



For training, visit mindshare.com 

MindShare Technology Series

x86 Instruction Set Architecture

Comprehensive 32- and 64-bit Coverage

Tom Shanley | MindShare, Inc.



x86 Instruction Set Architecture

*Comprehensive 32/64-bit Coverage
First Edition*



Also by Tom Shanley

HEAVEN'S FAVORITE

—A Novel of Genghis Khan—

Book 1, **ASCENT: THE RISE OF CHINGGIS KHAN**

Book 2, **DOMINION: DAWN OF THE MONGOL EMPIRE**

MINDSHARE TECHNICAL TRAINING

Please visit www.mindshare.com for a complete description of Mind-Share's technical offerings:

- Books
- eBooks
- eLearning modules
- Public courses
- On-site course
- On-line courses

Intel Core 2 Processor (Penryn)
Intel Nehalem Processor
Intel Atom Processor
AMD Opteron Processor (Barcelona)
Intel 32/64-bit x86 Software Architecture
AMD 32/64-bit x86 Software Architecture
x86 Assembly Language Programming
Protected Mode Programming
PC Virtualization
IO Virtualization (IOV)
Computer Architectures with Intel Chipsets
Intel QuickPath Interconnect (QPI)
PCI Express 2.0
USB 2.0
USB 3.0
Embedded USB 2.0 Workshop
PCI
PCI-X
Modern DRAM Architecture
SAS
Serial ATA
High Speed Design
EMI / EMC
Bluetooth Wireless Product Development
SMT Manufacturing
SMT Testing

x86 Instruction Set Architecture

Comprehensive 32/64-bit Coverage
First Edition

MINDSHARE, INC.

TOM SHANLEY

MindShare Press
Colorado Springs, USA

Refer to “Trademarks” on page 5 for trademark information.

The author and publisher have taken care in preparation of this book but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

ISBN: 0-9770878-5-3

Copyright © 2009 by MindShare, Inc.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher.
Printed in the United States of America.

Cover Design: Michelle Petrie
Set in 10 point Palatino by MindShare, Inc.

First Printing, December 2009

MindShare Press books are available for bulk purchases by corporations, institutions, and other organizations. For more information please contact the Special Sales Department at (575)-373-0336.

Find MindShare Press on the World Wide Web at:
<http://www.mindshare.com/>

To Nancy, the strongest person I know.

With Love,

Tom

P. S. It's done. I'm back.

At-a-Glance

Table of Contents

Part 1: Introduction, intended as a back-drop to the detailed discussions that follow, consists of the following chapters:

- Chapter 1, "Basic Terms and Concepts," on page 11.
- Chapter 2, "Mode/SubMode Introduction," on page 21.
- Chapter 3, "A (very) Brief History," on page 41.
- Chapter 4, "State After Reset," on page 63.

Part 2: IA-32 Mode provides a detailed description of two IA-32 Mode sub-modes—Real Mode and Protected Mode—and consists of the following chapters:

- Chapter 5, "Intro to the IA-32 Ecosystem," on page 79.
 - Chapter 6, "Instruction Set Expansion," on page 109.
 - Chapter 7, "32-bit Machine Language Instruction Format," on page 155.
 - Chapter 8, "Real Mode (8086 Emulation)," on page 227.
 - Chapter 9, "Legacy x87 FP Support," on page 339.
 - Chapter 10, "Introduction to Multitasking," on page 361.
 - Chapter 11, "Multitasking-Related Issues," on page 367.
 - Chapter 12, "Summary of the Protection Mechanisms," on page 377.
 - Chapter 13, "Protected Mode Memory Addressing," on page 383.
 - Chapter 14, "Code, Calls and Privilege Checks," on page 415.
 - Chapter 15, "Data and Stack Segments," on page 479.
 - Chapter 16, "IA-32 Address Translation Mechanisms," on page 493.
 - Chapter 17, "Memory Type Configuration," on page 599.
 - Chapter 18, "Task Switching," on page 629.
 - Chapter 19, "Protected Mode Interrupts and Exceptions," on page 681.
 - Chapter 20, "Virtual 8086 Mode," on page 783.
 - Chapter 21, "The MMX Facilities," on page 835.
 - Chapter 22, "The SSE Facilities," on page 851.
-

Part 3: IA-32e OS Kernel Environment provides a detailed description of the IA-32e OS kernel environment and consists of the following chapters:

- Chapter 23, "IA-32e OS Environment," on page 913.
- Chapter 24, "IA-32e Address Translation," on page 983.

Part 4: Compatibility Mode provides a detailed description of the Compatibility submode of IA-32e Mode and consist of the following chapter:

- Chapter 25, "Compatibility Mode," on page 1009.

Part 5: 64-bit Mode provides a detailed description of the 64-bit submode of IA-32e Mode and consists of the following chapters:

- Chapter 26, "64-bit Register Overview," on page 1023.
- Chapter 27, "64-bit Operands and Addressing," on page 1041.
- Chapter 28, "64-bit Odds and Ends," on page 1075.

Part 6: Mode Switching Detail provides a detailed description of:

- Switching from Real Mode to Protected Mode. This topic is covered in Chapter 29, "Transitioning to Protected Mode," on page 1113.
- Switching from Protected Mode to IA-32e Mode. This topic is covered in Chapter 30, "Transitioning to IA-32e Mode," on page 1139.

Part 7: Other Topics provides detailed descriptions of the following topics:

- Chapter 31, "Introduction to Virtualization Technology," on page 1147.
- Chapter 32, "System Management Mode (SMM)," on page 1167.
- Chapter 33, "Machine Check Architecture (MCA)," on page 1207.
- Chapter 34, "The Local and IO APICs," on page 1239.

