region of operation, the drain current is given by

$$I_D = \bar{\mu}_n C_O (w/L) (V_{GS} - V_{TN} - V_{DS}/2) V_{DS} \quad (9.1)$$

for $V_{GS} \geq V_{TN}$ and $V_{DS} \leq V_{GS} - V_{TN}$. $C_O = K_O \epsilon_O / X_O$ is the oxide capacitance per unit area, $\bar{\mu}_n$ is the average majority-carrier mobility in the inversion layer, and V_{TN} is the threshold voltage. W and L represent the width and length of the channel, respectively.

One of the first specifications required is the circuit power-supply voltages, which set the maximum value of V_{GS} and V_{DS} that the devices must withstand. Once this choice is made, the only variables in Eq. (9.1) that a circuit designer may adjust are the width and length of the transistor. Thus, the circuit designer varies the circuit topology and horizontal geometry to achieve the desired circuit function.

Other device parameters are fixed by the process designer, who must determine the process sequence, times, temperatures, etc., which ultimately determine the device

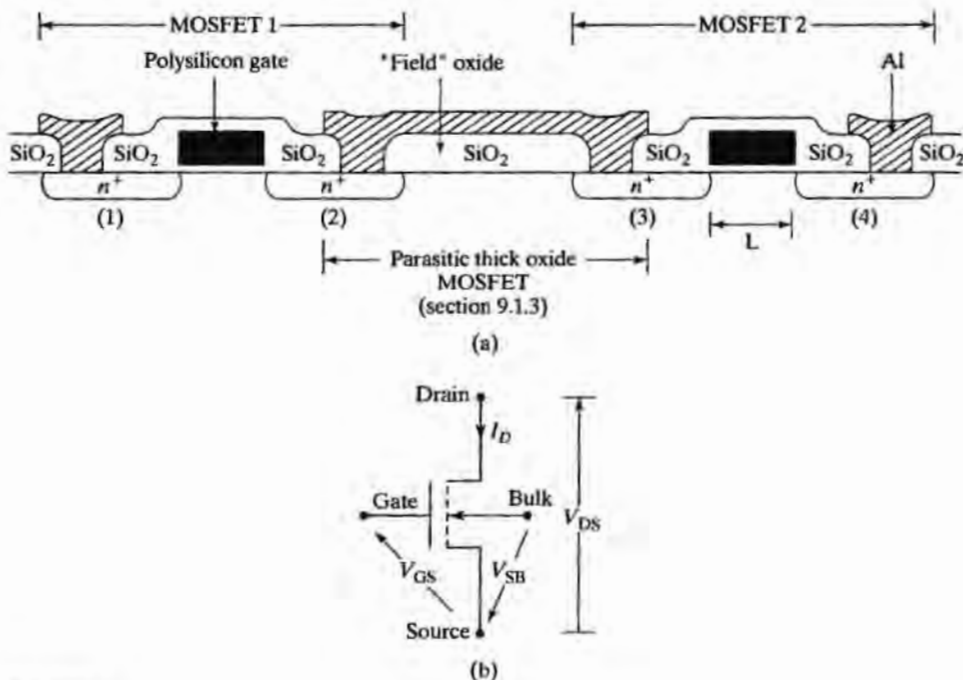


FIGURE 9.1

(a) Cross section of an integrated circuit showing two adjacent NMOS transistors. A parasitic NMOS device is formed by the aluminum interconnection over the field oxide with diffused regions (2) and (3) acting as source and drain. (b) An NMOS transistor with gate-to-source (V_{GS}), drain-to-source (V_{DS}), and source-to-bulk (V_{SB}) voltages defined.

structure and hence its characteristics. These include specifying the gate-oxide thickness, field-oxide thickness, substrate doping, and field and threshold-adjustment implantations. The process designer also supplies a set of design rules, or ground rules that must be obeyed during circuit layout. These include minimum channel length and width, spacings between features on the same and different mask levels, and overlaps between features on different mask levels. A mask alignment sequence and tolerances must also be developed for the process.

9.1.1 Gate-Oxide Thickness

Current flow in the MOS transistor, for a given set of terminal voltages, is inversely proportional to the gate-oxide thickness. The gate oxide will generally be made as thin as possible, commensurate with oxide breakdown and reliability considerations. High-quality silicon dioxide will typically break down at electric fields of 5 to 10 MV/cm, corresponding to 5 to 10 V across a 10-nm oxide. Present processes are using oxide thicknesses between 2 and 10 nm. Below 10 nm, current starts to flow by tunneling, and the oxide begins to lose its insulating qualities. The choice of oxide thickness is also related to hot electron injection into the oxide, a problem beyond the scope of this text [2–4]. Various alternative gate oxide materials are being investigated. Oxynitrides are formed by adding nitrogen to the silicon dioxide system either during or after formation of the gate oxide. Researchers are exploring a number of high-dielectric constant

gate materials that will permit the use of somewhat thicker oxides without reducing the oxide capacitance and transconductance of the transistor.

9.1.2 Substrate Doping and Threshold Voltage

Threshold voltage is an important parameter which determines the gate voltage necessary to initiate conduction in the MOS device. The threshold voltage [1] for a device with a uniformly doped substrate is given by

$$\text{NMOS: } V_{TN} = \Phi_M - \chi - \frac{E_g}{2q} + |\Phi_F| + \left[\sqrt{2K_s \epsilon_0 q N_B (2|\Phi_F| + V_{SB})} \right] / C_O - Q_{tot} / C_O \quad (9.2)$$

$$\text{PMOS: } V_{TP} = \Phi_M - \chi - \frac{E_g}{2q} - |\Phi_F| - \left[\sqrt{2K_s \epsilon_0 q N_B (2|\Phi_F| - V_{BS})} \right] / C_O - Q_{tot} / C_O$$

$$|\Phi_F| = (kT/q) \ln(N_B/n_i)$$

where N_B is the substrate doping, $\Phi_M - \chi = -0.11$ for an aluminum gate, $\Phi_M - \chi = 0$ for an n^+ -doped polysilicon gate, and $\Phi_M - \chi = +1.12$ for a p^+ -doped polysilicon gate.

Q_{tot} represents the total oxide and interface charge per cm^2 and adds a parallel shift of the curves in Fig. 9.2 to more negative values of threshold. This charge contribution to the threshold voltage had an extremely important influence on early MOS device fabrication. Q_{tot} tends to be positive, which makes the MOS transistor threshold more negative; n -channel transistors become depletion-mode devices ($V_{TN} < 0$), whereas p -channel transistors remain enhancement-mode devices ($V_{TP} < 0$). During early days of MOS technology, Q_{tot} was high, and the only successful MOS processing was done using PMOS technology. After the industry gained an understanding of the origin of oxide and interface charges, and following the advent of ion implantation, NMOS technology became dominant, because of the mobility advantage of electrons over holes. Total charge levels have been reduced to less than 5×10^{10} charges/ cm^2 in good MOS processes, and the oxide charge contribution to threshold voltage is minimal.

Substrate doping enters the threshold-voltage expression through both the $|\Phi_F|$ term and the square-root term. A plot of threshold voltage versus substrate doping for n - and p -channel, n^+ polysilicon-gate devices with 10-nm gate oxides is given in Fig. 9.2 for $Q_{tot} = 0$. The choice of substrate doping is complicated by other considerations, including drain-to-substrate breakdown voltage, drain-to-source punch-through voltage, source-to-substrate and drain-to-substrate capacitances, and substrate sensitivity or body effect.

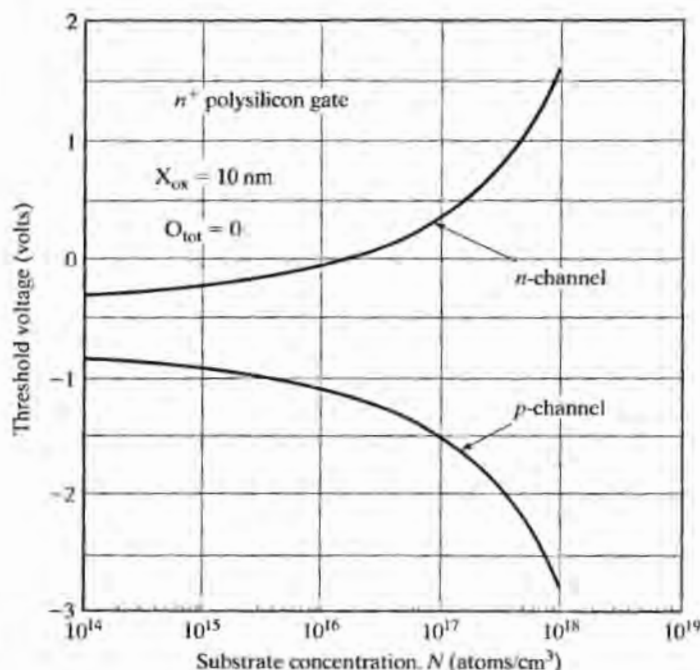


FIGURE 9.2

Threshold voltages for n - and p -channel polysilicon-gate transistors with 10-nm gate oxides, calculated from Eq. (9.2).

9.1.3 Junction Breakdown

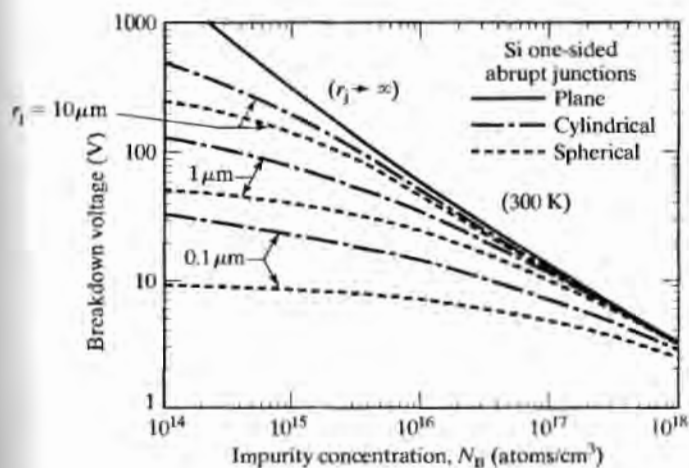
The source and drain regions are usually heavily doped to minimize their resistance and are essentially one-sided junctions in which the depletion region extends entirely into the substrate. Figure 9.3(a) gives the breakdown voltage of a one-sided pn junction as a function of the doping concentration on the lightly doped side of the junction [5]. Junction breakdown voltage decreases as doping level increases. Breakdown voltage is also a function of the radius of curvature of the junction space-charge region. Junction curvature enhances the electric field in the curved region of the depletion layer and reduces the breakdown voltage below that predicted by one-dimensional junction theory. A rectangular diffused area has regions with both cylindrical and spherical curvature, as shown in Fig. 9.3(b). It is worth noting that very shallow spherical junctions break down at voltages of less than 10 V, regardless of doping level.

9.1.4 Punch-through

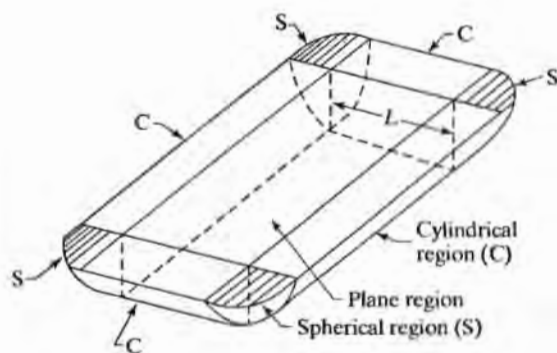
Punch-through occurs when the drain depletion region contacts the source depletion region, and substrate doping must be chosen to prevent the merging of these depletion regions when the MOSFET is off. Punch-through will not occur if the channel length exceeds the sum of the depletion-layer widths of the source-to-substrate and drain-to-substrate junctions. For a transistor used as a load device in a logic circuit, the source-to-substrate and drain-to-substrate junctions must both support a voltage equal to the drain supply voltage plus the substrate supply voltage. The depletion-layer widths can be estimated using the formula for the width of a one-sided step junction

$$W_d = \sqrt{(2K_s\epsilon_0(V_A + \Phi_{bi}))/qN_B}.$$

$$\Phi_{bi} = 0.56 + (kT/q) \ln(N_B/n_i), \quad (9.3)$$



(a)



(b)

FIGURE 9.3

(a) Abrupt pn junction breakdown voltage versus impurity concentration on the lightly doped side of the junction for both cylindrical and spherical structures. r_1 is the radius of curvature. (b) Formation of cylindrical and spherical regions by diffusion through a rectangular window. Copyright 1985, John Wiley & Sons, Inc. Reprinted with permission from Ref. [5].

where V_A is the total applied voltage and Φ_{bi} is the built-in potential of the junction. If the channel length is greater than $2W_d$ punch-through should not occur. Figure 9.4 gives the depletion-layer width of pn junctions as a function of doping and applied voltage. Punch-through is not a limiting factor for most doping levels, except for very short-channel transistors. Ion implantation has been used to enhance the doping concentration below the channel region of short-channel devices to increase the punch-through voltage.

9.1.5 Junction Capacitance

The capacitance per unit area associated with a diffused junction is given by the parallel-plate capacitance formula with a plate spacing of W_d :

$$C_j = K_s \epsilon_0 / W_d$$

The larger the doping, the larger the capacitance. Zero bias and a doping concentration of $10^{16}/\text{cm}^3$ result in a junction capacitance of approximately $10 \text{ nf}/\text{cm}^2$.

Eq. (9.2) shows that the threshold voltage depends on the source-to-substrate voltage, V_{SB} . This variation is known as "substrate sensitivity," or "body effect," and it becomes worse as the substrate doping level increases.

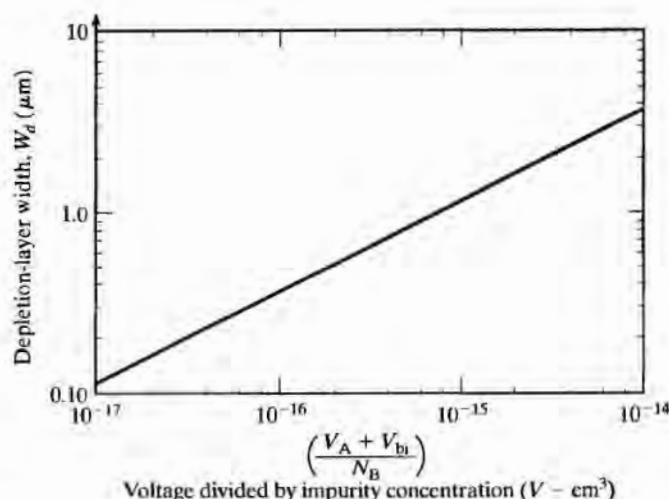


FIGURE 9.4

Depletion-layer width of a one-sided step junction as a function of doping and applied voltage calculated from eq. (9.3).

From the preceding discussion, one can see that there are trade-offs involved in the choice of substrate doping. Substrate doping is directly related to threshold voltage. It is desirable to reduce substrate doping to minimize junction capacitance and substrate sensitivity and to maximize breakdown voltage. Mobility also tends to be higher for lower doping levels. On the other hand, a heavily doped substrate will increase the punch-through voltage.

9.1.6 Threshold Adjustment

Ion implantation is routinely used to separate threshold-voltage design from the other factors involved in the choice of substrate doping. Substrate doping can be chosen based on a combination of breakdown, punch-through, capacitance, and substrate sensitivity considerations, and the threshold voltage is then adjusted to the desired value by adding a shallow ion-implantation step to the process. Figure 9.5 shows a step approximation to an implanted profile used to adjust the impurity concentration near the surface. These additional impurities cause a shift in threshold voltage given approximately by

$$\Delta V_{TN} = (1/C_O) (qQ_i) (1 - x_i/2x_d), \quad x_i \ll x_d, \quad x_d = \sqrt{qN_B/4K_s\epsilon_0|\Phi_F|} \quad (9.4)$$

where $Q_i = x_i N_i$ represents the implanted dose and x_d represents the depletion-layer width beneath the gate. For shallow implants, the threshold-voltage shift is proportional to the implanted dose. The threshold-voltage shift is positive for acceptor impurities and negative for donor impurities.

Example 9.1:

An NMOS transistor with an n^+ polysilicon gate is fabricated with a 10-nm gate oxide, a substrate doping of $10^{16}/\text{cm}^3$, and source-drain junction depths of 0.25 μm .

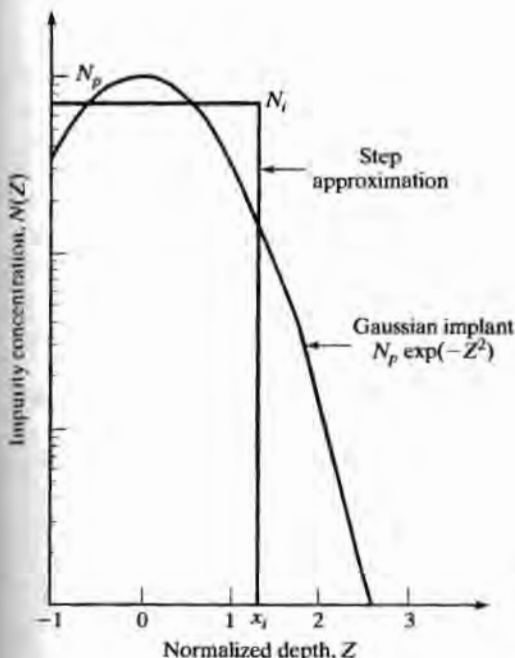


FIGURE 9.5

Step approximation to a Gaussian impurity profile used to estimate the threshold-voltage shift achieved using ion implantation.

Determine the threshold voltage and drain-to-substrate breakdown voltages for this device. What is the punch-through voltage for a channel length of $1\text{ }\mu\text{m}$ if the substrate bias is 0 V ? A shallow boron implantation is to be used to adjust the threshold to 0.7 V . What is the dose of this implant? (Assume that $V_{SB} = 0$ and $Q_{tot} = 0$.)

Solution : For the n^+ polysilicon-gate transistor, $\Phi_M - \chi - E_g/2q = -0.56\text{ V}$ and $|\Phi_F| = 0.36\text{ volts}$ (for $n_i = 1 \times 10^{10}/\text{cm}^3$ and $kT/q = 0.026\text{ V}$). For $V_{SB} = 0$, the threshold voltage expression yields $V_{TN} = -0.56 + 0.36 + 0.14\text{ V} = -0.06\text{ V}$. Interpolating Fig. 9.3 for spherical breakdown with a substrate doping of $10^{16}/\text{cm}^3$ and a radius of curvature of $0.25\text{ }\mu\text{m}$ gives an estimated drain-to-substrate breakdown voltage of 20 V . To estimate the punch-through voltage, we use Eq. (9.3) with $2W = 1\text{ }\mu\text{m}$ and $V_A = V_D$, where V_D is the drain voltage. Evaluating this expression yields $V_D = 1.01\text{ V}$.

For a shallow implant, the threshold-voltage shift is approximately $\Delta V_T = qQ/C_{ox}$. A voltage shift of 0.76 V with an oxide thickness of 10 nm yields $\Delta Q = 1.64 \times 10^{12}/\text{cm}^2$.

Thin gate oxides mentioned earlier in this chapter also have potential problems with impurity diffusion from the polysilicon gates through the oxide and into the substrates. Any doping that makes it into the substrate is directly in the MOS channel region and will shift the threshold of the devices.

In the past, NMOS depletion-mode ($V_{TN} < 0$) transistors were routinely used in processes designed for high-performance logic applications. To reduce the NMOS threshold voltage, n -type impurities can be implanted to form a built-in channel connecting the source and drain regions of the transistor, as shown in Fig. 9.6. The device characteristics of a depletion-mode transistor are similar, although not identical, to those of an enhancement-mode NMOS transistor, and the dose needed to shift the threshold voltage may be estimated using Eq. (9.4).

9.1.7 Field-Region Considerations

The region between the two transistors in Fig. 9.1 is called the *field* region and must be designed to provide isolation between adjacent MOS devices. Several factors must be considered. The metal line over the field region can act as the gate of a "parasitic NMOS transistor" with diffused regions (2) and (3) acting as its source and drain. To ensure that this parasitic device is never turned on, the magnitude of the threshold voltage in this region must be much higher than that in the normal gate region. Referring to Eq. (9.2), we find that the threshold voltage may be made higher by increasing the oxide thickness in the field region and by increasing the doping below the field oxide. The field oxide is typically made three to ten times thicker than the gate oxide of the transistors.

Another problem occurs for NMOS transistors. The substrate for NMOS transistors is *p*-type, usually doped with boron. We know that thermal oxidation results in depletion of boron from the surface of the silicon, and looking at Eq. (9.2) we see that boron depletion will lower the threshold voltage of the transistors in the field region. A field implant step is often added to processes to increase the threshold voltage and compensate for the boron depletion during field-oxide growth.

For PMOS devices, the substrate is typically doped with phosphorus. During oxidation, phosphorus pileup at the surface tends to increase the threshold voltage in the field region. Thus, phosphorus pileup helps to keep the parasitic field devices turned off.

9.1.8 MOS Transistor Isolation

When two properly biased MOS transistors are placed near each other, they are isolated by reverse-biased source-substrate and drain-substrate junctions, as shown in Fig. 9.7. The MOS devices are referred to as *self-isolated*. No additional structure is required to achieve isolation, and this fact gives an inherent size advantage to MOS technology over the junction isolated bipolar structures discussed in the next chapter. However, to maintain this isolation, the depletion layers surrounding the various

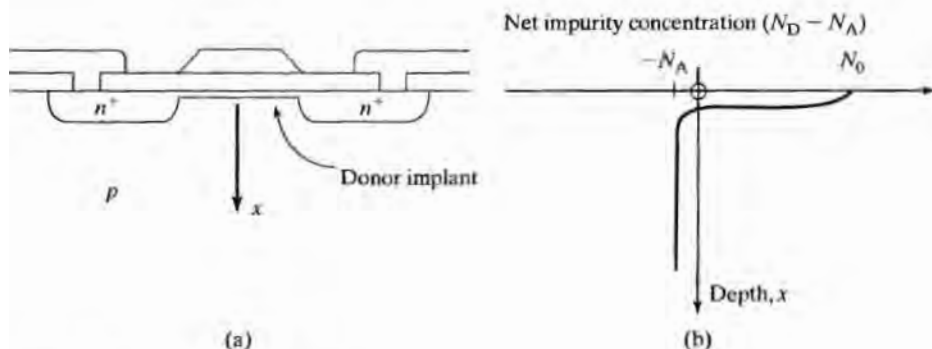


FIGURE 9.6

(a) Formation of a depletion-mode NMOS transistor using a shallow ion-implanted layer; (b) net impurity profile under the gate of the depletion-mode MOSFET.

source-drain diodes must not merge, and this requirement limits the minimum spacing between the devices. The spacing between adjacent transistors must be greater than twice the maximum depletion-layer width.

Example 9.2:

Use Eq. (9.3) to estimate the minimum spacing between the drains of two adjacent NMOS transistors if the substrate doping is $3 \times 10^{16}/\text{cm}^3$ and the maximum drain-substrate voltage is 5 V.

Solution: The n^+p drain-substrate junctions correspond to one-sided step junctions, so the use of Eq. (9.3) is appropriate. The built-in potential is equal to

$$V_{bi} = 0.56 \text{ V} + (0.0258 \text{ V}) \ln(3 \times 10^{16}/10^{10}) = 0.94 \text{ V}$$

and the depletion layer width is

$$W_d = \sqrt{2(11.7)(8.854 \times 10^{-14} \text{ F/cm})(5 \text{ V} + 0.94 \text{ V}) / (1.6 \times 10^{-19} \text{ C})(3 \times 10^{16} / \text{cm}^3)} \\ = 0.51 \mu\text{m}$$

Each transistor has a depletion layer around its drain, so the devices must be separated by at least twice this distance, and the minimum spacing between transistors (with no safety margin) is $1.02 \mu\text{m}$.

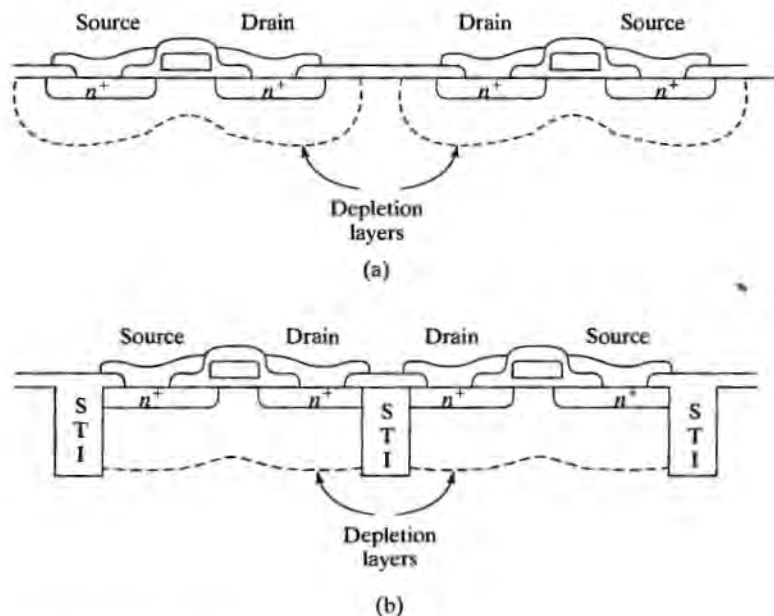


FIGURE 9.7

Isolation strategies (a) Intrinsic (b) Shallow trench isolation

In early technology based upon 2- μm or greater lithography, the spacing calculated in Example 9.2 does not represent a problem. However, for advanced processes with deep submicron feature sizes, this form of isolation is not satisfactory. This led first to the use of recessed oxide technology depicted in Fig. 3.12, and more recently to the pervasive utilization of shallow trench isolation (STI) in both MOS and bipolar technologies. For the STI depicted in Fig. 9.7(b), the depletion layers are effectively cut off and separated by the oxide. The minimum space between drain diffusions is now set by the minimum width of the STI region and can ideally approach a minimum feature size in the technology. Highly planar STI regions are produced using the CMP process described in Chapter 3. Note that the pn junction boundary intersects the STI oxide in Fig. 9.7(b). This interface represents a potential junction leakage site, but it is not a problem with well controlled processing. Note that pn junctions in MOS and bipolar transistors have always intersected the oxide at the surface of the silicon. (See Fig. 9.7(a).)

9.1.9 Lightly Doped Drain Structures

As devices are scaled to smaller dimensions, the substrate doping level tends to be increased, very shallow junctions with a high curvature are used, and the applied electric fields tend to increase. All of these factors tend to cause breakdown problems with the drain-substrate junction. A number of lightly doped drain (LDD) structures have evolved to control the breakdown problem. The concept is depicted in Fig. 9.8. After defining the polysilicon gate, we use an n -type implantation to form the LDD extension that ultimately defines the extent of the channel. An oxide or nitride "spacer" is formed on the edges of the gate by thermal oxidation, or CVD process, and then the highly doped source and drain contact regions are implanted. The reduced doping in the LDD region enhances the breakdown voltage of the transistor. A wide array of different process have been developed to achieve structures similar to that in Fig. 9.8.

9.1.10 MOS Transistor Scaling

The phenomenal increase in IC density and complexity has been driven by our ability to aggressively scale the physical dimensions (W , L , X_O , x_p , etc.) of the MOS transistor. A theoretical framework for MOSFET miniaturization was first provided by Dennard.

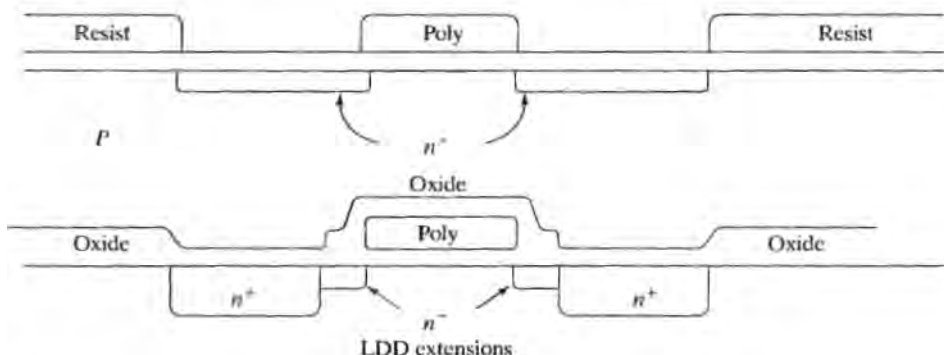


FIGURE 9.8

Self aligned polysilicon-gate transistor with lightly-doped source/drain regions

Gaensslen, Kuhn and Yu [22]. The basic tenant of the theory requires the electrical fields to be maintained constant within the device as the geometry is changed. Thus, if a physical dimension is reduced by a factor of α , then the voltage applied across that dimension must also be decreased by the same factor.

These rules are applied to the transconductance and linear region drain current for the MOSFET in Eq. (9.5) in which the three physical dimensions— W , L and X_O —are all reduced by the factor α , and each of the voltages including the threshold voltage is reduced by the same factor. For the n -channel MOSFET, we have

$$I_D^* = \bar{\mu}_n \frac{K_O \epsilon_O}{\left(\frac{X_O}{\alpha}\right)} \left(\frac{W}{\alpha}\right) \left(\frac{L}{\alpha}\right) \left(\frac{V_{GS}}{\alpha} - \frac{V_{TN}}{\alpha} - \frac{V_{DS}}{2\alpha}\right) \frac{V_{DS}}{\alpha} = \frac{I_D}{\alpha^3} \quad (9.5)$$

$$K_n^* = \bar{\mu}_n \frac{K_O \epsilon_O}{\left(\frac{X_O}{\alpha}\right)} \left(\frac{W}{\alpha}\right) \left(\frac{L}{\alpha}\right) = \alpha \bar{\mu}_n \frac{K_O \epsilon_O}{X_O} \frac{W}{L} = \alpha K_n$$

We see that the scaled drain current is actually reduced from the original value by the scale factor α , whereas the scaled transconductance parameter K_n^* is increased by the scale factor. In a similar manner, the total gate-channel capacitance of the device is also found to be reduced by α :

$$C_{GC}^* = (C_{OX})^* W^* L^* = \frac{K_O \epsilon_O}{\left(\frac{X_O}{\alpha}\right)} \left(\frac{W}{\alpha}\right) \left(\frac{L}{\alpha}\right) = \frac{C_{GC}}{\alpha} \quad (9.6)$$

We know that the delay of logic gates is limited by the transistor's ability to charge and discharge the capacitance associated with the circuit. Based upon $i = Cdv/dt$, an estimate of the delay of a scaled logic circuit is

$$\tau^* = C_{GC}^* \frac{\Delta V^*}{I_D^*} = \frac{C_{GC}}{\alpha} \frac{\frac{\Delta V}{\alpha}}{\frac{I_D}{\alpha^3}} = \frac{\tau}{\alpha} \quad (9.7)$$

We find that circuit delay is also improved by the scale factor α .

As we scale down the dimensions by α , the number of circuits in a given area will increase by a factor of α^2 . An important concern in scaling is therefore what happens to the power per circuit, and hence the power per unit area (power density) as dimensions are reduced. The total power supplied to a transistor circuit will be equal to the product of the supply voltage and the transistor drain current:

$$P^* = V_{DD}^* I_D^* = \left(\frac{V_{DD}}{\alpha}\right) \left(\frac{I_D}{\alpha}\right) = \frac{P}{\alpha^2} \text{ and } \frac{P^*}{A^*} = \frac{P^*}{W^* L^*} = \frac{\frac{P}{\alpha^2}}{\left(\frac{W}{\alpha}\right) \left(\frac{L}{\alpha}\right)} = \frac{P}{WL} = \frac{P}{A} \quad (9.8)$$

This equation is extremely important. It indicates that the power per unit area remains constant if a technology is properly scaled. Even though we are increasing the number of circuits by α^2 , the total power for a given size integrated circuit die will remain constant. Violation of the scaling theory over many years, by maintaining a constant 5-V power supply as dimensions were reduced, led to almost unmanageable power levels in many of today's integrated circuits. The problem could only be resolved by moving away from NMOS technology and into CMOS technology!

A useful figure of merit for comparing logic families is the power-delay product (PDP). The product of power and delay time represents energy, and the PDP represents a measure of the energy required to perform a simple logic operation:

$$\text{PDP}^* = P^* \tau^* = \frac{P}{\alpha^2} \frac{\tau}{\alpha} = \frac{\text{PDP}}{\alpha^3} \quad (9.9)$$

The PDP figure of merit shows the full power of technology scaling. The PDP is reduced by the cube of the scaling factor!

Each generation of lithography corresponds to a scale factor $\alpha = 1/\sqrt{2}$, so each new technology generation increases the number of circuits by a factor of 2 and improves the PDP by a factor of almost 3. Table 9.1 summarizes the performance changes achieved with constant field scaling.

9.2 MOS TRANSISTOR LAYOUT AND DESIGN RULES

Design of the layout for transistors and circuits is constrained by a set of rules called the design rules, or ground rules. These rules are technology specific and specify minimum sizes, spacings, and overlaps for the various shapes that define transistors. Processes are designed around a *minimum feature size*, which is the width of the smallest line or space that can be reliably transferred to the surface of the wafer using a given generation of lithography.

To produce a basic set of ground rules, we must also know the maximum misalignment that can occur between two mask levels. Figure 9.9(a) shows the nominal position of a metal line aligned over a contact window. The metal overlaps the contact window by at least one *alignment tolerance* in all directions. During the fabrication process, the alignment will not be perfect, and the actual structure may have misalignment in both the x - and y -directions. Figures 9.9(b)–(d) show the result of worst-case misalignment of the patterns in the x -, y -, and both directions simultaneously. Our set of design rules will assume that this alignment tolerance is the same in both directions.

TABLE 9.1 Constant Electric Field Scaling Results

Performance Measure	Scale Factor
Area/Circuit	$1/\alpha^2$
Transconductance Parameter	α
Current	$1/\alpha$
Capacitance	$1/\alpha$
Circuit Delay	$1/\alpha$
Power/Circuit	$1/\alpha^2$
Power/Unit Area (Power Density)	1
Power-Delay Product (PDP)	$1/\alpha^3$

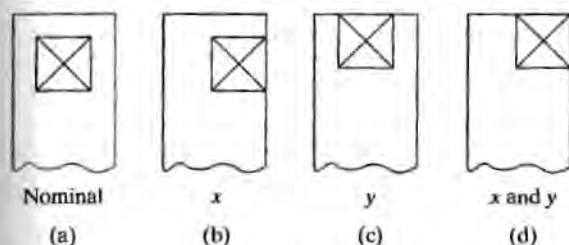


FIGURE 9.9

(a) Nominal alignment of the contact and metal masks; (b) worst-case misalignment in the x -direction, (c) in the y -direction, and (d) in both directions.

9.2.1 Metal-Gate Transistor Layout

The first successful MOS technologies utilized aluminum for the gate material. Although these metal-gate devices are seldom used in today's silicon processes, an understanding of their layout provides significant insight into the limitations of the metal gate process and the importance of the self-aligned silicon-gate technologies that replace them. The low melting temperature of aluminum greatly limits the type of processing steps that can be used following the metal deposition step. Refractory metals such as tungsten, which can withstand very high temperatures, have been used in experimental self-aligned metal-gate MOS processes.

Figure 9.10 shows the process sequence for a basic metal-gate process. The first mask defines the position of the source and drain diffusions. Following diffusion, the second mask is used to define a window for growth of the thin gate oxide. The third and fourth masks delineate the contact openings and metal pattern. The metal-gate mask sequence, omitting the final passivation layer mask, is as follows:

- | | |
|--------------------------------|------------------|
| 1. Source/drain diffusion mask | First mask |
| 2. Thin oxide mask | Align to level 1 |
| 3. Contact window mask | Align to level 1 |
| 4. Metal mask | Align to level 2 |

An alignment sequence must be specified in order to properly account for alignment tolerances in the ground rules. In this metal-gate example, mask levels two and three are aligned to the first level, and level four is aligned to level two.

We will first look at a set of design rules for metal-gate transistors similar in concept to the rules developed by Mead and Conway [6]. These ground rules were designed to permit easy movement of a design from one generation of technology to another by simply changing the size of a single parameter λ . In order to achieve this goal, the rules are quite loose in terms of level-to-level alignment tolerance. We will explore tighter ground rules later in this chapter.

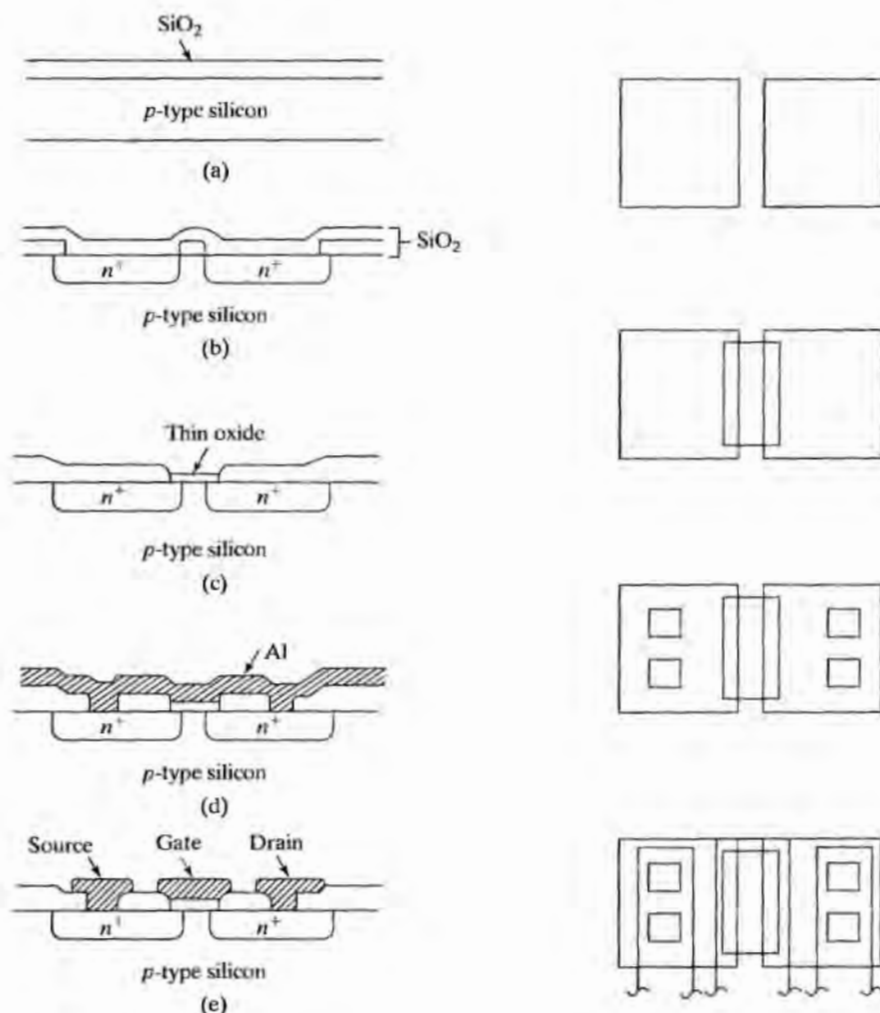


FIGURE 9.10

Mask steps and device cross sections in a metal-gate process. (a) Substrate ready for first mask step; (b) substrate following source/drain diffusion and oxide regrowth, (c) following gate-oxide growth, (d) following contact window mask and aluminum deposition, and (e) following metal delineation.

A set of metal-gate rules is shown in Fig. 9.11. The minimum feature size $F = 2\lambda$, and the alignment tolerance $T = \lambda$. The parameter λ could be $1\text{ }\mu\text{m}$, $.25\text{ }\mu\text{m}$, or $0.1\text{ }\mu\text{m}$, for example. Transistors designed using our ground rules will fail to operate properly if the misalignment exceeds the specified alignment tolerance T .

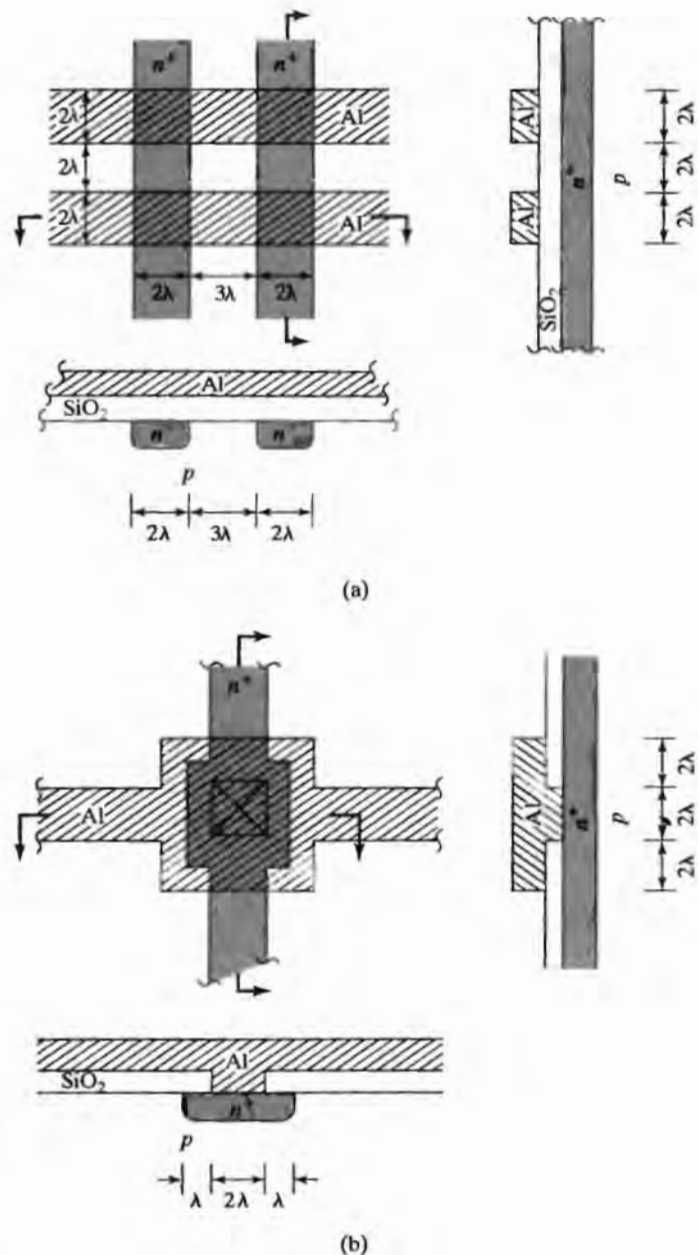


FIGURE 9.11

A simple " λ -based" set of "design rules" or "ground rules" based on an alignment sequence in which levels 2 and 3 are aligned to level 1 and level 4 is aligned to level 2. (a) Rules for metal and diffused interconnection lines; (b) rules for contacts between metal and diffusion.

On the metal level, minimum line widths and spaces are equal to 2λ . In some processes, the metal widths are made larger, because the metal level encounters the most mountainous topology of any level.

On the diffusion level, the minimum linewidth is 2λ . The minimum space between diffusions is increased to 3λ to ensure that the depletion layers of adjacent lines do not merge together. However, the spacing between the source-drain diffusions of a transistor may be 2λ .

In this set of rules, the alignment tolerance between two mask levels is assumed to be 1λ , which represents the maximum shift of one level away from its nominal position, relative to the level to which it is being aligned. A 1λ shift can occur in both the x - and y -directions.

Square contacts are a minimum feature size of 2λ in each dimension. It is normal practice to ensure that the contact is completely covered by metal even for worst-case alignment. Depending on the alignment sequence, a 1λ or 2λ metal border will be required around the contact window. Likewise, a contact window must be completely surrounded by a 1λ or 2λ border of the diffused region beneath the contact.

For our metal-gate transistors, the thin oxide region will be aligned to diffusion, so it requires a 1λ overlap over the source-drain diffusions in the length direction. The source-drain regions must also extend past the thin oxide by at least 1λ in the width direction. Contacts must be inside the diffusions by 1λ . The metal level is aligned to the thin oxide level, whereas the contacts are aligned to the diffusion level. A worst-case layout therefore requires a 2λ border of metal around contact windows, but only a 1λ border around the thin oxide regions.

Figure 9.12 shows the horizontal layout and vertical cross section of a minimum-size NMOS metal-gate transistor with $W/L = 10\lambda/2\lambda = 5/1$ at the mask level. The two diffusions are spaced by a minimum feature size of 2λ . Thin oxide must overlap the diffusions by 1λ in the length direction and underlap the diffusions by 1λ in the width direction. Metal must overlap thin oxide by 1λ . Accumulated alignment tolerances cause the minimum width of the gate metal to be 6λ . The spacing between metal lines must be 2λ . The metal over the contact holes must be 8λ wide, because of the alignment sequence used, and the contact hole must be 1λ inside the edge of the diffusion. The resulting minimum transistor is 26λ in the length direction and 16λ in the width direction.

A new design rule has been introduced into this layout. The gate metal is spaced 1λ from the diffusion to prevent the edge of a metal line from falling directly on top of the edge of the diffusion in the nominal layout.

Several observations can be made by looking at this structure. First, note that the transistor is $416\lambda^2$ in total area, whereas the active channel area of the device is $20\lambda^2$. The rest of the area is required in order to make contacts to the various regions, within the constraints of the minimum feature size and alignment tolerance rules. Second, there is a substantial area of thin and thick oxide in which the gate metal overlaps the source and drain regions of the transistor. This increases the gate-to-source and gate-to-drain capacitance of the transistor. In this metal-gate transistor layout, the channel is defined by the junction edges in the length direction and by the thin oxide region in the width direction.

It should also be noted that there are several small contact windows in the source and drain regions. The usual practice is to make all the contact windows the same size throughout the wafer. From a processing point of view, equal-size contact windows will all tend to open at the same time during the etching process. The uniform size of the contacts also facilitates modeling of the contact resistance as the area of the diffusion is changed.

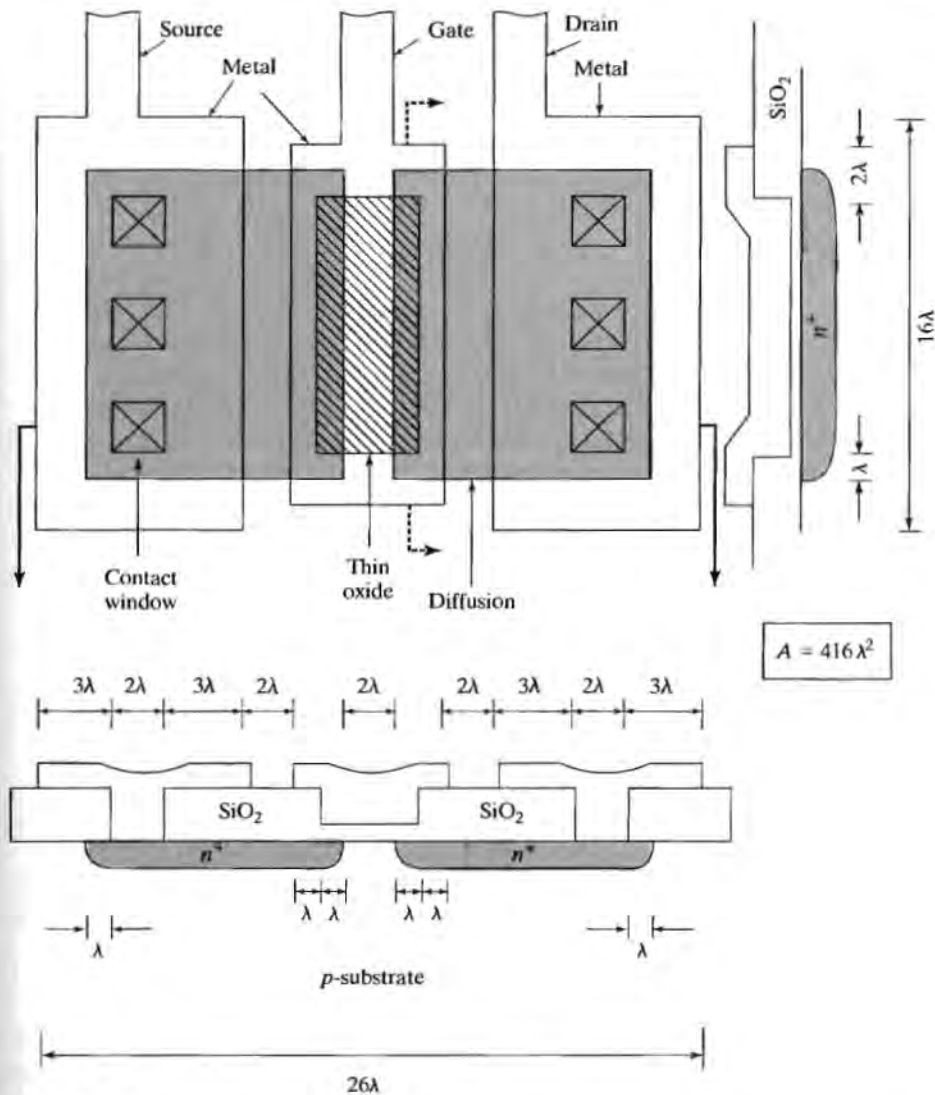


FIGURE 9.12

Minimum-size metal-gate transistor with W/L ratio of 5/1 using the design rules of Fig. 9.11. The active gate region is less than 5% of the total device area.

9.2.2 Polysilicon-Gate Transistor Layout

Transistors fabricated using polysilicon-gate technology have a number of important advantages over those built using metal-gate processes. The polysilicon gate can withstand high-temperature processing following its deposition, and this significantly improves the flexibility of the process. The silicon gate can be directly oxidized at high temperature to form an insulating layer over the gate. The heavily doped polysilicon represents an additional interconnect layer that other metal layers can easily cross, because of the oxide isolation. However, the most significant advantages are in layout

and parasitic capacitance reduction, and we will discover some of these advantages by looking at the layout and structure of the polysilicon-gate transistor.

The mask sequence for the basic polysilicon-gate process from Chapter 1 is (again without passivation layer) as follows:

- | | |
|------------------------------------|------------------|
| 1. Active region (thin oxide) mask | First mask |
| 2. Polysilicon mask | Align to level 1 |
| 3. Contact window mask | Align to level 2 |
| 4. Metal mask | Align to level 3 |

Some new design rules must be introduced for this process. Polysilicon lines and spaces will both be a minimum feature size of 2λ . The polysilicon gate must overlap the thin oxide region by an alignment tolerance λ . The preceding alignment sequence requires 1λ polysilicon and 1λ metal borders around contacts. However, contact holes should have a 2λ border of thin oxide due to tolerance accumulation.

Figure 9.13 shows the layout of the polysilicon-gate device with $WL = 5/1$ using these design rules. The total area is $168\lambda^2$. The active channel region now represents 12% of the total area, compared with less than 5% for the metal-gate device. The polysilicon gate acts as a barrier material during source-drain implantation and results in the self-alignment of the edge of the gate to the edge of the source-drain regions. Self-alignment of the gate to the channel reduces the size of the transistor and eliminates the overlap region between the gate and the source-drain regions. In addition, the size of the transistor is reduced, because the source-drain metallization can be placed nearer to the gate. In the polysilicon-gate layout, the channel is defined by the polysilicon gate in the length direction and by the thin oxide in the width direction.

A very important side benefit resulting from this process is the third level of interconnection provided by the polysilicon. Circuit wiring may be accomplished on the diffusion, metal, and polysilicon levels in the polysilicon-gate technology.

A design rule concerning edges has again been introduced into this layout. Metal lines are spaced 1λ from the polysilicon gate to prevent the edge of the metal line from falling directly on top of the edge of the polysilicon line in the nominal layout.

9.2.3 More-Aggressive Design Rules

The design rules discussed so far have focused on minimum feature size and alignment tolerance. F and T are determined primarily by the type of lithography being practiced. However, linewidth expansion and shrinkage throughout the process also strongly affect the ground rules. Expansion or shrinkage may occur during mask fabrication, resist exposure, resist development, etching, or diffusion. These linewidth changes are normally factored into the design rules.

In addition, alignment variation is a statistical process. Worst-case misalignments occur only a very small percentage of the time. (For a Gaussian distribution, a 3σ misalignment occurs only 2% of the time.) Our set of rules based on worst-case alignment tolerances is very pessimistic. For example, assuming that contacts are misaligned by λ in one direction, at the same time that the metal level is misaligned in the opposite direction by λ , results in an accumulated tolerance of 2λ . However, this situation would most probably never occur.

Let us consider the impact of tightening two design rules in the polysilicon-gate process. First, we will let the edge of one layer align with the edge of another layer.

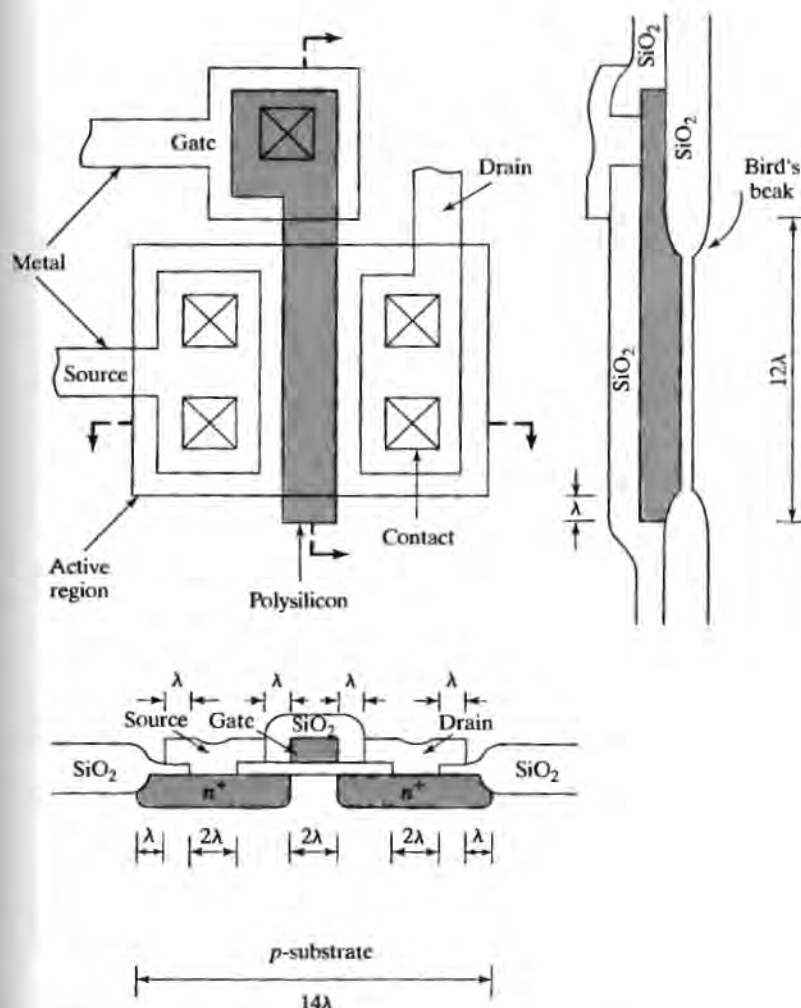


FIGURE 9.13

Minimum-size polysilicon-gate transistor layout for $W/L = 5/1$. The active gate region occupies 12% of the transistor area, and parasitic gate capacitance is minimized.

Second, a contact window will be allowed to run over onto the field oxide by 1λ . The resulting layout using our polysilicon-gate alignment sequence is shown in Fig. 9.14. The total area of the device has been reduced 25% to $120\lambda^2$, and the active channel region now represents 17% of the total transistor area. We see how ground rule changes can have a substantial effect on device area.

9.2.4 Channel Length and Width Biases

Figure 9.15 presents another example of the interaction of the process with design-rule definitions. Here we assume a polysilicon-gate process in which the source-drain junction depth is equal to $\lambda/2$ and lateral diffusion equals vertical diffusion. Since we know

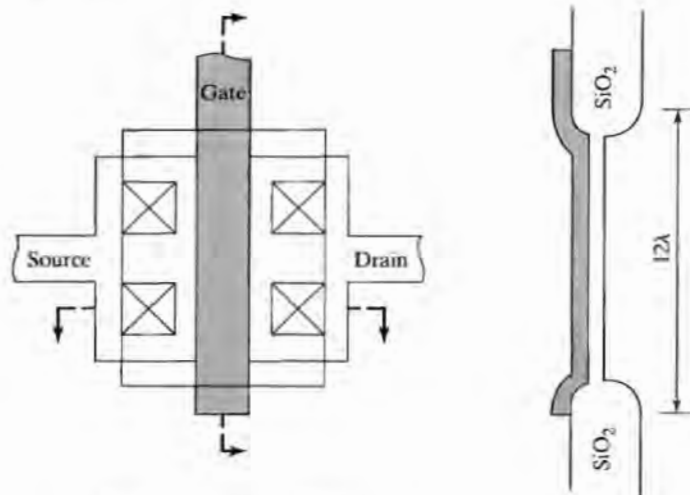


FIGURE 9.14

More aggressive layout of the polysilicon-gate transistor in which two ground rules have been relaxed. Active gate area is now 17% of total device area.

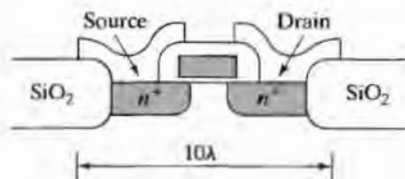
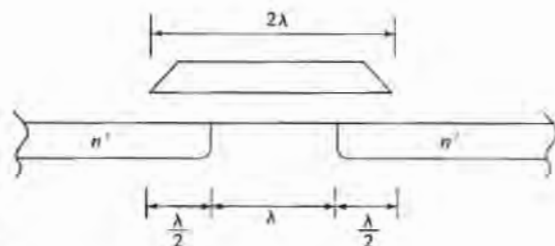


FIGURE 9.15

Channel-length bias in a polysilicon-gate NMOS device caused by lateral diffusion under the edge of the diffusion window. The transistor has $L = 2\lambda$ at the mask level but ends up with an actual $L = \lambda$ after the device is fabricated.



that the source-drain diffusions will move laterally under the edge of the oxide openings, the contact windows can be aligned within $\lambda/2$ of the edge of the diffusions at the mask level, but will still be 1λ within the border of the diffusion in the final structure.

However, lateral diffusion requires the length of channel at the mask level to be increased by λ to achieve the same electrical channel length in the device. The actual channel length $L = L_m - \Delta L$, where L_m is the channel length as originally drawn on the mask and ΔL is the channel-length shrinkage that occurs during processing. This is an important area where the process must be controlled. For devices with short channel lengths, ΔL may be so severe that the devices become unusable. For the layout of Fig. 9.15 width $W_m = 10\lambda$, $W_m/L_m = 10\lambda/2\lambda = 5/1$ at the mask level, whereas $W/L = 10\lambda/\lambda = 10/1$ in the fabricated transistor.

The development of self-aligned polysilicon-gate technology with ion-implanted source-drain regions was a major improvement. The polysilicon-gate process significantly reduces both the channel shrinkage caused by lateral diffusion and the overlap capacitance resulting from alignment tolerances in the metal-gate process.

In Fig. 9.14, one can see another source of channel bias. The "bird's beak" reduces the size of the active region to below that defined by the active region mask, and it introduces a process bias into the channel width of the polysilicon-gate transistor. $W = W_m - \Delta W$, where W_m is the width at the mask level and ΔW is the channel-width shrinkage during processing.

In sets of very tight design rules developed for high-volume-production ICs such as dynamic memories, all critical dimensions are adjusted to account for the processing and alignment sequences. This often results in a layout that must conform to a set of 50 to 100 design rules [7]. Such a set of design rules is highly technology-specific and cannot be transferred from one generation of lithography to the next. The Mead-and-Conway-style rules [6] reach a compromise between a set of rules that is overly pessimistic and wastes a lot of silicon area, and one that is extremely complex, but squeezes out all excess area. The Mead-and-Conway-style design rules are used for ICs in which design time, and not silicon area, is of dominant importance.

9.3 COMPLEMENTARY MOS (CMOS) TECHNOLOGY

Complementary MOS (CMOS) technology is arguably the most commercially important silicon technology. It came to the forefront in the mid 1980s when its low-power benefits finally outweighed the perceived increase in process complexity. Today, scaling of CMOS to submicron dimensions has made the technology highly competitive not only in term of power, but also in raw speed.

9.3.1 The *n*-Well Process

The basic CMOS process of Fig. 1.8 requires a *n*-well diffusion and formation of both NMOS and PMOS transistors. Substrate resistivity is chosen to give the desired NMOS characteristics, and an additional implant step may be introduced to adjust the NMOS threshold separately. The *n*-well-to-substrate junction may range from a few microns to as much as 20 microns in depth. The net surface concentration of the *n*-well must be high enough above the substrate concentration to provide adequate process control without severely degrading the mobility and threshold voltage of the PMOS transistors. The surface concentration of the *n*-well typically ranges between 3 and 10 times the substrate impurity concentration. An additional implant step is often introduced to adjust the PMOS threshold voltage.

9.3.2 *p*-Well and Twin-Well Processes

The first successful CMOS technologies actually utilized *p*-well processes whose structures are simply a mirror image of Fig. 1.8. However, the drive toward ever higher performance led to the development of the *n*-well processes in which the NMOS transistors are placed in the lightly doped substrate region where the *n*-channel

mobility will be the highest. More recently, twin-well processes, such as in Fig. 9.16, have been developed that permit individual optimization of the characteristics of both the n - and p -channel devices [10].

A lightly doped n - or p -type epitaxial layer is grown on a heavily doped n - or p -type substrate. (Lightly doped n - and p -type regions are often referred to as ν and π regions, respectively.) Separate implantations and diffusions are used to form wells for both the NMOS and PMOS transistors. The low-resistivity substrate substantially reduces the substrate resistance R_s and improves latchup resistance as discussed later.

9.3.3 Gate Doping

Early polysilicon gate CMOS processes used n^+ polysilicon gates for both transistors as assumed in the graph of Fig. 9.2. In many processes, it was found that use of an n^+ gate on a PMOS transistor led to formation of a buried channel rather than a surface channel device that causes problems with subthreshold turn-off of the device. With the push toward optimization of both devices with the advent of twin-well processes, p^+ doped polysilicon gates were introduced. With use of the p^+ gate, the PMOS device characteristics become more symmetrical to those of the NMOS devices, except for the inherent mobility differences. The threshold voltages also become symmetrical (See Prob. 9.4). Note that the twin-well process depicted in Fig. 9.16 produces p^+ and n^+ polysilicon gates. The p^+ gate is protected by the photoresist layer during the n^+ implant.

Example 9.3

An n^+ polysilicon gate CMOS process uses an n -type substrate with a doping of $10^{16}/\text{cm}^3$. An implant/drive-in schedule will be used to form a p -well with a net surface concentration of $10^{17}/\text{cm}^3$ and a junction depth of $3\text{ }\mu\text{m}$. (a) What is the drive-in time at 1150°C ? (b) Solve for the implanted dose in silicon. (c) What are the threshold voltages of the n - and p -channel transistors, if the oxide thickness is 10 nm ?

Solution: The $3\text{-}\mu\text{m}$ junction depth and low surface concentration suggest that the well has a Gaussian profile resulting from a two-step diffusion or implant/diffusion process. A final surface concentration of $1.1 \times 10^{17}/\text{cm}^3$ is required to produce a net concentration of $1 \times 10^{17}/\text{cm}^3$ at the surface. Solving for the Dt product yields

$$Dt = x_j^2/4 \ln(N_0/N_B) = 9.38 \times 10^{-9} \text{ cm}^2.$$

At 1150°C , $D = 8.87 \times 10^{-13} \text{ cm}^2/\text{sec}$, which gives $t = 2.94 \text{ h}$. The dose in silicon is given by $Q = N_0\sqrt{\pi Dt} = 1.89 \times 10^{13}/\text{cm}^2$. The p -channel devices reside in the n -type substrate with a doping concentration of $10^{16}/\text{cm}^3$. From Fig. 9.2, the threshold voltage will be -1.1 V . The deep well diffusion will be almost constant near the surface with a value of $10^{17}/\text{cm}^3$. Figure 9.2 yields an n -channel threshold of 0.4 V . A threshold adjustment implant would be needed in this process to increase the n -channel threshold voltage.

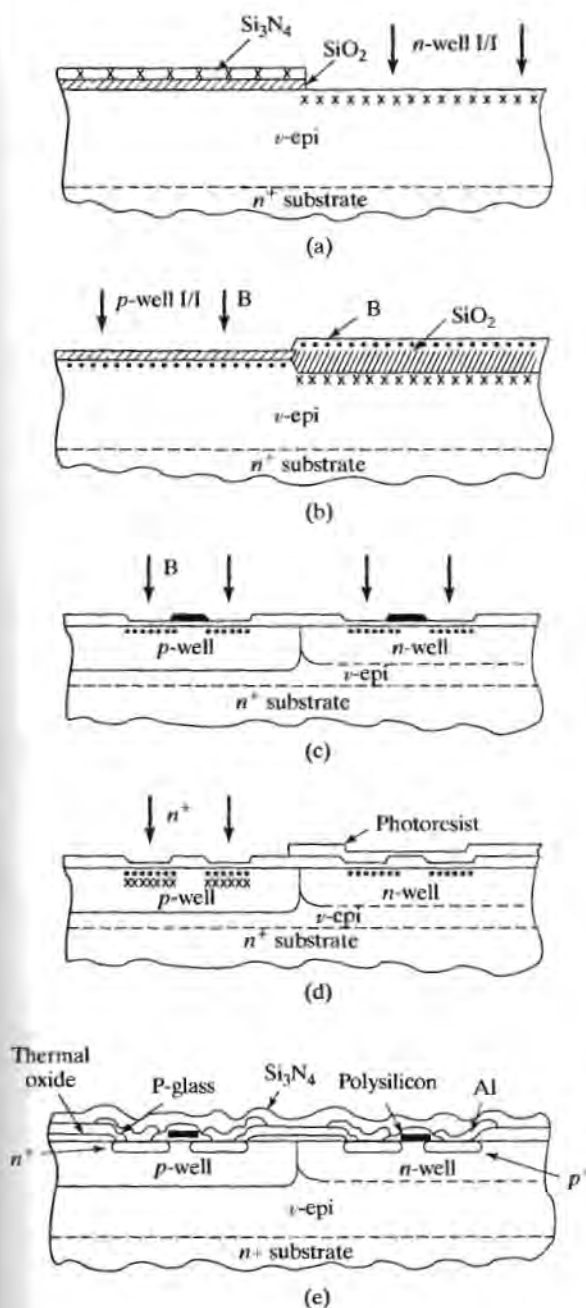


FIGURE 9.16

Twin-well CMOS structure at several stages of the process. (a) *n*-well ion implant; (b) *p*-well implant; (c) nonselective p^+ source/drain implant; (d) selective n^+ source/drain implant using photoresist mask; (e) final structure. Copyright 1980, IEEE. Reprinted with permission from Ref. [10].

In addition to the depletion layer extents, alignment tolerances must be added to the total spacing to ensure that minimum spacing is maintained under worst-case alignment errors. Also, diffusion of deep wells leads to significant lateral diffusion of the well boundary that must be taken into account in the CMOS layout.

9.3.5 CMOS Latchup

Parasitic bipolar devices are formed in the CMOS process in which merged *pnp* and *nnp* transistors form a four-layer (*pnpn*) lateral SCR, as shown in Fig. 9.18. If this SCR is turned on, the device may destroy itself via a condition called *latchup* [8, 9]. The *n*-well depth and the spacings between the source-drain regions and the edge of the *n*-well must be carefully chosen to minimize the current gain of the bipolar transistors and the size of the shunting resistors R_s and R_w . A CMOS process will have a number of additional ground rules not present in an NMOS or PMOS process. A more detailed discussion of the design of bipolar transistors will be given in Chapter 10.

To reduce the resistance of the two shunting resistors, "guard ring" diffusions are sometimes added to the process, as shown in Fig. 9.18. Guard rings can be formed using the source-drain diffusions of the PMOS and NMOS transistors or can be added as separate diffusion steps.

9.3.6 Shallow Trench Isolation

Advanced processes with deep submicron feature sizes make use of shallow trench isolation, as depicted in Fig. 9.19, in which a twin-well process is shown [20]. Both the NMOS and PMOS devices are bounded by the STI oxide region. The STI in combina-

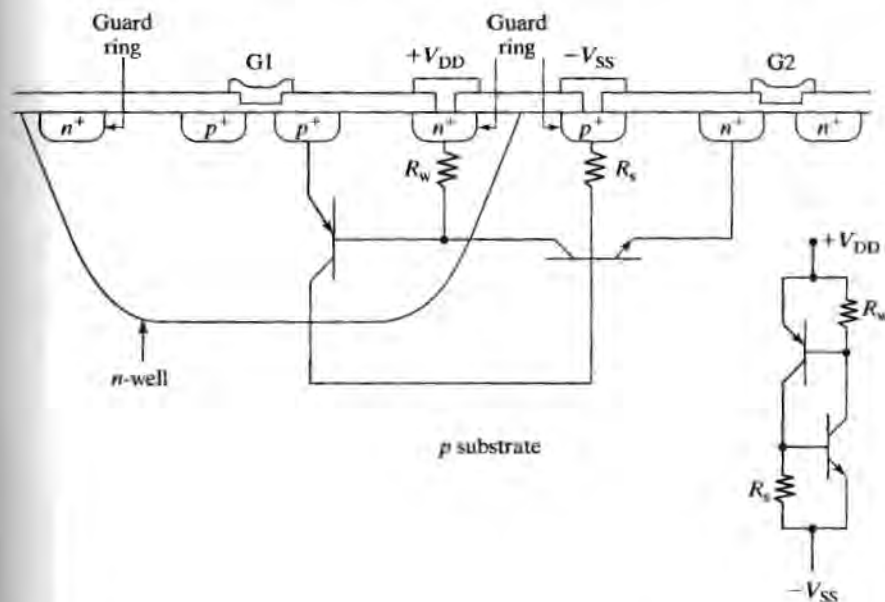


FIGURE 9.18

Cross-section of a CMOS structure, showing the existence of a parasitic lateral *pnpn* SCR and the use of guard rings to reduce the value of R_s and R_w .

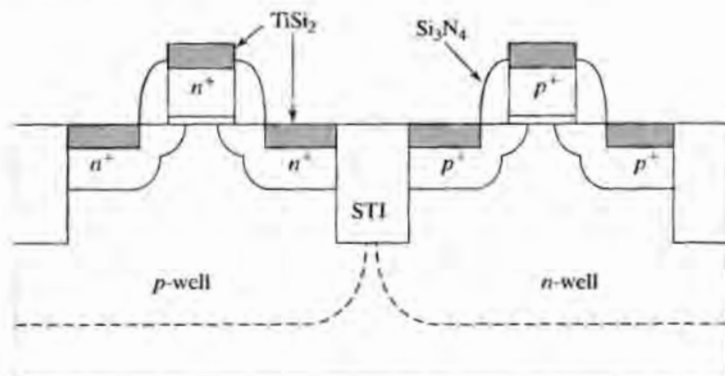


FIGURE 9.19

Application of shallow trench isolation to a twin-well CMOS technology. Copyright IEEE 1998. Reprinted with permission from Reference [20].

tion with a heavily doped substrate eliminates the need for guard rings by substantially reducing the current gain of the bipolar devices and decreasing the value of the shunt resistances. LDD extensions can be noted on both devices, as well as the silicon nitride spacers around the perimeter of the gate. Self-aligned tantalum silicide layers are used to reduce the effective sheet resistance of the polysilicon gate, as well as the source and drain regions. This type of CMOS process is being used for 0.18 μm devices and below.

9.4 SILICON ON INSULATOR

Insulating substrates provide the ultimate in device isolation and freedom from latchup problems. The earliest efforts to achieve an insulating substrate grew thin layers ($< 10 \mu\text{m}$) of single crystal silicon on a sapphire substrate that provides a reasonable match to the silicon crystal lattice. NMOS and PMOS devices were fabricated in the silicon film to produce a CMOS technology. This technology was termed silicon-on-sapphire (SOS) technology. The early attempts were plagued by problems at the silicon-sapphire interface, but the problems were eventually controlled well enough to produce a usable technology.

Our ability to produce a highly controlled silicon-silicon dioxide interface has led to newer forms of silicon-on-insulator (SOI) processes. High-energy ion-implantation can be used to place oxygen atoms in a layer well below the surface of a lightly-doped silicon wafer. Following implantation, the wafers are annealed at elevated temperature to produce a buried oxide layer well below the silicon surface, as depicted in Fig. 9.20. This technology is often referred to as Separation by Implanted Oxygen or

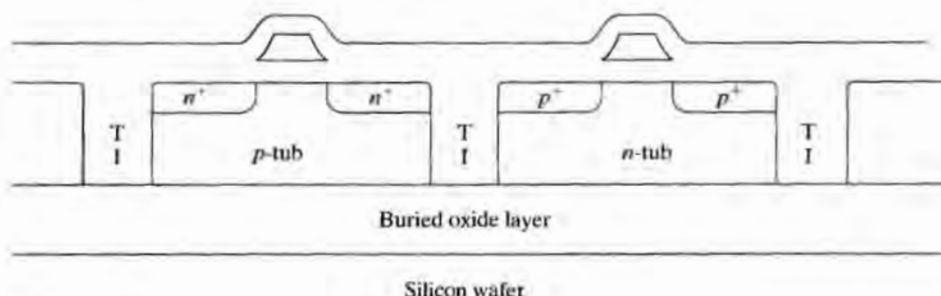


FIGURE 9.20

Trench isolated silicon-on-insulator technology.

SIMOX [19]. Twin-well tubs can then be formed in the lightly doped substrate and separated by trench isolation. NMOS and PMOS devices are then fabricated in the tubs to complete the SOI CMOS process.

Silicon wafer-to-wafer bonding, originally developed for use with MEMs, has also been used to fabricate SOI substrates. A silicon wafer is oxidized to form an insulating SiO_2 layer. A second silicon wafer is brought in contact with the oxidized surface and annealed at elevated temperature to form a bond between the two wafers. Obviously, surface cleanliness is of extreme importance to the success of this process, as indicated in Fig. 9.21. After the bonding is completed, the upper silicon layer is thinned by chemical etching until the desired silicon layer thickness is achieved. An alternative is to use mechanical lapping and polishing processes to thin the silicon wafer.

SUMMARY

In this chapter, we explored the interaction of process design with MOS device characteristics and transistor layout, including the relationships between processing parameters and breakdown voltage, punch-through voltage, threshold voltage, and junction capacitance. A low value of substrate doping is desired to minimize junction capacitance, substrate sensitivity, and junction breakdown voltage, whereas a high substrate doping is needed to maximize punch-through voltage. The use of ion implantation permits the designer to separately tailor the threshold voltage of the transistor.

We have developed basic ideas relating minimum feature size and alignment tolerances and have discussed simple sets of layout design rules. The strong relation between layout design rules and the size of transistors has been demonstrated. Polysilicon-gate technology has been shown to result in a much smaller device area than metal-gate technology for a given transistor W/L ratio, as well as to minimize the parasitic gate capacitance of the device. In addition, the polysilicon-gate process substantially reduces channel-length bias caused by lateral diffusion.

A combination of ion implantation and diffusion is commonly used to form the p - or n -well required for CMOS technology. VLSI CMOS often uses twin-well processes which permit separate optimization of both the n - and p -channel devices.

To achieve a high packing density for submicron processes, trench isolation, which provides excellent isolation between devices, is utilized. The source and drain regions of the transistors can be abutted with the oxide isolation regions. The combination of trench isolation and twin-well processes on a heavily doped substrate suppresses latchup and eliminates the need for guard ring diffusions. The ultimate in

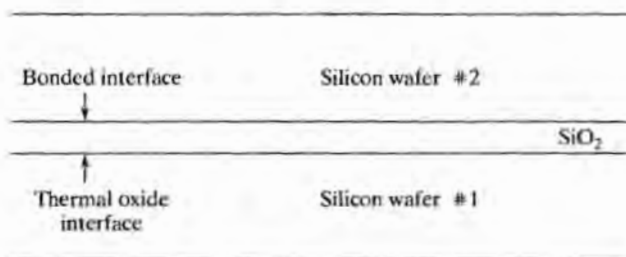


FIGURE 9.21

Formation of bonded wafer SOI

isolation and capacitance reduction is achieved with silicon-on-insulator or SOI substrates. The earliest versions, termed SOS, grew thin silicon layers on sapphire substrates. Today's SOI substrates are formed by high-energy implantation of oxygen or direct wafer-to-wafer bonding followed by chemical etching.

REFERENCES

- [1] R.F. Pierret, *Field Effect Devices*, Volume IV in the Modular Series on Solid State Devices, Addison-Wesley, Reading, MA, 1983.
- [2] S.A. Abbas and R.C. Dockerty, "N-channel Design Limitations due to Hot Electron Trapping," *IEEE IEDM Digest*, pp. 35–38, 1975.
- [3] T.H. Ning, C.M. Osburn, and H.N. Yu, "Threshold Instability in IGFETs due to Emission of Leakage Electrons from Silicon Substrate into Silicon Dioxide," *Applied Physics Letters*, 29, 198–199, 1976.
- [4] P.E. Cottrell and E.M. Buturla, "Steady State Analysis of Field Effect Transistors via the Finite Element Method," *IEEE IEDM Digest*, pp. 51–54, 1975.
- [5] S.M. Sze, *Semiconductor Devices—Physics and Technology*, John Wiley & Sons, New York, 1985.
- [6] C.A. Mead and L. Conway, *VLSI Design*, Addison-Wesley, Reading, MA, 1980.
- [7] Brian Spinks, *Introduction to Integrated Circuit Layout*, Chapter 7, Prentice-Hall, Englewood Cliffs, NJ, 1985.
- [8] A. Ochoa, W. Dawes, and D. Estreich, "Latchup Control in CMOS Integrated Circuits," *IEEE Transactions on Nuclear Science*, NS-26, 5065–5068, December 1979.
- [9] R.S. Payne, W.N. Grant, and W.J. Bertram, "The Elimination of Latchup in Bulk CMOS," *IEEE IEDM Digest*, p. 248–251, December 1980.
- [10] L.C. Parrillo, R.S. Payne, R.E. Davis, G.W. Reutlinger, and R.L. Field, "Twin-Tub CMOS—A Technology for VLSI Circuits," *IEEE IEDM Digest*, p. 752–755, December 1980.
- [11] B.J. Baliga and D.Y. Chen, *Power Transistors: Device Design and Applications*, IEEE Press, New York, 1984.
- [12] K.P. Roenker and L.W. Linholm, "An NMOS Test Chip for a Course in Semiconductor Parameter Measurements," *National Bureau of Standards Internal Report* 84–2822, April 1984.
- [13] T.J. Russell, T.F. Leedy, and R.L. Mattis, "A Comparison of Electrical and Visual Alignment Test Structures for Evaluating Photomask Alignment in Integrated Circuit Manufacturing," *IEEE IEDM Digest*, p. 7A–7F, December 1977.
- [14] D.S. Perloff, "A Four-Point Electrical Measurement Technique for Characterizing Mask Superposition Errors on Semiconductor Wafers," *IEEE Journal of Solid-State Circuits*, SC-13, 436–444, August 1978.
- [15] A. Hori et al., "High Speed 0.1 μm Dual Gate CMOS with Low Energy Phosphorus/Boron Implantation and Cobalt Silicide," *IEEE IEDM Technical Digest*, pp. 575–578, December 1996.
- [16] H. Hwang, D-H Lee and J.M. Hwang, "Degradation of MOSFETs Drive Current Due to Halo Ion Implantation," *IEEE IEDM Technical Digest*, pp. 567–570, December 1996.
- [17] A.J. Auberton-Hervé, "SOI: Materials to Systems," *IEEE IEDM Technical Digest*, pp. 3–10, December 1996.
- [18] T. Hashimoto et al., "A 0.2- μm Bipolar-CMOS Technology on Bonded SOI with Copper Metallization for Ultra High-speed Processors," *IEEE IEDM Technical Digest*, pp. 209–212, December 1998.

- [19] CMOS: From Bulk to SOI, IBIS Technology Corporation *Technical Note*, 1999, <http://www.ibis.com>.
- [20] S. Yang et al., "A High Performance 180 nm Generation Logic Technology," *IEEE IEDM Digest*, pp. 197–200, December 1998.
- [21] W. Maly, *Atlas of IC Technologies: An Introduction to VLSI Processes*, The Benjamin/Cummings Publishing Company, Inc., Menlo Park, CA: 1987.
- [22] R. H. Dennard, F. H. Gaensslen, L. Kuhn and H. N. Yu, "Design of Micron MOS Switching Devices," *IEEE IEDM Digest*, pp. 168, December 1972.

PROBLEMS

- 9.1 What is the maximum gate-to-source voltage that a MOSFET with a 10-nm gate oxide can withstand? Assume that the oxide breaks down at 5 MV/cm and that the substrate voltage is zero.
- 9.2 Two n^+ diffused lines are running parallel in a substrate doped with 10^{15} boron atoms/cm³. The substrate is biased to -2 V, and both lines are connected to $+3$ V. Using one-dimensional junction theory, calculate the minimum spacing needed between the lines to prevent their depletion regions from merging. (b) Repeat for a 3×10^{16} /cm³ doping level.
- 9.3 Use one-dimensional junction theory to estimate the punch-through voltage of a MOSFET with a channel length of 1 μ m. Assume a substrate doping of 3×10^{16} /cm³ and a substrate bias of 0 V.
- 9.4 Plot a graph of thresholds similar to Fig. 9.2, but assume that a p^+ gate is used for the PMOS transistors.
- 9.5 What is the minimum substrate doping required to realize an enhancement-mode NMOS device ($V_{TN} > 0$) with a 10-nm gate oxide?
- 9.6 Calculate the threshold voltage for the NMOS transistor with the doping profile shown in Fig. P9.6. Assume an n^+ polysilicon-gate transistor with a gate-oxide thickness of 20 nm.
- 9.7 An implant with its peak concentration at the silicon surface is used to adjust the threshold of an NMOS transistor. We desire to model this implant by a rectangular approximation similar to that of Figure 9.5. Show that $N_I = N_p \pi/4$ and that $x_I = \Delta R_p \sqrt{8/\pi}$ by matching the first two moments of the two impurity distributions.
- 9.8 A MOS technology is scaled from a 1- μ m feature size to 0.25 μ m. What is the increase in the number of circuits/cm²? What is the improvement in the power-delay product?

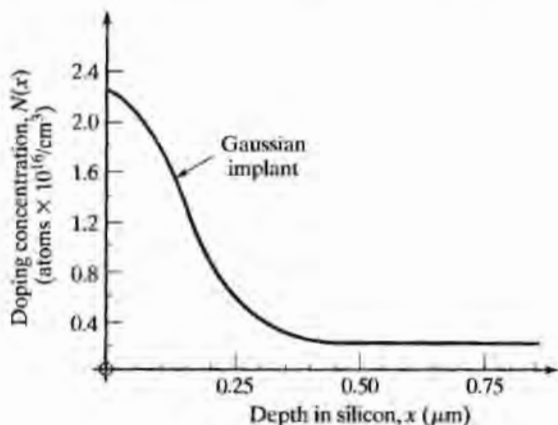


FIGURE P9.6

- 9.9** Suppose that the voltages are not scaled as the dimensions are reduced by a factor of α ? How does the drain current of the transistor change? How do the power/circuit and power density scale?
- 9.10** High-performance NMOS logic processes used depletion-mode NMOS transistors for load devices. This requires a negative threshold, which can be obtained by implanting a shallow arsenic or phosphorus dose into the channel region. Calculate the arsenic dose needed to achieve a -3 -V threshold in an n^+ polysilicon-gate NMOS transistor that has a substrate doping of $3 \times 10^{16}/\text{cm}^3$ and a gate-oxide thickness of 50 nm.
- 9.11** Draw the p -well version of the CMOS process in Fig. 1.8.
- 9.12** Our design rule examples used an alignment tolerance that was one-half the feature size. This ratio represents a very loose alignment capability. Develop a new set of design rules similar to those of Fig. 9.13 for $T = \alpha$ and $F = 4\alpha$. Draw the new minimum-size polysilicon-gate transistor using your rules. Compare the area of your transistor with the area of the transistor of Fig. 9.13 if $\lambda = 2\alpha$.
- 9.13** An n -well CMOS process starts with a substrate doping of $3 \times 10^{15}/\text{cm}^3$. The well doping near the surface is approximately constant at a level of $3 \times 10^{16}/\text{cm}^3$. The gate-oxide thicknesses are both 15 nm.
- Calculate the thresholds of the n - and p -channel transistors using Eqs. (9.2). Assume n^+ polysilicon gates.
 - Calculate the boron doses needed to shift the NMOS threshold to $+1$ V and the PMOS threshold to -1 V. Assume that the threshold shifts are achieved through shallow ion implantations. Neglect oxide charge.
- 9.14** Early CMOS logic circuits operated from power supplies of 8 V or more. Estimate the minimum spacing between the drains of adjacent NMOS and PMOS transistors in a CMOS process if the substrate doping is $3 \times 10^{15}/\text{cm}^3$, the well doping is $5 \times 10^{16}/\text{cm}^3$, and the maximum drain-substrate voltage is 8 V. Assume that the well is also reverse biased by 8 V.
- 9.15** A twin-well process starts with a 3- μm -thick, 10- $\Omega\text{-cm}$ v epi layer on an n^+ substrate.
- A p -well is to be formed by ion-implantation followed by a drive-in diffusion and is to have a surface concentration of $10^{16}/\text{cm}^3$ with a depth of 2 μm . What are the drive-in time at a temperature of 1075°C and impurity dose in silicon? What is the lateral diffusion distance of the well?
 - A phosphorus n -well is to be formed in the same substrate with a surface concentration of $5 \times 10^{16}/\text{cm}^3$ and a depth of 1.5 μm . What are the drive-in time at a temperature of 1075°C and impurity dose in silicon? What is the lateral diffusion distance of the n -well?
 - What is the total out-diffusion from the n^+ substrate following the formation of both wells if the substrate is arsenic doped with an arsenic concentration of $10^{20}/\text{cm}^3$?
- 9.16** Draw a composite view of the situation resulting from a worst-case misalignment of the masks for the MOSFET layout shown in Fig. 9.12. Assume that metal aligns to thin oxide and that thin oxide and contacts align to the diffusion.
- 9.17** Develop a new set of ground rules for the metal-gate transistor of Section 9.2, assuming that levels 2, 3, and 4 are all aligned to level 1. Redraw the transistor of Fig. 9.12 using your new rules. In what ways is this layout better or worse than that originally given in Fig. 9.12?
- 9.18** Draw a cross section of a metal-gate NMOS transistor and a composite view of its mask set, assuming an aggressive layout that takes into account all lateral diffusion. Assume a source-drain junction depth of 2.5 μm , and assume that lateral diffusion equals 80% of vertical diffusion. Assume λ is 2 μm and $W/L = 10/1$.
- 9.19** Draw the layout of a three-input NMOS NOR-gate with the dimensions given on the circuit schematic in Fig. P9.19. Be sure to merge diffusions wherever possible. Use the more aggressive ground rules developed for polysilicon-gate devices.

- 9.20** Draw the layout (top view) of the CMOS inverter in Fig. P9.20 for an n -well technology using the λ -based ground rules from Fig. 9.13 for the transistors. In addition, assume that source and drain regions must be a minimum of 8λ from the edge of the well. What is the total area of the CMOS inverter (in λ^2)? What is the total gate area?
- 9.21** Repeat Problem 9.20, but this time assume shallow trench isolation with a minimum width of 4λ . Assume that the source and drain regions can butt against the oxide, as in Fig. 9.19.
- 9.22** A high energy (4 MeV) is used to implant oxygen well below the silicon surface in order to form a buried SiO_2 layer. Assume that the SiO_2 layer is desired to be $0.25\ \mu\text{m}$ wide. (a) What is the oxygen dose required in silicon? (b) What beam current is required to achieve a throughput of five 200 mm wafers per hour? (c) How much power is being supplied to the ion beam?