About This Book

Is This the Book for You?	1
A Moving Target	1
x86 Instruction Set Architecture (ISA)	1
Glossary of Terms	2
32-/64-bit x86 Instruction Set Architecture Specification	2
The Specification Is the Final Word	2
Book Organization	3
Topics Outside the Scope of This Book	4
The CPUID Instruction	4
Detailed Description of Hyper-Threading	4
Detailed Description of Performance Monitoring	5
Documentation Conventions	5
Trademarks	5
Visit Our Web Site	6
We Want Your Feedback	7

Part 1: Introduction

Chapter 1: Basic Terms and Concepts

ISA Definition	11
This Book Focuses on the Common Intel/AMD ISA	11
For Simplicity, Intel Terminology Is Used Throughout	11
Some Terms in This Chapter May Be New To the Reader	12
Two x86 ISA Architectures	12
Processors, Cores and Logical Processors	13
Fundamental Processing Engine: Logical Processor	14
IA Instructions vs. Micro-ops	15
RISC Instructions Sets Are Simple	15
x86 Instruction Set Is Complex	15
But You Can't Leave It Behind	16
Complexity vs. Speed Dictated a Break With the Past	16
Why Not Publish a Micro-Op ISA?	16
Some Important Definitions	17
Virtual vs. Physical Memory	17
Other Important Terms	18

Chapter 2: Mode/SubMode Introduction

Basic Execution Modes	21
-----------------------------	----

Contents

IA-32 SubModes	25
IA-32e SubModes	28
Mode Switching Basics	30
Initial Switch from IA-32 to IA-32e Mode	30
IA-32e SubMode Selection	33
Protected/Compatibility 16-/32-bit SubModes	38

Chapter 3: A (very) Brief History

Major Evolutionary Developments	42
16-bit Mode Background	46
8086 and Real Mode.....	46
286 Introduced 16-bit Protected Mode.....	48
386 Supported Both 16- and 32-bit Protected Mode	51
The Intel Microarchitecture Families	55
A Brief Timeline	57

Chapter 4: State After Reset

State After Reset	64
Soft Reset	73
Boot Strap Processor (BSP) Selection	73
AP Discovery and Configuration.....	74
Initial Memory Reads.....	74

Part 2: IA-32 Mode

Chapter 5: Intro to the IA-32 Ecosystem

The Pre-386 Register Sets.....	80
8086 Register Set.....	80
286 Register Set.....	82
IA-32 Register Set Overview	84
Control Registers.....	85
Status/Control Register (Eflags)	88
Instruction Fetch Facilities	89
General.....	89
Branch Prediction Logic	90
General Purpose Data Registers.....	90
Defining Memory Regions/Characteristics.....	92
MTRRs	92
Segment Registers	92
Address Translation Facilities.....	93

Contents

Interrupt/Exception Facilities	93
Kernel Facilities.....	94
Real Mode Has No Memory Protection	95
Memory Protection in Protected Mode	95
Introduction.....	95
Segment Selection in Protected Mode	95
Access Rights Check.....	96
The Descriptor Tables	96
Descriptor Table Registers.....	96
Task Data Structure	97
Address Translation Facilities.....	97
Effective/Virtual/Linear/Physical Addresses.....	97
Introduction to Address Translation (Paging).....	98
RAM Is Finite and Can't Hold Everything	98
RAM and Mass Storage Are Managed on a Page Basis	99
This Requires a Series of Directories.....	99
Malloc Request.....	99
Problem: Non-Contiguous Memory Allocation.....	100
Malloc Returns a Virtual Address to the Application.....	100
IA-32 Applications Have a 4GB Virtual Address Space	101
Legacy FP Facilities.....	101
In the Beginning, FPU Was External and Optional	101
It Was Slow... ..	102
486DX Integrated It.....	102
x87 Register Set.....	102
x87 FP Instruction Set.....	102
General Purpose Instruction Set	102
MMX Facilities.....	102
Introduction	103
SIMD Programming Model.....	103
SSE Facilities.....	104
Introduction	104
Motivation.....	104
Instruction Set.....	105
Model-Specific Registers.....	105
General.....	105
Accessing the MSRs.....	106
Debug Facilities.....	106
Automatic Task Switching Mechanism.....	107

Chapter 6: Instruction Set Expansion

Why a Comprehensive Instruction Set Listing Isn't Included.....	110
--	------------

Contents

386 Instruction Set.....	111
Instruction Set (as of March, 2009).....	117

Chapter 7: 32-bit Machine Language Instruction Format

64-bit Machine Language Instruction Format	156
A Complex Instruction Set with Roots in the Past	156
Effective Operand Size	157
Introduction	157
Operand Size in 16- and 32-bit Code Segments	157
Operand Size in 64-bit Code Segments.....	158
Instruction Composition.....	160
Instruction Format Basics	162
Opcode (Instruction Identification)	168
In the Beginning	168
1-byte Opcodes.....	169
2-byte Opcodes Use 2-Level Lookup	172
2nd-Level Opcode Map Introduced in 286	172
Instructions with 2-byte Opcodes: Five Possible Forms	172
3-byte Opcodes Use 3-Level Lookup	176
3-Level Opcode Maps Introduced in Pentium 4 Prescott.....	176
Currently There Are Two 3rd-Level Maps Defined	176
Instructions with 3-byte Opcodes: Three Possible Forms.....	176
Special Use of Prefix Bytes	177
Opcode Micro-Maps (Groups).....	180
Micro-Maps Associated with 1-byte Opcodes.....	180
Some Opcodes Employ 2 x 8 Micro-Maps	180
Micro-Maps Associated with 2-byte Opcodes.....	183
3-byte Opcodes Don't Use Micro-Maps	187
x87 FP Opcodes Inhabit Opcode Mini-Maps	187
Special Opcode Fields	189
Operand Identification	194
General.....	194
Specifying Registers as Operands	195
Implicit Register Specification	196
Explicit Register Specification in Opcode	196
Explicit Register Specification in ModRM Byte	196
Addressing a Memory-Based Operand.....	198
Instruction Can Specify Only One Memory-Based Operand.....	198
Addressing Memory Using the ModRM Byte.....	200
When Effective Address Size = 16-Bits.....	200
When Effective Address Size = 32-Bits.....	202
Using the SIB Byte to Access a Data Structure	203