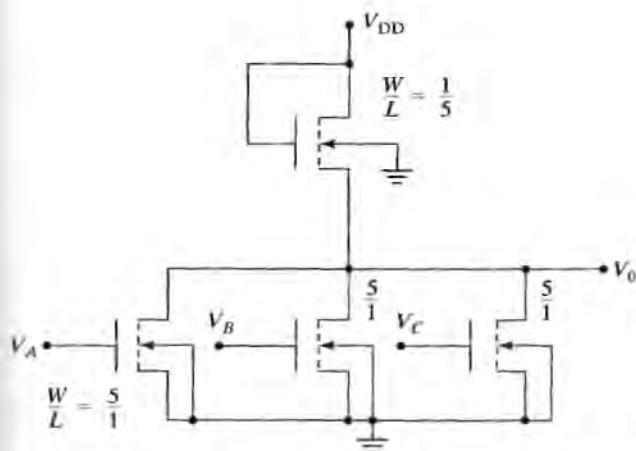


FIGURE P9.19

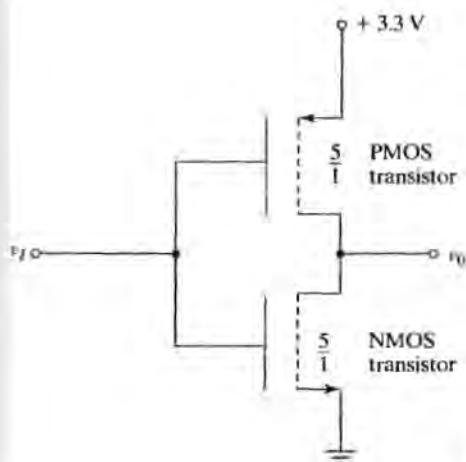


FIGURE P9.20

CMOS inverter with both $W/L = 5/1$

9.23 A number of types of alignment test structures have been developed [12, 13]. Figure P9.23 shows a simple test structure that can be used to measure the misregistration of the contact window mask relative to the diffusion mask [14]. Two linear potentiometers, one in the horizontal direction and one in the vertical direction, are fabricated using diffused resistors. The distance between contacts *A* and *C* is the same as that between *C* and *E*, and the contact from pad *D* is nominally one-half the distance between pads *C* and *E*. A current is injected between pads *B* and *F*, and the voltages between pads *C*–*D* and *D*–*E* are measured.

- (a) Show that the misregistration in the *y*-direction is given by $\Delta Y = \frac{1}{2} L (V_{DE} - V_{CD})/V_{AC}$.
- (b) Derive a similar relationship for misregistration in the *x*-direction.

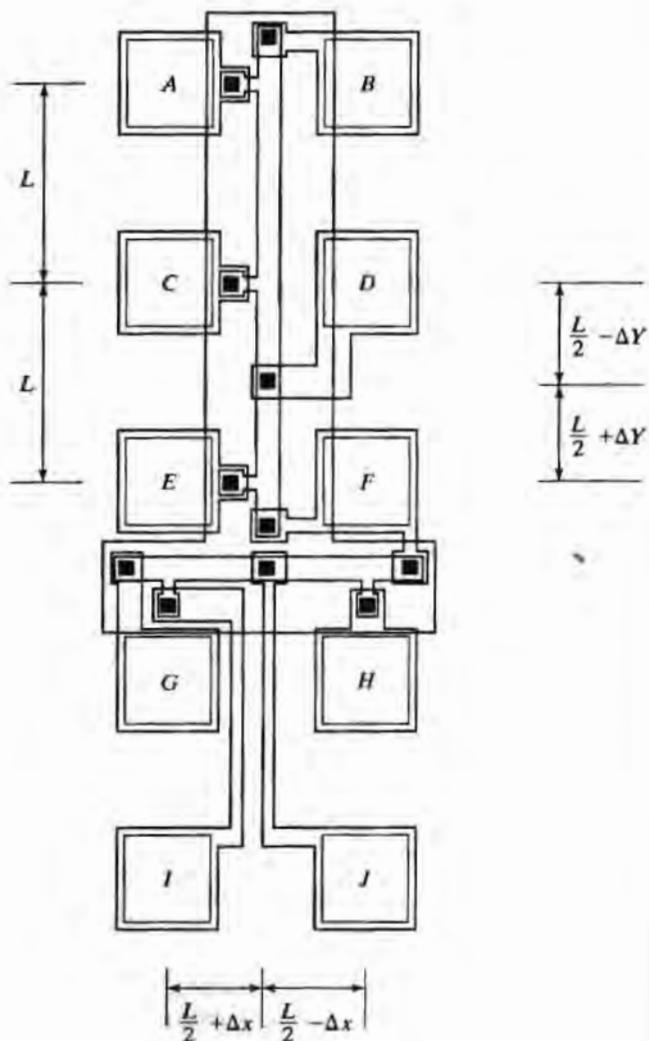


FIGURE P9.23

Bipolar Process Integration

In this chapter, interactions between fabrication processes and bipolar device design and layout will be explored. In particular, we will look closely at relationships between impurity profiles and device parameters such as current gain, transit time, and breakdown voltage. Basic design rules for bipolar structures are introduced. The use of recessed oxidation, deep and shallow trenches, polysilicon electrodes and self-aligned processes in the formation of high-performance bipolar transistors will be presented. Dielectric and collector-diffused isolation processes are discussed, as well as silicon-germanium epitaxial-base transistors and advanced BiCMOS technologies, which provide bipolar and CMOS devices.

10.1 THE JUNCTION-ISOLATED STRUCTURE

The classic SBC process provides a backdrop for understanding the limitations of the basic bipolar transistor, as well as the structure of various other devices that are fabricated in bipolar IC processes. The basic junction-isolated bipolar process of Fig. 10.1 has been used throughout the IC industry for many years and has become known as the *standard buried collector* (SBC) process. In this junction-isolated process, adjoining devices are separated by back-to-back pn junction diodes that must be reverse biased to ensure isolation. (See Fig. 10.1(b).) The SBC process remains the primary bipolar process for analog and power circuit applications with power supplies exceeding 15 V. Although the SBC process was also originally used for logic circuits, most digital technologies have evolved to self-aligned, oxide-isolated processes using polysilicon and other technology advances first developed for MOS processes. Wafers with a $\langle 111 \rangle$ surface orientation were used specifically for bipolar fabrication for many years. However, in the past few years, it has become common to find bipolar processes also using $\langle 100 \rangle$ substrate material, which facilitates transfer of processes from MOS technology. Certainly, all BiCMOS technologies utilize $\langle 100 \rangle$ material.

The process flow for the SBC structure of Fig. 10.1(b) was discussed in Section 1.4 and will only be outlined here. An n^+ buried layer is formed by selective diffusion into a $\langle 111 \rangle$ -oriented p -type substrate and is followed by growth of an n -type epitaxial

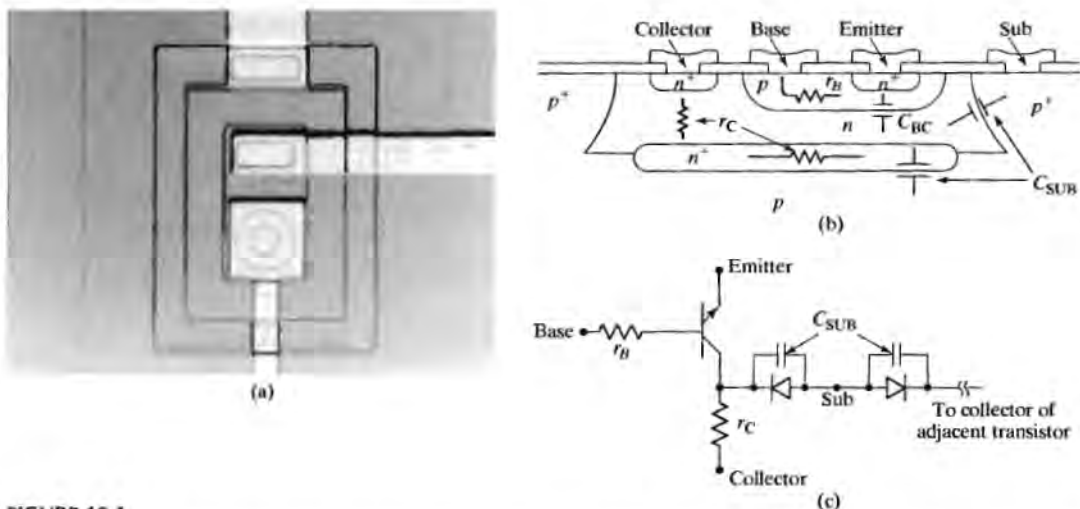


FIGURE 10.1

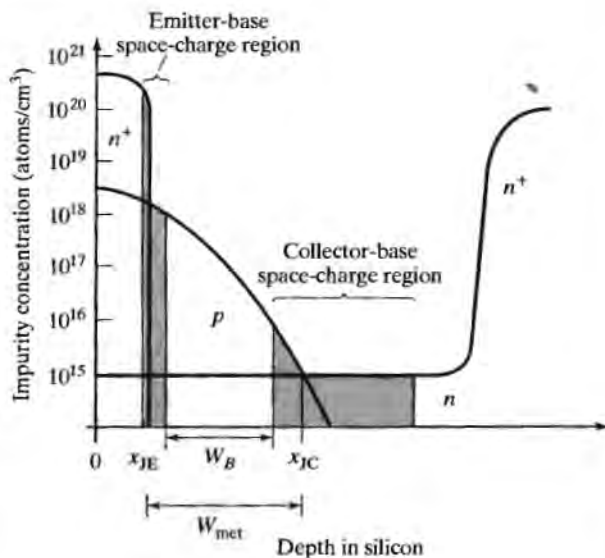
(a) Photo of an SBC transistor; (b) Cross section of a transistor fabricated with the SBC process showing the collector-base capacitances and the base and collector series resistances; (c) lumped circuit model for the transistor showing back-to-back diodes, which provide isolation between adjacent transistors.

layer. Isolated n -type collector islands are formed using a deep boron diffusion that surrounds the collector island. The base and emitter are formed by successive p - and n -type diffusions into the epitaxial layer. The structure is completed with contact window formation and metallization.

A cross section of the SBC impurity profile through the center of the device is shown in Fig. 10.2. In the next several sections, we will consider how the design of this profile is related to several important measures of device performance. An understanding of the basic profile design for the SBC process will help us see the advantages and disadvantages of other types of processes.

FIGURE 10.2

Vertical impurity profile in typical bipolar junction transistor. The shaded regions represent the emitter-base and collector-base space-charge regions. The metallurgical basewidth and electrical basewidth are indicated by W_{met} and W_B , respectively.



10.2 CURRENT GAIN

To be useful in circuits, the bipolar transistor must have a current gain of at least 10 to 20 for digital applications and an order of magnitude greater for analog applications. An expression for the current gain of the bipolar transistor is

$$\beta^{-1} = \frac{G_B}{G_E} + \frac{W_B^2}{\eta L_B^2} \quad (10.1)$$

The basewidth, W_B , is the width of the electrically neutral base region of the transistor. The constant η is determined by the shape of the impurity profile in the base and ranges from 2 to 20. L_B is the diffusion length for minority carriers in the base, and G_B and G_E are called the *Gummel numbers* in the base and emitter, respectively.

The *Gummel numbers* are defined by

$$G_B = \int_{\text{base}} \frac{N(x)}{D_B(x)} dx \quad \text{and} \quad G_E = \int_{\text{emitter}} \frac{N(x)}{D_E(x)} dx \quad (10.2)$$

where D_E and D_B are the minority-carrier diffusion constants in the emitter and base. Heavy doping effects in the emitter typically limit the value of G_E to 10^{13} to 10^{14} sec/cm⁴. The basewidth is defined by the distance between the edges of the two space-charge regions in the base. For wide-base transistors, this is approximately equal to the distance between the metallurgical junctions, as shown in Fig. 10.2. For narrow-base transistors, the space-charge regions must be subtracted from the metallurgical basewidth, as discussed further in Section 10.4.

For large current gain, Eq. (10.1) should be as small as possible. The ratio of the Gummel numbers in the base and emitter should be low, the width of the base region should be small, and L_B should be large. Figure 10.3 shows the dependence of the diffusion length on impurity concentration. As the doping level increases, L_B decreases, but is greater than 10 μm for typical base-doping concentrations. In modern high-frequency transistors, the basewidth W_B is typically far less than the diffusion length L_B , and the first term in Eq. (10.1) determines the current gain.

From Eqs. (10.1) and (10.2), the emitter must be heavily doped relative to the base in order to obtain high gain. In fabricating a bipolar transistor, each successive diffusion is heavier than the last, and the final n^+ diffusion naturally performs best as the emitter. Thus, the n^+ layer nearest the surface is used as the emitter.

Example 10.1

Estimate the current gain for a transistor with the following parameters: $N_E/D_E = 5 \times 10^{13}$ sec/cm⁴, $N_B/D_B = 10^{12}$ sec/cm⁴, $W_B = 1 \mu\text{m}$, $L_B = 20 \mu\text{m}$, and $\eta = 10$.

Solution: Plugging these parameters into Eq. (10.1) yields $\beta^{-1} = 0.02 + 0.00025$ and $\beta = 50$. In this transistor, the current gain is dominated by the ratio of the Gummel number terms.

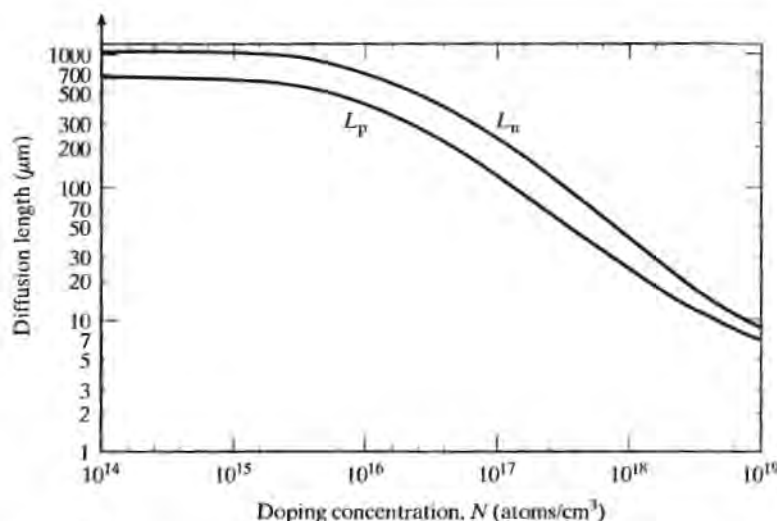


FIGURE 10.3

Calculated minority-carrier diffusion lengths as a function of doping concentration for bulk silicon using lifetime equations from Ref. [1].

10.3 TRANSIT TIME

Another important bipolar device parameter is the delay incurred during carrier propagation between the emitter and collector terminals of the transistor. Both logic switching speed and amplifier frequency response are limited by the *transit time*, which is defined by

$$\tau = r_E C_{BE} + W_B^2 / \eta D_B + (C_{JC} + C_{sub}) r_C + X_C / 2V_s \quad (10.3)$$

The unity-gain frequency of the transistor, f_T , is given approximately by

$$f_T = \frac{1}{2\pi\tau} \quad (10.4)$$

The first term is the product of the small-signal resistance of the emitter region r_E and the emitter-base capacitance C_{BE} . The second term in Eq. (10.3), called the *base transit time*, represents the time required for a carrier to move across the neutral base region W_B . The third term is the delay associated with charging the capacitances connected to the collector node through the collector series resistance r_C . The capacitances C_{JC} and C_{sub} are determined by the collector-base and collector-substrate junction areas and by the doping concentrations of the base, collector, and substrate regions. The last term is the delay time associated with a carrier crossing the depletion region of the collector-base junction. X_C is the width of the depletion layer and V_s is the saturation velocity of the carriers.

In order to minimize τ , the basewidth is made as narrow as possible, the buried layer is added to minimize the value of r_C , and light doping is used to minimize the capacitances. Estimates of the capacitance of the junctions can be made using the one-sided step-junction expression (Eq. (9.3)) in which the capacitance is determined by the concentration on the lightly-doped side of the junction.

Example 10.2

Calculate the transit time for a bipolar transistor with the following parameters: $r_E = 25 \Omega$, $C_{BE} = 10 \text{ pF}$, $W_B = 1 \mu\text{m}$, $\eta = 10$, $D_B = 20 \text{ cm}^2/\text{sec}$, $C_{JC} + C_{sub} = 2 \text{ pF}$, $r_C = 250 \text{ ohms}$, $X_C = 10 \mu\text{m}$, and $V_s = 10^7 \text{ cm/sec}$.

Solution: Substituting these values into Eq. (10.3) gives the following values for the four terms: $0.25 \times 10^{-9} \text{ sec}$; $0.05 \times 10^{-9} \text{ sec}$; $0.5 \times 10^{-9} \text{ sec}$; $0.05 \times 10^{-9} \text{ sec}$. The resulting value of transit time is $0.85 \times 10^{-9} \text{ sec}$. The unity-gain frequency f_T is equal to 188 MHz.

Another important measure of the high-frequency performance of bipolar transistors is the product of the base resistance and collector-base capacitance, $r_B \cdot C_{BC}$. This product can be shown to limit the gain-bandwidth product of single-stage amplifiers [6]. A narrow basewidth increases the value of r_B . Employing heavier base doping counteracts the increase in base resistance, but increases the values of C_{BC} and C_{BE} . Choosing the base profile for optimum $r_B C_{BC}$ product is a delicate design issue.

10.4 BASEWIDTH

Equations (10.1) through (10.3) indicate that device performance is improved by making the basewidth as narrow as possible. The primary restrictions on reducing the basewidth are set by breakdown-voltage requirements and by tolerances on the basewidth, due to variations in process control. For low-voltage logic devices, the metallurgical basewidth may be less than $1 \mu\text{m}$. For higher voltage devices used in analog circuit or power applications, the basewidth must be wide enough to support the collector-base depletion-layer width under large reverse bias.

The actual basewidth of the transistor is determined by reducing the metallurgical basewidth by the portions of the emitter-base and collector-base space-charge regions, which protrude into the base as shown in Fig. 10.2. The emitter and base are both heavily doped near the base-emitter junction, and although the space-charge-region width of the emitter-base junction is usually quite small, it does extend almost entirely into the base. Its width can be estimated from Fig. 9.4.

The collector-base space-charge-region width is dependent on the voltage across the junction and extends into both the base and collector regions. Figure 10.4 shows the depletion-layer width on either side of a pn junction formed by a Gaussian diffusion into a uniformly doped substrate, the normal situation for a bipolar transistor fabricated using the SBC process.

Total space-charge region width X_T as a function of the ratio of applied voltage to background concentration (V/N_B) appears in Fig. 10.4(a), whereas the division of the total between the heavily doped side (x_1) and the lightly doped side (x_2) appears in the second half of the same figure.

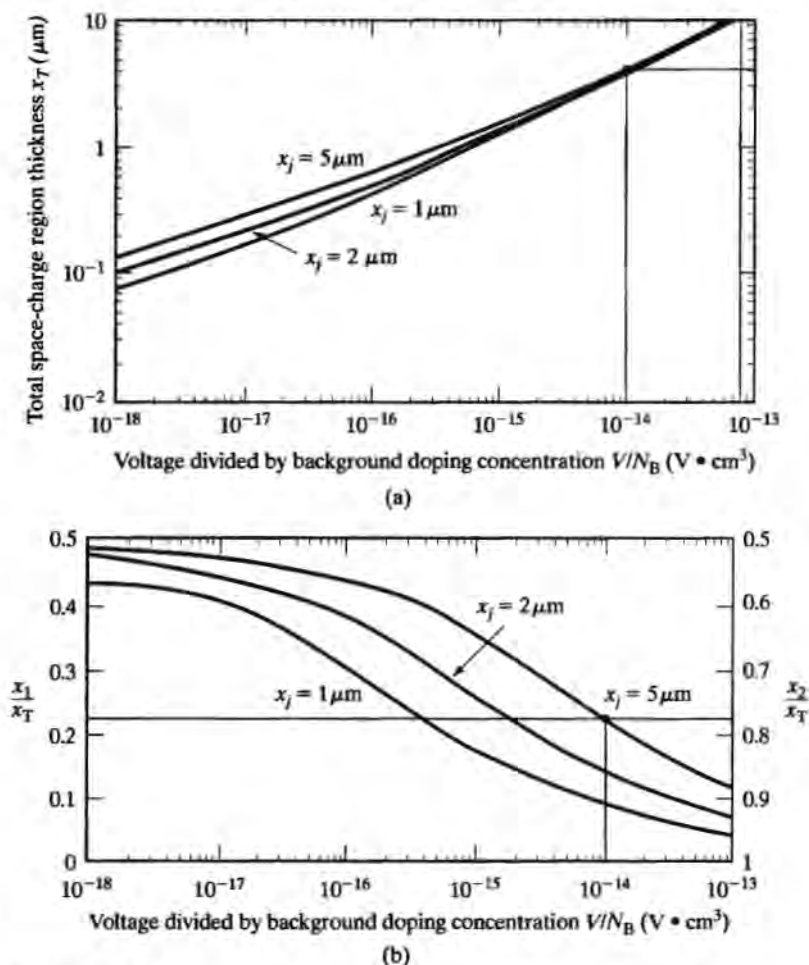


FIGURE 10.4

The space-charge region width as a function of voltage and doping for a pn junction formed by a Gaussian diffusion into a uniformly doped substrate. (a) Total space-charge region width x_T ; (b) fraction of total space-charge region width extending on the heavily doped side, x_1 , and on the lightly doped side, x_2 , respectively. After Ref. [4]. Reprinted with permission from the *AT&T Technical Journal*. Copyright 1960 AT&T.

Example 10.3

Estimate the space-charge region widths on each side of the collector-base junction of a bipolar transistor fabricated on a 1-ohm-cm n -type epitaxial layer. The reverse-bias voltage across the junction is 40 V, and the collector-base junction depth is 5 μm .

Solution: The doping of the epitaxial layer is $4 \times 10^{15} / \text{cm}^3$, giving a value of $V/N_B = 1 \times 10^{-14} \text{ V} \cdot \text{cm}^3$. From Fig. 10.4(a), the total depletion-layer width is approximately 4 μm . From Fig. 10.4(b), 77%, or 3.1 μm , extends into the collector region, and 23%, or 0.9 μm , extends into the base region.

Increased base doping reduces the size of the space-charge regions in the base, permitting a narrow-base design. However, heavy base doping tends to increase the Gummel number in the base, which reduces the current gain of the transistor. Heavy doping also increases the collector-junction capacitance, thus increasing the transit time in Eq. (10.3). This is another situation in which conflicts arise when trying to optimize several different device parameters simultaneously.

10.5 BREAKDOWN VOLTAGES

The process designer must understand the magnitude of the voltages that will be applied to the transistors in circuit applications. The device in Ex. 10.3 had to withstand 40 V and was probably designed for analog-circuit applications. On the other hand, transistors designed for logic applications must support only relatively low voltages. For example, the devices used in TTL circuits are designed to withstand only 7 V. The multi-gigahertz oxide-isolated devices described subsequently in Section 10.8 often have breakdown voltages of 2.5 V or less.

10.5.1 Emitter-Base Breakdown Voltage

The emitter-base breakdown voltage is determined by the doping concentration and radius of curvature of the junction, as was discussed previously in Section 9.1.2. Breakdown occurs first in the region of the junction where the electric field is the largest, usually corresponding to the portion of the junction where the doping levels and curvature are the highest. The actual breakdown voltage is then determined by the doping on the more lightly doped side of the junction.

To achieve high current gain, the emitter region is doped heavily, and the breakdown voltage of this junction will be determined by the impurity concentration of the more lightly doped base region. The base impurity concentration is highest at the surface, so the emitter-base junction will tend to break down first at the surface. The curvature of the junction enhances the electric field and reduces the breakdown voltage. Figure 10.5 gives the breakdown voltage of the emitter-base junction as a function of

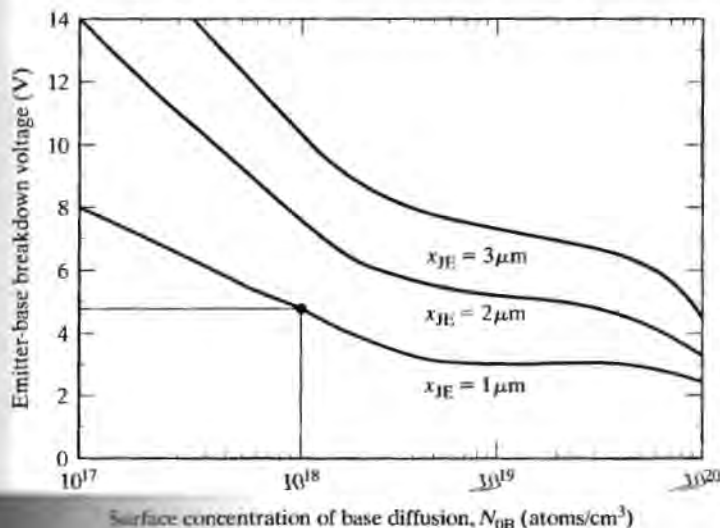


FIGURE 10.5

Emitter-base junction breakdown voltage as a function of base surface concentration with emitter-base junction depth as a parameter. After Ref. [4]. Reprinted with permission from *Solid-State Electronics*, Vol. 17, P. R. Wilson, "The Emitter-Base Breakdown Voltage of Planar Transistors." Copyright 1974, Pergamon Press, Ltd.[3]

the final surface concentration of the base region, with junction radius as a parameter. Emitter-base breakdown voltages are low because of the relatively large impurity concentrations on both sides of the junction.

Example 10.4

An *npn* transistor has a 1- μm -deep emitter-base junction and the base diffusion given in Example 4.2. What is the expected breakdown voltage of this junction?

Solution: The base-region surface concentration in Fig. 4.9 is $1.1 \times 10^{18}/\text{cm}^3$. Figure 10.5 predicts the breakdown voltage of a 1- μm -deep junction with this surface concentration to be approximately 4.8 V. (The emitter-base junctions of most common bipolar transistors will break down well below 10 V.)

10.5.2 Circular Emitters

Although the structures shown thus far have been drawn with square or rectangular emitter regions, circular emitters (see Fig. 10.6) commonly appear in technologies used for analog applications. Matching tends to be better with circular emitters, and the use of circular emitters increases the radius of curvature of the junction and therefore the breakdown voltage of the emitter-base junction. (See Figs. 9.3 and 10.5.) Figure 10.6 shows a quad of cross-connected *npn* transistors that are often used to improve matching in differential amplifiers.

10.5.3 Collector-Base Breakdown Voltage

The bipolar transistor can begin to conduct excessive collector current by two mechanisms. The first is Zener or avalanche breakdown of the collector-base junction. As previously discussed, breakdown is localized to the region where the doping concentrations

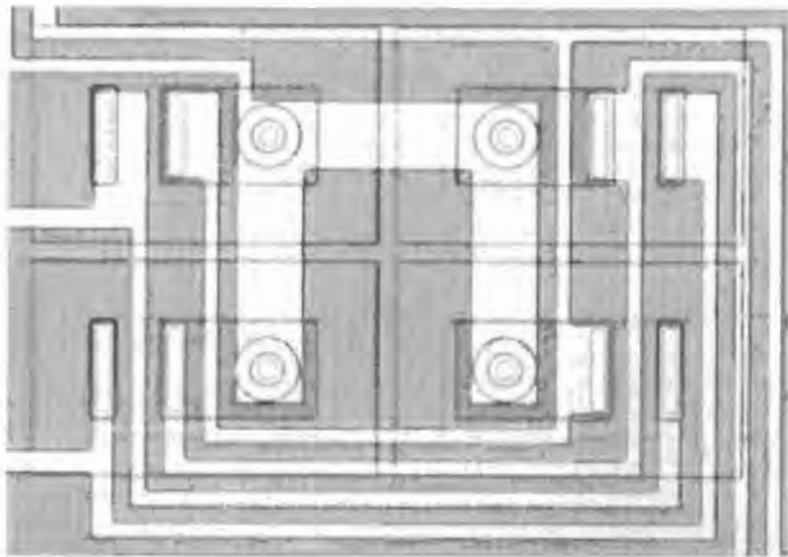


FIGURE 10.6

Cross-connected quad of transistors with circular emitters.

are the largest, but it is determined primarily by the doping concentration on the more lightly-doped side of the junction. The collector is formed in the uniformly doped epitaxial layer. The base region is diffused into the epitaxial layer and is the more heavily doped side of the junction. Since the collector is uniform, junction breakdown will occur first where the electric field is enhanced by junction curvature. The breakdown voltage for the collector-base junction as a function of epitaxial-layer impurity concentration and junction radius is given in Fig. 10.7.

The second breakdown mechanism is punch-through of the base region and is illustrated in Fig. 10.8. The epitaxial layer is more lightly doped than the base region, and the collector-base junction depletion layer extends predominantly into the epitaxial layer. As the collector-base voltage increases, the depletion layer expands further into the epitaxial layer and will eventually hit the n^+ buried layer. At this point, further depletion-layer expansion will occur in the base, and any increase in collector-base voltage will quickly punch through the remaining base region.

The second set of curves in Fig. 10.7 shows collector-base junction breakdown limitations set by punch-through. As already described, the primary parameters determining the punch-through voltage are the epitaxial-layer doping and the width, $X_{BL} - X_{BC}$, of the region between the collector-base junction and the n^+ buried layer.

Example 10.5

What is the collector-base breakdown voltage of a transistor with a 10- μm -thick epitaxial layer doped at a level of $10^{15}/\text{cm}^3$ if the collector-base junction depth is 5 μm ? Assume that the buried layer has diffused upward 2 μm .

Solution: First, determine the avalanche breakdown voltage of an isolated pn junction. Figure 10.7 gives a breakdown voltage of approximately 130 V for a doping of $10^{15}/\text{cm}^3$. Next, we must also check the punch-through limitations for $X_{BL} - X_{BC} = 3$ μm . From Fig. 10.7 the transistor will punch through at approximately 30 V. So the collector-base breakdown voltage is limited to 30 V by punch-through in this transistor.

As usual, different device requirements produce conflicting design constraints. High Zener breakdown voltage requires low epitaxial-layer doping. Low epitaxial-

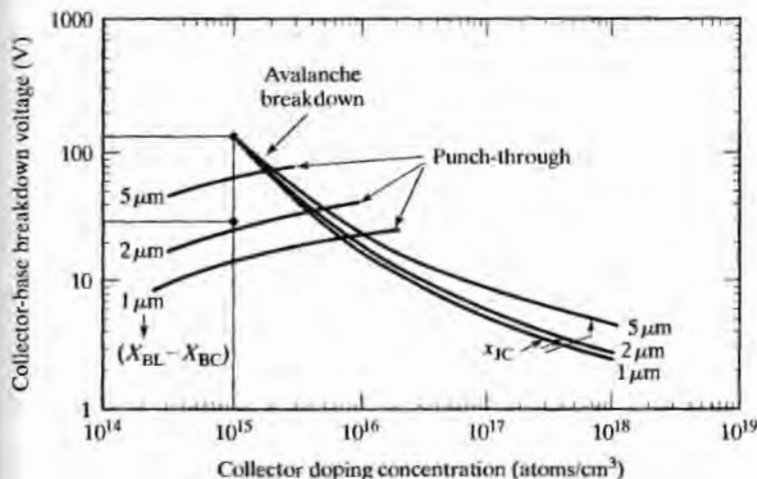


FIGURE 10.7

Collector-base junction breakdown voltage as a function of collector-doping concentration with collector-base junction depth and punch-through limits as parameters. After Ref. [4]. Reprinted from the *Journal of the Electrochemical Society*, Volume 113 (1966), pages 508-510, by permission of the publisher, The Electrochemical Society, Inc., [19] and Pergamon Press, Ltd. [20]

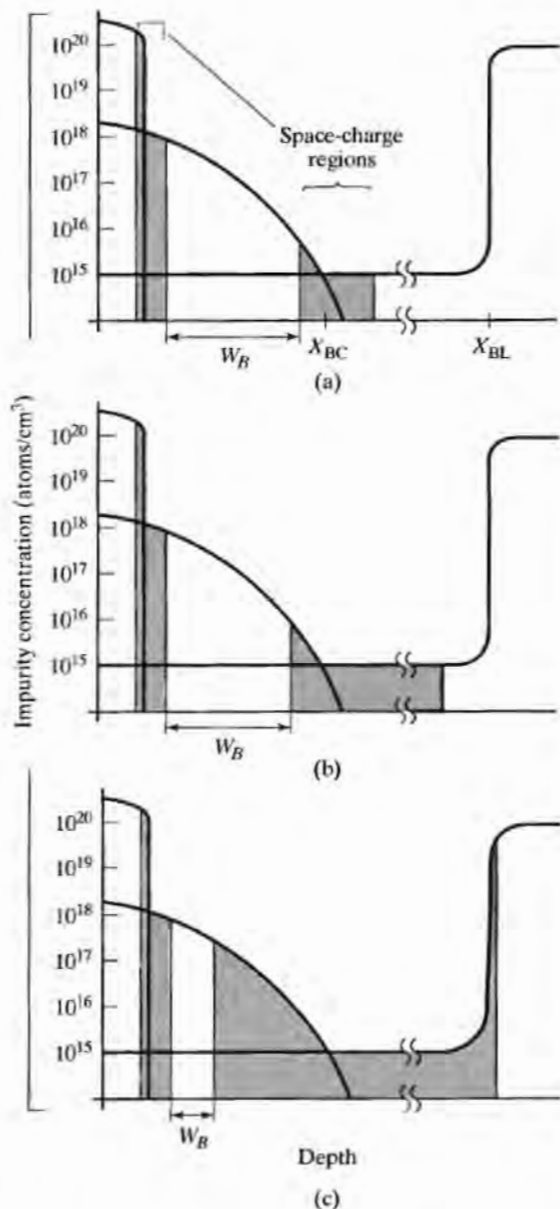


FIGURE 10.8

Collector-base space-charge region growth as the collector-base voltage is increased. (a) Zero bias; (b) intermediate collector-base voltage; (c) large collector-base voltage just below the punch-through voltage.

layer doping requires a relatively wide epitaxial-layer thickness in order to prevent punch-through. However, a wide depletion-layer width in the epitaxial layer increases the transit time and reduces the frequency response of the device.

10.6 OTHER ELEMENTS IN THE SBC TECHNOLOGY

Obviously, other electronic elements, such as resistors, diodes, and *pnp* transistors, are required to build circuits in bipolar technology. Several resistor structures are available with widely divergent values of sheet resistance. The SBC technology is optimized

TABLE 10.1 Resistors in the SBC Technology.

Resistor Layer	Sheet Resistance (Ω/\square)	Absolute Tolerance (%)	Matching (%)
Emitter Diffusion	5–20	20	2
Base Diffusion	100–200	20	0.2–2
Epitaxial Layer	1000–5000	30	5
Pinched Base	2000–10,000	50	10
Ion Implanted	100–1000	3	0.1–1

around the *npn* transistor whose characteristics are dominated by high-mobility electron transport across the base region. High-quality *pnp* transistors are more difficult to form in this technology. However, there are two types of useful *pnp* transistors that can readily be fabricated. The first is a *substrate pnp*, which has its collector permanently tied to the substrate potential. The second is the *lateral pnp*, which does provide uncommitted collector, base, and emitter terminals, but has much poorer current gain and frequency response compared with the vertical *npn* transistor.

10.6.1 Emitter Resistor

A resistor may be formed from the emitter diffusion, as depicted in Fig. 10.9. The emitter layer is a low-sheet-resistance region and is therefore useful only for relatively small-value resistors. A meandering resistor pattern is often used to achieve the

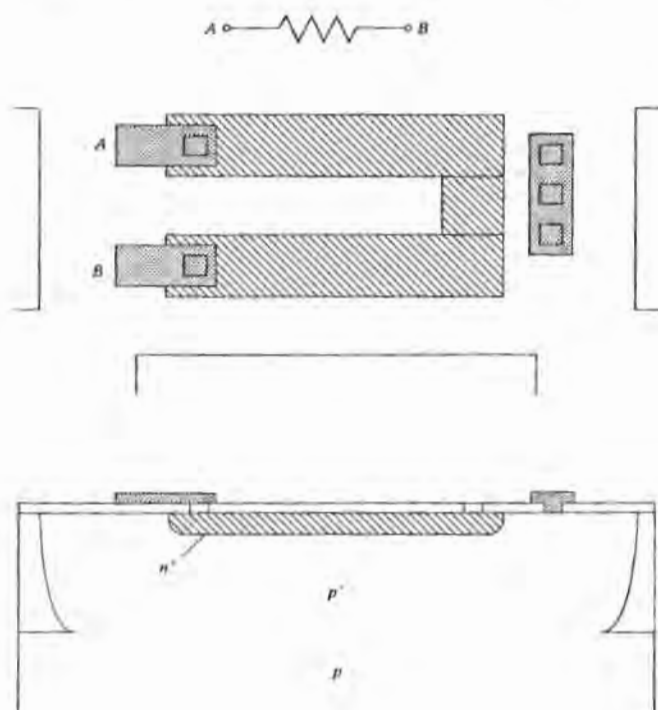


FIGURE 10.9

Resistor using the n^+ emitter diffusion into the isolation region.

required number of squares in the resistor. (See Prob. 4.10.) In Fig. 10.9, the emitter layer is placed directly in the relatively heavily doped p -type isolation region, which is normally connected directly to the most negative power supply. Resistor-substrate isolation will be maintained by a reverse bias across the pn junction. The resistor body represents an n^+-p^+ junction, which leads to a relatively large junction capacitance per unit area and low breakdown voltage. (See Fig. 7.8.)

10.6.2 Base Resistor

The most common resistor utilizes the base diffusion within an isolated n -type epitaxial region, as in Fig. 10.10. Here, again, a meandering resistor pattern can be used to achieve the desired resistance value. The resistor body must be kept reverse biased

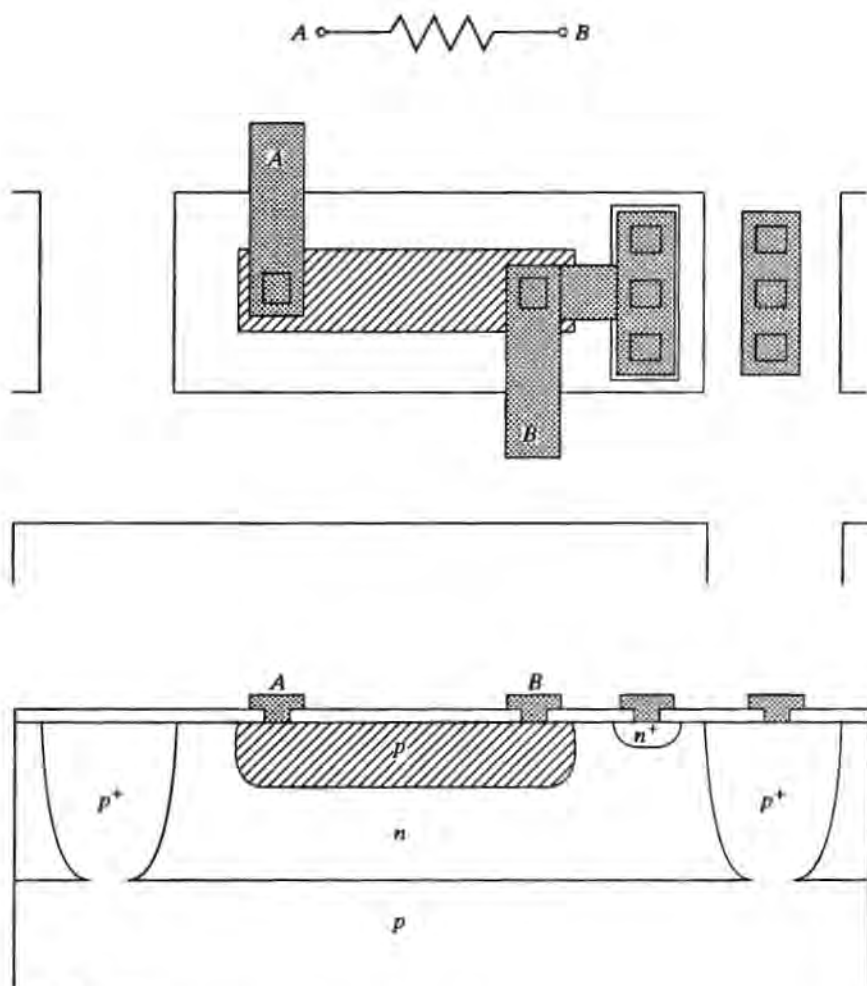


FIGURE 10.10

Resistor formed from the standard p -type base diffusion.

along its length, so the n -region contact is tied to either the most positive end of the resistor or directly to the most positive power supply in the circuit. An n^+ region is placed below the aluminum contact to ensure the formation of an ohmic contact and not a Schottky barrier diode.

10.6.3 Epitaxial Layer Resistor

The epitaxial layer itself provides a high-sheet-resistance layer, and high-value resistors can be formed with a small number of squares, but the absolute value is often poorly controlled. The epitaxial-layer resistor is formed by making two ohmic contacts to an isolated epitaxial region, as in Fig. 10.11.

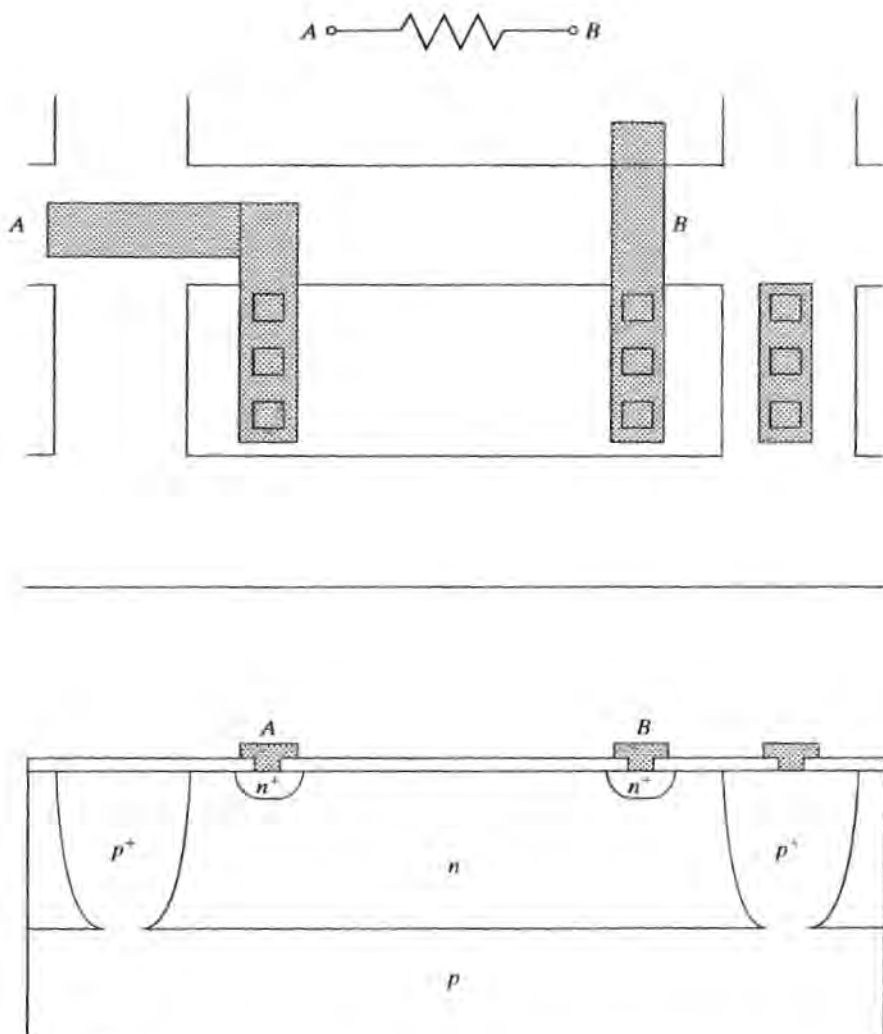


FIGURE 10.11

Resistor formed from the n -type epitaxial layer.

10.6.4 Pinch Resistor

An alternative technique used to obtain higher sheet resistance is to narrow the thickness of the base diffusion by crossing it with an emitter diffusion, as drawn in Fig. 10.12. The structure is very similar to the *n*pn transistor, except that the emitter diffusion merges with the epitaxial-layer island beyond the edges of the *p*-type region. The high-sheet-resistance region that is obtained corresponds to the active base region of the *n*pn transistor. The resistance value is highly dependent upon the thickness of the pinched base region, and the absolute tolerance is relatively poor. The pinched base resistor may also be used as a JFET with terminals *A* and *B* acting as source and drain and the *n* serving as the gate.

10.6.5 Substrate pnp Transistor

The so-called *substrate pnp* transistor is formed from the *p*-type diffusion (emitter), isolated epitaxial tub (base) and *p*-type substrate (collector), as shown in Fig. 10.13.