Near and Far Branch Target Addressing.....	206
Specifying an Immediate Value As an Operand	209
Instruction Prefixes.....	210
Operand Size Override Prefix (66h)	211
In 32-bit Mode	211
In 16-bit Mode	213
Special Usage of 66h Prefix	213
Address Size Override Prefix (67h)	214
In 32-Bit Mode	214
In 16-Bit Mode	215
Lock Prefix	215
Shared Resource Concept.....	215
Race Condition Can Present Problem.....	216
Guaranteeing Atomicity of Read/Modify/Write.....	216
Use Locked RMW to Obtain and Give Up Semaphore Ownership	217
Instructions That Accept Lock Prefix.....	218
Repeat Prefixes	218
Normal Usage.....	218
Special Usage.....	220
Segment Override Prefix.....	220
General	220
Usage In String Operations	221
Segment Override Use With MMX and SSE1 - 4 Instructions	221
Branch Hint Prefix	221
Summary of Instruction Set Formats	222

Chapter 8: Real Mode (8086 Emulation)

8086 Emulation.....	229
Unused Facilities	231
Real Mode OS Environment.....	232
Single-Task OS Environment Overview.....	232
Command Line Interface (CLI).....	232
Program Loader	233
OS Services.....	233
Direct IO Access	234
Application Memory Usage	234
Task Initiation, Execution and Termination.....	234
Running Real Mode Applications Under a Protected Mode OS.....	235
Real Mode Applications Aren't Supported in IA-32e Mode.....	235
Real Mode Register Set.....	235
Introduction	235

Contents

Control Registers.....	237
CR0.....	238
Address Translation (Paging) Control Registers	243
CR2.....	243
CR3.....	243
CR4 (Feature Control Register).....	244
XCR0 (XFEM)	249
Flags Register.....	251
General Purpose Registers (GPRs)	255
A, B, C and D Registers.....	255
General Usage	255
Special Usage Examples	255
EBP Register: Stack Frame Address Register	256
Index Registers.....	258
Stack Pointer (SP) Register	259
Instruction Pointer Register.....	259
Kernel Registers.....	260
x87/MMX FPU Register Set	260
SSE Register Set	262
Debug Address Breakpoint Register Set	262
General	262
Defining Trigger Address Range.....	263
Defining Access Type.....	264
Defining Scope (Current Task or All Tasks)	264
Special Notes	264
Local APIC Register Set	269
Architecturally-Defined MSRs.....	272
General	272
Determining MSR Support.....	272
Accessing the MSRs.....	272
IO Space versus Memory Space	281
IO Operations	281
IO Operations in IO Address Space	281
IN and OUT Instructions	281
Block (String) IO Operations.....	282
Block Transfer from IO Port to Memory	282
Block Transfer from Memory to an IO Port.....	283
IO Space is Limited and Crowded	284
Memory-Mapped IO (MMIO) Operations.....	284
Introduction.....	284
Know the Characteristics of Your Target.....	284
Why the Logical Processor Must Know the Memory Type	284

Uncacheable (UC) Memory	285
No IO Protection	286
Operand Size Selection.....	286
Address Size Selection	287
Real Mode Memory Addressing	288
No Address Translation.....	288
Introduction to Real Mode Segmentation	288
All Segments are 64KB in Size.....	292
Memory Address Representation.....	293
Accessing the Code Segment.....	293
Jumping Between Code Segments	294
Far Jumps and Calls	294
Near Jumps and Calls	295
IP-Relative Branches.....	296
Operations That Default to the Code Segment.....	296
Accessing the Stack Segment.....	297
Introduction.....	297
Stack Characteristics.....	298
Pushing Data Onto the Stack	298
Popping Data From the Stack	300
Stack Underflow/Overflow	300
Processor Stack Usage.....	301
Accessing Parameters Passed on the Stack	301
Operations That Default To the Stack Segment	302
Accessing the DS Data Segment	303
General	303
Operations That Default to the DS Data Segment	303
Accessing the ES/FS/GS Data Segments.....	304
General	304
Operations That Default to the ES Data Segment.....	305
Segment Override Prefixes.....	305
Example Segment Register Initialization.....	305
Accessing Extended Memory in Real Mode	307
Big Real Mode.....	310
286 DOS Extender Programs.....	311
Hot Reset and 286 DOS Extender Programs.....	311
Alternate (Fast) Hot Reset	312
286 DOS Extenders on Post-286 Processors	313
String Operations	315
Real Mode Interrupt/Exception Handling	316
Events and Event Handlers	316
Events Are Recognized on an Instruction Boundary	317

Contents

The IDT	317
Definition of the IDT	317
IDT and IDTR Initialization	324
Stack Initialization	325
Event (Interrupt and Exception) Handling	326
Software Event Types.....	327
Introduction.....	327
Software Exceptions	328
Definition of an Exception.....	328
Exception Handling.....	328
Three Categories of Software Exceptions.....	330
Software Interrupt Instructions	330
INT nn Instruction.....	330
BOUND Instruction.....	331
INTO Instruction	331
INT3 (Breakpoint) Instruction	332
Hardware Event Types	332
NMI.....	332
Definition of NMI and Delivery Mechanisms.....	332
External NMI Masking Mechanism.....	332
NMI Handling.....	333
SMI.....	334
Maskable Interrupts	334
Maskable Interrupts Are Originated by Devices	334
Enabling/Disabling Maskable Interrupt Recognition	334
Selective Masking of Maskable Interrupts.....	335
Maskable Interrupt Delivery Mechanisms	335
IDT Entries Associated with Maskable Interrupts	335
Handling Maskable Interrupts	335
Machine Check Exception	336
Summary of Real Mode Limitations	337
Transitioning to Protected Mode	337