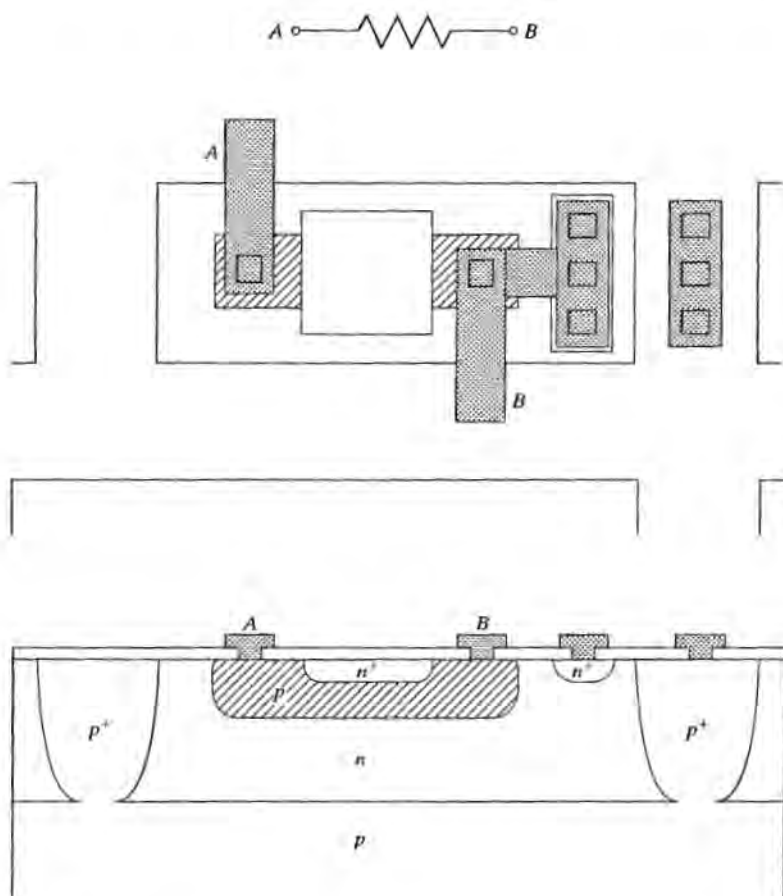


FIGURE 10.12
Pinched base resistor.

The basewidth W_B , set by the distance between the bottom of the p -type diffusion and the bottom of the n -type epitaxial layer, is typically much wider than that of the npn device. Thus, the current gain and frequency response are correspondingly less than those associated with the npn transistor. Also, the collector, formed from the p -type substrate, is inherently tied to the most negative power supply level in the circuit. Thus, only the base and emitter terminals are free for connection. Even so, the substrate pnp is often found as an emitter follower in the output stages of op-amps. The substrate pnp structure also appears in the n -well CMOS technology, where it is formed from the p -type source-drain region, the n -well, and the p -type substrate. In p -well CMOS technology, a substrate nnp transistor can be formed.

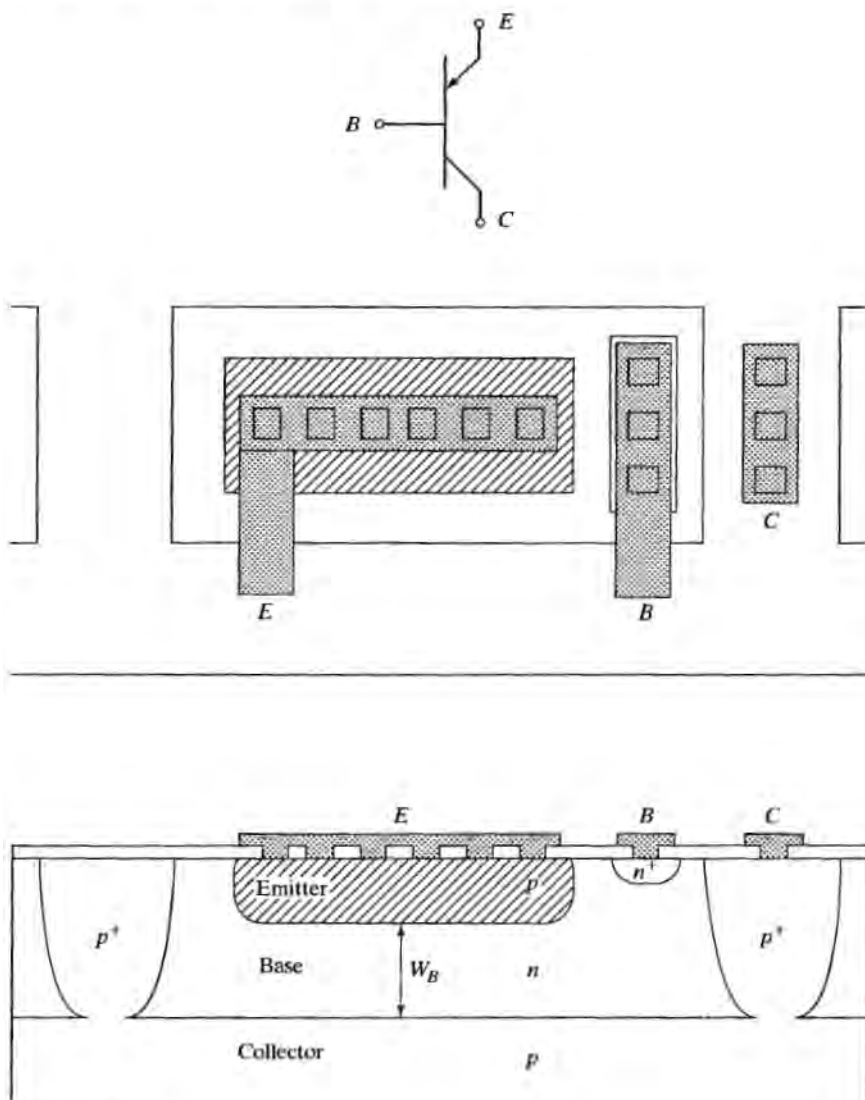


FIGURE 10.13

Substrate pnp transistor.

10.6.6 Lateral *pnp* Transistors

The lateral *pnp* structure in Fig. 10.14 provides a transistor with uncommitted collector, base, and emitter contacts. The emitter and collector are formed from the *p*-type region that is used for the base of the *npn* transistor. The basewidth is determined by the lithographic spacing between the two diffusions plus the degree of lateral diffusion

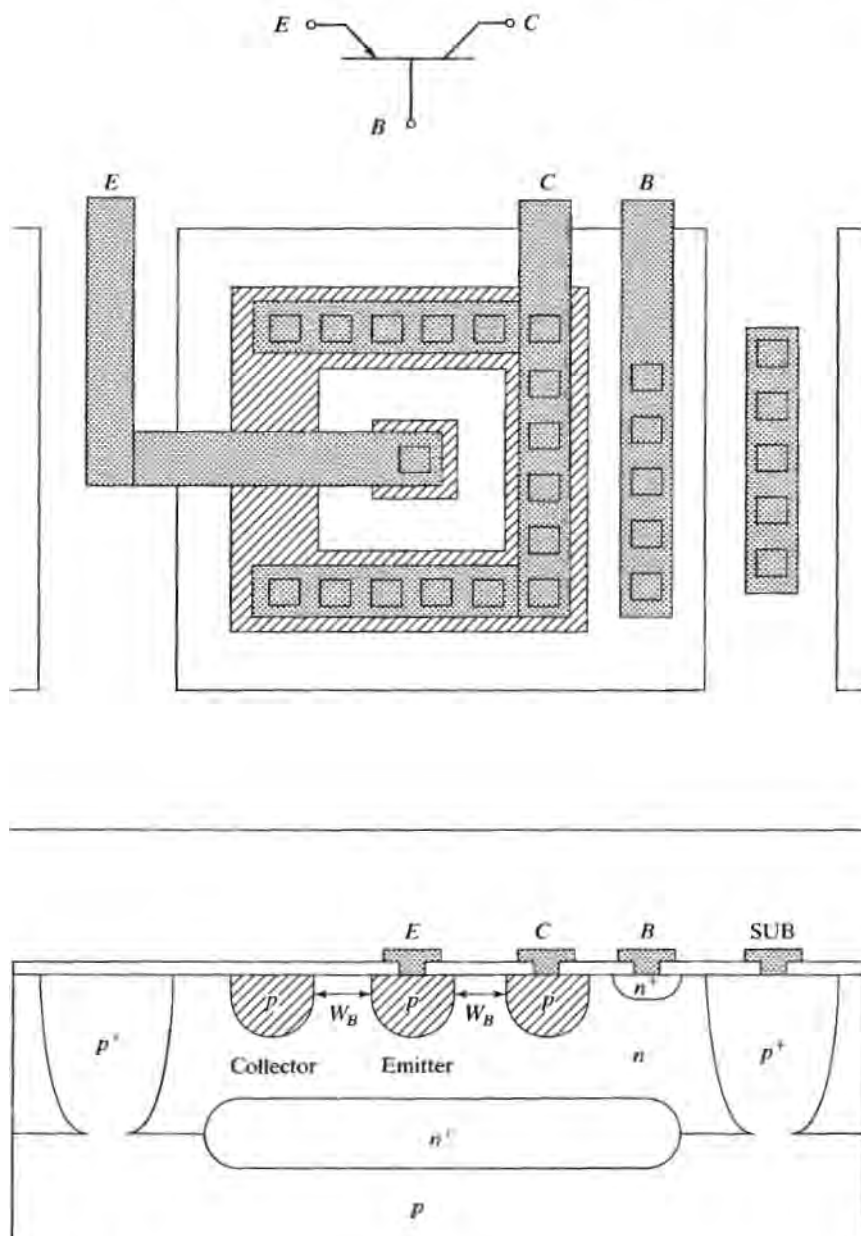


FIGURE 10.14

Enclosed lateral *pnp* transistor structure.

and is therefore much wider than that of the *nnp* transistor. This leads to low current gain and poor f_T . It is difficult to achieve a high current gain, because the emitter tends to inject carriers equally in all directions, but the collector only collects those carriers that cross the base region between the sides of the two *p*-type regions. Thus, *pnp* transistors are often enclosed structures with the emitter surrounded by the collector, as in Fig. 10.14. (See Prob. 10.9.)

10.6.7 Schottky Diodes

Bipolar logic circuits, for example, classic Schottky TTL, have used Schottky diodes to prevent the saturation of *nnp* transistors. A simple Schottky diode can be formed from an aluminum contact to the lightly doped *n*-type layer as shown in Fig. 10.15(a). A *p*-type ring is often used to minimize problems with high fields at the edges of the metal contact, as are circular diode layouts. Although the resulting structure places a Schottky diode in parallel with a *pn* junction diode, forward conduction is dominated by the low forward voltage drop of the Schottky diode.

Schottky diodes tend to have both higher leakage current and lower breakdown voltage in the reverse direction than a *pn* junction diode. A clever solution to this problem appears in Fig. 10.15(b), in which multiple *p*-type diffusions have been added to the interior of the Schottky diode region [7]. In the forward direction, multiple Schottky and *pn* junction diodes are in parallel, but the low voltage drop of the Schottky diode again dominates the forward region of the diode. However, under reverse bias, the depletion regions of the *pn* diodes merge, and the reverse breakdown appears as that of the *pn* junction, rather than the Schottky barrier diode. This structure shows the use of a detailed understanding of the behavior of the *pn* junction to achieve a novel result.

10.7 LAYOUT CONSIDERATIONS

This section will explore mask layout for a classic SBC transistor. The analysis will expand our understanding of the interaction of process design and layout. In particular, the top view of the mask set for a bipolar transistor often differs greatly from the final device structure, due to large lateral diffusions, although this is much less true of the high-performance digital technologies discussed in Section 10.8.

10.7.1 Buried-Layer and Isolation Diffusions

The spacing between adjacent buried layers and the width of the intervening isolation diffusion determines how closely two transistors can be spaced. To maintain electrical isolation, the substrate is tied to the most negative voltage present in the circuit. The collector of a transistor, on the other hand, is often connected to the most positive voltage in the same circuit. Thus, the collector-substrate junction must be designed to support a voltage equal to the sum of the positive and negative voltages supplying the circuit. For analog circuits, this may be more than 40 V. A typical design value of 60 V would provide an adequate safety margin.

Figure 10.16 shows the depletion layer in the *p* and *n* material near the isolation region of the transistor. The doping of the isolation diffusion is heavy at the surface and

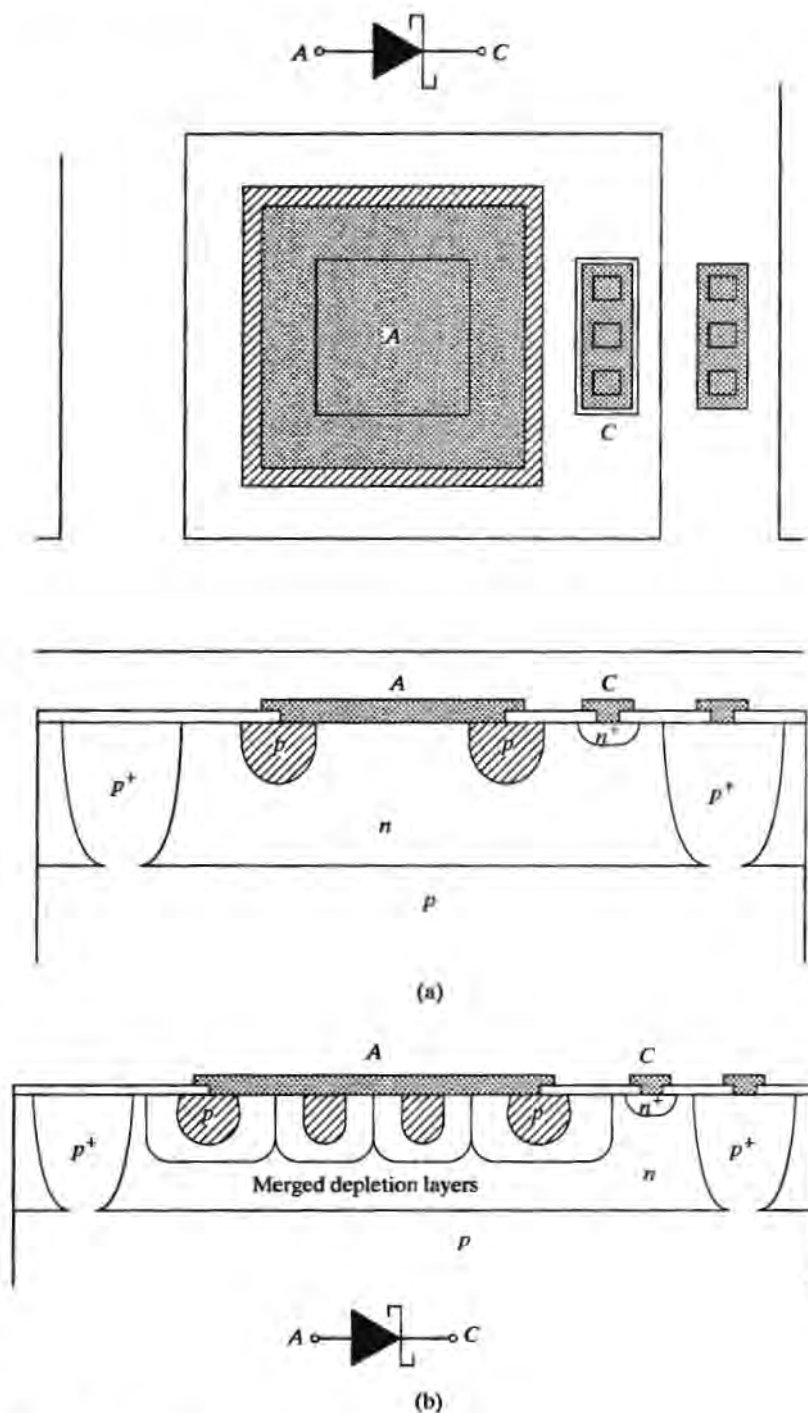


FIGURE 10.15

(a) Schottky barrier diode with p -type breakdown protection ring; (b) high reverse breakdown Schottky barrier diode.

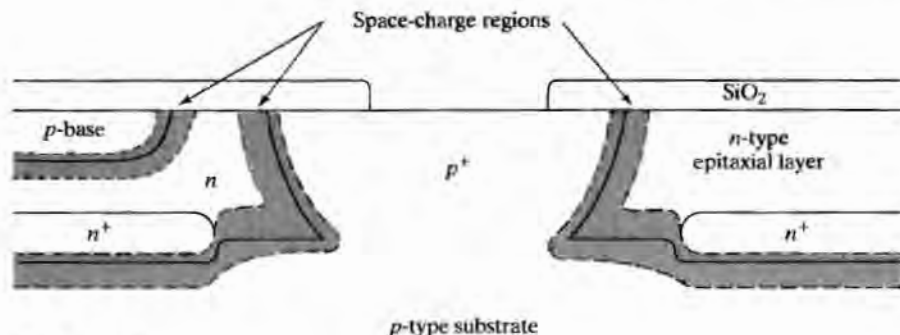


FIGURE 10.16

Isolation region between two bipolar transistors. The spacing must be large enough to ensure that the two space-charge regions do not merge together.

intersects the original substrate to produce isolation. At the surface, the window defining the isolation diffusion may be a minimum feature size, but the total width of the isolation region at the surface will be determined by lateral diffusion. For example, if the epitaxial layer is $15\text{ }\mu\text{m}$ thick and the minimum feature size is $10\text{ }\mu\text{m}$, the isolation region will approach $40\text{ }\mu\text{m}$ in width, assuming that lateral diffusion equals vertical diffusion. Obviously, this form of isolation is not acceptable for VLSI/ULSI logic circuits.

The n^+ buried-layer diffusion is not usually permitted to intersect the p^+ isolation diffusion. If the two diffusions meet, the breakdown voltage of the junction will decrease, and the capacitance of the junction will increase. Thus, there will be a layout design rule associated with the minimum spacing between the n^+ region and the isolation diffusion.

10.7.2 Base Diffusion to Isolation Diffusion Spacing

At the surface, the collector-base and collector-substrate depletion regions of Fig. 10.16 must not merge. The minimum spacing can be determined from a knowledge of the applied voltages and the epitaxial-layer impurity concentration. Additional spacing must be added to account for the alignment sequence and accumulated alignment tolerances.

Example 10.6

What is the minimum spacing between the edge of the base diffusion and the edge of the isolation diffusion at the surface of a bipolar transistor if the alignment tolerance is $1\text{ }\mu\text{m}$ and the epitaxial-layer resistivity is 10 ohm-cm ? Assume that the two junctions must each support 40 V . Use a collector-base junction depth of $5\text{ }\mu\text{m}$.

Solution: A 10-ohm-cm epitaxial layer has an impurity concentration of $5 \times 10^{14}/\text{cm}^3$, giving a value of $V/N_B = 8 \times 10^{-14}\text{ V-cm}^3$. From Fig. 10.7, the total depletion-layer width is approximately $10\text{ }\mu\text{m}$, with $8.7\text{ }\mu\text{m}$ in the epitaxial layer. The conditions at the isolation-collector junction are essentially the same, so the minimum spacing will be two times the depletion-layer width of $8.7\text{ }\mu\text{m}$ plus the alignment tolerance of $1\text{ }\mu\text{m}$ for a total of $18.4\text{ }\mu\text{m}$.

10.7.3 Emitter-Diffusion Design Rules

The minimum spacing between the edges of the emitter and base diffusions must be greater than the sum of the emitter and collector depletion-layer widths in the base, the accumulated alignment tolerance between the emitter and base masks, and the active base-region width.

In the basic SBC process, the n^+ emitter diffusion is also used to ensure the formation of a good ohmic contact to the collector, and this collector contact diffusion should not intersect the depletion layers associated with either the p -type base or isolation diffusions. If this occurs, the breakdown voltage of the junctions will be reduced and the junction capacitances increased. (See Prob. 10.7.)

10.7.4 A Layout Example

A set of design rules for a hypothetical bipolar transistor is given in Table 10.2, and Fig. 10.17 shows the layout of a minimum-size transistor based on these rules and making maximum use of lateral diffusion. It is important to note that the active area of the transistor—the region directly under the emitter—is a small fraction of the total device area of approximately $4536 \mu\text{m}^2$. The final emitter area is $11 \times 11 \mu\text{m}$, or $121 \mu\text{m}^2$, which represents only 2.67% of the total area of the transistor. (Remember, a similar area utilization problem occurred with the layout of the metal gate MOS transistor.)

The rest of the area is needed to make contacts, to support depletion layers, and to provide isolation between adjacent devices. In this layout, the isolation area is $2800 \mu\text{m}^2$, or more than 60% of the total area! Minimization of the isolation region represents an important issue in high-performance devices, not only for density improvement, but also for junction-capacitance reduction.

The solid lines in the top view of the transistor layout represent the edges of the masks used to fabricate the transistor. The various dotted lines represent the final positions of the emitter, base, and isolation diffusions. For the design rules of Table 10.2, the window for the emitter diffusion happens to coincide exactly with the emitter contact window. The lateral diffusion of $3 \mu\text{m}$ provides more than the required $2\text{-}\mu\text{m}$ alignment tolerance for the emitter contact window.

TABLE 10.2 Bipolar Transistor Design Rules for Fig. 10.17.

Minimum feature size	5 μm
Worst-case alignment tolerance between levels	2 μm
Epitaxial-layer thickness	10 μm
Collector-base junction depth	5 μm
Emitter-base junction depth	3 μm
Minimum emitter-to-collector spacing at surface	5 μm
Minimum base-to-isolation spacing at surface	5 μm
Minimum collector contact n^+ diffusion to isolation spacing	5 μm
Minimum collector contact n^+ diffusion to base spacing	5 μm
Buried-layer diffusion (both up and down)	2 μm
Buried layer to isolation spacing	5 μm
Buried layer pattern shift	10 μm
Lateral diffusion = vertical diffusion	

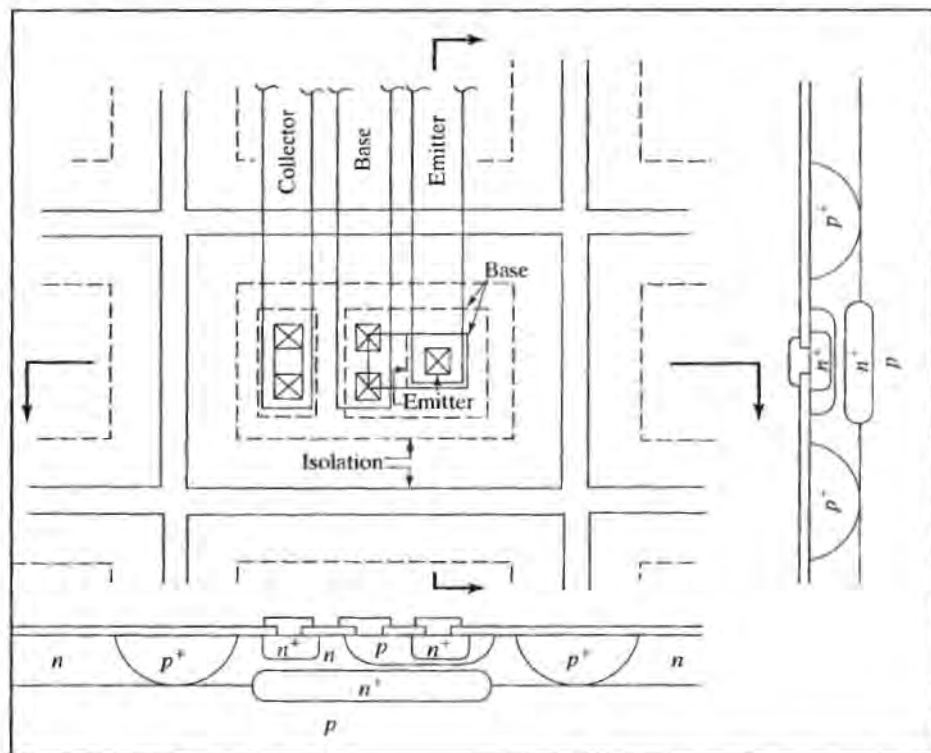


FIGURE 10.17

Minimum-size bipolar transistor layout based on the design rules of Table 10.2. The buried layer is not shown in the top view for reasons of clarity. Each square is $5\text{ }\mu\text{m} \times 5\text{ }\mu\text{m}$.

The base diffusion is $5\text{ }\mu\text{m}$ deep and is assumed to diffuse laterally $5\text{ }\mu\text{m}$. In the layout, the base contact windows actually extend outside the base region at the mask level, but are more than one alignment tolerance within the base region following diffusion. This is an excellent example of the interaction between processing and layout. However, only a most aggressive designer would consider a layout ground rule such as this.

The width of the base-contact metallization has been widened to more than a minimum feature size to help clarify the figure. This did not affect the size of the device, as space was available, because of other design rule limitations. Two base and two collector contact windows fit within the minimum base and n^+ collector contact regions. The collector contact windows align with the edges of the n^+ diffusion window, as was the case for the emitter.

The buried-layer mask has also been omitted from the figure for clarity. In this particular structure, the design rules relating to the buried layer are not limiting factors in the size of this layout.

10.8 ADVANCED BIPOLAR STRUCTURES

For digital logic circuits, structures are optimized to provide as short a transit time as possible. This requires minimizing the basewidth, eliminating as much capacitance as possi-

ble by minimizing total junction area, minimizing the width of the collector space-charge region, and reducing the collector and base series resistances. A reduced current gain is traded for a shorter transit time. As mentioned earlier, although $\langle 111 \rangle$ wafers were originally used for bipolar processing, silicon with a $\langle 100 \rangle$ surface orientation is now commonly used in most advanced bipolar processes.

10.8.1 Locos-Isolated Self-Aligned Contact Structure

Figure 10.18 shows a high-performance bipolar structure that attempts to achieve the aforementioned goals by using a very thin epitaxial layer and shallow ion-implanted base and emitter regions. As much pn junction area as possible is eliminated through the use of oxide isolation. The sides of the emitter and base regions are actually walled by the oxide isolation regions. The n^+ buried layer is relatively large to minimize r_C , and the total base region is minimized to reduce the base resistance. Self-aligned contacts are made to the base, emitter, and collector regions.

The formation of the transistor of Fig. 10.18 begins with the implantation and diffusion of the buried layer with a typical sheet resistance of 10 to 50 ohms per square. The masking oxide is removed, and a thin epitaxial layer is grown on the surface. A recessed oxide isolation process is used to form the isolation regions between devices and to eliminate the unnecessary junction area between the collector and emitter contacts. Prior to oxidation, part of the epitaxial layer is etched away so that the subsequent oxidation will extend completely through the epitaxial layer. An implantation is used to overcome boron depletion in the substrate during oxidation.

Next, the silicon nitride-oxide sandwich is removed, and an oxide is regrown on the surface. A boron implantation creates the shallow active base region. A mask is used to create windows for the emitter, and contacts to the base, emitter, and collector are all defined at the same time. Note that a single oxide strip defines both the emitter and base contact regions, eliminating alignment tolerances that would be needed if the regions were formed separately. The width of this strip is set by the metal-to-metal spacing plus accumulated alignment tolerances.

Photoresist is used as a barrier material during implantation of the base contact region. This p^+ implantation further reduces the base resistance of the device. Photoresist is also used as a barrier material during implantation of the emitter and collector contact regions. Note that these two mask steps use noncritical *blockout* masks similar to those used for threshold adjustment in a CMOS process.

Contacts are made through the same openings used for the base and emitter implantations. These contact areas are all cleared by a short wet or dry etch prior to metallization. Because of the very shallow junction depths involved in this structure, the metallization will be a multilayer sandwich structure including a barrier metal in the contact region. Interconnection of devices to form circuits will typically involve a multilevel metal process.

10.8.2 Dual Polysilicon Self-Aligned Processes

Figure 10.19 outlines an early dual-polysilicon self-aligned process whose steps remain representative of those used to form many of today's high-performance bipolar devices [8]. A nonselective n^+ buried layer and n -epitaxial layer are fabricated on the p -type substrate. Deep trenches are used to provide isolation between devices. After etching, the

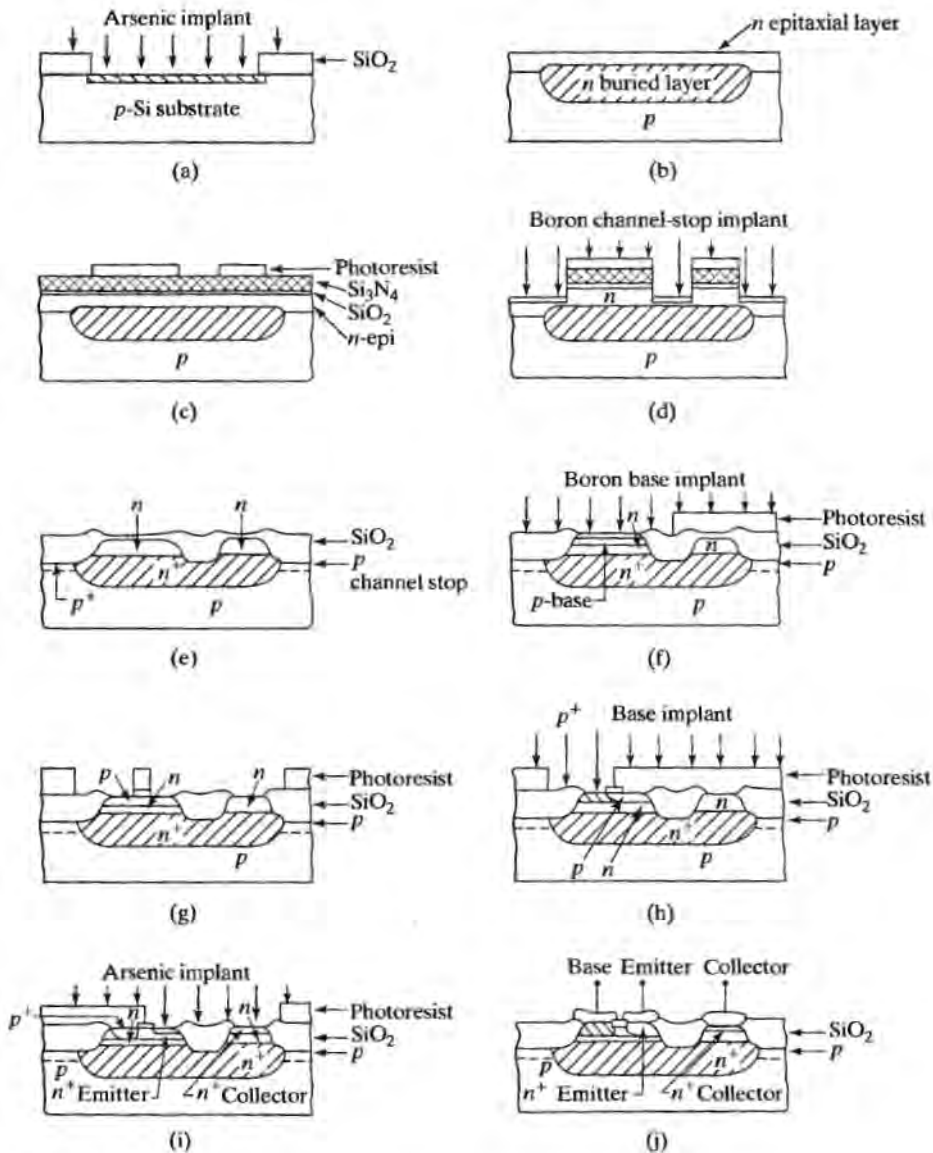


FIGURE 10.18

Process sequence for a high-performance oxide-isolated bipolar transistor. (a) Buried-layer formation; (b) epitaxial layer growth; (c) mask for selective oxidation; (d) boron implant prior to recessed oxide growth; (e) selective oxidation; (f) base mask and boron base implantation; (g) emitter, base contact, and collector contact mask; (h) p^+ base contact implantation; (i) arsenic implantation for emitter and collector contact; (j) structure completed with multilayer metallization. Copyright 1985, John Wiley & Sons, Inc. Reprinted with permission from Ref. [5].

trench walls are oxidized and then refilled with polysilicon. (CVD oxide is also used in some processes.) A p^+ implant is used to minimize the impact of boron depletion in the trench. Next, shallow trench regions are formed, yielding a planarized surface. An n^+ "sinker" diffusion is used to contact the n^+ buried collector layer.

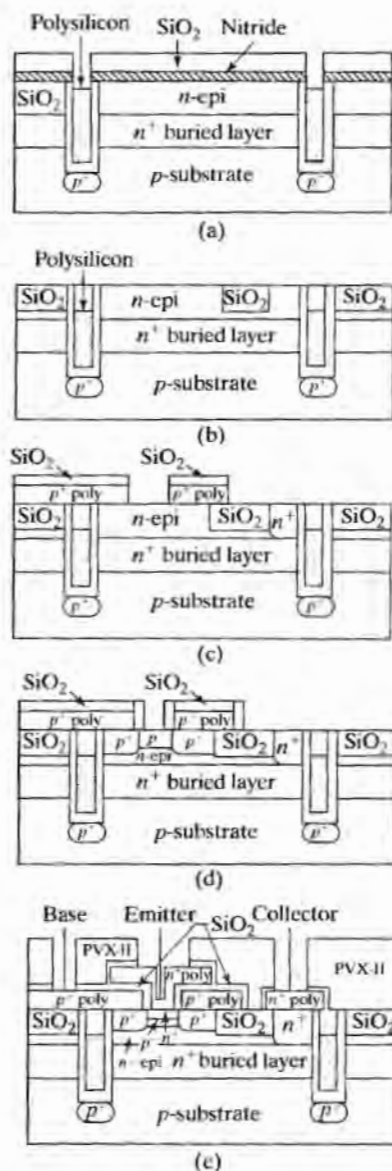


FIGURE 10.19

A high-performance bipolar transistor structure with an f_T of 10 GHz. (a) Isolation is achieved using deep-trench isolation with polysilicon and silicon dioxide refill; (b) structure following selective oxidation; (c) p^+ polysilicon deposited and patterned; (d) diffusion from doped polysilicon forms the extrinsic base region and base contacts; a self-aligned implantation forms the intrinsic base; (e) diffusion from n^+ polysilicon forms the emitter and the emitter and collector contacts of the transistor. Copyright 1997, IEEE. Reprinted with permission from Ref. [8].

Polysilicon is deposited and heavily doped with p -type impurities, commonly by ion implantation. The base opening is delineated, the intrinsic base region is implanted through the base window, and the polysilicon is oxidized. During subsequent annealing and processing steps, the extrinsic base region outdiffuses from the polysilicon forming p^+ side contacts to the more lightly doped intrinsic base region. The heavy poly doping minimizes the extrinsic base region's contribution to the base resistance. The oxide on the polysilicon provides an insulating layer that will separate the emitter contacts from the base contacts. Note the added formation of oxide on the sides of the polysilicon layer similar to that used for LDD processing in MOS devices.

The second layer of polysilicon, this time n^+ doped, is deposited and annealed, forming a very shallow emitter region in the p -type base. Arsenic doping is typically used, because of its smaller diffusion coefficient. The n^+ polysilicon simultaneously provides contacts to both the emitter and the n^+ collector sinker. Note that the intrinsic base, base contacts, and emitter are defined and formed as a result of a single lithographic step and are referred to as being self-aligned. The elimination of alignment tolerances between emitter, base, and base contact layers results in significant reduction in overall size of the transistor and attendant minimization of device capacitances.

In these processes, the devices have extremely shallow emitters and narrow base-widths in order to achieve high f_T . To constrain and control the depletion-layer extents and value of the base resistance, the doping of the base region must be increased, and both the narrow base and increased doping levels lead to relatively low breakdown voltage in the transistor. Many of these devices are designed to operate with supply voltages of 2.5 V and below. In addition, the overall "thermal budget" (Dt product) is very small.

10.8.3 The Silicon-Germanium Epitaxial Base Transistor

By adding germanium to silicon [9], one can modify the silicon bandgap and perform "bandgap engineering" in a manner similar to that achieved in III-V compound semiconductor materials such as GaAs, InP, and InAs. Germanium is a Column IV element, as is silicon, so it does not act as a dopant impurity in silicon. However, when germanium replaces silicon in the lattice, it modifies the bandgap of the composite material. The bandgap change enhances the current gain and can result in a built-in base field that increases the unity-gain frequency of the transistor. At the time of this writing, cut-off frequencies in excess of 200 GHz have been reported in such SiGe heterojunction bipolar transistors (SiGe HBTs).

In SiGe HBTs, germanium is introduced into the base region of the device, with peak concentrations of less than 15%, during an epitaxial growth step that is used to form the base region of the transistor. A representative cross section of the SiGe device structure appears in Fig. 10.20(a), with a microphotograph of the structure in Fig. 10.20(b). The structure uses a combination of deep and shallow trenches for isolation. (The deep trenches in the photograph are approximately 5 μm deep.) The base region is formed by a carefully controlled epitaxial growth under ultrahigh-vacuum (UHV) conditions, which produces a single crystal intrinsic base region in the silicon window, but results in the formation of polycrystalline silicon over the oxide. During the epitaxial growth, various germanium profiles can be introduced into the intrinsic base region to produce the SiGe base device. A sacrificial gate structure is used to protect the emitter region during a p^+ implant that heavily dopes the polysilicon and provides the side wall contacts to the intrinsic base region. (This implant is actually the S/D implant for a PMOS transistor in a BICMOS process; see Section 10.10.) The emitter is formed by outdiffusion following deposition of the heavily doped polysilicon emitter contact. Transient enhanced diffusion (TED, discussed in Section 5.6.3) of base impurity implants can result in loss of control of the base profile, and the use of an epitaxial base results in a more abrupt transition and narrower basewidth.

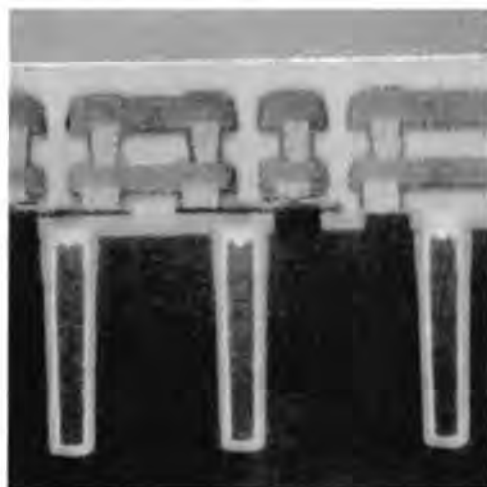
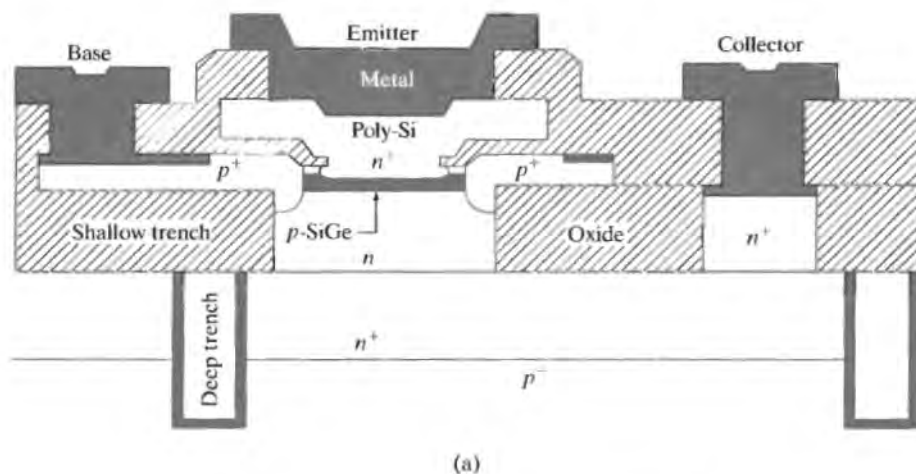


FIGURE 10.20

(a) Cross section of an epitaxial base silicon-germanium heterojunction bipolar transistor (SiGe HBT) (b) Photomicrograph of actual structure with 5- μm -deep trench isolation. Copyright 1985. IEEE. Reprinted with permission from Ref. [9].

A typical profile for an SiGe HBT is shown in Fig. 10.21 [10]. Without the addition of germanium, this profile is also characteristic of the oxide-isolated double-polysilicon structures, as depicted in Fig. 10.19, although TED may increase the basewidth in implanted base devices. In the profile in Fig. 10.21, the metallurgical basewidth is less than 100 nm, and the emitter-base junction is only 150 nm below the surface of the polysilicon. The desired germanium-doping profile is achieved during the base epitaxy. For film stability, the peak concentration is typically maintained below 15%. Also evident in the profile are the phosphorus collector and heavily doped arsenic subcollector regions.

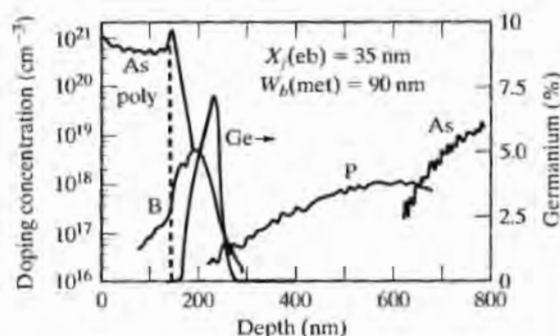


FIGURE 10.21

SiGe HBT impurity profile. The same profile applies to high-performance silicon BJTs if the Ge is eliminated. Copyright 1985, IEEE. Reprinted with permission from Ref. [9].

10.9 OTHER BIPOLAR ISOLATION TECHNIQUES

Several other interesting approaches to device isolation have been developed over the years, and two are surveyed here. Dielectric isolation processes are in use in high-performance analog circuits; the CDI process is replaced by oxide-isolated structures.

10.9.1 Collector-Diffusion Isolation (CDI)

The collector-diffusion-isolation (CDI) [8] structure, shown in Fig. 10.22, was developed primarily for digital applications. The process eliminates the p -type isolation diffusion, achieving reduced device area and process complexity.

The process starts with diffusion of a low-sheet-resistance buried layer, which will serve as the collector of the transistor. A thin p -type epitaxial layer forms the base region and is grown in the next step. Device isolation is achieved via an n^+ diffusion that completely encloses the transistor and also provides the collector contact area. Next, a shallow emitter is implanted or diffused into the device, followed by contacts and metallization. Typical parameters for the CDI process include a buried-layer sheet resistance of 15 to 30 ohms per square, a 2- μm -thick, 0.25-ohm-cm epitaxial layer, and an emitter depth of less than 1 μm .

This process produces high-performance, narrow-base transistors with minimum r_c , but with relatively large collector-base and collector-substrate capacitances. It has for the most part been replaced with advanced oxide-isolated structures that also minimize these capacitances, although at a cost of considerable process complexity.

10.9.2 Dielectric Isolation

One of the "holy grails" of bipolar processes has been to find a process that can provide $n\text{pn}$ and $p\text{np}$ devices with similar levels of performance. Such processes are referred to as complementary bipolar processes. However, parasitic coupling between

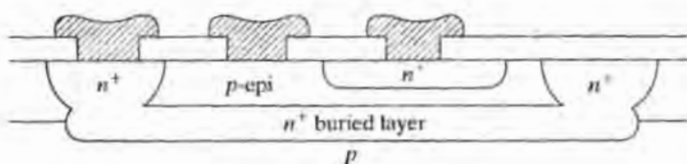


FIGURE 10.22

Cross section of a transistor fabricated in the CDI process.[18]

devices and susceptibility to latchup prevent the simple addition of diffusions to the *npn* process to produce vertical *pnp* transistors. Instead, some form of dielectric isolation has been used to isolate the two types of transistors.

The first successful dielectric-isolated process was developed by Harris Semiconductor [11–12], and the basic process flow is depicted in Fig. 10.23. Deep V-grooves are first etched in the surface of <100> oriented silicon wafers. (This etching process will be discussed in more detail in Chapter 11.) A nonselective n^+ diffusion is performed, and a silicon dioxide layer is grown on the surface. A thick layer of polycrystalline silicon is then deposited on the surface of the wafer. Silicon is removed from the back surface of the wafer by lapping until the silicon dioxide in the V-grooves is exposed, as indicated by the dotted line in Fig. 10.23(b). The surface is then mechanically and chemically polished. The wafer is turned over, yielding islands of silicon completely isolated from each other by the silicon dioxide dielectric layer. Standard processing is then used to form *npn* transistors. This process is expensive, but an important variation permits formation of complementary, vertical *npn* and *pnp* transistors. It is used in the fabrication of high-performance analog circuits. In addition, the structures are highly tolerant to radiation and are used in military applications.