Chapter 9: Legacy x87 FP Support

A Little History	340
x87 FP Instruction Format.....	341
FPU-Related CR0 Bit Fields	341
x87 FPU Register Set.....	343
The FP Data Registers.....	343
x87 FPU's Native Data Operand Format.....	344
32-bit SP FP Numeric Format.....	346
Background.....	346

A Brief IEEE FP Primer	346
The 32-bit SP FP Format.....	347
Representing Special Values	348
An Example	348
Another Example	349
DP FP Number Representation.....	350
FCW Register	350
FSW Register.....	352
FTW Register	355
Instruction Pointer Register.....	355
Data Pointer Register.....	355
Fopcode Register.....	356
General	356
Fopcode Compatibility Mode	356
FP Error Reporting	357
Precise Error Reporting.....	357
Imprecise (Deferred) Error Reporting.....	357
Why Deferred Error Reporting Is Used.....	358
The WAIT/FWAIT Instruction.....	358
CR0[NE].....	358
DOS-Compatible FP Error Reporting	359
FP Error Reporting Via Exception 16.....	359
Ignoring FP Errors	360

Chapter 10: Introduction to Multitasking

Concept.....	363
An Example—Timeslicing.....	364
Another Example—Awaiting an Event.....	364
1. Task Issues Call to OS for Disk Read	364
2. Device Driver Initiates Disk Read	364
3. OS Suspends Task	365
4. OS Makes Entry in Event Queue	365
5. OS Starts or Resumes Another Task.....	365
6. Disk-Generated Interrupt Causes Jump to OS	365
7. Interrupted Task Suspended.....	366
8. Task Queue Checked.....	366
9. OS Resumes Task	366

Chapter 11: Multitasking-Related Issues

Hardware-based Task Switching Is Slow!	368
Private (Local) and Global Memory	369

Contents

Preventing Unauthorized Use of OS Code	369
With Privilege Comes Access	370
Program Privilege Level	370
The CPL.....	370
Calling One of Your Equals.....	371
Calling a Procedure to Act as Your Surrogate.....	371
Data Segment Protection.....	371
Data Segment Privilege Level.....	371
Read-Only Data Areas	371
Some Code Segments Contain Data, Others Don't.....	372
IO Port Anarchy.....	374
No Interrupts, Please!	375
BIOS Calls	376

Chapter 12: Summary of the Protection Mechanisms

Protection-Related Mechanisms.....	378
------------------------------------	-----

Chapter 13: Protected Mode Memory Addressing

Real Mode Segment Limitations.....	384
An Important Reminder: Segment Base + Offset = Virtual Address	385
Descriptor Contains Detailed Segment Description.....	386
Segment Register—Selects Descriptor Table and Entry	386
Introduction to the Descriptor Tables.....	390
Segment Descriptors Reside in Memory	390
Global Descriptor Table (GDT)	393
GDT Description	393
Setting the GDT Base Address and Size.....	393
GDT Entry 0.....	394
Local Descriptor Tables (LDTs)	395
General	395
Creating and Selecting an LDT	396
General Segment Descriptor Format.....	399
Granularity Bit and the Segment Size	399
Segment Base Address Field	400
Default/Big Bit	400
In a Code Segment Descriptor, D/B = “Default” Bit.....	400
Override Prefixes	401
In a Stack Segment Descriptor, D/B = “Big” Bit.....	402
Segment Type Field	403
Introduction to the Type Field.....	403
Non-System Segment Types	403

Segment Present Bit	406
Descriptor Privilege Level (DPL) Field	406
System Bit	407
Available Bit	408
Goodbye to Segmentation	408
Introduction	408
IA-32 Flat Memory Model	409
No Protection? Paging Takes Care of It	412
A Reminder of Where We Are	412

Chapter 14: Code, Calls and Privilege Checks

Abbreviation Alert	416
Selecting the Active Code Segment	416
CS Descriptor	418
CS Descriptor Selector	418
Calculating the Descriptor's Memory Address	419
Descriptor Read and Privilege Checked	419
CS Descriptor Format	419
Accessing the Code Segment	423
In-Line Code Fetching	423
Short and Near Branches (Jumps and Calls)	423
General	423
Example Near Jump	424
Far Branches (Far Jumps and Calls)	424
General	424
Example Far Jump	425
Short/Near Jumps	427
General	427
No Privilege Check	427
Unconditional Short/Near Branches	427
Conditional Branches	428
General	428
Loop Instructions	431
Unconditional Far Jumps	434
The Privilege Check	434
Far Jump Targets	435
Far Jump Forms	435
Privilege Checking	436
No Check on Near Calls or Near Jumps	436
General	437
Definitions	437
Definition of a Task	437

Contents

Definition of a Procedure.....	437
CPL Definition.....	437
CS DPL Definition	438
Conforming and Non-Conforming Code Segments.....	438
Definition	438
Examples.....	438
RPL Definition.....	439
General	439
RPL Usage in Privilege Check	439
RPL Use on RET or IRET	440
Privilege Check on Far Call or Far Jmp	440
General	440
Example.....	440
Jumping from a Higher-to-Lesser Privileged Program.....	441
Direct Procedure Calls.....	442
Introduction	442
General.....	442
Near Calls/Returns	443
Description.....	443
Call/Ret Operand Size Matching.....	445
Near Call/Return Forms	445
Far Calls.....	447
General	447
Far Call Forms.....	447
Far Call, Same Privilege Level	451
Far Call to a More-Privileged Procedure.....	451
Far Call to a Procedure in a Different Task.....	452
Indirect Procedure Far Call Through a Call Gate	452
Example Scenario Defines the Problem	452
The Scenario.....	452
The Problem.....	453
The Solution—Different Gateways	454
The Call Gate Descriptor.....	454
Call Gate Example.....	456
Execution Begins	456
Call Gate Descriptor Read	457
Call Gate Contains Target Code Segment Selector	458
Target Code Segment Descriptor Read	459
The Big Picture	461
The Call Gate Privilege Check	461
Automatic Stack Switch	462
Background.....	462

A Potential Problem	462
The Solution: Pre-Allocated Stacks	462
Far Call From 32-bit CS to 16-bit CS.....	466
General.....	466
Method 1: Far Call with Operand Size Override Prefix	466
Method 2: Far Call Via 16-bit Call Gate	467
Method 3: Call 32-bit/16-bit Interface Procedure	470
Far Call From 16-bit CS to 32-bit CS.....	471
Method 1: Far Call With an Operand Size Prefix	472
Method 2: Far Call Via a 32-bit Call Gate	473
Far Returns	475
General.....	475
Far Return Forms	476

Chapter 15: Data and Stack Segments

A Note Regarding Stack Segments.....	480
Data Segments	481
General.....	481
Two-Step Permission Check.....	481
An Example.....	482
Selecting and Accessing a Stack Segment.....	484
Introduction	484
Expand-Up Stack.....	485
Expand-Down Stack	487
The Problem.....	487
Expand-Down Stack Description	489
An Example	490
Another Example.....	491

Chapter 16: IA-32 Address Translation Mechanisms

Three Generations.....	494
Demand Mode Paging Evolution.....	495
Background	496
Memory and Disk: Block-Oriented Devices.....	496
Definition of a Page	496
Example Scenario: Block Transfer from Disk to Memory.....	497
A Poor Memory Allocation Strategy	498
Applications Are Presented With a Simplified World-View	499
Introduction	499
Life Without Paging Would Be Chaotic	499
The Virtual World Is a Simple One	500

Contents

Virtual Address Space Partitioning.....	506
Example Virtual Buffer Allocation	508
Address Translation Advantages.....	509
Introduction	509
Simplifies Memory Management	509
Efficient Memory Usage.....	510
A Wasteful Approach.....	510
A Better Approach: Load On Demand	511
Attribute Assignment.....	511
Track Access History	512
Allows DOS Applications to Co-Exist	512
Problem: Running Multiple DOS Programs	512
Solution: Address Redirection	512
First-Generation Paging.....	513
Definition of First Generation Paging.....	513
Paging Logic's Interpretation of a Virtual Address	514
First-Generation Paging Overview	515
The Set-Up	515
Virtual-to-Physical Address Translation.....	516
The Goal	516
The Translation	516
Two Overhead Memory Reads Take a Toll	522
The TLBs.....	523
TLB Miss.....	523
TLB Hit	524
TLB Maintenance	524
TLBs Are Cleared on Task Switch or Page Directory Change	525
Updating a Single Page Table Entry	525
Global Pages	526
Problem	526
Global Page Feature.....	526
Enabling Paging	527
Detailed Description of PDE and PTE.....	529
PDE Layout.....	529
PTE Layout.....	531
Checking Page Access Permission.....	535
The Privilege Check.....	535
Segment Privilege Check Takes Precedence Over Page Check	535
U/S Bit in PDE and PTE Are Checked	536
Accesses with Special Privilege	537
The Read/Write Check	537

Missing Page or Page Table	538
Introduction	538
Page Table Not Present	538
Page Not Present	542
Page Faults	545
Page Fault Causes	545
Page Fault During a Task Switch	546
Page Fault while Changing to a Different Stack	547
Page Fault Error Code	547
Additional Page Fault Information	547
Access History.....	548
4MB Pages.....	550
Basic Concept.....	550
Enabling the PSE Feature.....	550
Simplifies Housekeeping	550
How To Set Up a 4MB Page	551
The Address Translation.....	552
Second-Generation Paging.....	553
First-Gen Problem: 4GB Physical Memory	553
The Solution: PAE-36 Mode	553
Enabling PAE-36 Mode	553
CR4[PSE] Is “Don’t Care”	554
Application Still Limited to a 4GB Virtual Address Space	554
Virtual Address Space Partitioning.....	555
First Generation Partitioning	555
Second Generation Partitioning.....	556
Second Generation Uses 3-Level Lookup Mechanism	557
CR3 Points to PDPT in Lower 4GB.....	558
Enlarged Physical Address Space.....	559
The Translation.....	560
Step 1: PDPT Lookup	560
Step 2: Page Directory Lookup	562
Step 2a: PDE Points to a Page Table.....	562
Step 2b: PDE Points to a 2MB Physical Page	563
Step 3: Page Table Lookup	565
Page Protection Mechanisms.....	567
General	567
Write-Protection.....	568
Example Usage: Unix Copy-on-Write Strategy	569
3-Level Lookup—Increased TLB Size	571
Microsoft PAE Support	572
Linux PAE Support.....	574

Contents

PSE-36 Mode (PAE-36 Mode's Poor Cousin)	574
PSE-36 Mode Background	575
Detecting PSE-36 Mode Capability	575
Enabling PSE-36 Mode	575
Per Application Virtual Memory Space = 4GB	576
First-Generation Lookup Mechanism	576
Selected PDE Can Point to 4KB Page Table or a 4MB Page	576
Virtual Address Maps to a 4MB Page in 64GB Space	577
Windows and PSE-36	578
AMD Enhanced PSE-36 to PSE-40	579
Execute Disable Feature	579
Problem: Malicious Code	579
The Overflow	579
The Exploit	581
The Fix: Intercept Code Fetches from Data Pages	582
Enabling the Execute-Disable Feature	582
Available in both IA-32 and IA-32e Mode	583
How It Works	583
Defining a Page's Caching Rules	585
Introduction	585
Translation Table Caching Rules	585
General	585
First-Generation Paging Tables	586
Second-Generation Paging Tables	586
Page Caching Rules	587
PAT Feature (Page Attribute Table)	587
What's the Problem?	587
Detecting PAT Support	588
PAT Allows More Memory Types	588
Default Contents of IA32_CR_PAT MSR	589
Memory Type When Page Definition and MTRR Disagree	590
General	590
The UC- Memory Type	590
Altering IA32_CR_PAT MSR	593
Ensuring IA32_CR_PAT and MTRR Consistency	593
Assigning Multiple Memory Types to a Single Physical Page	595
Compatibility with Earlier IA-32 Processors	596
Third Generation Paging	597