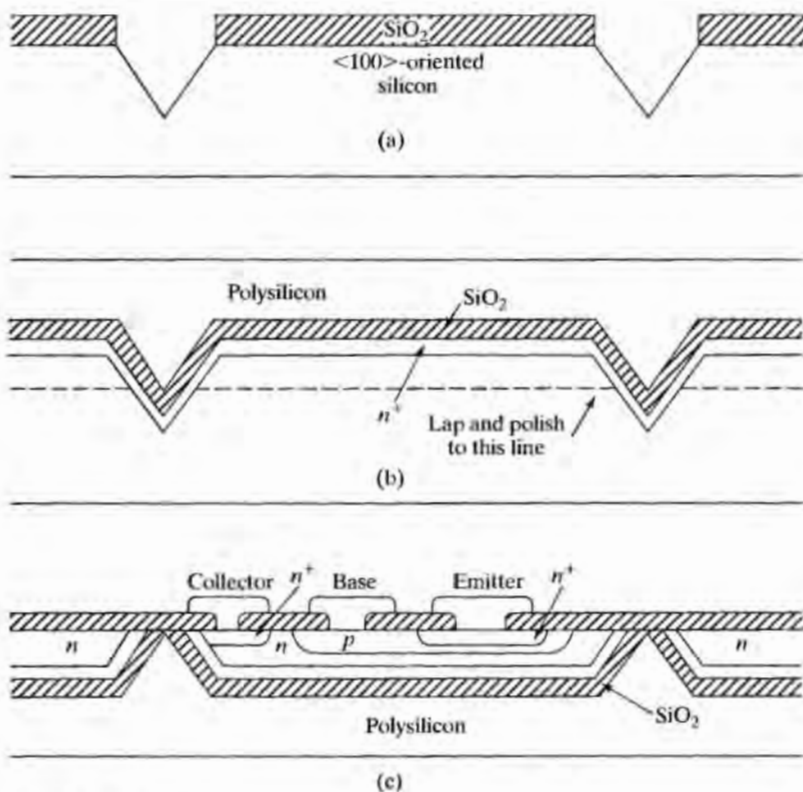
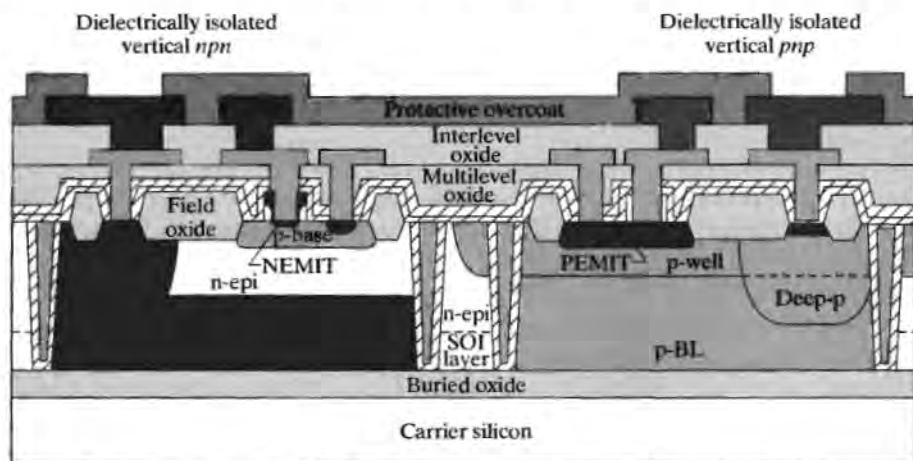


FIGURE 10.23

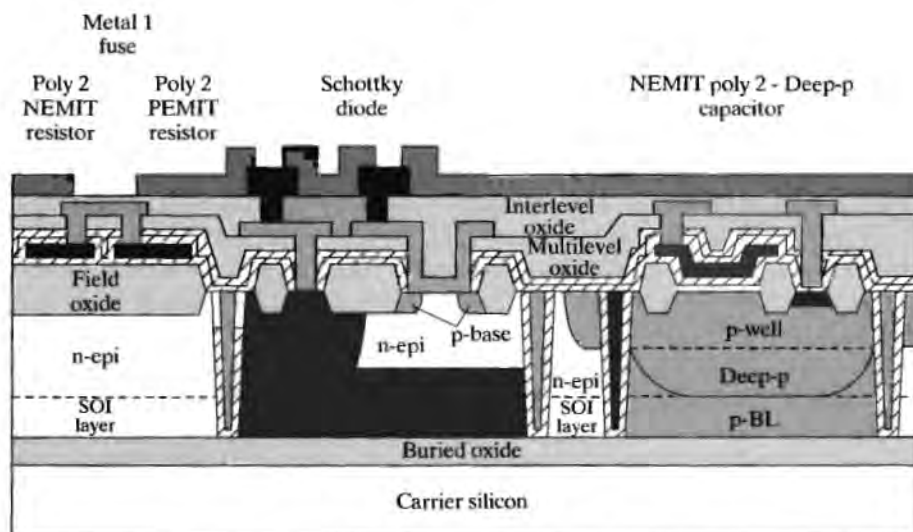
Several steps in the dielectric isolation process [11,12] (a) V-grooves anisotropically etched in the silicon substrate; (b) the structure following n^+ diffusion and oxidation; the wafer is turned over, lapped back and polished to the dotted line; (c) bipolar transistors are then fabricated in the isolated islands of silicon.

A myriad of recent processing developments, including buried-oxide and bonded-wafer SOI, deep-trench RIE with oxide and polysilicon refill, etc., have come together to permit the development of more straightforward processes for producing dielectrically isolated BJTs. Figure 10.24(a) depicts the cross sections resulting from one such process [13]. In this particular process, the starting substrates are bonded SOI wafers. However, SIMOX wafers could also be utilized. Heavily doped p - and n -type buried collectors are implanted and followed by the growth of an n -type epitaxial layer for the



npn and pnp transistors

(a)



Passive components

(b)

FIGURE 10.24

(a) High-performance bipolar transistors and (b) passive components fabricated on a SIMOX wafer using shallow trench isolation. Copyright 1997, IEEE. Reprinted with permission from Ref. [13].

nnp collector. A *p*-well forms the *pnp* collector region, and deep *n*- and *p*-type sinker diffusions are added to contact the buried collectors of both transistors. Isolation between devices is achieved by etching deep trenches down to the buried-oxide layer. The trenches are then refilled with oxide and polysilicon. LOCOS isolation at the surface is followed by separate base region implants for the *nnp* and *pnp* transistors. An undoped polysilicon layer is deposited and then selectively implanted with *n*- and *p*-type impurities. Outdiffusion from the implanted polysilicon layer forms the emitters of the two types of transistors. This process results in complementary transistors with f_T 's of 2.5 GHz. Figure 10.24(b) shows the realization of resistors, capacitors, and Schottky diodes in the same process.

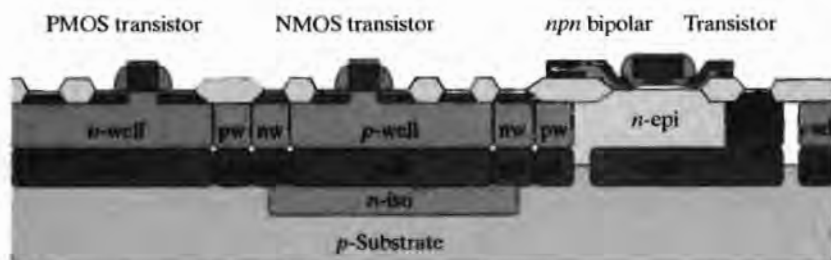
10.10 BICMOS

Achieving a process that can provide the benefits of both bipolar and MOS devices has been another of the long-term goals in fabrication. Over time, the complexity of NMOS, CMOS, and bipolar processes have all grown dramatically to the point that the processes are essentially indistinguishable in terms of complexity. For example, a state-of-the-art CMOS process may have 20–25 mask steps. At the same time, the basic process steps have tended to converge, and it has become relatively straightforward, although certainly more expensive, to combine bipolar and CMOS structures into a technology most often termed BiCMOS. (The inclusion of JFETs in the bipolar process is often called a BiFET process.) A good review can be found in [14].

Three approaches to realizing a BiCMOS technology are presented in Fig. 10.25. Because of the presence of CMOS devices, <100> silicon is the substrate of choice for BiCMOS processing. The first case [15] forms heavily doped *n*-type and *p*-type buried layers in the *p*-type substrate. Following epitaxial layer growth, *n*-well and *p*-well regions are formed by implantation and diffusion, and then NMOS and PMOS devices are formed in the two wells. The *nnp* transistor, actually a SiGe HBT, has an epitaxial base with polysilicon side-wall contacts and a polysilicon emitter. Standard lateral *pnp* transistors can be formed in *n* regions similar to those used for the *nnp* devices.

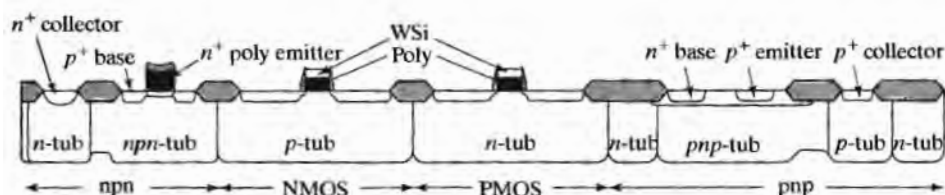
The process illustrated in Fig. 10.25(b) [16] starts with LOCOS isolation and utilizes ion implantation and diffusion to form *n*-tubs and *p*-tubs in the substrate for NMOS, PMOS, *nnp*, and *pnp* transistors. In this process, high-energy boron and phosphorus implantations are used to form the buried layers needed to reduce the bipolar transistor collector resistance. The *nnp* has an ion-implanted base and an arsenic-doped *n*⁺ polysilicon emitter. Here, a vertical *pnp* transistor is formed in one of the *p*-tub regions using a phosphorus base implantation and the *p*⁺ S/D implant as its emitter. The f_T of the *nnp* and *pnp* transistors exceed 40 GHz and 10 GHz, respectively. The NMOS and PMOS devices are silicided polysilicon gate LDD devices with 5-nm gate oxides.

The third BiCMOS structure [17] is a portion of an IC process that also includes vertical power devices. Here, *n*⁺ buried layers are formed below all the CMOS and bipolar devices. Following epitaxial layer growth, *n*-tub formation, and *p*⁺ isolation, straightforward processing is used to form the vertical *nnp*, lateral *pnp*, and PMOS transistors. A *p*-well is added to the *n*-tubs to act as the substrate for the NMOS device. Note the common connection of the source, *p*-well and *n*-tub by the source contact of the NMOS transistor. The *n*-epitaxial layer with *n*⁺ backside contact is required by vertical power device structures, which are also included in the process.



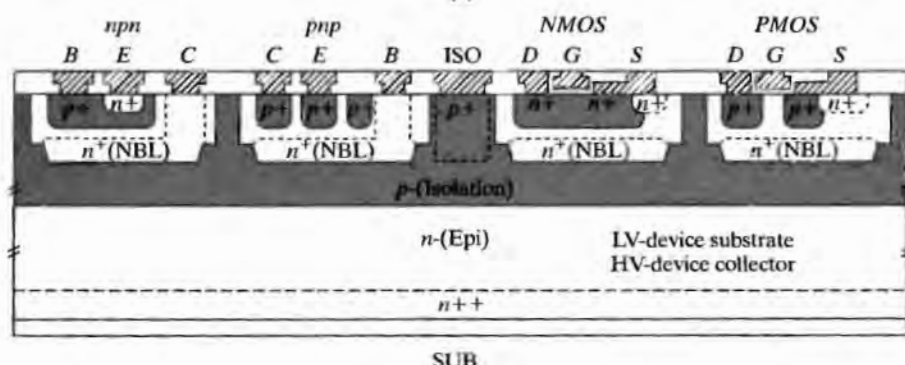
Schematic cross section of the BiCMOS6G technology

(a)



Schematic cross-sectional view of HECIBiC technology.

(b)



(c)

FIGURE 10.25

BiCMOS technology cross sections. (a) A dual-well CMOS + *npn* process. Well formation is preceded by *n*- and *p*-type buried-layer diffusion; (b) A twin-well process, which utilizes high-energy implantation to reduce collector resistance of the BJTs. Vertical *npn* and *pnp* transistors are available in this process; (c) BiCMOS devices in a process designed for power semiconductor applications. Copyright 1999, 1998, 1997, IEEE. Reprinted with permission from Refs. [15, 16, 17].

SUMMARY

The *standard buried collector* (SBC) process is widely used throughout the IC industry for analog and power circuit applications. More recently developed digital bipolar technologies have benefited greatly from process advances originally developed for use in MOS dynamic RAMs. These include *locos*-isolation, deep- and shallow-trench formation, the use of polysilicon, and the introduction of ion implantation.

Bipolar technologies for analog applications are typically designed to yield current gains of several hundred, with breakdown voltages of up to 50 V. The resulting devices have cutoff frequencies of less than 500 MHz. Dielectric isolation provides the ultimate in isolation between devices and can permit formation of truly complementary bipolar devices, as well as devices that can operate at high voltages.

Devices for digital applications can operate with lower current gains and supply voltages than traditional analog circuits. The modern VLSI bipolar process utilizes self-aligned epitaxial base region with a polysilicon emitter and base contact regions. Shallow- and deep-trench isolation are combined to minimize the area and capacitance of the bipolar transistor. Very narrow base transistors can also be formed by ion implantation. These factors permit device designs with cutoff frequencies exceeding 40 GHz. Germanium can be introduced into the base region of the epitaxial base transistor to produce an SiGe heterojunction transistor (HBT). These devices have already exhibited cutoff frequencies in excess of 200 GHz.

Bipolar and MOS processes and complexity have largely converged to the point that it is possible to fabricate both types of transistors in a single composite technology. Although there are widely varying approaches, these processes are termed BiCMOS technology.

REFERENCES

- [1] J. G. Fossum, "Computer-Aided Numerical Analysis of Silicon Solar Cells," *Solid-State Electronics*, 19, 269–277, April 1976.
- [2] H. Lawrence and R. M. Warner, Jr., *Bell System Technical Journal*, 39, 389–403, March 1960.
- [3] P. R. Wilson, "The Emitter-base Breakdown Voltage of Planar Transistors," *Solid-State Electronics*, 17, 465–467, May 1974.
- [4] R. A. Colclaser, *Microelectronics Processing and Device Design*, John Wiley & Sons, New York, 1980.
- [5] S. M. Sze, *Semiconductor Devices—Physics and Technology*, John Wiley & Sons, New York, 1985.
- [6] R. C. Jaeger, *Microelectronic Circuit Design*, McGraw-Hill Book Company, Burr Ridge, IL, 1997.
- [7] B. M. Wilamowski, "Schottky Diodes with High Breakdown Voltage", *Solid-State Electronics*, Vol. 26, no. 5, pp 491–493, May 1983.
- [8] M. Vora, Y. L. Ho, S. Bhamre, F. Chien, G. Bakker, H. Hingarh, and C. Schmitz, "A Sub-100 Picosecond Bipolar ECL Technology," *IEEE IEDM Technical Digest*, pp. 34–37, December 1985.
- [9] John D. Cressler, "Reengineering Silicon: Si-Ge Heterojunction Bipolar Technology," *IEEE Spectrum*, pp. 49–55, March 1995.
- [10] E. F. Crabbe, et al., "Vertical Profile Optimization of Very High-Frequency Epitaxial Si and SiGe-Base Bipolar Transistors," *1993 IEEE IEDM Digest*, pp. 83–86, December 1993.
- [11] C. J. Sanders, W. R. Morcom, and C. S. Kim, "An Improved Dielectric-junction Combination Isolation Technique for Integrated Circuits," *IEEE IEDM Technical Digest*, pp. 38–40, December 1973.
- [12] J. L. Davidson and D. R. Mason, Method of Etching Silicon Crystals, U.S. Patent #3,728,179, issued April 17, 1973.
- [13] R. Patel, W. Milam, G. Cooley, M. Corsi, J. Erdeljic and L. Hunter, "A 30-V Complementary Bipolar Technology on SOI for High-Speed Precision Analog Circuits," *1997 IEEE BCTM Digest*, pp. 48–51, September 1997.
- [14] D. Hareme, "High-performance BiCMOS Process Integration Trends, Issues and Future Direction," *IEEE BCTM Digest*, pp. 36–43, September 1997.

- [15] A. Monroy, et al., BiCMOS6G: A High-performance, 0.35- μm SiGe BiCMOS Technology for Wireless Applications," *IEEE BCTM Digest*, pp. 121–134, September 1999.
- [16] Y-F Chyan, et al., "A 50-GHz, 0.25- μm Implanted-Base High-energy Implanted Collector Complementary Modular BiCMOS (HEICBiC) Technology for Low-power Wireless-communication VLSIs," *IEEE BCTM Digest*, pp. 128–131, September 1998.
- [17] A. Leone, N. Speciale, S. Graffi, G. Masetti and V. Graziano, "Modeling Parasitic Bipolar Devices in Advanced Smart-power Technologies," *IEEE BCTM Digest*, pp. 127–130, September 1997.
- [18] B. T. Murphy, V. J. Glinski, P. A. Gary, and R. A. Pederson, "Collector Diffusion Isolated Integrated Circuits," *Proceedings of the IEEE*, 57, 1523–1527, September 1969.
- [19] C. C. Allen, L. H. Clevenger, and D. C. Gupta, "A Point Contact Method of Evaluating Epitaxial Layer Resistivity," *Journal of the Electrochemical Society*, 113, 508–510, May 1966.
- [20] H. F. Wolf, *Silicon Semiconductor Data*, Pergamon Press, Oxford, 1969.

PROBLEMS

- 10.1** Evaluate the Gummel number expressions for a uniformly doped transistor with impurity concentrations of N_E and N_B in the emitter and base, respectively. The effective width of the emitter is L_E , and W_B is the base width. What is the current gain for a device with $N_E = 10^{20}/\text{cm}^3$, $N_B = 10^{18}/\text{cm}^3$, $W_B = 4 \mu\text{m}$, $L_E = 20 \mu\text{m}$, $L_B = 50 \mu\text{m}$, $D_B = 20 \text{ cm}^2/\text{sec}$, and $D_E = 5 \text{ cm}^2/\text{sec}$? Assume that $\eta = 1$.
- 10.2** Using Eq. (10.1), estimate the current gain of the transistor with the impurity profiles given in Fig. 10.2. Assume that $W_B \ll L_B$, and $D_B = 20 \text{ cm}^2/\text{sec}$. Use $N_E/D_E = 5 \times 10^{13} \text{ sec}/\text{cm}^4$, $x_{JE} = 1.5 \mu\text{m}$, $x_{JC} = 4 \mu\text{m}$.
- 10.3** What is the maximum collector-base breakdown voltage of a transistor with $X_{BL} - X_{BC} = 5 \mu\text{m}$? What range of epitaxial layer dopings may be used to achieve this breakdown voltage?
- 10.4** A Zener reference diode is often formed using breakdown of the emitter-base junction.
- What base surface concentration is required to produce a breakdown voltage of 6 V for a 1- μm -deep junction?
 - If the base surface concentration is too small by a factor of two, what is the actual breakdown voltage of the diode?
- 10.5** Calculate the collector-base depletion-layer width for the transistor of Example 10.3 using the expression for a one-sided step junction given in Eq. (9.3). How does this compare with the width derived from Fig. 10.4?
- 10.6** The effective Gummel number in the emitter is substantially reduced from that calculated from the profile by bandgap narrowing in the emitter. Calculate G_E using the expressions in Chapter 4 for an As emitter with a sheet resistance of 10 ohms per square. Compare G_E to the values stated in the text.
- 10.7** In Section 10.7.3, the importance of correct positioning of the n^+ collector contact diffusion was discussed. To illustrate this point, three simple npn transistors are fabricated in an n -type substrate using the structure drawn in Fig. P10.7. An n^+ collector ring is used to reduce the collector series resistance R_c . Three different spacings—0 μm , 3 μm , and 5 μm —are used between the edge of the n^+ ring and the edge of the base diffusion. Explain why the breakdown voltage will be different for these three cases, and estimate the collector-base breakdown voltage for the three devices. Assume a base surface concentration of $10^{18}/\text{cm}^3$ and a junction depth of 5 μm .
- 10.8** Surface conditions can degrade the breakdown voltage of Zener diodes. Subsurface breakdown can be achieved using an ion-implanted process with the profile shown in Fig. P10.8. Estimate the breakdown voltage of this diode.

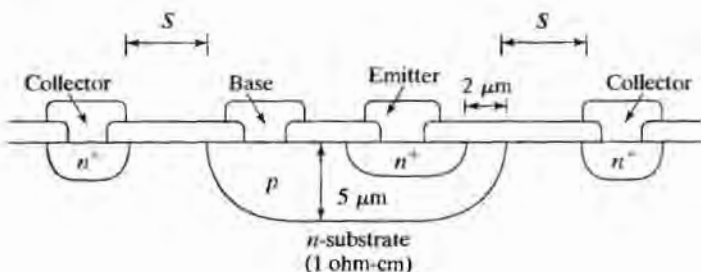


FIGURE P10.7

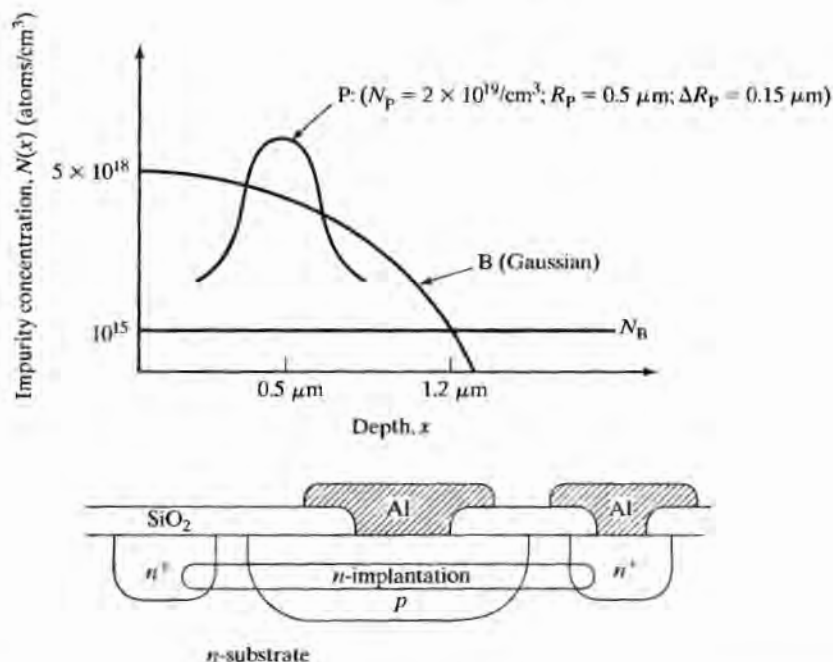


FIGURE P10.8

10.9 A simple lateral *pnp* structure is shown in Fig. P10.9(a). The current gain of this transistor is collection limited and does not obey Eq. (10.1). Assume that the emitter injects current uniformly in all directions and that the collector collects all the current coming its way.

- Under these assumptions, what is the value of the common base current gain = I_C/I_E ? What is the common emitter current gain β ?
- Derive an expression for β as a function of the length and width of the device. For a given area, what relationship between the length and width maximizes the gain?
- What geometry would be used to optimize the current gain?
- Repeat Prob. 10.9(a) for the circular device structure in Fig. 10.9(b).

10.10 Determine a reasonable diffusion schedule for the isolation diffusion of a junction-isolated structure with a 15- μm -thick epitaxial layer with the geometry of Fig. P10.10.

- Assume that the width of the isolation at the bottom is to be 10 μm , that there is no up-diffusion from substrate, and that lateral diffusion equals vertical diffusion.
- Modify your diffusion time in part (a) to account for up-diffusion of boron from the substrate. Assume that the substrate represents an infinite supply of boron impurities with a constant concentration of $10^{18}/\text{cm}^3$.

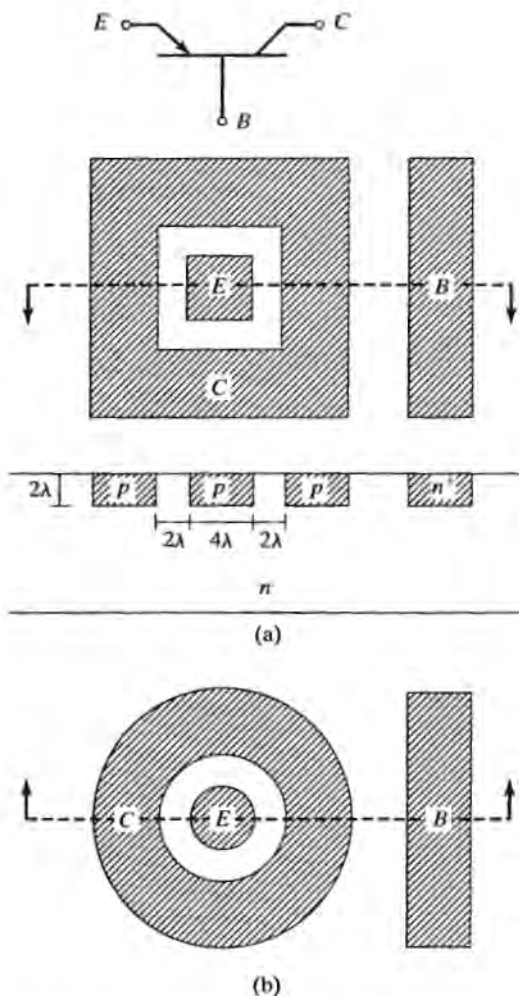


FIGURE P10.9

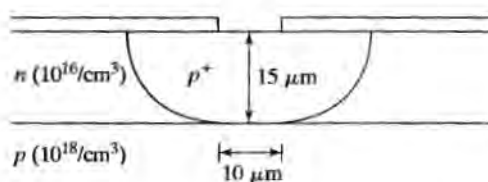


FIGURE P10.10

- 10.11** The distance $X_{BL} - X_{BC} = 5 \mu\text{m}$ in a certain bipolar SBC process with an epitaxial layer doping of $10^{15}/\text{cm}^3$. Use the one-sided step-junction formula [Eq. (9.3)] to estimate the punch-through voltage for this transistor. Compare with Fig. 10.7.
- 10.12** What epitaxial-layer doping would give approximately the same value for the Zener and punch-through voltage limits in Example 10.5.
- 10.13** Suppose the resistor in Fig. P4.10 is formed as an epilayer resistor. What is its resistance, based upon the values in Table 10.2 if lateral diffusion is ignored? (b) How about if the resistor is formed from the base layer? (c) How about for an emitter layer resistor?
- 10.14** Draw the cross section of (a) a substrate *n*p*n* transistor in a *p*-well CMOS process, and (b) a substrate *p*n*p* transistor in an *n*-well CMOS process.

- 10.15** Redraw the layout in Fig. 10.17 using an alignment tolerance of $1\text{ }\mu\text{m}$.
- 10.16** A Schottky clamped bipolar transistor is formed by placing a Schottky barrier diode in parallel with the collector-base junction by simply overlapping the base contact metal over the collector region. Draw a cross section of this structure, and include a guard ring on the edge of the diode.
- 10.17** List the mask steps required for the oxide-isolated bipolar transistor of Fig. 10.19. Which are noncritical alignment steps?
- 10.18** Estimate the emitter-base and collector-base breakdown voltages and punch-through voltage, for the impurity profiles in Fig. 10.21. Ignore the presence of germanium in the base region.
- 10.19** Estimate the space-charge region extents in the base for the impurity profiles in Fig. 10.21. Ignore the presence of germanium in the base region. The metallurgical base-width was estimated to be 100 nm in the text. What is the actual base region?
- 10.20** Estimate the permissible Dt products associated with the arsenic and boron profiles in Fig. 10.21. (a) Assume that the arsenic follows an erfc profile with a concentration of $10^{21}/\text{cm}^3$ at the surface ($x' = 0$) and intersects the boron profile at a level of $5 \times 10^{16}/\text{cm}^3$ at $x' = 50\text{ nm}$. What is the Dt product? If the As were diffused at 1000°C , what would be the diffusion time? (b) Assume that the boron follows an Gaussian profile with a concentration of $5 \times 10^{16}/\text{cm}^3$ at the surface ($x' = 0$) and intersects the phosphorus profile at a level of $5 \times 10^{16}/\text{cm}^3$ at $x' = 100\text{ nm}$. What is the Dt product? If the boron were diffused at 1000°C , what would be the diffusion time?
- 10.21** Draw the cross section of a complementary bipolar process that simply adds an additional p diffusion to an $n\text{pn}$ process. Show the required biasing of the various diffusions, and indicate any problems you encounter.
- 10.22** A silicon bipolar transistor is to be designed to have an f_T of 50 GHz . What is its transit time? Suppose the transistor has $r_C = 40\text{ }\Omega$. What is the largest possible value for $(C_{JC} + C_{sub})$? What is an upper bound on the width of the collector depletion layer? What is an upper bound on the basewidth? If $r_E = 25\text{ }\Omega$, what is an upper bound on the emitter-base capacitance? (Assume that $V_s = 10^7\text{ cm/sec}$, $\eta = 10$, and $D_B = 20\text{ cm}^2/\text{sec}$.)
- 10.23** Identify potential latchup paths in the BiCMOS structures in Figs. 10.25(a)–(c).
- 10.24** A high energy (5 MeV) is used to implant oxygen well below the silicon surface to form a buried SiO_2 layer. Assume that the SiO_2 layer is desired to be $0.50\text{ }\mu\text{m}$ wide. (a) What is the oxygen dose required in silicon? (b) What beam current is required to be able to implant at least four 200-mm wafers per hour? (c) How much power is being supplied to the ion beam?
- 10.25** It was noted that the CDI process is used mainly for digital applications. What characteristics of the structure make this true?
- 10.26** (a) How many masks are required for the CDI process?
(b) Design a good mask-alignment sequence for this process.
- 10.27** A CDI process uses a 0.25-ohm-cm epitaxial base layer and a 5-ohm-cm substrate. Estimate the breakdown voltages of the emitter-base and collector-base junctions. The emitter junction depth is $1\text{ }\mu\text{m}$, and the epitaxial layer thickness is $2\text{ }\mu\text{m}$.

Processes for MicroElectroMechanical Systems: MEMS

MEMS represent one of today's most exciting areas of microelectronics activity. MEMS technology has brought together innovations from many areas of microelectronics only to develop rapidly into a discipline of its own. Today's micromachined systems combine the signal processing and computational capability of analog and digital integrated circuits with a wide variety of nonelectrical elements, including pressure, temperature and chemical sensors, mechanical gears, and actuators, 3D mirror structures, etc., and we have only begun to scratch the surface of biomedical applications. As a brief introduction, this chapter attempts to provide a flavor of the creativity and wide variety of devices and structures that are being conceived as you read this text. A wealth of greater detail on the subject can be found in the books by Madou [1] and Kovacs [2], the paper compendia edited by Muller [3] and Trimmer [4], and many other publications [5–8]. The latest research in the field is presented at the biennial International Conference on Solid-State Sensors and Actuators, the Solid-State Sensors and Actuator Workshop held on intervening years, and the yearly IEEE International Electron Devices Meeting.

MEMS structures are based upon our ability to sculpt or machine silicon on a microelectronic scale. These micromachining technologies can be broken into three groups: (i) Bulk micromachining, (ii) Surface micromachining, and (iii) High-aspect-ratio electroplated structures. The first, bulk micromachining dates back to the 1960s, when techniques for wet anisotropic etching of various forms of trenches, grooves, and membranes in silicon wafers were first developed. Advances in bulk micromachining continued rapidly through the 1970s with development of impurity-dependent etch stops, wafer-dissolution processes, and wafer fusion bonding. During the next two decades, surface micromachining, which makes use of the full lithography capability of IC processing, emerged to form a wide range of new beam, comb, microactuator, and rotary structures. The application of circuits to improve the sensor characteristics

advanced greatly during this time as well. In the 1990s, processes capable of producing high-aspect-ratio structures were developed using thick polymer molds and electroplating, and a variety of methods for merging MEMS fabrication with standard CMOS and bipolar processes were demonstrated. The sections that follow provide a basic introduction to the forms of processing previously mentioned.

11.1 MECHANICAL PROPERTIES OF SILICON

Before we explore the fabrication of mechanical devices in silicon, we will look at some of the mechanical properties of crystalline silicon. For example, these properties are needed to model the electromechanical behavior of resonant beams, comb actuators, and pressure diaphragms. Table 11.1 contains a comparative list of the mechanical properties of several materials [9, 2].

Most students in the laboratory quickly discover (upon dropping their first wafer) that silicon is a brittle material that yields catastrophically. This behavior is in contrast to our view of most common metals that tend to deform in a plastic manner. However, if we explore the table, we find that silicon is not as delicate as one might think. In fact, silicon is quite strong, with an intrinsic yield strength that exceeds high-strength steel and tungsten. Its Knoop hardness is similarly large, and Young's modulus for silicon is similar to that of iron and steel.

Wafers, however, tend to be large in diameter and relatively thin. In this state, they are fragile. Single-crystal material tends to cleave or fracture along crystal planes, and breakage can often be traced to flaws in the bulk, surface, and edges of the wafers. Pieces that have been sawed from wafers often have microscopic chips and defects along the edges that can also lead to premature failure. On the other hand, once the silicon wafer has been sawed into final form, even a large 2.5×2.5 cm die is quite robust. Note also the high value of Knoop hardness for silicon nitride, which is one reason nitride makes an excellent protective passivation layer for finished wafers. In addition, silicon nitride forms a nearly hermetic seal that resists moisture penetration.

TABLE 11.1 Mechanical Properties of Selected Materials

Material	Density (g/cm ³)	Yield Strength (GPa)	Knoop Hardness (kg/mm ²)	Young's Modulus (GPa)	Thermal Expansion Coefficient (10 ⁻⁶ /K)	Thermal Conductivity (W/cm-K)
Si	2.3	7.0	850	190	2.3	1.6
SiO ₂ (Fibers)	2.5	8.4	820	73	0.55	0.014
Si ₃ N ₄	3.1	14	3500	390	0.8	0.19
Diamond	3.5	53	7000	1000	1.0	20
Aluminum	2.7	0.17	130	70	25	2.4
Tungsten	19	4.0	490	410	4.5	1.8
Steel (max.)	7.9	4.2	1500	210	12	0.97
Polysilicon		1.2		170		

11.2. BULK MICROMACHINING

Bulk silicon etching, or bulk micromachining [10], is the selective removal of significant regions of silicon from the silicon substrate, ranging from simple cavity formation to almost complete dissolution of the wafer. Wet chemistry was used in most of the original micromachining processes, whereas vapor and plasma etching are in widespread use today. With careful design, the etching processes can be made to be compatible with on-chip CMOS circuitry, which has facilitated application of bulk micromachining to a wide range of sensors and actuators.

11.2.1. Isotropic and Anisotropic Etching

A summary of a few of the possible bulk micromachined profiles are outlined in Figs. 11.1–11.3. The isotropic etching case of Fig. 11.1 was also discussed in Chapter 2. Here the etch progresses equally in all directions and will undercut the masking material at the surface. Agitation of the fluid during etching creates more rounded features, but does not significantly change the undercutting at the surface.

In contrast, anisotropic etches attack certain crystal planes much more rapidly than others and produce etched cavities with flat surfaces that intersect each other at sharp well-defined angles. For the case of silicon, $\langle 100 \rangle$ and $\langle 110 \rangle$ planes etch at much higher rates than $\langle 111 \rangle$ planes. For example, a solution of KOH, water, and alcohol [11, 12] etches the $\langle 100 \rangle$, $\langle 110 \rangle$, and $\langle 111 \rangle$ crystal planes at relative rates of 40:30:1.

This etch rate selectivity can be used to create various cavity and groove structures, as shown in Fig. 11.2. On $\langle 100 \rangle$ silicon, the $\langle 111 \rangle$ planes intersect the surface at an angle of 54.74° , forming pyramidal cavities and v-grooves bounded by these planes. The etching front is defined by the oxide opening and proceeds vertically downward as the $\langle 100 \rangle$ surface is etched away, exposing the slow-etching $\langle 111 \rangle$ planes. The cavity thus formed has a flat bottom with sloping walls on all four sides of the cavity, since the $\langle 100 \rangle$ surface is symmetrical under 90° rotation. If the etch is permitted to continue, the

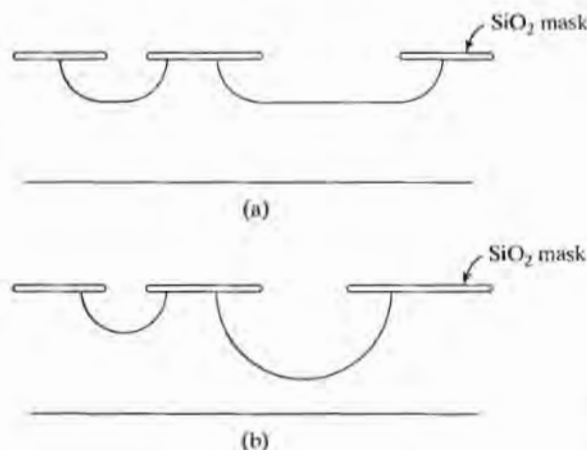


FIGURE 11.1

(a) Isotropic etching without agitation (b) with agitation

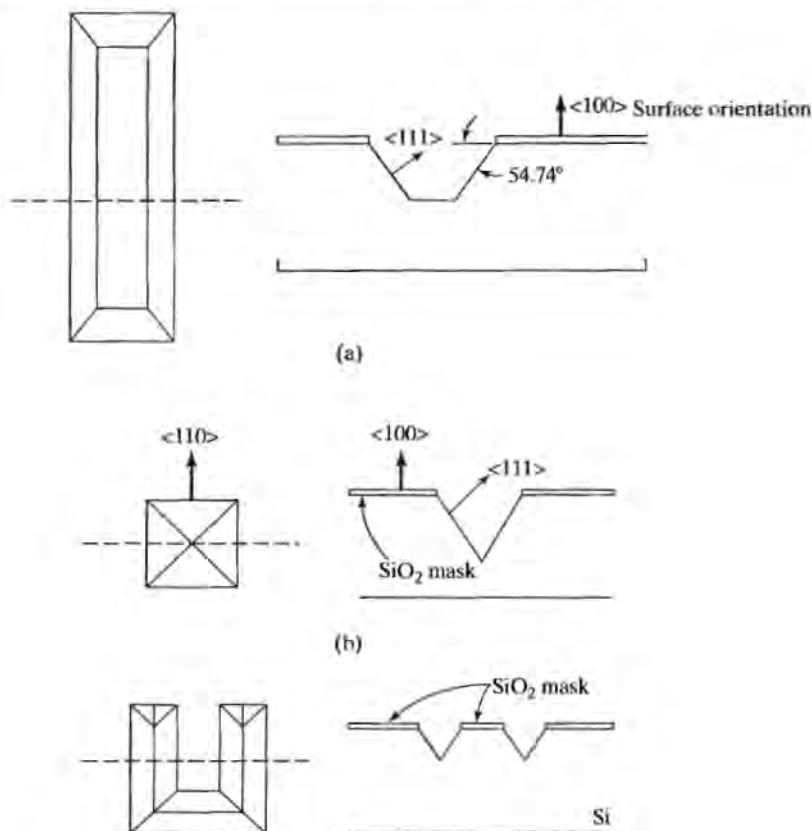


FIGURE 11.2

Anisotropic etching of $\langle 100 \rangle$ silicon exposes slow etching $\langle 111 \rangle$ planes which form a 54.74° degree angle with the surface. (a) Formation of a flat-bottomed trench (b) A cavity in the form of an inverted pyramid (c) Cavity in the form of a "v-groove".

$\langle 111 \rangle$ planes ultimately intersect each other, forming an inverted pyramid, or v-groove, as indicated in the figure. The depth of flat-bottomed cavities is determined by the etch rate and etching time. On the other hand, the v-groove is essentially a self-terminating structure, because of the slow etch rate of the $\langle 111 \rangle$ planes, and the groove depth is defined by the width of the mask opening at the surface. Inkjet printer nozzles [11] formed by etching cavities completely through the wafer represent an early application of anisotropic etching techniques. It is important to realize that the sidewalls are truly $\langle 111 \rangle$ planes only for an etchant with infinite selectivity. The actual angle of the sidewalls will be approximately 54.6° for 400:1 etch selectivity and 52° for a 20:1 selectivity.

Cavities with vertical side walls can be etched in wafers with $\langle 110 \rangle$ surface orientations as indicated in Fig. 11.3(a). For this surface, the slow-etching $\langle 111 \rangle$ planes produce vertical side walls. This technique was used by Tuckerman and Pease [13] to fabricate an experimental microchannel heat sink with very low thermal resistance as

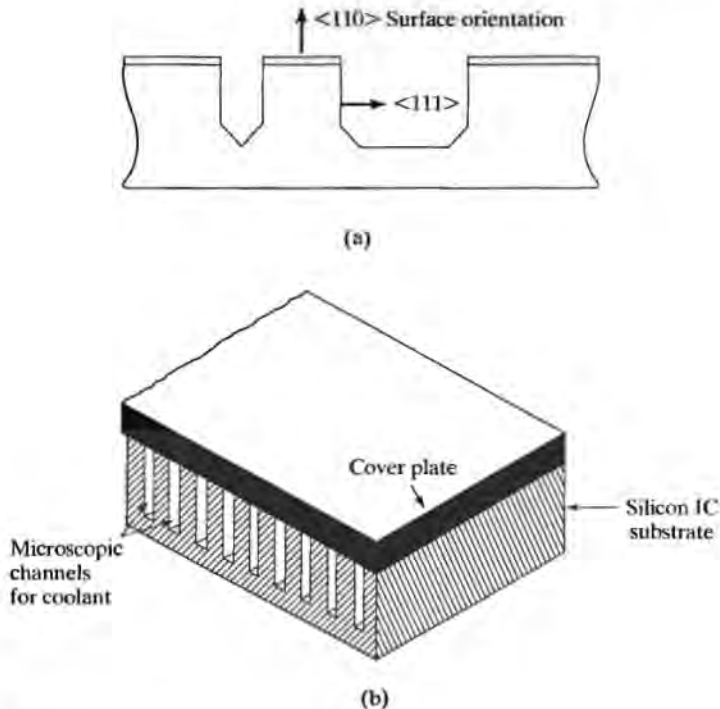


FIGURE 11.3

(a) Anisotropic etching of $\langle 110 \rangle$ silicon can form cavities with vertical walls (b) Microchannel heat sink developed by Tuckerman and Pease. Thermal resistances below 0.1°C/W were measured. Copyright 1981, IEEE. Reprinted with permission from Ref. [13].

depicted in Fig. 11.3(b). Here is an additional note of caution. There are additional pairs of $\langle 111 \rangle$ planes that become exposed in the corners of the trenches forming angles of approximately 22° relative to the surface. The etch barrier openings need to be long and relatively narrow to produce deep cavities similar to those indicated in the figure.

11.2.2 Diaphragm Formation

Silicon diaphragms similar to those in Fig. 11.4 have been used for thermal isolation and in various forms of pressure sensors for many years. Here, a relatively large cavity is etched from the backside of the wafer and is terminated with a few-micron-thick silicon diaphragm remaining at the surface. Double-polished wafers are utilized in double-sided processing, and special two-sided alignment systems employing special split-field optics or infrared through-wafer viewing have been developed especially for MEMS applications. A layer of wax can be used to protect the topside of the wafer during etching. For production, however, more complex mechanical holders are used to protect one side of the wafer during wet anisotropic etching.