Chapter 17: Memory Type Configuration

Characteristics of Memory Targets	600
Introduction	600

Example Problem: Caching from MMIO	600
Early Processors Implemented Primitive Mechanism.....	601
Solution/Problem: Chipset Memory Type Registers	602
Solution: Memory Type Register Set	602
MTRR Feature Determination.....	603
MTRRs Are Divided Into Four Categories	604
MTRRDefType Register.....	604
State of the MTRRs after Reset	605
Fixed-Range MTRRs.....	605
The Problem: Legacy Issues	605
Enabling the Fixed-Range MTRRs.....	605
Defining Memory Types in Lower 1MB.....	606
Variable-Range MTRRs.....	607
How Many Variable-Range Register Pairs?	607
Variable-Range Register Pair Format.....	607
MTRRPhysBasen Register	608
MTRRPhysMaskn Register.....	608
Programming Variable-Range Register Pairs	609
Enabling Variable-Range Register Pairs	609
Memory Types	609
Memory Type Defines Processor Aggressiveness	609
Five Memory Types.....	610
Uncacheable (UC) Memory	610
Uncacheable Write-Combining (WC) Memory	611
Description.....	611
Weakly-Ordered Writes.....	612
Cacheable Write-Protect (WP) Memory	612
Cacheable Write-Through (WT) Memory	613
Cacheable Write-Back (WB) Memory	614
The Definition of a Speculatively Executed Load	615
Rules as Defined by MTRRs	616
Memory Type Provided in Memory Transaction	617
Paging Also Defines Memory Type	617
In an MP System, MTRRs Must Be Synchronized.....	618
Posted-Write Related Issues.....	618
General.....	618
Synchronizing Events.....	618
PMWB and WCBs Aren't Snooped	619
WCB Usage	620
An Example.....	620
All WCBs in Use.....	627
Draining the WCBs	628

Contents

Chapter 18: Task Switching

Hardware- vs. Software-Based Task Switching.....	630
A Condensed Conceptual Overview	631
A More Comprehensive Overview	631
The Scheduler and the Task Queue.....	631
Setting Up a Task	632
The Task Data Structure.....	632
The LDT	633
The TSS	633
The Address Translation Tables	634
The GDT and GDTR Register.....	635
The LDTR Register.....	636
The Task Register (TR)	637
Starting a Task	638
Suspend Task and Resume Scheduler Execution.....	639
Hardware-Based Task Switching.....	641
It's Slow	641
Why Didn't OSs Use It?	642
Why Wasn't It Improved?	642
Why Does It Still Exist?	642
Introduction to the Key Elements.....	642
The Trigger Events.....	646
The Descriptors	648
TSS Descriptor	648
Task Gate Descriptor	649
Task Gate Selected by a Far Call/Jump	650
Gate Selected by Hardware Interrupt/Software Exception.....	650
Task Gate Selected by a Software Interrupt Instruction	650
The Task Register	652
General	652
TR Instruction Pair.....	652
STR Instruction	652
LTR Instruction	653
TSS Data Structure Format	654
General	654
Required Fields	655
Optional Fields	655
Register Snapshot Area.....	656
LDT Selector Field	657
Segment Register Fields.....	657
General Purpose Register Fields	657

SS:ESP Register Pair Fields.....	657
Extended Flags (Eflags) Register Field	657
CS:EIP Register Pair Fields.....	658
Control Register 3 (CR3) Field	658
Debug Trap Bit (T)	659
IO Port Access Protection	659
IO Protection in Real Mode	659
Definition of IO Privilege Level (IOPL).....	659
IO Permission Check in Protected Mode	660
IO Permission Check in VM86 Mode	661
IO Permission Bit Map	661
Required or Optional?	661
The Bitmap Offset Field	662
The Permission Check.....	662
Interrupt Redirection Bit Map.....	664
OS-Specific Data Structures.....	664
Privilege Level 0 - 2 Stack Definition Fields	664
Link Field (to Old TSS Selector).....	665
Comprehensive Task Switch Description	665
Calling Another Task	670
An Overview	670
A Comprehensive Example.....	671
LTR Instruction and the Busy Bit	677
When Is Busy Cleared?	677
Critical Error: Switching to a Busy Task.....	677
Busy Toggle Is a Locked Operation	678
Linkage Modification	678
Task Switching and Address Translation.....	678
One GDT to Serve Them All	678
Each Task Can Have Different Virtual-to-Physical Mapping.....	679
TSS Mapping Must Remain the Same for All Tasks.....	679
Placement of a TSS Within a Page(s).....	680
Switch from More-Privileged Code to Lower	680
Software-Based Task Switching	680

Chapter 19: Protected Mode Interrupts and Exceptions

Handler vs. ISR.....	682
Real Mode Interrupt/Exception Handling	683
The IDT	684
General.....	684
Protected Mode IDT and the IDTR.....	685