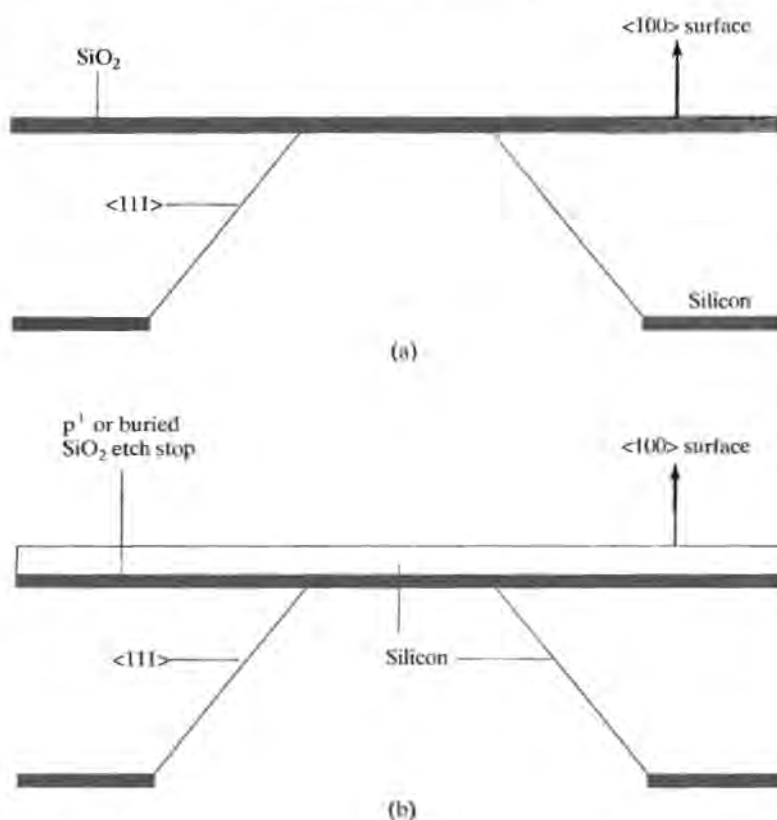


FIGURE 11.4

Diaphragms formed by anisotropic backside etching of the silicon wafer (a) SiO_2 layer/diaphragm used as an etch stop (b) Buried SiO_2 or p^+ layer can be used as an etch stop to form thin diaphragms.

Etching through a 500- μm wafer, but stopping with only 25 μm or less remaining, requires tight etching control. However, another characteristic of the anisotropic etching process was found to solve this control problem. It was discovered that the etch will terminate when it encounters a heavily doped boron region. Thus, a buried p^+ layer can be used as an etch stop. A Boron doping level of $10^{20}/\text{cm}^3$ reduces the etch rate between 10 and 100 times depending upon the etchant. Another possibility is to terminate the etch on a buried oxide layer, such as is available in either SIMOX or fusion-bonded SOI wafers. In fact, the wafer-to-wafer bonding techniques that are utilized today to produce some SOI wafers were originally developed for MEMS applications.

A more flexible etch stop technique employs an electrochemical method in which voltage biases are applied to different n - and p -type regions in the semiconductor while the sample is immersed in the etchant. With the electrochemical approach, the etch can be controlled to terminate on either n - or p -type regions.

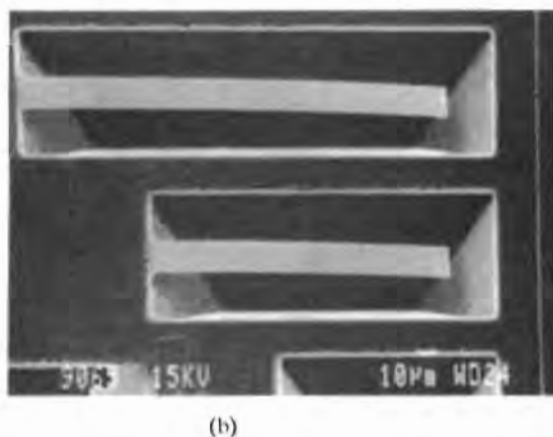
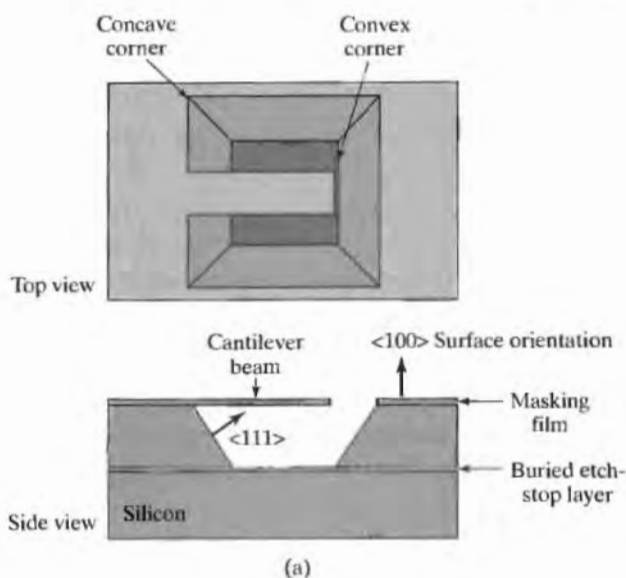


FIGURE 11.5

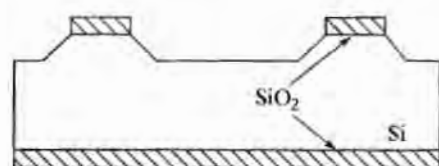
(a) Cavity anisotropically etched into $\langle 100 \rangle$ silicon using a buried etch stop and undercutting to form a free cantilever beam (b) SEM photograph of cantilever beams formed of SiO_2 . Reprinted with permission from Ref. [14].

11.2.3 Cantilever Beams and Released Structures

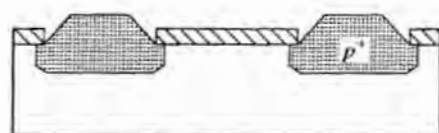
Figure 11.5 depicts formation of a cantilever beam structure using bulk micromachining. In this case, a cavity is etched in the surface of an epitaxial silicon wafer and terminates on a p^+ or oxide etch-stop layer. At the same time, the etchant removes the silicon completely from below the oxide region. Upon the completion of the process, a

cantilever beam overhangs the silicon cavity. Silicon dioxide, silicon nitride, and metal coated beams represent various process options. As one application, electrostatic deflection can be used to modulate the position of such cantilever beams to form the basis of a simple mechanical light modulator.

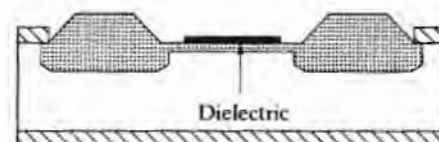
Machined structures can be completely released from the wafer using the p^+ etch stop technique. Fabrication of a capacitive pressure sensor is shown in Fig. 11.6 [14]. The wafer is first anisotropically etched to form a cavity topology on the surface. A heavy deep boron diffusion is then used to define a region that will become a physical support rim for the structure. A shallow boron diffusion is added, and an insulating dielectric film is deposited over the shallow diffusion and patterned. The bulk of the wafer is then sacrificed to a silicon etch, producing a fully sculptured silicon structure that is bounded by the diffused p -type regions. This process is referred to as a *dissolved-wafer* process. In the application, the released structure forms one plate of a parallel-plate capacitor and is anodically bonded (see Section 11.6.3) to a glass substrate containing the other capacitor electrode and interconnections.



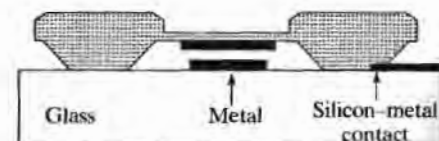
(a) KOH etch



(b) Deep boron diffusion



(c) Shallow boron diffusion; dielectric deposition



(d) Electrostatic bonding; final wafer dissolution

FIGURE 11.6

Dissolved wafer process (a) Etch cavities to form desired surface topology (b) Selective diffusion to form p^+ etch stop regions (c) Second p region defines thin diaphragm and insulating dielectric layer (d) Structure following bonding and removal of most of wafer. Copyright 1990, IEEE. Reprinted with permission from Ref. [14].

11.3 SILICON ETCHANTS

The structures described in the previous section require the use of well-characterized etchants, and a variety of chemistries have been formulated for micromachining silicon. A summary of the characteristics of the important etchant categories appears in Table 11.2 [10]. Much of the fabrication is moving from wet chemistry to plasma-based dry etching, because of environmental concerns and clean room compatibility issues.

11.3.1 Isotropic Etching

HNA, a mixture of hydrofluoric acid, nitric acid, and acetic acid is the most common isotropic wet etchant for silicon. Rounding and undercutting can be controlled to some degree by solution agitation, as depicted in Fig. 11.2. Xenon difluoride, XeF_2 , a solid at

TABLE 11.2 Comparison of Bulk Silicon Etchants [Reprinted with permission from Kovacs et al. [10].

	HNA (HF + HNO_3 + Acetic Acid)	Alkali-OH	EDP (ethylene diamine pyrochatechol)	TMAH (tetramethyl- ammonium hydroxide)	XeF_2	SF_6 Plasma	DRIE (Deep Reactive Ion Etch)
Etch Type	wet	wet	wet	wet	dry ²	dry	dry
Anisotropic?	no	yes	yes	yes	no	varies	yes
Availability	common	common	moderate	moderate	limited	common	limited
Si Etch Rate $\mu\text{m}/\text{min}$	1 to 3	1 to 2	0.02 to 1	≈ 1	1 to 3	≈ 1	> 1
Si Roughness	low	low	low	variable ³	high ⁴	variable	low
Nitride Etch	low	low	low	1 to 10 nm/min	low ⁵	low	low
Oxide Etch	10 to 30 nm/min	1 to 10 nm/min	1 to 80 nm/min	$\approx 1 \text{ nm}/\text{min}$	very low	low	low
Al Selective	no	no	no ⁶	yes ⁷	yes	yes	yes
Au Selective	likely	yes	yes	yes	yes	yes	yes
p++ Etch Stop?	no (n slows)	yes	yes	yes	no	no (some dopant effects)	no
Electrochemical Stop?		yes	yes	yes	no	no	no
CMOS Compatible? ⁸	no	no	yes	yes	yes	yes	yes
Cost ⁹	low	low	moderate	moderate	moderate	high	high
Disposal	low	easy	difficult	moderate	N/A	N/A	N/A
Safety	moderate	moderate	low	high	moderate?	high	high

²Sublimation from solid source.

³Varies with wt% TMAH, can be controlled to yield very low roughness.

⁴Addition of Xe to vary stoichiometry in F or Br etch systems can yield optically smooth surfaces.

⁵Low for low hydrogen content film.

⁶Some formulations do not attack Al, but are not common.

⁷With added Si, polysilic acid or pH control.

⁸Defined as (1) allowing wafer to be immersed directly with no special measures and (2) no alkali ions.

⁹Includes cost of equipment.

room temperature and pressure, sublimates at room temperature in a vacuum and provides one form of dry isotropic etchant for silicon. XeF_2 exhibits high selectivity for photoresist, silicon dioxide, silicon nitride, and aluminum. [15]

11.3.2 Anisotropic Etching

Anisotropic etching has most often been performed using wet chemistry with one of the more common etchants being potassium hydroxide (KOH). Various KOH solutions can achieve $\langle 110 \rangle$: $\langle 100 \rangle$: $\langle 111 \rangle$ plane selectivity as high as 600:400:1 and provide etch rates of up to 2 $\mu\text{m}/\text{min}$ with good selectivity for oxide, nitride, and heavily doped boron etch stop regions. However, KOH is corrosive and attacks aluminum. Remember that alkali ions—particularly sodium—are a potentially serious contaminant in MOS gate oxides. Potassium ions represent a similar problem, and KOH is not permitted in most IC clean rooms.

Although quite hazardous, EDP (a solution of ethylene diamine, pyrochatechol and water¹) is one of the useful dopant-dependent anisotropic etches for silicon, with a $\langle 100 \rangle$: $\langle 111 \rangle$ selectivity of approximately 35:1 and an etch rate approaching 1 $\mu\text{m}/\text{min}$. Although the selectivity and etch rate are lower than KOH, the selectivity for *p*-type doping is considerably higher. EDP, like KOH, attacks aluminum metallization and limits its use with CMOS integrated circuits. Unfortunately, EDP mixtures are extremely corrosive, potentially carcinogenic, and seldom, if ever, allowed in commercial IC clean rooms. Safety considerations also require the use of a reflux condenser.

TMAH, a quaternary ammonium compound ($(\text{CH}_3)_4\text{NOH}$), is one of the more useful clean room compatible wet etchants for silicon. TMAH does not contain alkali ions and is available as clean-room grade solutions. The etch rate is approximately 1 $\mu\text{m}/\text{min}$, although the $\langle 100 \rangle$: $\langle 111 \rangle$ selectivity of 10–35:1 is the lowest of the three etchants mentioned here. The etch rate falls rapidly for boron concentrations of $10^{20}/\text{cm}^3$ and higher. Etch rates for silicon dioxide and nitride, as well as aluminum, can be made quite low with proper formulation of the TMAH solution.

Dry-plasma and reactive-ion etching of bulk silicon provides a full range of isotropic and anisotropic etching possibilities. Sulfur hexafluoride (SF_6) is a gas source commonly utilized in plasma systems. Fluorine free radicals formed by disassociation in the plasma produce etch rates comparable to wet chemistry with isotropic profiles. CFCs are used in some plasma systems to simultaneously deposit polymers on the side walls of etched regions. These polymers substantially reduce the side-wall etch rate and result in anisotropic plasma etching.

Deep reactive-ion etching (DRIE) can be used to produce high-aspect-ratio structures. The Bosch process [16] involves a sequence of alternative etching and polymer deposition phases within the plasma-etching system. High-aspect-ratio, vertical-side-wall structures can be produced with high selectivity for photoresist, and silicon-dioxide-masking layers. Precise etch stops can be provided by buried SiO_2 layers formed by SIMOX, or wafer-bonding processes. DRIE systems have also been developed that can rapidly etch vertical profiles completely through silicon wafers. Dry-plasma techniques are also now available for thinning of silicon wafers.

¹Also referred to as EPW.

A structure produced by deep RIE appears in Fig. 11.7. A cavity is first etched in a silicon wafer. A second wafer is then bonded on top of the first and thinned to the desired thickness by etching. DRIE is then used to etch a high-aspect-ratio spring structure into the thinned silicon material. The DRIE process scales well geometrically with a useful aspect ratio of 20–25:1.

11.4 SURFACE MICROMACHINING

Bulk micromachining techniques, other than plasma processes, are most often used to produce relatively large sculpturing of the surface or backside of silicon wafers. Surface micromachining on the other hand, builds up microstructures on the surface of the silicon wafer by making use of the full lithography capability of integrated circuit processing. The processes commonly involve deposition, patterning and removal of various layers of silicon dioxide, polysilicon, phosphosilicate glass (PSG), and silicon nitride. PSG is widely used as a “sacrificial layer,” because of its relatively high etch rate in hydrofluoric acid solutions.

11.4.1 Cantilever Beams, Bridges, and Sealed Cavities

Typical steps in surface micromachining may best be described through examples of the formation of such structures. Figure 11.8 depicts the formation of a free-standing polysilicon bridge and a cantilever beam [17, 18]. First, a sacrificial silicon-dioxide layer is formed by a combination of thermal and CVD oxides and is patterned to open an anchor point for the polysilicon layer, which is then deposited and patterned to form the desired beam and bridge geometries. Polysilicon layers provide excellent step coverage. The wafer is then immersed in a 49% HF solution to etch away the sacrificial oxide, leaving the cantilever beam and bridge structures anchored to the silicon substrate.

If the bridge structure is tapered on all four sides, it may be turned into a sealed cavity, as shown in Fig. 11.9. Such cavities can be used as fluid channels or sealed-pressure references. Access holes are etched through the polysilicon layer, or formed as gaps in the anchor region of the polysilicon. Cavities can be sealed with thin films of materials deposited by sputtering, evaporation, or CVD, or a reactive seal may be formed by oxidizing the polysilicon layer. Materials such as photoresist, epoxy, or polyimide can also seal the cavity. Processes have been developed to achieve a variety of sealing conditions ranging from vacuum to atmospheric pressure.

11.4.2 Movable In-Plane Structures

A second surface micromachining example demonstrating the formation of a basic rotating structure [20] that may form the basis for a turbine, gear train, or micromotor is depicted in Fig. 11.10. A layer of thermal oxide is grown on the silicon surface and patterned. A second mask is used to create the center step in the oxide shown in Fig. 11.10(a). A second layer of polysilicon is deposited and covered with CVD oxide that is patterned and used as a barrier layer during reactive-ion etching of the polysilicon. An additional CVD oxide is deposited over the complete structure and patterned using plasma etching to open the hole in the center that becomes a substrate anchor

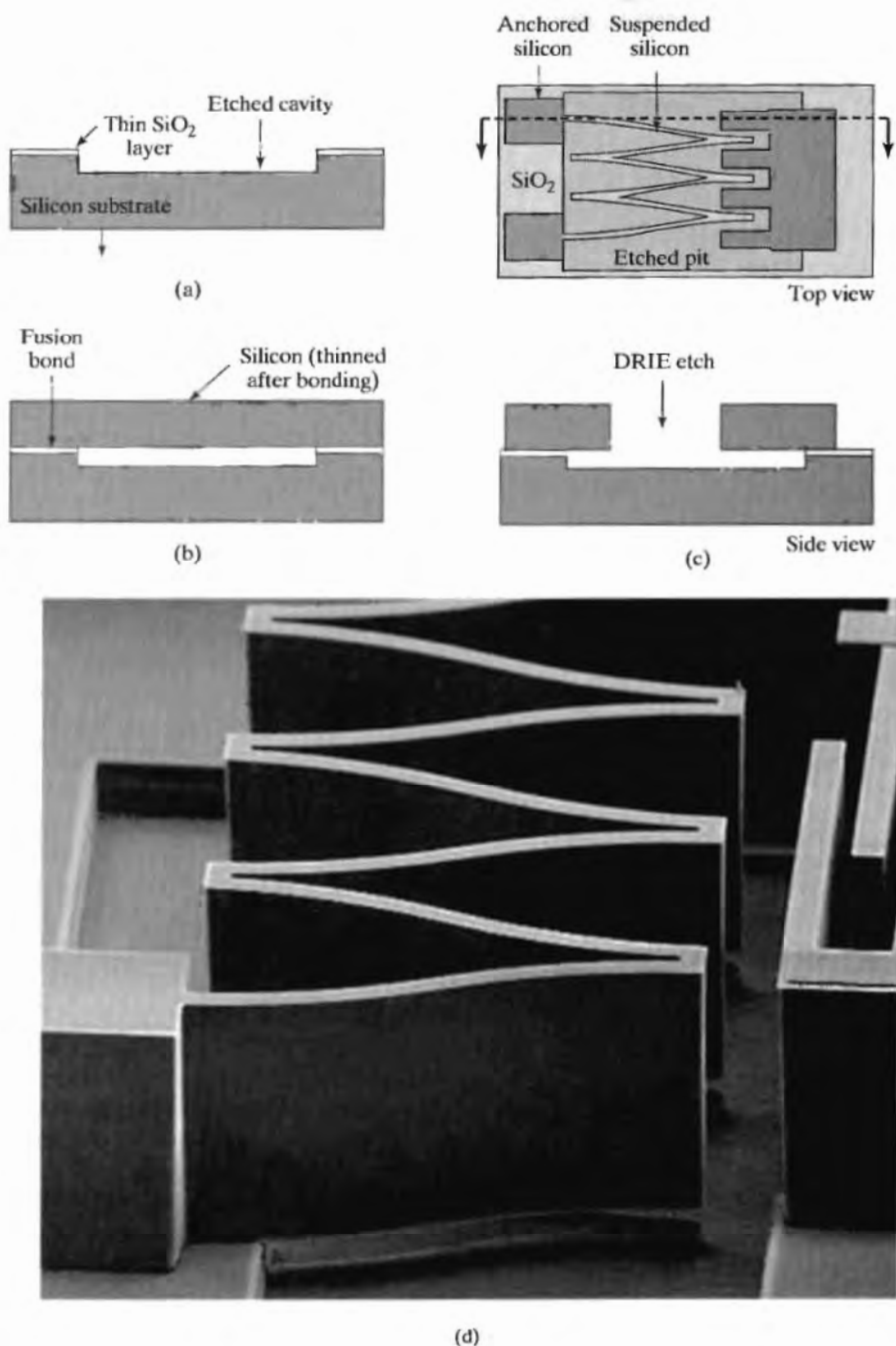


FIGURE 11.7

Deep RIE etching applied to form spring structure. (a) Cavity etched in silicon (b) Second wafer bonded over etched cavity and thinned (c) Spring structure following DRIE processing. (d) Photomicrograph of completed spring structure. Copyright 1998, IEEE. Reprinted with permission from Ref. [10].

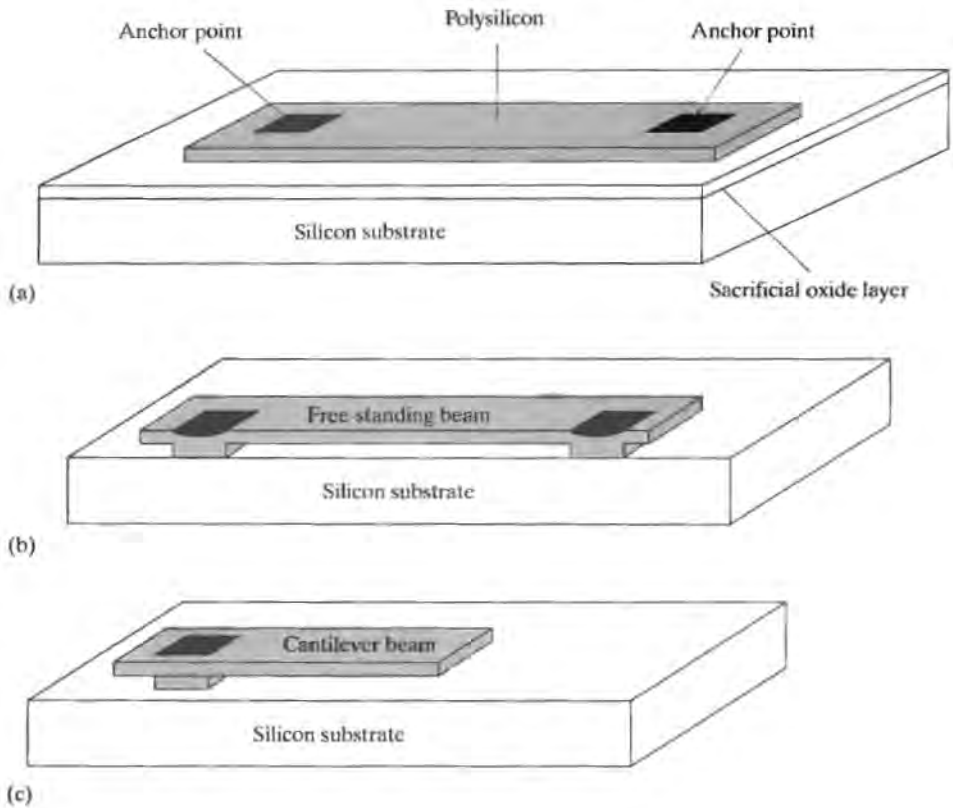


FIGURE 11.8

Surface micromachining (a) Polysilicon bridge following deposition and patterning (b) Bridge following removal of the sacrificial oxide (c) Cantilever beam can be formed if one bridge support is eliminated.

point for the rotor shaft. A second polysilicon layer is deposited and etched to become the shaft for the rotor. The oxide layers represent sacrificial layers that are all removed by etching in a final release step in hydrofluoric acid. Once the oxide is etched away, the rotor is free to turn around the shaft. The various clearances between the rotor, shaft, and substrate are all determined by differences in thicknesses of the various deposited layers and are of the order of $1\text{--}3\text{ }\mu\text{m}$ in this structure.

Figure 11.11 presents a resonant beam structure with electrostatic combs utilized for input drive and output signal sensing. The beams in the center of the structure are anchored to the substrate at two points and are otherwise free to vibrate. Excitation is applied electrostatically [21] to the comb at one end, and the output signal is removed from the other end. Signal transfer is frequency dependent with high Q resonances determined by the physical dimensions and mechanical properties of the beams. Resonant comb structures can be used to induce motion in microactuators and are being studied as alternatives to crystal oscillators and a variety of filters used in communications receivers and transmitters.

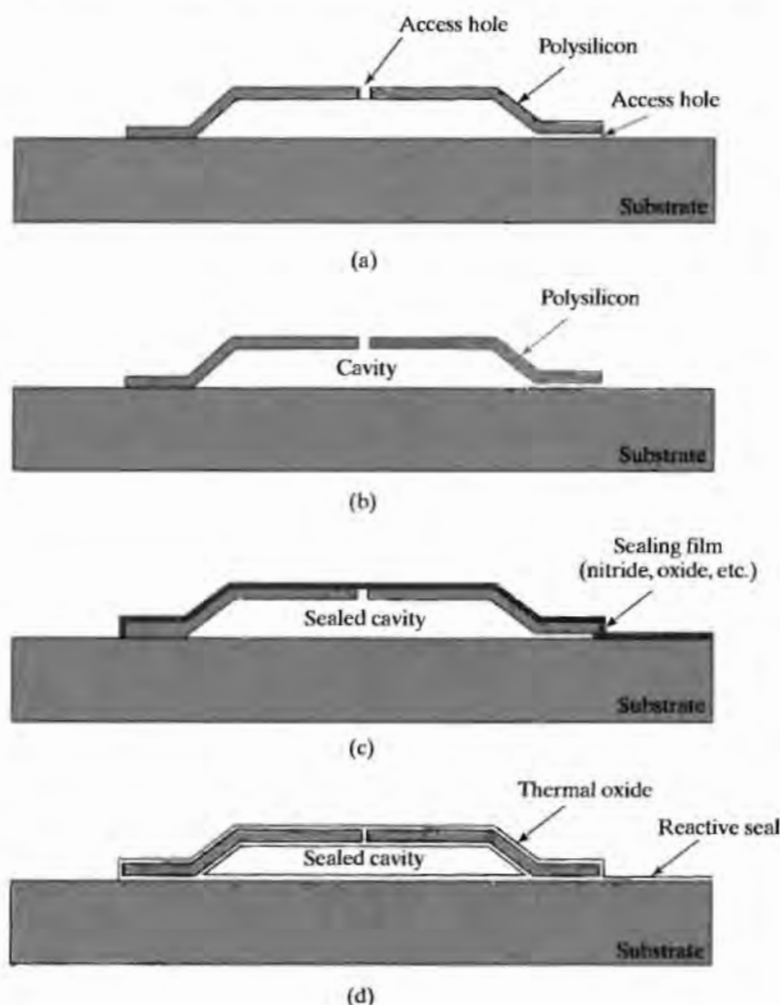


FIGURE 11.9

Scaled cavity formation (a) Polysilicon cap over sacrificial oxide with holes for etchant access (b) Covered cavity is formed upon removal of the oxide (c) Cavity can be sealed with a variety of deposited films (d) A reactive seal can be formed by thermal oxidation.

11.4.3 Out-of-Plane Motion

The various surface micromachined structures outlined thus far are essentially two-dimensional in nature. To achieve operation up out of the wafer plane and into the third dimension, some form of actuator and hinge structure must be fabricated. Figure 11.12 depicts three possible hinge structures [22] that may be formed using the fabrication sequence outlined in Fig. 11.13.

The formation of these hinge structures begins with the deposition of a trilayer film consisting of a 1–2- μm thick sacrificial layer of PSG on the surface of the silicon wafer, an

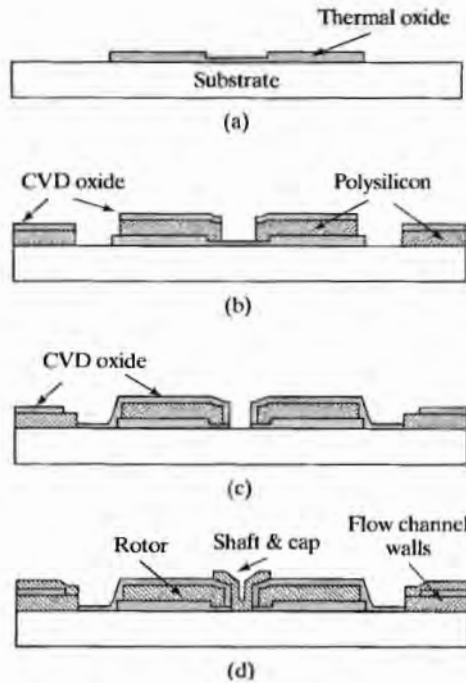


FIGURE 11.10

Formation of rotary structures by surface micromachining (a) Sacrificial oxide definition (b) Structure following deposition of polysilicon and CVD oxide. Lithography has defined the rotor of the structure in the center as well as polysilicon anchored to the silicon surface at the sides. (c) Deposition and patterning of a second sacrificial oxide layer (d) Definition of rotor shaft in a second layer of polysilicon (e) SEM of partially released 240 μm diameter gear that is sawed in half. Copyright 1988, IEEE. Reprinted with permission from Ref. [20].

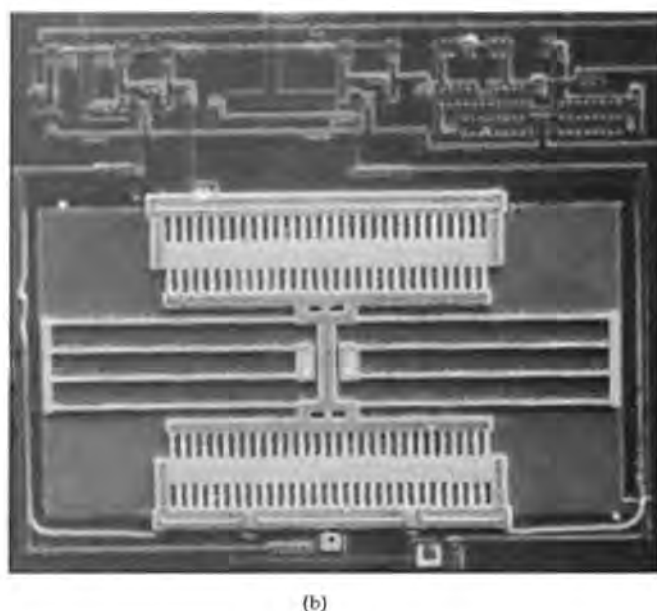
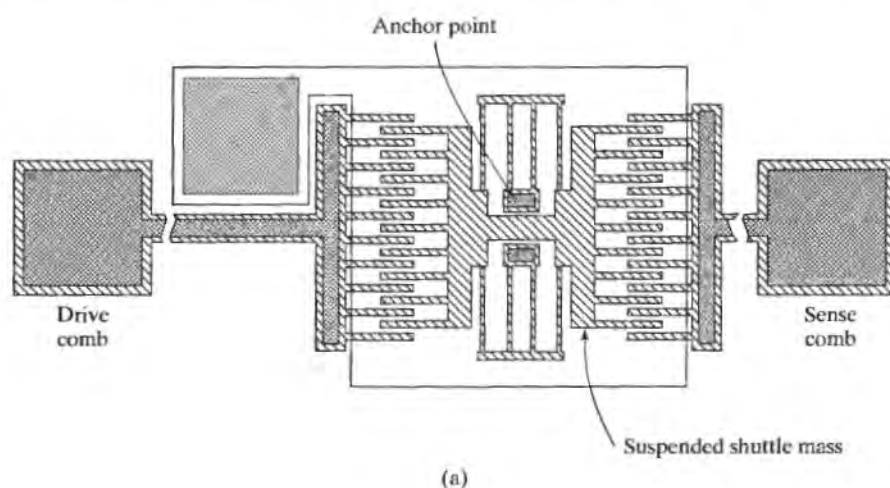


FIGURE 11.11

(a) Layout of a polysilicon resonant beam structure with electrostatic comb drive. (b) Microphotograph of a fabricated structure. Copyright 1993, IEEE. Reprinted with permission from Ref. [21].

LPCVD deposition of an undoped polysilicon layer of similar thickness, and, finally, a second PSG layer. An anneal in nitrogen at 950°C is used to dope the polysilicon with phosphorus from the PSG layer and to minimize mechanical stress in the polysilicon.

Next, the top PSG layer is completely removed, and the polysilicon layer is patterned to achieve the desired hinge plate and pin geometry from Fig. 11.12. The second sacrificial layer of PSG is deposited and patterned to open contact points to either the substrate or the first layer of polysilicon. A second sandwich of polysilicon and PSG is deposited and undergoes another combined doping-annealing step. The top PSG is removed, and the second layer of polysilicon is patterned to define staples that anchor

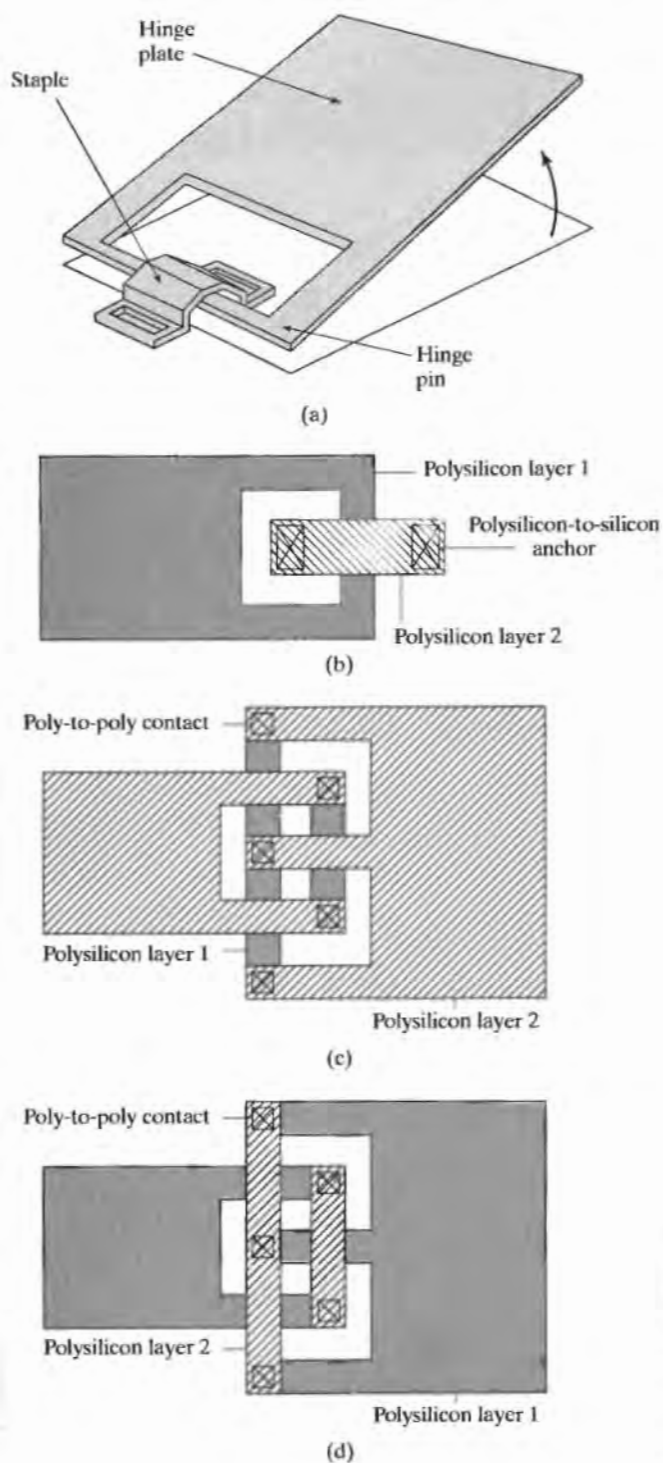


FIGURE 11.12

Three-dimensional hinges (a) Artists' concept of raised hinge (b) Layout of a hinge similar to that in (a) which is anchored to the silicon surface; (c) and (d) represent polysilicon plates that are hinged to each other using two forms of hinges, but are free from the silicon. Copyright 1992, Pergamon Press, *Sensors and Actuators*, reprinted with permission from Ref. [22].

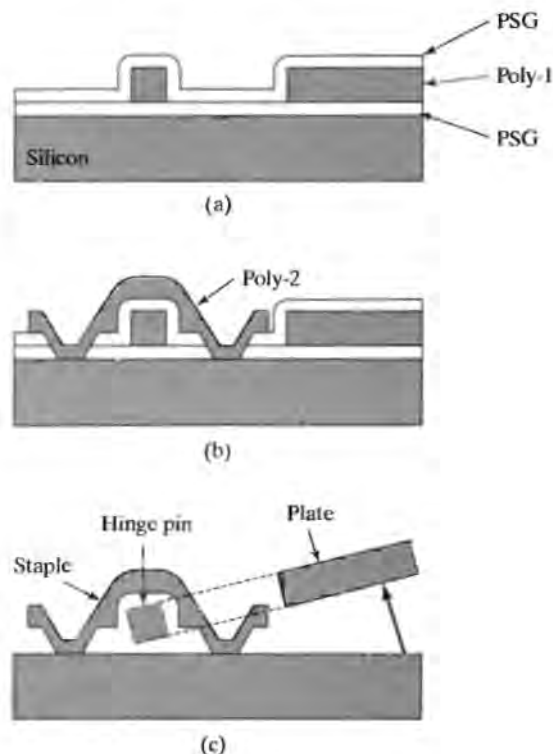


FIGURE 11.13

Formation of the hinge in Fig. 11.12a. (a) Definition of the hinge plate in the first polysilicon layer with two sacrificial oxide layers (b) Formation of the anchored "staple" in the second poly layer (c) Anchor and plate following removal of the sacrificial oxide layers.

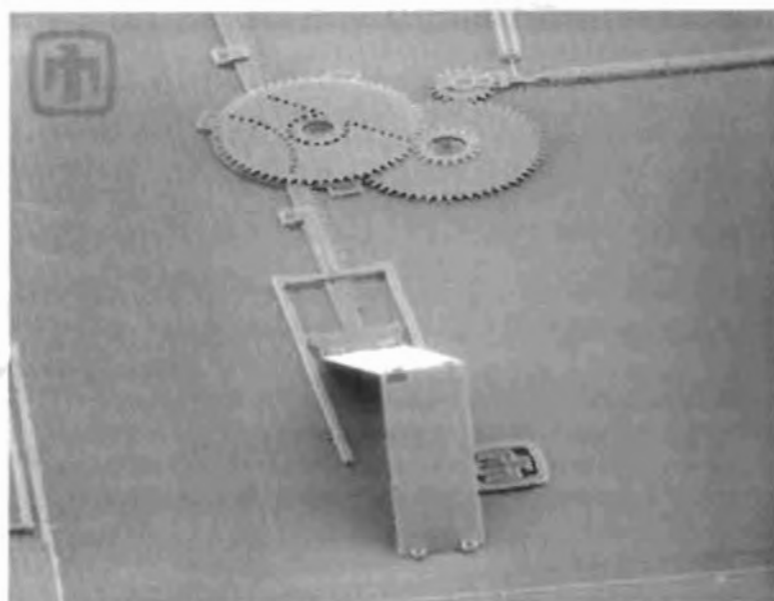
Copyright 1992, Pergamon Press, Sensors and Actuators, reprinted with permission from Ref. [22].

the structure to the substrate or connect the two polysilicon layers to form the hinges. The final step is to release the structures by etching away the sacrificial PSG oxides in HF, rinsing, and drying.

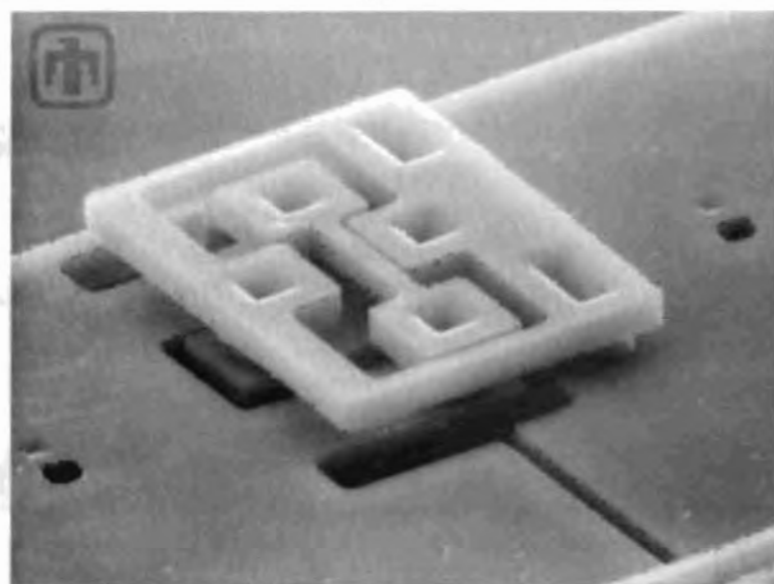
The hinged devices are now free to be moved up off the surface. In the first demonstration devices, movement was accomplished in a tedious manual fashion. More recently, structures have been fabricated that include mechanisms necessary to raise the hinge plate out of the wafer plane and control its position. Figure 11.14 provides an example of such a three-dimensional structure fabricated at Sandia National Laboratories (Courtesy of Sandia National Laboratories, [23]).

11.4.4 Release Problems

A subtle problem with wet-released structures is the tendency of the surfaces to stick together. Capillary forces tend to pull the surfaces toward each other, and once in contact, the surfaces remain stuck together. Hydrogen-oxygen bonding and van der Waals forces have been proposed as two mechanisms for the sticking problem. One anti-sticking approach uses critical-point drying of carbon dioxide. Water is replaced by methanol and then by liquid CO_2 under high pressure. As the temperature is raised, the CO_2 changes to gas and is removed. A phase-change release technique using *t*-butyl alcohol is a second method around the problem [24]. Water in the structure is



(a)



(b)

FIGURE 11.14

(a) SEM photograph of hinged structure raised with gear drive. (b) Close-up of one of the hinges. (Courtesy of Sandia National Laboratories. [23])

replaced by the *t*-butyl alcohol, which is then frozen by reducing the temperature to slightly below room temperature. The alcohol is subsequently removed by sublimation in a vacuum, leaving an inherently dry structure.

11.5 HIGH-ASPECT-RATIO MICROMACHINING: THE LIGA MOLDING PROCESS

The aspect ratio of devices produced by surface micromachining is limited by the relatively thin layers that can be deposited. An alternative process, called LIGA [lithography, galvanoforming (electroplating) and abforming (molding)], described in this section uses thick photoresists as molds that are subsequently filled via metal plating processes. The LIGA process produces micromachined devices with very high aspect ratios with thicknesses that range from 50 to as much as 500 μm [25].

The conceptual process depicted in Fig. 11.15 represents the basic LIGA process enhanced with surface micromachining (SLIGA). A sacrificial layer is deposited on

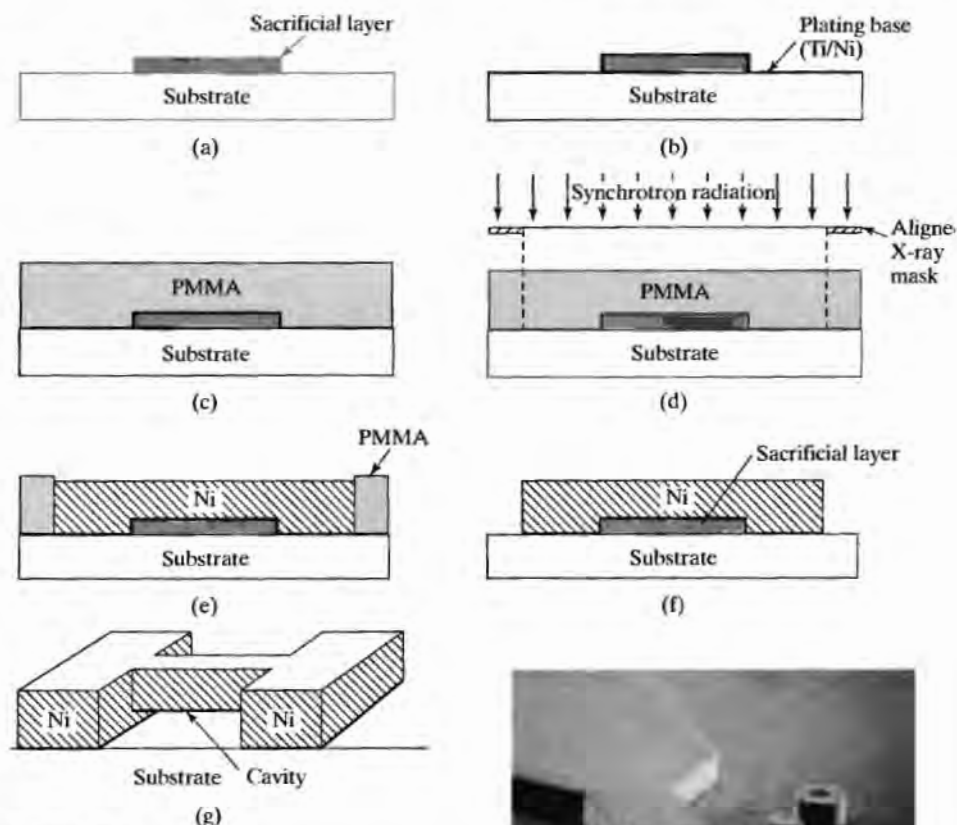


FIGURE 11.15

LIGA process including sacrificial oxide layer. (a) Patterned sacrificial oxide (b) Titanium nitride plating base applied (c) Thick PMMA layer has been deposited (d) Pattern is defined in PMMA using x-ray lithography (e) After development, the open cavities are filled with electroplated nickel (f) Structure following removal of the PMMA and plating base (g) Final molded nickel structure after etching of sacrificial oxide layer. (h) SEM photograph. Copyright 1998, IEEE. Reprinted with permission from Ref. [25].