Contents

The Gates.....	688
Introduction.....	688
Interrupt Gate.....	691
Trap Gate.....	693
Actions Taken When Interrupt or Trap Gate Selected	694
Actions Taken When Task Gate Selected	696
Interrupt/Exception Event Categories	697
General Event Handling.....	699
State Saved on Stack (but which stack?)	704
Return to the Interrupted Program	708
General.....	708
The IRET Instruction	709
Maskable Hardware Interrupts	713
General.....	713
Maskable Interrupt Vector Delivery	713
PC-Compatible Vector Assignment	714
Actions Performed by the Handler	719
Effect of CLI/STI Execution	720
General	720
Other Events That Affect Interrupt Flag Bit.....	722
Protected Mode Virtual Interrupt Feature	723
Non-Maskable Interrupt (NMI) Requests	723
PC-Compatible NMI Logic.....	723
NMI Description	723
More Detailed Coverage of Hardware Interrupt Handling.....	724
Machine Check Exception	724
SMI (System Management Interrupt)	725
Software Interrupts.....	725
INT nn Instruction	726
INTO Instruction.....	726
BOUND Instruction.....	727
INT3 (Breakpoint) Instruction.....	727
Software Exceptions	728
General.....	728
Faults, Traps, and Aborts.....	728
Instruction Restart After a Fault	735
Exception Error Codes	735
Interrupt/Exception Priority.....	739
Detailed Description of Software Exceptions	743
Divide-by-Zero Exception (0).....	743
Processor Introduced In.....	743
Exception Class	743

Description.....	743
Error Code.....	743
Saved Instruction Pointer	743
Processor State.....	743
Debug Exception (1)	744
Processor Introduced In.....	744
Exception Class	744
Description.....	744
Error Code.....	745
Saved Instruction Pointer	745
Processor State.....	745
The Resume Flag Prevents Multiple Debug Exceptions	745
NMI (2)	746
Processor Introduced In.....	746
Exception Class	746
Error Code.....	746
Saved Instruction Pointer	746
Processor State.....	746
Breakpoint Exception (3).....	746
Processor Introduced In.....	746
Exception Class	747
Description.....	747
Error Code.....	747
Saved Instruction Pointer	747
Processor State.....	747
Overflow Exception (4)	748
Processor Introduced In.....	748
Exception Class	748
Description.....	748
Error Code.....	748
Saved Instruction Pointer	748
Processor State.....	748
Array Bounds Check Exception (5)	748
Processor Introduced In.....	748
Exception Class	749
Description.....	749
Error Code.....	749
Saved Instruction Pointer	749
Processor State.....	749
Invalid OpCode Exception (6).....	749
Processor Introduced In.....	749
Exception Class	749

Contents

Description.....	749
Error Code.....	750
Saved Instruction Pointer	750
Processor State.....	751
Device Not Available (DNA) Exception (7)	751
Processor Introduced In.....	751
Exception Class	751
Description.....	751
General	751
X87 FPU Emulation	751
CR0[TS]: Task Switch, But FP/SSE Registers Not Saved	751
CR0[MP].....	752
Error Code.....	752
Saved Instruction Pointer	752
Processor State.....	752
Double Fault Exception (8)	752
Processor Introduced In.....	752
Exception Class	752
Description.....	753
Shutdown Mode.....	755
Error Code.....	756
Saved Instruction Pointer	756
Processor State.....	756
Coprocessor Segment Overrun Exception (9).....	756
Processor Introduced In.....	756
Exception Class	756
Description.....	756
Error Code.....	756
Saved Instruction Pointer	756
Processor State.....	757
Invalid TSS Exception (10).....	757
Processor Introduced In.....	757
Exception Class	757
Description.....	757
Error Code.....	758
Saved Instruction Pointer	759
Processor State.....	759
Segment Not Present Exception (11)	759
Processor Introduced In.....	759
Exception Class	759
Description.....	760
Error Code.....	760

Saved Instruction Pointer	760
Processor State.....	761
Stack Exception (12).....	761
Processor Introduced In.....	761
Exception Class	761
Description.....	761
Error Code.....	762
Saved Instruction Pointer	762
Processor State.....	763
General Protection (GP) Exception (13)	763
Processor Introduced In.....	763
Exception Class	763
Description.....	763
Error Code.....	765
Saved Instruction Pointer	765
Processor State.....	766
Page Fault Exception (14).....	766
Processor Introduced In.....	766
Exception Class	766
Description.....	766
Error Code.....	767
CR2.....	768
Saved Instruction Pointer	768
Processor State.....	768
The More Common Case	768
Page Fault During a Task Switch	769
Page Fault During a Stack Switch	769
Vector (Exception) 15.....	770
FPU Exception (16)	770
Processor Introduced In.....	770
Exception Class	770
Description.....	770
Handling of Masked Errors.....	771
Handling of Unmasked Errors	772
Error Code.....	773
Saved Instruction Pointer	773
Processor State.....	773
Alignment Check Exception (17)	774
Processor Introduced In.....	774
Background: Misaligned Transfers Affect Performance	774
Alignment Is Important!	774
Exception Class	775

Contents

Description.....	775
Implicit Privilege Level 0 Accesses	776
Storing GDTR, LDTR, IDTR or TR	776
FP/MMX/SSE Save and Restore Accesses	777
MOVUPS and MOVUPD Accesses	777
FSAVE and FRSTOR Accesses	777
Error Code.....	777
Saved Instruction Pointer	777
Processor State.....	778
Machine Check Exception (18)	778
Processor Introduced In.....	778
Exception Class	778
Description.....	778
Error Code.....	778
Saved Instruction Pointer	779
Processor State.....	779
SIMD Floating-Point Exception (19).....	779
Processor Introduced In.....	779
Exception Class	779
Description.....	779
Exception Error Code	782
Saved Instruction Pointer	782
Processor State.....	782
Legacy Problem: 2-Step SS:ESP Update	782
Problem Description	782
The Solution	782

Chapter 20: Virtual 8086 Mode

A Special Note	784
Real Mode Applications Are Dangerous.....	784
Solution: a Watchdog	785
Real Mode Applications Run at Privilege Level 3.....	787
Switching Between Protected Mode and VM86 Mode	787
Eflags[VM] = 1 Switches Processor into VM86 Mode	787
But Software Cannot Directly Access Eflags[VM]	788
Scheduler Activates VM86 Mode	788
Exiting VM86 Mode.....	789
Determining Interrupted Task Is a Real Mode Task.....	789
Returning to VM86 Mode from VMM.....	790
VMM Passes Control to Real Mode Interrupt/Exception Handler.....	790
Real Mode Application's World View	790
The DOS World	790