(h)

the substrate, patterned, and then covered by a thin metal plating base, such as sputtered Titanium/Nickel. A thick, 50–500 μm , layer of PMMA (polymethylmethacrylate), an X-ray sensitive photoresist material, is deposited on the surface. Photon illumination from a synchrotron X-ray source is used to expose the desired pattern through the full thickness of the PMMA layer. After development, a microscopic mold of the desired part is formed. Then the resulting cavities in the PMMA are filled with electroplated Nickel. Permalloy (nickel/iron), gold, copper, and many other materials are commonly electroplated as well. Fabrication is completed by the removal of the PMMA and base-plating electrodes, as well as the sacrificial layer beneath the plated Ni. Structures ranging from 50 μm to several millimeters in height have been fabricated by using the LIGA process.

An example of a completed test device appears in Fig. 11.15(h) [25]. Dimensional variations in such a structure can be held to as little as 0.1 μm per 100 μm of height for a 100- μm -diameter gear. However, very narrow-width devices are not possible because the PMMA mold is not mechanically stable. Obviously a primary limitation of the process is the requirement for an X-ray exposure system. In addition, X-ray masks are also considerably more difficult to fabricate than masks for standard optical lithography.

11.6 SILICON WAFER BONDING

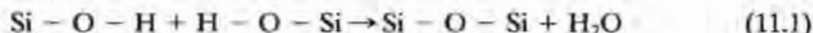
Bonding of silicon to itself or other materials has long been used to create MEMS devices. Bonding permits the joining of silicon devices that must be fabricated by incompatible silicon processes, as well as devices that must involve materials or substrates other than silicon, particularly glass. Bonding is also a highly important process for the packaging of many sensors.

11.6.1 Adhesive Bonding

The simplest form of bonding, and one that satisfies many requirements, utilizes an adhesive layer such as a commercial epoxy or other glue. Photoresist, polyimide, and PMMA have been used successfully as adhesive and protective layers. If necessary, oxygen plasmas can remove (or ash) these compounds, although residual inorganic contaminants may remain.

11.6.2 Silicon Fusion Bonding

Direct wafer-to-wafer bonding is used to fuse silicon substrates to each other [26]. Although the mechanism is not totally understood, the bonding is assumed to involve formation of a silicon–oxygen–silicon bond between the two wafers with liberation of a water molecule:



The surfaces to be bonded must be extremely smooth, flat, and clean, and as part of the final cleaning procedure, the surfaces are normally hydrated to insure presence

of an abundance of O-H groups. The centers of the wafers are brought into contact. Once initiated, a contact wave front sweeps across the wafers joining the entire surfaces. A high-temperature-annealing step significantly increases the strength of the bond. At temperatures above 1000°C, the bond strength is said to approach that of the silicon wafer itself. Obviously, any particulate matter on either wafer surface will cause voids to form in the bonding.

An example of the formation of an experimental reentrant cavity heat sink for liquid cooling applications using anisotropic etching and silicon fusion bonding is shown in Fig. 11.16. Here, a two-step anisotropic etching process has been used to form cavities that reach through the silicon wafer. The cavity wafer is then fusion bonded to the wafer containing high-heat-flux ICs.

An alternative fusion-bonding technique first forms an oxide on the surface of each wafer. Bonding then proceeds in a manner similar to that just described. This second method is one of the techniques used to form the silicon-on-insulator substrates required for both the SOI CMOS and dielectrically isolated bipolar technologies mentioned in Chapters 9 and 10. One additional bonding technique utilizes thin intermedi-

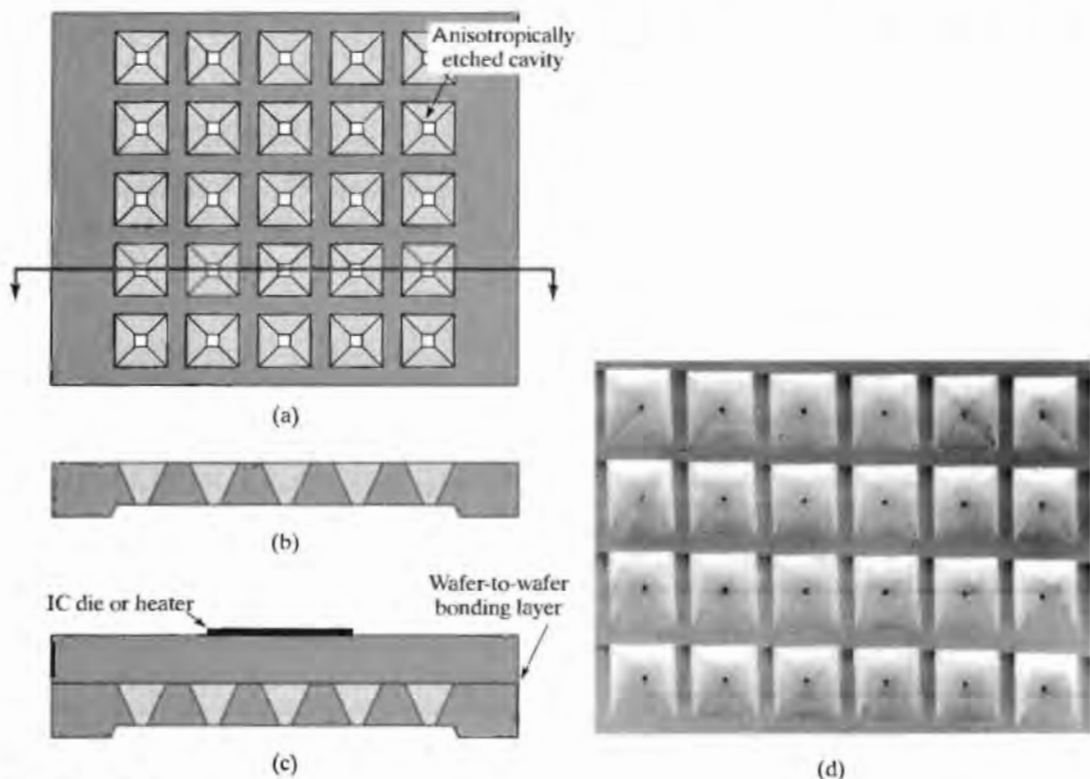


FIGURE 11.16

Reentrant cavity heat sink structure (a) Top view of silicon wafer with etched cavities (b) cavity cross section following second etching step (c) Cavity wafer fusion bonded to second silicon wafer with thermal source on top (d) SEM photograph of cavity array with 420 cavities/cm². Copyright 1993, IEEE. Reprinted with permission from Ref. [27].

ate layers of evaporated gold can be utilized to form a conductive silicon–gold eutectic bond between two wafers.

11.6.3 Anodic Bonding

Silicon (as well as metals) can be bonded to glass using a low-temperature field-assisted anodic bonding process as depicted in Fig. 11.17. The silicon wafer is biased as the electrical anode and the glass substrate as the cathode with an applied voltage as high as 1200 V. The sandwich is heated to 300–400°C, and mobile sodium ions in the glass migrate to the cathode leaving fixed charge and a high electric field at the glass–silicon interface supported by image charges in the silicon. The combination of elevated temperature and high field causes chemical bonds to form between the silicon wafer and glass substrate. Glasses with thermal expansion coefficients similar to silicon are utilized to minimize thermal stresses. Corning 7740 (Pyrex™) & 7707, Schott 8329 and 8330, and Iwaki 7570 are widely used for MEMS applications.

Anodic bonding can also be used to join silicon wafers. A thin (a few μm) intermediate layer of glass is sputtered on one of the silicon wafers. The second wafer is brought into contact with the glass layer, and the two wafers are anodically bonded. Since the glass layer is quite thin, much lower ($< 75\text{V}$) voltages are required for the anodic bonding process. Spin-on and evaporated-glass layers have also been used successfully in anodic bonding.

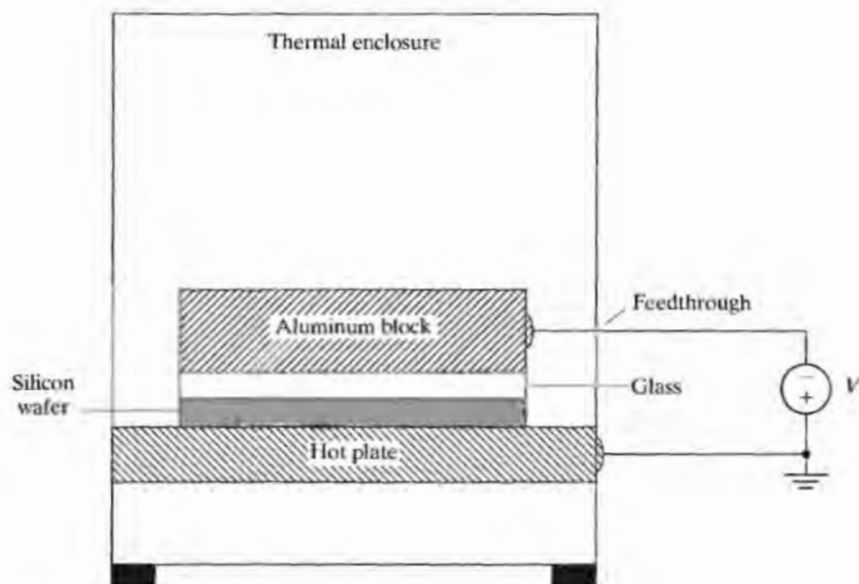


FIGURE 11.17

Electrostatic bonding of glass to a silicon wafer in a small oven

Both anodic and fusion bonding can be utilized to produce sealed and open cavity structures using combinations of silicon and glass substrates. Anodic bonding can be done at either atmospheric or reduced pressure, and the sealed cavity pressure will reflect the ambient pressure and temperature of the bonding conditions.

11.7 IC PROCESS COMPATIBILITY

To combine the power of VLSI signal processing with MEMS devices, MEMS fabrication steps must somehow be merged with high-density IC processes. However, a major issue in has always been process incompatibility between MEMS and standard CMOS or bipolar processes. Many of the original processes and chemicals specific to the formation of MEMS structures are not compatible with clean room technology. This has led to movement away from wet chemistry and into dry processing. Still, most attempts to merge MEMS with IC processes tend to separate the two process sequences. In pre-processing, for example, the MEMS structures are completely fabricated before any CMOS processing begins. After completion of the CMOS processing, the final release step is performed. In postprocessing, the CMOS process flow is completed first, and then the MEMS process steps are added to the end. Although considerable research is currently being conducted into monolithic integration of MEMS, hybrid fabrication, and packaging techniques still often represent the most economical solution to combining MEMS devices and standard integrated circuits.

11.7.1 Preprocessing

An example of the preprocessed approach is the Sandia CMOS technology [29] shown in Fig. 11.18. Cavities are etched in the wafer and passivated with a nitride insulator wherever MEMS devices are to be formed. Surface micromachining processes are utilized to form the MEMS devices in the cavities including the first polysilicon interconnect layer that is used to make electrical connections between the mechanical structures and polysilicon contact studs. The studs subsequently provide the interface path to the CMOS circuitry. Following fabrication, the MEMS cavities are completely filled with a sacrificial oxide and planarized using CMP. Next, a standard CMOS process flow forms the CMOS circuitry connected to the polysilicon studs using the first CMOS metal level. The completed CMOS circuitry is then protected by a nitride layer, and the sacrificial oxide is removed from the MEMS cavities to complete the fully released MEMS device. In this preprocess sequence, the CMOS circuitry does not see any detrimental additions to its thermal budget from high-temperature processing steps related to the MEMS fabrication.

11.7.2 Postprocessing

Figure 11.19 depicts a process developed at the University of California at Berkeley [30] in which the mechanical structures are formed following completion of the CMOS circuitry. A standard *p*-well CMOS process forms the CMOS circuitry, except for the use of tungsten metallization, which is a refractory metal that can withstand high-temperature processing. A nitride passivation layer protects the CMOS devices during MEMS machining.

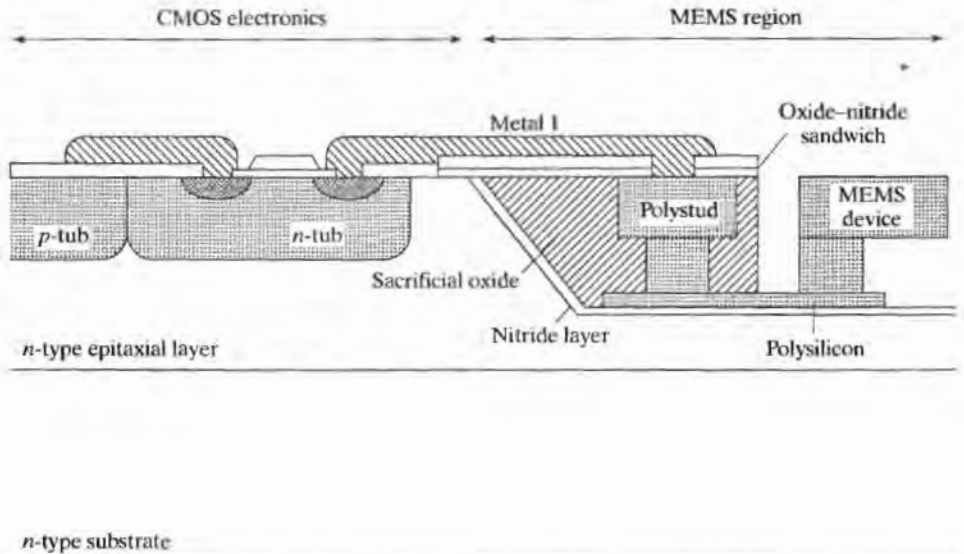


FIGURE 11.18

A cross section of CMOS electronics combined with preprocessed MEMS devices. The MEMS elements are fabricated in a cavity that is refilled with a sacrificial oxide and planarized via CMP prior to CMOS processing. After J. H. Smith et al., Ref. [29].

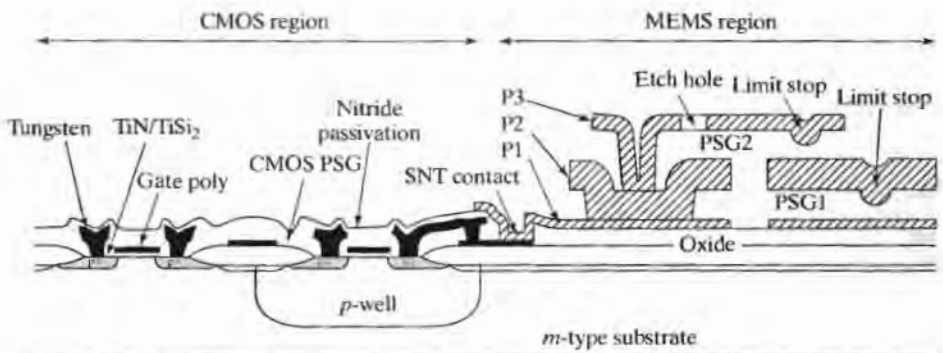


FIGURE 11.19

CMOS electronics combined with postprocessed MEMS devices. After J. M. Bustillo et al., Ref. [30].

The MEMS devices are fabricated over the nitride, and oxide-insulated region. Interconnection between the electromechanical devices and the active CMOS circuitry is accomplished with the polysilicon gate layer of the integrated circuit and the first polysilicon layer of the MEMS region.

A recent approach to postprocessing [31, 32] utilizes deposition of sacrificial layers of polycrystalline Germanium layers that can be etched in H_2O_2 . Structural layers are formed from poly-SiGe. These depositions can be done below 450°C and are thus compatible with aluminum layers that are present in CMOS ICs.

11.7.3 Merged Processes

Various approaches have been developed in which the MEMS processes are fully merged with standard CMOS and bipolar fabrication steps. Three possibilities are outlined next. In the first case, (see Fig. 11.20), dry-processing steps are added to a CMOS process to release MEMS structures [33]. The third level of metal in the CMOS process is used as an etch mask in the definition of the MEMS devices. An anisotropic RIE step etches vertical walled trenches through the oxide layers and is followed by a similar anisotropic RIE of the silicon. Release is achieved by isotropic etching of the material between the trenches.

Figure 11.21 shows a merged process flow developed at the University of Michigan [34] that forms diaphragms that can be utilized in a variety of applications. The process begins with the definition and implantation of the boron dose for the CMOS p -well. A barrier oxide is deposited, and windows are opened and followed by a deep ($15\text{-}\mu\text{m}$) p + etch stop diffusion. At the same time, the p -well is driven in to a depth of approximately $5\text{ }\mu\text{m}$. Next, a $2\text{-}\mu\text{m}$ -deep p -type diffusion is used to define the pressure diaphragm thickness. Recessed oxidation is used for CMOS isolation, followed by deposition and patterning of oxide and nitride dielectric layers on the diaphragms. A standard CMOS process flow continues from this point, and polysilicon and metal layers can also be patterned as needed in the diaphragm regions. At the final process step, the diaphragm is etched from the backside using EDP. Gold on chromium metallization is used to prevent interconnect corrosion by the EDP etch.

Analog Devices has developed a merged bipolar-MEMS process (BiMEMS) [35]. In this case, a shallow n + layer is used as the under pass from the metal layer in the

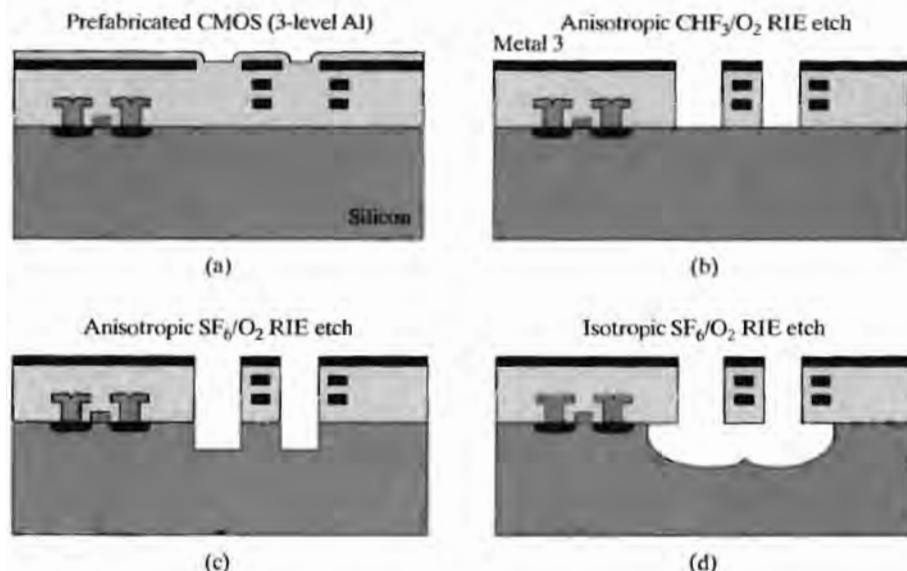


FIGURE 11.20

Example of the use of variable anisotropy dry etching to form MEMS devices within a CMOS IC technology. (a) CMOS with the third level of metal patterned as a mask for subsequent silicon etching steps (b) Structure following RIE of dielectric layers (c) Anisotropic etch of the silicon with SF_6 plasma (d) Beam release completed with isotropic etch in SF_6 plasma. Copyright 1998, IEEE. Reprinted with permission from Ref. [10]. After G. K. Fedder et al., Ref. [33].

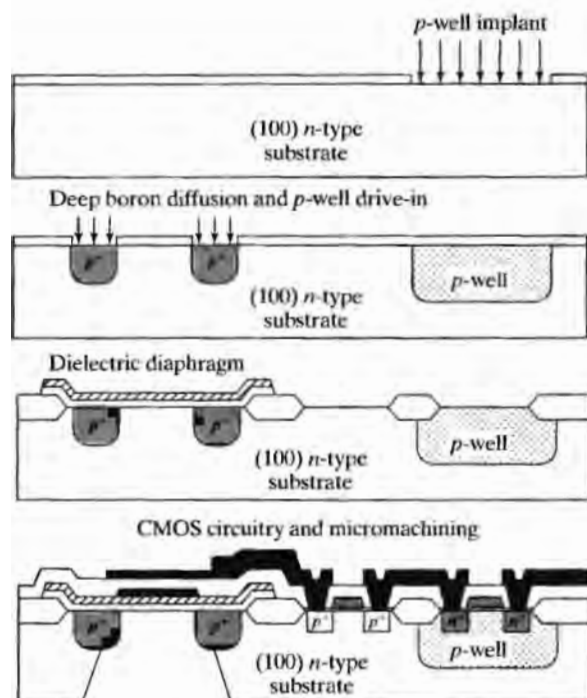


FIGURE 11.21

Process flow which merges standard p -well CMOS with p^+ etch stops and micromachined diaphragms. (Vertical dimensions not to scale.) Copyright 1998, IEEE. Reprinted with permission from Ref. [34].

bipolar circuitry to the polysilicon layer in the MEMS regions. The BiMEMS process is one of the few merged processes actually used in fabricating commercial products.

SUMMARY

MEMS technology development has brought together innovations from many areas of microelectronics to rapidly become a discipline of its own, and the MEMS area represents one of today's most exciting areas of microelectronics activity. Today's micromachined systems combine the signal processing and computational capability of analog and digital integrated circuits with a wide variety of nonelectrical elements, including pressure, temperature and chemical sensors, mechanical gears and actuators, 3D mirror structures, etc., and we have only begun to scratch the surface of biomedical applications.

MEMS structures are based upon our ability to sculpt or machine silicon on a microelectronic scale. Bulk micromachining dates back to the 1960s, when techniques for wet anisotropic etching of various forms of trenches, grooves, and membranes in silicon wafers were first developed. Advances in this area continued rapidly through the 1970s with the development of impurity-dependent etch stops, wafer-dissolution processes, and wafer fusion bonding. More recently, surface micromachining, which makes use of the full lithography capability of IC processing, came to the forefront, and many new beam, comb, microactuator, and rotary structures were conceived. Processes capable of producing high-aspect-ratio structures have been developed to use thick polymer molds and electroplating, or deep RIE.

A variety of methods for merging MEMS fabrication with standard CMOS and bipolar processes have been developed. Pre- and postprocessing techniques both attempt to maintain a separation between MEMS fabrication steps and standard CMOS process flows. Fully merged processes have also been developed. Present fabrication is moving from wet chemistry to plasma-based dry etching, because of environmental and clean room compatibility issues.

REFERENCES

- [1] Marc Madou, *Fundamentals of Microfabrication*, CRC Press, 1997.
- [2] G. T. A. Kovacs, *Micromachined Transducers Sourcebook*, The McGraw-Hill Companies, Boston: 1998.
- [3] R. S. Muller and R. T. Howe, eds., *Microsensors*, IEEE Press, New York: 1991.
- [4] W. S. Trimmer, Ed., *Micromechanics and MEMS: Classic and Seminal Papers to 1990*, IEEE Press, New York: 1997.
- [5] M. Elwenspoek and H. Jansen, *Silicon Micromachining*, Cambridge University Press, Cambridge, UK: 1998.
- [6] R. F. Wolffenbuttel, Ed., *Silicon Sensors and Circuits: On-Chip Compatibility*, Chapman and Hall, London: 1996.
- [7] J. B. Angell, S. C. Terry and P. W. Barth, "Silicon Micromechanical Devices," *Scientific American*, pp 44-55, April 1983.
- [8] K. Najafi, "Micromachined Micro Systems: Miniaturization Beyond Microelectronics," *Digest of the Symposium on VLSI Circuits*, pp. 6-13, June 2000.
- [9] K. E. Peterson, "Silicon as a Mechanical Material," *Proceedings of the IEEE*, no. 5, pp. 420-457, May 1982.
- [10] G. T. Kovacs, N. I. Maluf and K. E. Peterson, "Bulk Micromachining of Silicon," *Proceedings of the IEEE*, vol. 86, no. 8, pp. 1536-1551, August 1998.
- [11] K. E. Bean, "Anisotropic Etching of Silicon," *IEEE Trans. Electron Devices*, vol. ED-25, no. 10, pp. 1185-1193, October 1978.
- [12] E. Bassous, "Fabrication of Novel Three-dimensional Microstructures by the Anisotropic Etching of (100) and (110) Silicon," *IEEE Trans. Electron Devices*, vol. ED-25, pp. 1178-1185, October 1978.
- [13] D. B. Tuckerman and R. F. W. Pease, "High Performance Heat Sinking for VLSI," *IEEE Electron Device Letters*, vol. EDL-2, no. 5, pp. 126-129, May 1981.
- [14] S. T. Cho, K. Najafi and K. D. Wise, "Scaling and Dielectric Stress Compensation of Ultrasensitive Boron-Doped Silicon Microstructures," *Proceedings of IEEE Micro Electro Mechanical Systems*, pp. 50-55, February 1990.
- [15] P. B. Chu et al., "Controlled Pulse-etching with Xenon Difluoride," *Proceedings of Transducers '97*, Chicago, IL, June 1997. (<http://robotics.eecs.berkeley.edu/~pister/publications-/1997/ChuXeF2T97.pdf>)
- [16] A. A. Ayon, R. Braff, C. C. Lin, H. H. Sawin, and M. A. Schmidt, "Characterization of a Time Multiplexed Inductively Coupled Plasma Etcher," *Journal of the Electrochemical Society*, vol. 146, no. 1, pp. 339-349, January 1999.
- [17] D. Moser, M. Parameswaran and H. Baltes, "Field Oxide Microbridges, Cantilever Beams, Coils and Suspended Membranes in SACMOS Technology," *Transducers 89*, vol. 2, pp. 1019-1022, June 1990.

- [18] R. T. Howe and R. S. Muller, "Polycrystalline-Silicon Micromechanical Beams," *Journal of the Electrochemical Society*, vol. 130, pp. 1420-1423, June 1983.
- [19] J. M. Bustillo, R. T. Howe and R. S. Muller, "Surface Micromachining for Microelectromechanical Systems," *Proceedings of the IEEE*, vol. 86, no. 8, pp. 1552-1574, August 1998.
- [20] M. Mehregany, K. J. Gabriel and W. S. N. Trimmer, "Integrated Fabrication of Polysilicon Mechanisms," *IEEE Trans. Electron Devices*, vol. 35, no. 6, pp. 719-723, June 1988.
- [21] C. T. C. Nguyen and R. T. Howe, "CMOS Micromechanical Resonator Oscillator," *IEEE IEDM Digest*, pp. 199-202, December 1993.
- [22] K. S. J. Pister, M. W. Judy, S. R. Burgett and R. S. Fearing, "Microfabricated Hinges," *Sensors and Actuators A*, vol. 33, pp. 249-256, 1992.
- [23] <http://www.mdl.sandia.gov/micromachine/images.html>
- [24] N. Takeshima et al., "Electrostatic Parallelogram Actuators," *Proceedings of Transducers '91, the 1991 International Conference on Solid-State Sensors and Actuators*, IEEE Press, pp. 63-66, June 1991.
- [25] H. Guckel, "High-aspect-ratio Micromachining via Deep X-ray Lithography," *Proceedings of the IEEE*, vol. 86, no. 8, pp. 1586-1593, August 1998.
- [26] M. A. Schmidt, "Wafer-to-wafer bonding for microstructure formation," *Proceedings of the IEEE*, vol. 86, no. 8, pp. 1575-1585, August 1998.
- [27] A. Goyal, S. H. Bhavnani, R. C. Jaeger, C. D. Ellis, J. S. Goodling and M. Azimi-Rashti, "Formation of Silicon Re-entrant Cavity Heat Sinks Using Anisotropic Etching and Direct Wafer Bonding," *IEEE Electron Device Letters*, vol. EDL-14, no. 1, pp. 29-32, January 1993.
- [28] R. D. Nasby et al., "Application of Chemical-Mechanical Polishing to Planarization of Surface-Micromachined Devices," *Proceedings of the 1996 Solid-State Sensors and Actuator Workshop*, pp. 48-53, June 1996.
- [29] H. Smith, et al., "Embedded Micromechanical Devices for the Monolithic Integration of MEMS with CMOS," *IEEE IEDM*, pp. 609-612, December 1995.
- [30] J. M. Bustillo, G. K. Fedder, C. T. C. Nguyen and R. T. Howe, "Process Technology for the Modular Integration of CMOS and Polysilicon Microstructures," *Microsystems Technology*, vol. 1, pp. 30-41, 1994.
- [31] J. M. Heck et al., "High Aspect Ratio Poly-Silicon-Germanium Microstructures," *Proceedings of Transducers '99, Sendai, Japan*, June 7-10, 1999.
- [32] A. E. Franke et al., "Optimization of Poly-Silicon-Germanium as a Microstructural Material," *Proceedings of Transducers '99, Sendai, Japan*, June 7-10, 1999.
- [33] G. K. Fedder et al., "Laminated High-aspect-ratio Microstructures in a Conventional CMOS Process," *Proceedings of the IEEE International Workshop on Micro Electro Mechanical Systems*, pp. 13-18, February 1996.
- [34] E. Yoon and K. D. Wise, "A Multi-element Monolithic Mass Flowmeter with On-chip CMOS Readout Electronics," *Digest of the IEEE Solid-State Sensor and Actuator Workshop*, pp. 161-164, June 1990.
- [35] T. A. Core, W. K. Tsang and S. J. Sherman, "Fabrication Technology for an Integrated Surface-micromachined Sensor," *Solid State Technology*, vol. 36, no. 10, pp. 39-40, 42, 44, 46-47, October 1993.
- [36] T. J. Rodgers et al., "VMOS Memory Technology," *IEEE ISSCC Digest*, pp. 74-75, February 1977.
- [37] K. Hoffman and R. Losehand, "VMOS Technology Applied to Dynamic RAMs," *IEEE Journal of Solid-State Circuits*, vol. SC-13, pp. 617-622, October 1978.

PROBLEMS

- 11.1** Demonstrate that a $\langle 111 \rangle$ plane makes an angle of 54.74° with a $\langle 100 \rangle$ plane.
- 11.2** What is the angle of the planes in the bottom of the groove on the $\langle 110 \rangle$ surface in Fig. 11.3(b)?
- 11.3** A 10-micron-deep v-groove is to be formed in the surface of $\langle 100 \rangle$ silicon. What is the width of the mask opening?
- 11.4** Draw a possible cross section of the hinge structure in Fig. 11.12(b).
- 11.5** VMOS was a technology [36, 37] that made use of grooves to reduce the channel length and increase the W/L ratio of the MOS transistor. The channel is formed on the four sides of a pyramid, and its length is determined by the thickness of the epitaxial layer and the diffusions. At the time VMOS was invented, channel lengths available with this technology were much shorter than those that could be achieved with normal planar lithography. Figure P11.5 shows a vertical MOS transistor (VMOS) fabricated using anisotropic etching of $\langle 100 \rangle$ silicon. If the p -epitaxial layer is $4\text{ }\mu\text{m}$ thick and the n^+ buried layer and diffusions consume a total of $2\text{ }\mu\text{m}$ in depth, what is the length of the channel of the transistor?

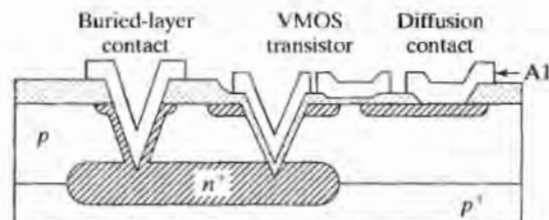


FIGURE P 11.5

- 11.6** Figure P11.6 shows a bipolar transistor fabricated on $\langle 100 \rangle$ silicon using v-groove isolation. What is the minimum isolation groove width at the surface if the epitaxial layer is $5\text{ }\mu\text{m}$ thick and a $1\text{-}\mu\text{m}$ minimum isolation width is required at the bottom of the groove? Does this seem competitive from an area point of view with other isolation processes? Which ones?

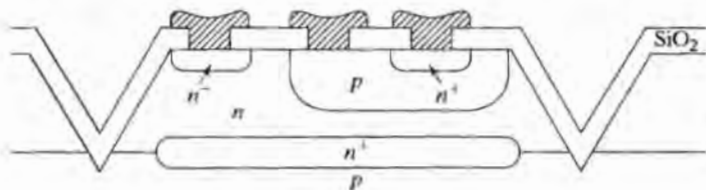


FIGURE P 11.6

- 11.7** Draw a top view of the mask set needed to produce a square sealed cavity similar to that in Fig. 11.9. Assume that the cavity is open along one edge.
- 11.8** A cavity is sealed by anodic bonding at atmospheric pressure at a temperature of 400°C . What is the cavity pressure at 27°C ? (The cavity dimensions at the surface are $250\text{ }\mu\text{m}$ on a side and the cavity is $25\text{ }\mu\text{m}$ deep.)
- 11.9** A sealed cavity is to have a pressure of 1 psi at room temperature (27°C). If the seal is formed by anodic bonding at a temperature of 350°C , what must be the ambient pressure during bonding? Assume that the cavity dimensions at the surface are $150\text{ }\mu\text{m}$ on a side and the cavity is $15\text{ }\mu\text{m}$ deep.

- 11.10** A 20×20 array of cavities with $10 \times 10 \mu\text{m}$ square openings is to be etched through a $500\text{-}\mu\text{m}$ -thick $\langle 100 \rangle$ silicon wafer using anisotropic wet etching. What is the closest center-to-center spacing of the $10\text{-}\mu\text{m}$ cavity openings? What is the area of the array?
- 11.11** An estimate for the lateral resonant frequency of the structure in Fig. 11.11 is given by

$$f_o = \frac{1}{2\pi} \sqrt{\frac{4E_y t W^3}{ML^3}},$$

assuming that there is no residual stress in the polysilicon film and where E_y is Young's modulus, W , L , and t are the width, length, and thickness, respectively, of the vertical supports connected to the anchor points, and M is the mass of the suspended shuttle. Suppose the width $W = 2.5 \mu\text{m}$ and $t = 2 \mu\text{m}$. Estimate f_o for the layout in Fig. 11.11(a). Assume that the density of polysilicon is similar to that of silicon.

- 11.12** A $\langle 100 \rangle$ silicon wafer is anisotropically etched to expose the $\langle 111 \rangle$ planes through a large square opening. As mentioned in the text, infinite selectivity is required to achieve the ideal slope of 54.74° . (a) Calculate the actual sidewall slope if the etchant has a selectivity of 400:1. (b) Repeat for selectivity of 20:1.

Answers to Selected Problems

Chapter 1

- 1.3 177, 148
 1.5 1.5 years, 5.1 years
 1.7 2.0 years, 6.7 years
 1.9 150 MW, 680 kA
 1.11 277 dice
 1.13 10^6 , 4×10^8 , 2.5×10^9

Chapter 2

- 2.1 95.3%, 98.6%
 2.5 0.536, 0.403 μm
 2.7 96.5 nm \equiv 0.1 μm

Chapter 3

- 3.1 9 min., 2.3 hrs.
 3.5 9.9 min, Eq. 3.8 is not valid for dry oxidation with $X_i < 25$ nm!
 3.7 2.7 hrs.
 3.9 58 nm, 1.54 μm
 3.11 0.430 μm
 3.13 3.55 hrs.
 3.15 0.15 μm
 3.17 carnation pink, carnation pink

Chapter 4

- 4.1 5.8 μm , 5.3 μm , 47 Ω/\square , 60 Ω/\square
 4.3 5.14 hr, 83 Ω/\square , $6.7 \times 10^{14}/\text{cm}^2$
 4.5 1150 for 10.7 hrs, $5 \times 10^{16}/\text{cm}^3$, $1.66 \times 10^{13}/\text{cm}^2$,
 4.7 2.18 μm
 4.9 5.8 \square , 4.7 \square , 390 Ω
 4.13 Drive in for 9 hrs at 1125°C, Dose is too low for preposition by solid-solubility limited diffusion at even 900°C.
 4.15 14.5 Ω/\square , 34.9 Ω/\square , 203 Ω/\square , 3770 Ω/\square , 9.72 Ω/\square , 12.5 Ω/\square

- 4.19 5 min
 4.23 100 ppm, 30 minutes exposure is life threatening – evacuate quickly!
 Arsine is a factor of 10 more toxic and is immediately life threatening – evacuate immediately!

Chapter 5

- 5.1 $3.33 \times 10^{18}/\text{cm}^3$, $1.1 \times 10^{13}/\text{cm}^2$, 0.40 μm
 5.3 900 keV
 5.5 $5.06 \times 10^{14}/\text{cm}^2$, 10.1 hrs.
 5.7 $4.50 \times 10^{14}/\text{cm}^2$, 66.7 sec.
 5.9 245 Ω/\square , $2.76 \times 10^{14}/\text{cm}^2$, 470 keV, 0.53 μm , 1.47 μm
 5.11 2.98×10^5 m/sec, 1.98×10^5 m/sec, 2.69×10^4 m/sec
 5.13 90 nm, 0.82 μm , $8.77 \times 10^{13}/\text{cm}^2$, $1.13 \times 10^{14}/\text{cm}^2$,
 5.15 $4 \times 10^{11}/\text{cm}^2$,
 5.19 $5.60 \times 10^{-12} \text{ cm}^2$, $5.48 \times 10^{-12} \text{ cm}^2$
 5.21 $1.82 \times 10^{-12} \text{ cm}^2$, $1.79 \times 10^{-12} \text{ cm}^2$

Chapter 6

- 6.1 120 μsec
 6.3 2.68×10^{14} molecules/sec, 7.5×10^{-7} Torr
 6.5 4.16×10^{-12} Pa
 6.7 8.0 nm
 6.9 71.1 cm
 6.11 0.14 $\mu\text{m}/\text{min}$, 0.02 $\mu\text{m}/\text{min}$, 1245°C, 0.35 eV, 1.52 eV
 6.15 1.13, 1.56, 2.44
 6.17 $10^{20} \operatorname{erfc}\left(x/2\sqrt{D_s t}\right)$,
 $D_s = 6.45 \times 10^{-15} \text{ cm}^2/\text{sec}$

Chapter 7

- 7.1 32 mΩ/□, 1.6 Ω, 175 fF, 0.28 psec
 7.3 7.8 Ω/□, 5.0 Ω/□, 40 Ω/□,
 7.5 0.5 μΩ-cm², 50 Ω
 7.7 40 time units, 430 units
 7.9 20 mA
 7.11 0.80 Ω

Chapter 8

- 8.3 392, 492, 3626
 8.5 Depends upon die placement:
 best case 2/26 or 7.7%, worst case 0,
 approximately 3.1
 8.7 \$7.65 vs. \$9.79, \$13.70 vs. \$16.40
 8.9 $\exp(-D_0A)$
 8.11 17.7, 31.4, 70.7
 8.13 17%, 0
 8.15 44, 22, 7, 2, 0
 8.17 0.0737/cm², 0.0412/cm²

Chapter 9

- 9.1 5 V
 9.3 14 V
 9.5 $2 \times 10^{16}/\text{cm}^3$

- 9.9 increases by α , increases by α ,
 increases by α^3 !
 9.13 4.4 μm
 9.19 $448 \lambda^2$, $40 \lambda^2$
 9.21 $2.5 \times 10^{18}/\text{cm}^3$, 175 mA, 699 kW!

Chapter 10

- 10.1 $4 \times 10^{16} \text{ sec/cm}^4$, $2 \times 10^{13} \text{ sec/cm}^4$, 145
 10.3 Approximately 70 V
 10.5 3.64 μm
 10.7 8 V, 50 V, 50 V
 10.9 0.5, 1.0, $2d(W+L)/WL$ where d is the
 diffusion depth, a square with $L = W$
 and $\beta = 4d/L$
 10.11 18 V vs. 60 V
 10.13 24.9 kΩ, 1.25 kΩ, 104 Ω
 10.19 $\approx 40 \text{ nm}$
 10.27 6–7 V, 20 V

Chapter 11

- 11.3 14.2 μm
 11.5 2.45 μm
 11.9 2.1 psi

Index

A

- Abrupt pn junction breakdown voltage
 - vs. impurity concentration, 205
- Additive metal liftoff, 165
- Additive process, 164–166
- Adhesion improvement, 22
- Adhesive bonding
 - silicon wafer bonding, 289
- Advanced bipolar structures, 253–258
 - dual polysilicon self-aligned process, 254–257
 - locos isolated self-aligned contact structure, 254
 - silicon germanium epitaxial base transistor, 257–258
- Aggressive design rules
 - MOS transistor layout and design rules, 218–219
- Alignment marks
 - example, 23
- Alignment tolerance, 212
- Alignment variation
 - design rules, 218
- Aluminum
 - bulk resistivity, 153
 - mechanical properties, 270
- Aluminum contacts to silicon
 - types, 154
- Aluminum interconnection failure
 - scanning electron micrographs, 158
- Aluminum-silicon alloying
 - aluminum spiking, 156
- Aluminum-silicon eutectic behavior
 - metal interconnections and contact technology, 154–155
- Aluminum-silicon system
 - contact resistivity, 157
 - phase diagram, 155
- Aluminum spiking
 - aluminum-silicon alloying, 156
 - and junction penetration
 - metal interconnections and contact technology, 155–156
- Amorphous layer
 - dose requirement, 120
- Analytical microscopes, 37
- Angle-lap method
 - junction-depth measurement, 90–91
- Anisotropic backside etching
 - diaphragms, 274
- Anisotropic etching, 272, 273
 - bulk micromachining, 271–273
 - silicon etchants, 278–279
- Anisotropy dry etching
 - MEMS device formation, 294

- Annealing
 - ion implantation, 118–121
 - rapid thermal
 - shallow ion implantation, 123
- Anodic bonding
 - silicon wafer bonding, 291–292
- Antimony
 - diffusion system, 100
 - during oxidation, 51
- Area array of pads
 - die layout, 180
- Arsenic
 - during oxidation, 51
- Arsenic diffusion
 - properties, 81
 - system, 99–100
- Arsine
 - threshold limit recommendations, 99
- Asymmetrical reactive ion etching system, 27
- Atmospheric-pressure reactor
 - CVD, 137–138
- Atomic diffusion
 - two-dimensional lattice, 68

B

- Background concentrations
 - vs. peak concentrations, 116
- Ball grid array (BGA), 190
 - cross-section, 190
- Barrel VPE reactor, 143, 144
- Barrier layer formation
 - photolithographic process, 21
- Barrier metals, 164
- Base diffusion to isolation diffusion spacing
 - layout considerations, 251
- Base resistor
 - elements in SBC technology, 244–245
- Basewidth
 - bipolar process integration, 237–239
- Basic multilevel metallization, 166–167
- BGA, 190
- BHF, 25
- BiCMOS
 - bipolar process integration, 262–263
- Bipolar isolation techniques, 259–262
 - CDI, 259
 - dielectric isolation, 259–262
- Bipolar junction transistor
 - vertical impurity profile, 234
- Bipolar process, 10–11
 - cross-sectional views, 12

- flowchart, 13
- integration, 233–264
- Bipolar process, *continued*
 - advanced bipolar structures, 253–258
 - basewidth, 237–239
 - BICMOS, 262–263
 - bipolar isolation techniques, 259–262
 - breakdown voltages, 239–242
 - current gain, 235
 - elements in SBC technology, 243–249
 - exercise problems related to, 265–268
 - junction-isolated structure, 233–235
 - layout considerations, 249–253
 - transit time, 236–237
- Bipolar transistors
 - isolation region, 251
- Body effect, 205
- BOE, 25
- Boron
 - during oxidation, 51
 - p*-type dopant, 97
- Boron diffusion
 - deep, 276
 - properties, 81
 - shallow, 276
 - system, 97–98
- Boron impurity distributions
 - four-moment distribution functions, 121
- Breakdown voltages, 239–242
 - circular emitters, 240
 - collector-base breakdown voltage, 240–242
 - emitter-base breakdown voltage, 239–240
- Bridges
 - surface micromachining, 279
- Buffered oxide etch (BOE or BHF), 25
- Bulk micromachining, 271–276
 - cantilever beams, 275–276
 - diaphragm formation, 273–274
 - isotropic and anisotropic etching, 271–273
 - released structures, 275–276
- Buried contacts, 152, 160
 - polysilicon interconnections, 160
 - structure, 161
- Buried layer, 10
 - epitaxy, 145–147
 - and isolation diffusions
 - layout considerations, 249–251
- Butted contacts, 152
 - polysilicon interconnections, 162
 - structure, 161
- C**
- Cantilever beams
 - bulk micromachining, 275–276
 - surface micromachining, 279
- Cavity anisotropic etching, 275
- Cavity formation
 - sealed, 282
- CDI
 - bipolar isolation techniques, 259
- Ceramic leadless chip carriers, 186
- Channeling
 - ion implantation, 118–121
- Channel length
 - MOS transistor layout and design rules, 219–221
- Chemical mechanical polishing (CMP), 56–57
 - multilevel metallization fabrication, 59
 - process, 167
 - Damascene plating, 168
 - technique, 58
 - thermal oxidation of silicon, 56
- Chemical vapor deposition (CVD), 5, 136–141
 - metal, 141
 - polysilicon, 138–139
 - reactors, 137–138
 - silicon dioxide, 139–140
 - silicon nitride, 140–141
- Chip-on-board packaging, 193
- Chip Scale Packages (CSPs), 184, 193
- Circular emitters
 - breakdown voltages, 240
 - cross connected quad, 240
- Circular TO-style
 - packages, 184
- Class ratings
 - filtration effectiveness, 18
- Clean rooms
 - filtration effectiveness
 - class rating, 18
- CMOS. *See* Complementary MOS (CMOS)
- CMP. *See* Chemical mechanical polishing (CMP)
- Cobalt
 - bulk resistivity, 153
- Collector-base breakdown voltage
 - breakdown voltages, 240–242
- Collector-based junction breakdown voltage, 241
- Collector-based space-charge region growth, 242
- Collector-diffusion isolation (CDI)
 - bipolar isolation techniques, 259
- Color chart
 - comparisons, 59
 - thermally grown silicon dioxide, 60
- Complementary error function, 69–71
 - barrier layer, 80
 - vs. Gaussian distribution, 71
- Complementary MOS (CMOS)
 - electronics
 - cross-section, 293
 - process, 9–10
 - illustrated views, 11
 - structure
 - cross-section, 225

- twin-well processes, 223
- VLSI, 141
- Complementary MOS (CMOS) technology, 221–226
 - gate doping, 222–223
 - isolation, 224–225
 - latchup, 225
 - n*-well process, 221
 - minimum spacing requirements, 224
 - p*-well process, 221–222
 - shallow trench isolation, 225–226
 - twin well process, 221–222
- Concentration-dependent diffusion
 - diffusion profiles, 81
 - junction formation, 79–81
- Constant electric field scaling results, 213
- Constant-source diffusions
 - shallow phosphorus diffusion profiles, 82
- Contact, 151–172
- Contact mask
 - nominal alignment, 213
- Contact printing, 28
 - illustrated, 31
- Contact resistance
 - metal interconnections and contact technology, 156
- Contamination levels
 - filament evaporation, 132
- Copper
 - bulk resistivity, 153
 - IC processing, 151
- Copper Damascene process, 168–171
 - Damascene plating, 168
 - dual Damascene structures, 169–171
 - electroplated copper interconnect, 168
 - steps, 169
- Copper interconnects process, 168–171
- Copper metallization
 - electromigration performance improvement, 158
- CSPs, 184
- Current gain, 235
 - bipolar process integration, 235
- CVD. *See* Chemical vapor deposition (CVD)
- D**
- Damascene plating
 - copper interconnects and Damascene process, 168
- Dc sputtering system, 136
 - illustration, 135
- Dc tests, 177
- Deep boron diffusion, 276
- Deep reactive-ion etching (DRIE), 278
 - application, 280
- Deep submicron MOS devices, 121
- Deep ultraviolet (DUV) region, 34
- Defect densities
 - yield curves, 197
- Defect probability density functions, 196
- Deionized (DI) water
 - wafer cleaning, 20
- Depletion-layer width
 - one-sided step junction, 206
- Depletion-mode NMOS transistor
 - formation, 208
- Deposition rate
 - monitoring, 134
- Depth of focus, 33
- Diamond
 - mechanical properties, 270
- Diamond saws
 - die separation, 178
- Diaphragm formation
 - anisotropic backside etching, 274
 - bulk micromachining, 273–274
- Diborane
 - threshold limit recommendations, 99
- Diborane oxidation, 98
- Dichlorosilane, 140
 - threshold limit recommendations, 99
- Dicing saw
 - wafer, 179
- Die attachment, 178–179
 - epoxy, 178
 - eutectic, 179
- Dielectric isolation
 - bipolar isolation techniques, 259–262
 - process steps, 260
- Die
 - number per wafer, 3
- Die separation, 178
 - diamond saws, 178
- Diffused interconnections
 - interconnections and contacts, 158–159
- Diffused region, 10
- Diffusion, 5, 67–101
 - coefficient, 72–74
 - exercise problems related to, 103–108
 - gettering, 100–101
 - junction-depth and impurity profile
 - measurement, 90–92
 - junction formation and characterization, 76–81
 - mathematical model, 68–72
 - constant-source, 69
 - limited-source, 70
 - two-step, 71–72
 - process, 67–68
 - sheet resistance, 81–90
 - simulation, 93–94
 - solid-solubility limits, 74–75
 - successive, 74
 - systems, 95–100
 - antimony, 100
 - arsenic, 99–100
 - boron, 97–98
 - phosphorous, 98–99

- Diffusion coefficient values
 - for impurities, 74
 - Diffusion constants
 - vs. temperature is, 72–73
 - DIP, 184, 185
 - Dishing, 56
 - Dissolved wafer process, 276
 - D1 water
 - wafer cleaning, 20
 - Dopant redistribution
 - oxidation, 50–51
 - Doping layers
 - epitaxy, 145
 - Dose, 69
 - DRIE, 278, 280
 - Drive-in step, 71
 - Dry etching plasma systems
 - techniques, 26
 - Dry-plasma etching, 278
 - Dry silicon dioxide
 - growth, 49
 - Dual Damascene process flow, 170
 - Dual Damascene structures
 - copper interconnects and Damascene process, 169–171
 - Dual-in-line packages (DIP), 184, 185
 - Dual polysilicon self-aligned process
 - advanced bipolar structures, 254–257
 - DUV region, 34
 - Dynamic memory density
 - microprocessor
 - yearly change, 4
 - Dynamic memory feature size
 - yearly change, 6
- E**
- EDP, 278
 - Electrically active impurity-concentrations limits
 - in silicon, 75
 - Electromigration
 - metal interconnections and contact technology, 157
 - Electromigration performance improvement
 - copper metallization, 158
 - Electron-beam evaporation
 - evaporation, 132–134
 - Electron-beam source evaporation
 - illustration, 132
 - photograph, 134
 - Electron microscopy, 37–38
 - Electroplated copper interconnect
 - copper interconnects and Damascene process, 168
 - Electrostatic bonding, 276
 - glass to silicon wafer, 291
 - Ellipsometer, 59
 - Emitter-base breakdown voltage
 - breakdown voltages, 239–240
 - Emitter-base junction breakdown, 239
 - Emitter-diffusion design rules
 - layout considerations, 252
 - Emitter resistor
 - elements in SBC technology, 243–244
 - Enclosed lateral pnp transistor structure, 248
 - Epitaxial growth process
 - geometrical model, 146
 - model, 142
 - pattern shift, 147
 - Epitaxial layer resistor
 - elements in SBC technology, 245
 - Epitaxy, 10, 141–148
 - buried layers, 145–147
 - doping layers, 145
 - liquid-phase and molecular-beam epitaxy, 148
 - vapor-phase epitaxy, 142–145
 - Epoxy
 - die attachment, 178
 - Etch bias, 25
 - Etching pressure ranges, 27
 - Etching profiles
 - wet vs. dry, 25
 - Etching techniques, 25–28
 - dry plasma systems, 26
 - metrology and critical dimension control, 28
 - photoresist removal, 27
 - wet chemical, 25
 - Ethylene diamine pyrochatechol (EDP), 278
 - Eutectic
 - die attachment, 179
 - Eutectic temperature, 154
 - silicide and silicon, 163
 - Evaporation, 129–135
 - electron-beam evaporation, 132–134
 - filament evaporation, 132
 - flash evaporation, 134
 - kinetic gas theory, 130–131
 - shadowing and step coverage, 134–135
 - Evaporation sources
 - forms, 132
 - Exposure sources
 - lithography, 34
 - Exposure systems
 - lithography, 28–34
- F**
- Fairchild Semiconductor, 1
 - Fick's first law of diffusion, 68–69
 - Fick's second law of diffusion, 69
 - Field-region considerations
 - MOS device considerations, 208
 - Filament evaporation, 132
 - illustration, 132
 - Film deposition, 129–148
 - CVD, 136–141
 - epitaxy, 141–148

- evaporation, 129–135
- exercise problems related to, 149–150
- parameter
 - mean free path, 131
- sputtering, 135–136
- Film thickness measurement, 59, 61
- Flash evaporation
 - evaporation, 134
- Flip-chip technology, 188–189
- Four-moment distribution functions
 - boron impurity distributions, 121
- Four-point probe
 - correction factors
 - wafer correction, 89
 - sheet resistance, 88
- Furnaces
 - oxidation and diffusion, 53–54

G

- Gallium
 - during oxidation, 51
- Gas-phase autodoping, 146
- Gas-source
 - diffusion, 97
 - threshold limit recommendations, 99
- Gate doping
 - CMOS technology, 222–223
- Gate metal
 - spacing, 216
- Gate-oxide thickness
 - MOS device considerations, 202–203
- Gaussian density
 - triangular approximation, 196
- Gaussian diffusion
 - barrier layer, 80
- Gaussian distribution
 - vs. complementary error function, 71
 - Irvin's curves, 117–118
 - limited-source diffusion, 70
- Gaussian impurity profile
 - threshold-voltage shift, 207
- Gaussian theory
 - deviations
 - ion implantation, 121
- Germanium, 1
- Gettering agent, 98
- Gettering diffusion, 100–101
- Glass to silicon wafer
 - electrostatic bonding, 291
- Gold
 - bulk resistivity, 153
- Gold ball bonding
 - SEM micrograph, 181
- Gold bumps
 - process sequence
 - aluminum metallurgy devices, 191

- Gold-silicon eutectic point, 179
- Gold wire
 - thermosonic ball-wedge bonding, 182
- Groove-and-stain method
 - junction-depth measurement, 90–91
- Guard ring diffusions, 225
- Gull-wing surface-mount, 187
- Gummel numbers, 235

H

- Hard baking
 - photolithographic process, 25
- Hexamethyldisilazane (HMDS), 22
- High aspect ratio micromachining
 - MEMS process, 288–289
- High-density chip interconnection
 - ITRS plan, 188
- High-density IC process, 292–295
- High-density VLSI processes
 - junction penetration, 156
- High-performance oxide-isolated bipolar transistor
 - process sequence, 255
 - SIMOX wafer, 261
 - structure illustrated, 256
- High-pressure mercury, 34
 - lamp, 34
- High-resolution VLSI lithography systems, 29
- High-voltage accelerator, 110
- Hillocks, 157
- Hinges
 - formation, 285
 - SEM photograph, 287
 - three-dimensional, 285
- HMDS, 22
- Horizontal furnace
 - oxidation and diffusion, 54
- Horizontal VPE reactor, 143, 144
- Hot-wall LPCVD system, 137–138
- Hydrofluoric acid
 - wafer cleaning, 20
- Hydrofluoric acid, 13
- Hydrogen
 - diffusivities, 44

I

- IC. *See* Interconnection (IC) process
- Ideal ohmic contact
 - characteristics, 154
- Implantation
 - argon atoms, 100
- Implantation damage, 120
- Implanted impurity profile, 115
- Impulse
 - defect probability density function, 196

- Impurity concentration
 - vs. abrupt pn junction breakdown voltage, 205
 - vs. temperature, 75
 - Impurity depletion, 51
 - Impurity implantations
 - silicon
 - junction formation, 118
 - Impurity profile
 - measurement
 - junction-depth, 91–92
 - oxidation effects, 52
 - Initial functional testing
 - wafer
 - scribing, 178
 - Integrated circuit
 - cross-section, 202
 - Intel microprocessor packaging, 190
 - Interconnecting polysilicon and n+ diffusion techniques, 161
 - Interconnection formation
 - comparison, 165
 - Interconnection (IC) process
 - compatibility, 292–295
 - merged processes, 294–295
 - postprocessing, 292–293
 - preprocessing, 292
 - in integrated circuits, 151–152
 - Interconnections and contacts, 151–172
 - copper interconnects and Damascene process, 168–171
 - diffused interconnections, 158–159
 - exercise problems related to, 174
 - interconnections in integrated circuits, 151–152
 - liftoff process, 164–166
 - metal interconnections and contact technology, 153–157
 - multilevel metallization, 166–168
 - polysilicon interconnections, 159–162
 - silicides and multilayer-contact technology, 162–164
 - International technology road map for semiconductors (ITRS), 6–7
 - lithography projections, 36
 - plan
 - high-density chip interconnection, 188
 - Interstitial diffusion, 67
 - Interstitial diffusers, 73
 - Interstitial diffusion, 67
 - Interstitial space, 67
 - Intrinsic isolation strategy
 - illustration, 209
 - Ion concentrations
 - contours, 115
 - Ion energy
 - vs. sputtering heated
 - dc sputtering system, 136
 - Ion implantation, 5, 109–128
 - channeling, lattice damage, and annealing, 118–121
 - exercise problems related to, 126–128
 - junction depth, 117–118
 - mathematical model, 111–114
 - selective, 114–117
 - shallow, 121–124
 - low-energy, 122
 - rapid thermal annealing, 123
 - transient enhanced diffusion (TED), 124
 - sheet resistance, 117–118
 - technology, 109–111
 - Ion milling, 26
 - Ion source, 110
 - Irvin's curves
 - Gaussian distribution, 117–118
 - sheet resistance, 85–87
 - Isolation
 - CMOS technology, 224–225
 - Isolation diffusions
 - layout considerations, 249–251
 - Isotropic etching
 - agitation, 271
 - bulk micromachining, 271–273
 - silicon etchants, 277–278
 - Isotropic process, 25
 - ITRS. *See* International technology road map for semiconductors (ITRS)
- ## J
- J-lead surface-mount, 187
 - Junction breakdown
 - MOS device considerations, 204
 - Junction capacitance
 - MOS device considerations, 205–206
 - Junction depth
 - calculations, 76–77
 - ion implantation, 117–118
 - measurement, 90–95
 - angle-lap and stain method, 92
 - angle-lap method, 90–91
 - groove-and-stain method, 90–91
 - impurity-profile, 91–92
 - Junction formation
 - characterization, 76–81
 - concentration-dependent diffusion, 79–81
 - lateral diffusion, 78
 - vertical diffusion, 76–78
 - silicon impurity implantations, 118
 - Junction-isolated bipolar transistor structure, 8
 - Junction-isolated structure, 233–235
 - bipolar process integration, 233–235
 - Junction penetration
 - high-density VLSI processes, 156
 - metal interconnections and contact technology, 155–156

K

- Kinetic gas theory
 - evaporation, 130–131
- Kink effect, 79–80

L

- λ -based set of design rules, 215
- Latchup, 10, 141
 - CMOS technology, 225
- Lateral diffusion
 - junction formation, 78
- Lateral pnp transistors
 - elements in SBC technology, 248–249
 - enclosed structure, 248
- Lattice damage
 - ion implantation, 118–121
- Layout considerations, 249–253
 - base diffusion to isolation diffusion spacing, 251
 - buried-layer and isolation diffusions, 249–251
 - emitter-diffusion design rules, 252
 - layout example, 252–253
- LCCs, 184, 186
- Leadless chip carriers (LCCs), 184
 - packages, 186
- Lens system
 - wafer stepper concept, 32
- Liftoff, 151
 - process
 - interconnections and contacts, 164–166
- LIGA process, 288–289
- Lightly doped drain structures
 - MOS device considerations, 210
- Lindhard, Scharff, and Schiott (LSS) theory, 118
 - ion implantation, 112
- Liquid-phase epitaxy (LPE), 141
 - epitaxy, 148
- Liquid-source systems
 - diffusion, 97
- Lithography, 17–39
 - etching techniques, 25–28
 - exercise problems related to, 40–42
 - exposure sources, 34
 - exposure systems, 28–34
 - optical and electron microscopy, 37–38
 - photolithographic process, 17–25
 - photomask fabrication, 28
- Lithography, galvanofarming, and abforming (LIGA) process
 - illustrated, 288
 - MEMS, 288–289
- Lithography process
 - ultraclean conditions, 17
- Locos isolated self-aligned contact structure
 - advanced bipolar structures, 254
- Low dielectric constant interlevel dielectrics

- multilevel metallization, 167–168
- Low-energy arsenic implantations
 - shallow amorphous layers, 120–121
- Low-pressure CVD (LPCVD), 7, 138
- Low-pressure vacuum deposition
 - shadowing problem, 134–135
- Low-resistance ohmic contacts, 151
- Low-temperature processing
 - VLSI fabrication, 111
- LPCVD, 7, 138
- LPE, 141, 148
- LSS theory, 112, 118
- Lumped circuit model, 159

M

- Mask alignment
 - photolithographic process, 23
- Mask fabrication
 - illustrations, 30
 - process steps, 29
- Mass spectrometer, 110
- Mathematical model
 - diffusion, 68–72
 - ion implantation, 111–114
- MBE, 141, 148
- Mean free path
 - film-deposition parameter, 131
- Mean time to failure (MTF)
 - conductor, 157
- Mechanical properties of silicon
 - processes MEMS, 270
- Mechanical surface profiler, 59
- MEMS. *See* Microelectromechanical systems (MEMS)
- Mercury-rare gas discharge lamps, 34
- Metal deposition
 - CVD, 141
- Metal-gate mask sequence
 - listed, 213
- Metal-gate process
 - mask steps and device cross-sections, 214
- Metal-gate transistor layout
 - MOS transistor layout and design rules, 213–217
- Metal interconnections and contact
 - technology, 153–157
 - aluminum-silicon eutectic behavior, 154–155
 - aluminum spiking and junction penetration, 155–156
 - contact resistance, 156
 - electromigration, 157
 - ohmic contact formation, 153–154
- Metallization
 - following, 177
 - six-level
 - microphotographs, 171
- Metallurgical junction depth, 76

- Metal mask
 - nominal alignment, 213
 - Metal-oxide-semiconductor FET (MOSFET), 204, 210–211
 - Metal-oxide-semiconductor (MOS)
 - device considerations, 201–211
 - field-region considerations, 208
 - gate-oxide thickness, 202–203
 - junction breakdown, 204
 - junction capacitance, 205–206
 - lightly doped drain structures, 210
 - punch-through, 204–205
 - substrate doping, 203
 - threshold adjustment, 206–207
 - threshold voltage, 203
 - transistor isolation, 208–210
 - transistor scaling, 210–212
 - device formation, 164
 - logic circuit
 - interconnections, 152
 - process, 7–10
 - process integration, 201–228
 - CMOS technology, 221–226
 - exercise problems related to, 229–232
 - MOS device considerations, 201–212
 - MOS transistor layout and design rules, 212–221
 - silicon on insulator, 226–227
 - structure
 - showing polycide, 163
 - TEM images, 39
 - transistor layout and design rules, 212–221
 - aggressive design rules, 218–219
 - channel length, 219–221
 - metal-gate transistor layout, 213–217
 - polysilicon-gate transistor layout, 217–218
 - width biases, 219–221
 - Metal pileup, 157
 - Metals
 - bulk resistivity, 153
 - Metrology and critical dimension control
 - etching techniques, 28
 - Microelectromechanical systems (MEMS)
 - device
 - anisotropy dry etching, 294
 - cross-section, 293
 - VLSI signal processing, 292
 - processes, 269–296
 - bulk micromachining, 271–276
 - high aspect ratio micromachining, 288–289
 - IC process compatibility, 292–295
 - LIGA molding process, 288–289
 - mechanical properties of silicon, 270
 - problem, 298–299
 - silicon etchants, 277–279
 - silicon wafer bonding, 289–295
 - surface micromachining, 279–287
 - Microelectronic fabrication, 1–16
 - safety, 12–14
 - Minimum feature size, 212
 - Minimum melting temperature, 154
 - Minimum-size metal-gate transistor, 217
 - Minority-carrier diffusion lengths
 - calculations, 236
 - Misalignments
 - design rules, 218
 - Molecular-beam epitaxy (MBE), 141
 - epitaxy, 148
 - Molybdenum
 - bulk resistivity, 153
 - CVD metal deposition, 141
 - Monolithic fabrication processes
 - basic steps, 5
 - and structures, 5–7
 - MOS. *See* Metal-oxide-semiconductor (MOS)
 - MOSFET, 204, 210–211
 - Movable in-plane structures
 - surface micromachining, 279–281
 - MTF
 - conductor, 157
 - Multilayer contacts, 164
 - Multilevel aluminum metallization, 167
 - Multilevel metallization, 166–168
 - basic, 166–167
 - low dielectric constant interlevel dielectrics, 167–168
 - planarized, 167
 - Multilevel metallization fabrication
 - CMP, 59
- N**
- n^+ emitter diffusion
 - resistor, 243
 - n -channel metal-oxide-semiconductor (NMOS), 7–9
 - basic process flowchart, 10
 - NOR-gate layout, 160
 - structure, 7
 - transistors
 - two-adjacent, 202
 - n -channel polysilicon-gate transistors
 - threshold voltages, 204
 - n -type epitaxial layer
 - resistor, 245
 - n -type silicon
 - gettering, 100
 - room-temperature resistivity, 77
 - n -well CMOS technology
 - minimum spacing requirements, 224
 - n -well process
 - CMOS technology, 221
 - Negative photoresists, 24
 - Negative resists, 24
 - Nickel
 - bulk resistivity, 153

NMOS. *See* *n*-channel metal-oxide-semiconductor (NMOS)

Nonuniform defect densities
yield, 195–197

Normal aluminum link, 161

*n*pn transistor, 225

O

Ohmic contact formation
metal interconnections and contact
technology, 153–154

One-sided step junction
depletion-layer width, 206

Open-furnace-tube diffusion systems, 96

Optical focal plane, 33

Optical microscopy, 37

Out-diffusion, 146

Out-of-plane motion
surface micromachining, 282–286

Oxidation, 5
dopant redistribution, 50–51
shallow trench formation, 55–56
technology, 52–53

Oxidation modeling, 44–46

Oxidation process, 43

Oxidation rate
factors influencing, 46–50

Oxides
properties, 140
quality, 53–54
thickness characterization, 57–60

Oxygen
diffusivities, 44
wet and dry
coefficient *D* and activation energy values, 47

P

p⁺ etch stops
p-well CMOS
process flow, 295
p-channel polysilicon-gate transistors
threshold voltages, 204

p-glass reflow, 139

p-type base diffusion
resistor, 244

p-type dopant
boron, 97

p-type silicon
room-temperature resistivity, 77

p-well and twin well processes
CMOS technology, 221–222

p-well CMOS
p⁺ etch stops
process flow, 295

Packages, 184–186

circular TO-style, 184

DIPs, 184

LCCs, 186

PGAs, 185

surface mounting, 186

Packaging and yield, 177–198

die attachment, 178–179

die separation, 178

exercise problems related to, 199–200

flip-chip and tape-automated-bonding processes, 187–193

packages, 184–186

testing, 177–178

wafer thinning, 178

wire bonding, 179–184

yield, 194–197

Paladium

bulk resistivity, 153

Pancake susceptors, 144

Parallel plate plasma-enhanced CVD system, 137–138

Parallel plate plasma etcher, 27

Parametric test dice, 177

Parasitic bipolar devices, 225

Parasitic lateral pnpn SCR, 225

Parasitic NMOS transistor, 208

Passivation-layer processing
following, 177

Pattern shift
epitaxial growth process, 147

Pattern transfer, 35

PDP, 212

Peak concentrations
vs. background concentrations, 116

Pearson IV
boron impurity distributions, 121

PECVD system, 137–138

Peripheral pads

die layout, 180

PGA, 184–186

Phase-shifting mask technology, 34

Phosphine
threshold limit recommendations, 99

Phosphorus
during oxidation, 51

Phosphorus diffusion system, 98–99

Phosphorus impurity profiles, 119

Phosphorus MOS (PMOS) devices, 208

Photolithographic process, 17–25

barrier layer formation, 21

hard baking, 25

mask alignment, 23

photoresist application, 21–22

photoresist exposure and development, 23–24

soft baking or prebaking, 22

steps, 18

steps illustrated, 19

wafer and wafer cleaning, 19–21

Photomask fabrication

- lithography, 28
 - Photoresist, 21
 - description, 21
 - Photoresist application
 - photolithographic process, 21–22
 - Photoresist deposition system
 - photograph, 22
 - Photoresist development
 - photolithographic process, 23–24
 - Photoresist exposure
 - photolithographic process, 23–24
 - Photoresist removal
 - etching techniques, 27
 - Physical evaporation, 129–135
 - Pile up, 51
 - Pinched base resistor, 246
 - Pinch resistor
 - elements in SBC technology, 246
 - Pin-grid array (PGA), 184, 186
 - packages, 185
 - Planarized multilevel metallization, 167
 - Planetary substrate holder
 - geometry for evaporation, 133
 - Plasma-enhanced CVD (PECVD) system, 137–138
 - Plasma-etching sources, 27
 - Plasma oxide
 - properties, 140
 - Plastic-leaded chip carrier (PLCC), 187
 - PLA structure, 159
 - Plated copper, 169
 - Platinum
 - bulk resistivity, 153
 - PLCC, 187
 - PMMA, 289
 - PMOS devices, 208
 - pn* junction breakdown voltage
 - vs. impurity concentration, 205
 - pn* junction by diffusion
 - formation, 76
 - pnpn* transistor, 225
 - cross-section, 225
 - npn* transistor, 225
 - Poisson distribution, 195
 - Polycides, 162–163, 163
 - Polymethylmethacrylate (PMMA), 289
 - Polysilicon, 21
 - CVD process, 8
 - mechanical properties, 270
 - Polysilicon deposition
 - chemical vapor deposition, 138–139
 - Polysilicon gate process
 - mask sequence list, 218
 - Polysilicon gate transistor
 - aggressive layout, 220
 - channel-length bias, 220
 - layout
 - MOS transistor layout and design rules, 217–218
 - minimum-size layout, 218
 - Polysilicon interconnections, 159–162
 - buried contacts, 160
 - butted contacts, 162
 - Polysilicon resonant beam
 - layout, 284
 - Positive resist, 23
 - Postprocessing
 - IC process compatibility, 292–293
 - Potassium oxide etch, 276
 - Power-delay product (PDP), 212
 - Practical nonlinear ohmic contact
 - characteristics, 154
 - Prebaking
 - photolithographic process, 22
 - Predeposition step, 71
 - Preprocessing
 - IC process compatibility, 292
 - Production-level ion implanter
 - cost, 111
 - Programmable-logic array (PLA) structure, 159
 - Projected range
 - Gaussian distribution, 112
 - LSS theory, 113
 - Projection printing, 28
 - illustrated, 31
 - Proximity
 - illustrated, 31
 - Proximity printing
 - illustrated, 31
 - Punch-through
 - MOS device considerations, 204–205
- ## Q
- Quad flat pack (QFP), 187
 - Quad J-leaded pack (QJP), 187
- ## R
- Rapid thermal annealing (RTA), 123
 - Rapid thermal nitridation (RTN), 123
 - Rapid thermal oxidation (RTO), 123
 - Rapid thermal processing (RTP)
 - concept illustrated, 123
 - Reactive-ion etching (RIE), 26, 278
 - Reactors
 - CVD, 137–138
 - Recessed oxide, 55
 - Rectifying contact
 - characteristics, 154
 - Reentrant cavity heat sink structure, 290
 - Released structures
 - bulk micromachining, 275–276
 - Release problems
 - surface micromachining, 286–287

- Resist, 23
 - ashing, 26
 - patterns, 24
 - Resistance
 - and size, 83
 - and thickness, 82
 - Resistors
 - effective square contributions, 84
 - n^+ emitter diffusion, 243
 - n -type epitaxial layer, 245
 - pinched base, 246
 - p -type base diffusion, 244
 - Reticle, 28
 - Retrograde
 - ion implantation, 118
 - Reverse-biased pn junction, 158
 - RJE, 26, 278
 - Room-temperature resistivity
 - n - and p -type silicon, 77
 - Rotary structures
 - formation, 283
 - Round TO-style
 - packages, 184
 - RTA, 123
 - RTN, 123
 - RTO, 123
 - RTP
 - concept illustrated, 123
 - Rubylith, 28
- S**
- Safety
 - microelectronic fabrication, 12–14
 - Salicides, 162–164
 - SBC. *See* Standard buried collector (SBC)
 - Scanning electron microscopy (SEM), 37
 - aluminum interconnection failure, 158
 - Scanning system, 110
 - Schottky barrier, 154
 - diode, 250
 - Scribing wafer
 - initial functional testing, 178
 - Sealed cavities
 - surface micromachining, 279
 - Sealed cavity formation, 282
 - Secondary Ion Mass Spectroscopy (SIMS), 91–94
 - analysis concept, 94
 - impurity profile measurement example, 94
 - silicon, 92
 - Selective ion implantation, 114–117
 - Self aligned polysilicon-gate transistor, 210
 - Self aligned silicides, 163, 164
 - Self isolated
 - MOS devices, 208
 - SEM, 37, 158
 - Semirecessed oxide NMOS process
 - sequence, 9
 - Semirecessed oxide structure, 55
 - Shadowing and step coverage
 - evaporation, 134–135
 - Shadowing problem
 - low-pressure vacuum deposition, 134–135
 - Shallow amorphous layers
 - low-energy arsenic implantations, 120–121
 - Shallow boron diffusion, 276
 - Shallow ion implantation, 121–124
 - low-energy, 122
 - rapid thermal annealing, 123
 - TED, 124
 - Shallow phosphorus diffusion profiles
 - constant-source diffusions, 82
 - Shallow trench formation
 - oxidation, 55–56
 - Shallow trench isolation
 - application, 226
 - CMOS technology, 225–226
 - illustration, 209
 - Sheet resistance, 81–90
 - definition, 82–84
 - four-point probe, 88
 - ion implantation, 117–118
 - Irvin's curves, 85–87
 - van der Pauw's method, 88–90
 - van der Pauw's test structure
 - defused layer, 90
 - Sheet resistance-junction depth
 - vs. surface impurity concentration, 86–87
 - Silane, 140
 - threshold limit recommendations, 99
 - Silicide, 162–163
 - contacts
 - device cross-section, 165
 - eutectic temperature, 163
 - layers
 - feature, 163
 - and multilayer-contact technology, 162–164
 - properties, 162
 - Silicon, 1
 - eutectic temperature, 163
 - mechanical properties, 270
 - SIMS analysis, 92
 - Silicon dioxide, 21
 - chemical vapor deposition, 139–140
 - layer formation, 44
 - masking properties, 51–52
 - patterns, 24
 - thermally grown
 - color chart, 60
 - thickness, 53
 - Silicon dopants
 - diffusion, 95
 - Silicon epitaxial growth process
 - mole fraction, 144

- temperature, 143
- Silicon etchants, 277–279
 - anisotropic etching, 278–279
 - comparisons, 277
 - isotropic etching, 277–278
- Silicon fusion bonding
 - silicon wafer bonding, 289–291
- Silicon germanium epitaxial base transistor
 - advanced bipolar structures, 257–258
- Silicon germanium heterojunction bipolar transistor
 - epitaxial base, 258
 - impurity profile, 259
- Silicon glass
 - diffusivities, 44
- Silicon lattice view, 119
- Silicon nitride, 21
 - chemical vapor deposition, 140–141
- Silicon on insulator
 - MOS process integration, 226–227
- Silicon wafer
 - bonding, 289–295
 - adhesive, 289
 - anodic, 291–292
 - silicon fusion, 289–291
 - cleaning procedure, 21
 - covered with barrier layer, 21
 - identification, 20
 - relative size, 3
- Silver
 - bulk resistivity, 153
- SIMS. *See* Secondary Ion Mass Spectroscopy (SIMS)
- Single-in-line packages (SIPs), 184
- Sintering step, 163
- SIOC, 187
- SIPs, 184
- Six-level metallization
 - microphotographs, 171
- Small outline integrated circuit (SIOC), 187
- Small outline J-leaded (SOJ), 187
- Small outline transistor (SOT), 187
- Sodium
 - diffusivities, 44
- Sodium-ion contamination, 54
- Soft baking or prebaking
 - photolithographic process, 22
- SOJ, 187
- Solder ball
 - cross-section, 188
 - footprints, 189
- Solid-solubility limits
 - diffusion, 74–75
- Solid-source system
 - diffusion, 95–97
- Solvent removal, 22
- SOT, 187
- Source-drain diffusion
 - PMOS and NMOS transistors, 225
- Space-charge region width
 - voltage and doping, 238
- Spreading-resistance measurements, 91
- Spreading resistance method
 - impurity profile measurement example, 93
- Sputtering, 5, 131
 - film deposition, 135–136
- Sputtering healed
 - vs. ion energy
 - dc sputtering system, 136
- Standard buried collector (SBC)
 - process, 233
 - transistor
 - photo and cross-section, 234
- Standard buried collector (SBC) technology
 - elements in, 243–249
 - base resistor, 244–245
 - emitter resistor, 243–244
 - epitaxial layer resistor, 245
 - lateral pnp transistors, 248–249
 - pinch resistor, 246
 - Schottky diodes, 249
 - substrate pnp transistor, 246–247
- Stanford University Process Engineering Modeling
 - program (SUPREM), 61–62
 - simulation results
 - two-step boron diffusion, 95
- Steel
 - mechanical properties, 270
- Step-and-scan method, 31
- Straggle
 - Gaussian distribution, 112
 - LSS theory, 113
- Straight edges, 20
- Subcollector, 10
- Substitutional diffusers, 72
- Substitutional diffusion, 67
- Substrate doping
 - MOS device considerations, 203
- Substrate pnp transistor, 247
 - elements in SBC technology, 246–247
- Substrate sensitivity, 205
- Subsurface peaks
 - ion implantation, 118
- Subtractive etching, 165
- Subtractive processes, 166
- SUPREM, 61–62, 95
- Surface impurity concentration
 - vs. sheet resistance-junction depth, 86–87
- Surface micromachining, 279–287
 - bridges, 279
 - cantilever beams, 279
 - movable in-plane structures, 279–281
 - out-of-plane motion, 282–286
 - release problems, 286–287

- rotary structure formation, 283
- sealed cavities, 279
- Surface mounting packages, 186

T

- TAB process, 191–192
- Tantalum
 - CVD metal deposition, 141
- Tape-automated-bonding (TAB) process, 191–192
 - illustration, 192
- Target chamber, 110
- TED, 257–258
 - shallow ion implantation, 124
- TEM, 38
- TEOS, 139, 140
- Testing, 177–178
- Tetraethylorthosilicate (TEOS)
 - LPCVD system, 139
 - oxide
 - properties, 140
- Texas Instruments, 1
- Thermally grown silicon dioxide
 - color chart, 60
- Thermal oxidation of silicon, 43–66
 - dopant redistribution during oxidation, 50–51
 - factors influencing oxidation rate, 46–50
 - model, 45
 - oxidation and shallow trench formation, 55–56
 - oxidation modeling, 44–46
 - oxidation process, 43
 - oxidation technology, 52–53
 - oxide quality, 53–54
 - oxide thickness characterization, 57–60
 - process simulation, 61
 - silicon dioxide masking properties, 51–52
 - temperature
 - linear rate constant, 48
 - parabolic rate constant, 47–48
- Thermal treatment, 163
- Thermocompression
 - wire bonding, 182–183
- Thermosonic
 - wire bonding, 184
- Thermosonic ball-wedge bonding
 - gold wire, 182
- Thickness measurement, 59, 61
- Thin small outline package (TSOP), 187
- Three-dimensional hinges, 285
- Threshold adjustment
 - MOS device considerations, 206–207
- Threshold voltage
 - MOS device considerations, 203
 - step approximation

- Gaussian impurity profile, 207
- Titanium
 - bulk resistivity, 153
 - CVD metal deposition, 141
- Titanium-tungsten
 - barrier metal, 164
- TMAH, 278
- TO-style packages, 184, 185
- Transient enhanced diffusion
 - examples, 122
- Transient enhanced diffusion (TED), 257–258
 - shallow ion implantation, 124
- Transistor
 - cross-section
 - CDI process, 259
 - microprocessor
 - yearly change, 4
- Transistor isolation
 - MOS device considerations, 208–210
- Transistor scaling
 - MOS device considerations, 210–212
- Transistor structure
 - enclosed lateral pnp, 248
- Transit time, 236–237
 - bipolar process integration, 236–237
- Transmission electron microscopy (TEM), 38
- Transverse straggle
 - selected implantation, 114
- Trench isolation
 - structures, 57
 - thermal oxidation of silicon, 56
- Triangular approximation
 - Gaussian density, 196
- TSOP, 187
- Tungsten
 - bulk resistivity, 153
 - CVD, 136, 141
 - mechanical properties, 270
- Twin-well CMOS structure
 - processes illustrated, 223
- Two-dimensional lattice
 - atomic diffusion, 68
- Two-level metallization processes, 166
- Two-step boron diffusion
 - SUPREM simulation results, 95

U

- Ultrasonic
 - wire bonding, 183
- Uniform defect densities
 - yield, 194–195
- Uniform density function, 196
- Unity-gain frequency
 - transistor, 236

V

- Vacuum system
 - illustration, 130
- van der Pauw's method
 - sheet resistance, 88–90
- van der Pauw's test structure
 - sheet resistance
 - defused layer, 90
- Vapor-phase epitaxy (VPE), 141
 - epitaxy, 142–145
 - reactors
 - types, 143
 - vertical reactor, 143
- Vapor-phase mass-transfer coefficient, 142
- Vertical diffusion
 - junction formation, 76–78
- Vertical furnace
 - oxidation and diffusion, 54
- Vertical laminar-flow hoods, 17
- Vertical VPE reactor, 143
- Vibrations
 - ultrasonic bonding, 183
- VLSI
 - chips, 166
 - CMOS, 141
 - devices, 79
 - fabrication
 - Low-temperature processing, 111
 - pattern transfer, 37
 - high-density process
 - junction penetration, 156
 - high-resolution
 - lithography systems, 29
 - signal processing, 292
 - MEMS devices, 292
 - structures
 - shallow amorphous layers, 120–121
 - ULSI, 17
- VPE. *See* Vapor-phase epitaxy (VPE)

W

- Wafer
 - diameter, 2
 - dicing saw, 179
 - die size vs. yield, 194
 - drawing, 19
 - initial functional testing
 - scribing, 178

- photolithographic process, 19–21
- relative size, 3
- size, 2
- Wafer cleaning
 - photolithographic process, 19–21
- Wafer correction
 - four-point-probe correction factors, 89
- Wafer flats, 20
 - illustration, 20
- Wafer SOI
 - bonded
 - formation, 227
- Wafer stepper
 - concept
 - lens system, 32
 - system drawing, 33
- Wafer thinning, 178
- Water vapor
 - diffusivities, 44
- Wet chemical
 - etching techniques, 25
- Wet silicon dioxide
 - growth, 49
 - pressures, 50
- Width biases
 - MOS transistor layout and design rules, 219–221
- Wire bonding, 179–184
 - process
 - bumps formed, 189
 - thermocompression, 182–183
 - thermosonic, 184
 - ultrasonic, 183

X

- Xenon-Hg lamp
 - spectral content, 35
- X-ray radiation, 134

Y

- Yield, 177–198
 - nonuniform defect densities, 195–197
 - uniform defect densities, 194–195
- Yield curves
 - defect densities, 197

Z

- Zig-zag-in-line packages (ZIPs), 184

USEFUL TABLES AND FIGURES

Table Number	Table Title	Page
Table 1.1	ITRS Projections	6
Table 2.2	Clean Room Specifications	18
Table 3.1	Linear and Parabolic Rate Constants	47
Figures 3.6/3.7	Silicon Dioxide Growth Rate	49
Figure 3.10	Oxide Mask Thickness	53
Table 3.2	Oxide Color Chart	60
Figure 4.4	Gaussian and Complementary Error Functions	71
Table 4.1	Diffusion Coefficient Data	74
Figure 4.6	Solid-Solubility Limits in Silicon	75
Figure 4.8	Silicon Resistivity	77
Figure 4.10	Lateral Diffusion under a Mask Edge	80
Table 4.2	High-Concentration Boron and Arsenic Diffusions	81
Figure 4.12	High-Concentration Phosphorus Diffusion	82
Figure 4.15	Resistor Terminals	84
Figure 4.16	Irvin's Curves	86-87
Figure 4.18	Four-Point-Probe Correction Factors	89
Table 4.4	Toxic Gas Characteristics	99
Figure 5.3	Ion Implantation Range and Straggle	113
Figure 6.10	Silicon Epitaxial Growth Rate	143
Table 7.1	Bulk of Resistivity of Metals	153
Figure 7.4	Aluminum-Silicon Phase Diagram	155
Table 7.2	Silicide Resistivity	162
Figure 8.18	Normalized Yield for Several Models	197
Figure 9.2	Threshold Voltage versus Substrate Concentration	204
Figure 9.3	Junction Breakdown Voltage	205
Figure 9.4	One-Sided Step Junction Depletion-Layer Width	206
Table 9.1	Constant Field Scaling Results	213
Figure 10.3	Diffusion Lengths in Silicon	236
Figure 10.4	Depletion-Layer Widths in Collector-Base Junction	238
Figure 10.5	Emitter-Base Breakdown Voltage	239
Figure 10.6	Collector-Base Breakdown Voltage	241
Table 11.1	Mechanical Properties of Selected Materials	270

		IIIA	IVA	VA	VIA
IIB	5 10.811	B Boron	6 12.01115	7 14.0067	8 15.9994
	13 26.9815	Al Aluminum	14 28.086 4	15 30.9738	16 32.064
30 65.37	31 69.72	32 72.59	33 74.922	34 78.96	
Zn Zinc	Ga Gallium	Ge Germanium	As Arsenic	Se Selenium	
48 112.40	49 114.82	50 118.69	51 121.75	52 127.60	
Cd Cadmium	In Indium	Sn Tin	Sb Antimony	Te Tellurium	
80 200.59	81 204.37	82 207.19	83 208.980	84 (210)	
Hg Mercury	Tl Thallium	Pb Lead	Bi Bismuth	Po Polonium	

Modular Series on Solid State Devices, Volume V

G.W. Neudeck and R.F. Pierret, both of Purdue University

Introduction to Microelectronic Fabrication, Second Edition, by Richard C. Jaeger, is a concise survey of the most up-to-date techniques in the field. It is devoted exclusively to processing and is highlighted by careful explanations, clean, simple language, and numerous fully solved example problems.

The second edition includes an entirely new chapter on MEMS, as well as substantial modifications to the chapters on MOS and bipolar process integration. Also included is new or expanded coverage of lithography and exposure systems, trench isolation, chemical mechanical polishing, shallow junctions, transient enhanced diffusion, copper Damascene processes and process simulation.

Volumes in Series:

Volume I: *Semiconductor Fundamentals, Second Edition* (R.F. Pierret)
0-201-12295-2

Volume II: *The PN Junction Diode, Second Edition* (G.W. Neudeck)
0-201-12296-0

Volume III: *The Bipolar Junction Transistor, Second Edition* (G.W. Neudeck)
0-201-12297-9

Volume IV: *Field Effect Devices, Second Edition* (R.F. Pierret)
0-201-12298-7

Volume V: *Introduction to Microelectronic Fabrication, Second Edition* (R.C. Jaeger)
0-201-44494-1

**Pearson
Education**

Prentice Hall
Upper Saddle River, NJ 07458
www.prenhall.com