Memory Address Formation in VM86 Mode	792
Multiple DOS Domains in Separate 1MB Areas.....	793
VMM Should Not Reside in the HMA.....	795
Dealing with Segment Wraparound	796
8088/8086 Processor	796
286 and Later Processors.....	796
Solutions.....	796
Using the Address Size Override Prefix.....	797
Sensitive Instructions.....	797
Problematic Instructions	797
CLI (Clear Interrupt Enable) Instruction.....	798
STI (Set Interrupt Enable) Instruction.....	798
PUSHF (Push Flags) Instruction.....	799
POPF (Pop Flags) Instruction.....	799
INT nn (Software Interrupt) Instruction	799
IRET (Interrupt Return) Instruction.....	799
Solution: IOPL Sensitive Instructions	800
Handling Direct IO	800
The Problem.....	800
IO-Mapped IO	800
IO Permission in Protected Mode	800
IO Permission in VM86 Mode.....	801
Memory-Mapped IO	802
To Permit an Access	803
To Deny an Access.....	803
For Finer Control	803
Handling Video Frame Buffer Updates.....	803
Handling Exceptions in VM86 Mode.....	804
Processor Actions.....	804
Option 1: Protected Mode Handler Services Exception	806
Option 2: Handler Passes Exception to VMM for Servicing.....	806
Option 3: Exception Handled by Another Task	808
Hardware Interrupt Handling in VM86 Mode.....	809
NMI, SMI, and Maskable Interrupts	809
Real Mode Application's Unreal Reality	811
VM86 Task Executes CLI When VME = 0	812
CLI Handling.....	812
Subsequent High-Priority Interrupt Detected.....	813
Servicing of Lower-Priority Interrupt Deferred.....	813
STI/POPF/PUSHF/IRET Handling When VME = 0.....	816
Attempted Execution of STI Instruction (VME = 0)	816
Attempted Execution of PUSHF Instruction (VME = 0)	816

Contents

Attempted Execution of POPF Instruction (VME = 0)	817
Attempted Execution of IRET Instruction (VME = 0)	817
CLI/STI/POPF/PUSHF Handling When VME = 1	818
VM86 Extensions	818
Background.....	819
When VME = 1 and IOPL = 3, Task Can Control Eflags[IF]	819
When VME = 1 and IOPL < 3, Task Controls VIF, Not IF	820
Eflags[VIP] Is Controlled by the VMM.....	820
Software Cannot Directly Access Eflags[VIP]	820
CLI Followed by a Maskable Interrupt	820
Subsequent STI Effect Depends on Eflags[VIP]	822
A Special Case	825
POPF/PUSHF Handling	825
Software Interrupt Instruction Handling.....	825
Software Interrupt Handling in Protected Mode	825
Software Interrupt Handling in VM86 Mode	826
INT3 Is Special.....	826
VMM Passes Control To Real Mode Handler.....	829
Halt Instruction in VM86 Mode	832
Protected Mode Virtual Interrupt Feature.....	832
General.....	832
1. Task executes CLI, Clears VIF	833
2. Maskable Interrupt Occurs and Is deferred	833
3. Task Executes STI.....	833
Registers Accessible in Real/VM86 Mode.....	834
Instructions Usable in Real/VM86 Mode	834

Chapter 21: The MMX Facilities

Introduction.....	836
Detecting MMX Capability	837
The Basic Problem	837
Assumptions	837
The Operation.....	837
Example: Processing One Pixel Per Iteration	838
Example: Processing Four Pixels Per Iteration	838
MMX SIMD Solution.....	840
Dealing with Unpacked Data	841
Dealing with Math Underflows and Overflows.....	842
Elimination of Conditional Branches	843
Introduction	843
Non-MMX Chroma-Key/Blue Screen Compositing Example.....	844
MMX Chroma-Keying/Blue Screen Compositing Example	845

Changes To the Programming Environment	847
Handling a Task Switch.....	848
MMX Instruction Set Syntax.....	848

Chapter 22: The SSE Facilities

Chapter Objectives	852
SSE: MMX on Steroids.....	852
Streaming SIMD Extensions (SSE).....	857
The Motivation Behind SSE.....	857
Detecting SSE Support	858
The SSE Elements.....	858
SSE Data Types.....	859
The MXCSR.....	860
MXCSR Description.....	860
Loading and Storing the MXCSR	863
SIMD (Packed) Operations.....	863
Scalar Operations	864
Cache-Related Instructions.....	864
Overlapping Data Prefetch with Program Execution	865
Streaming Store Instructions.....	868
Introduction.....	868
The MOVNTPS Instruction.....	871
MOVNTQ Instruction.....	871
MASKMOVQ Instruction.....	872
Ensuring Delivery of Writes Before Proceeding	873
An Example Scenario	873
SFENCE Instruction	874
Elimination of Mispredicted Branches.....	877
Background.....	877
SSE Misprediction Enhancements.....	877
Comparisons and Bit Masks	877
Min/Max Determination.....	878
The Masked Move Operation	878
Reciprocal and Reciprocal Square Root Operations	878
MPEG-2 Motion Compensation.....	879
Optimizing 3D Rasterization Performance	880
Optimizing Motion-Estimation Performance	880
Accuracy vs. Fast Real-Time 3D Processing (FTZ)	880
SSE Alignment Checking.....	881
The SIMD FP Exception	881
Saving and Restoring x87/MMX/SSE Registers.....	881
General	881

Contents

MXCSR Mask Field.....	882
OS Support for SSE	883
General	883
Enable SSE Instruction Sets and Register Set Save/Restore.....	884
Enable the SSE SIMD FP Exception	884
SSE Setup.....	885
Summary of the SSE Instruction Set.....	885
The SSE2 Instruction Set	887
General.....	887
DP FP Number Representation.....	887
SSE2 Packed and Scalar DP FP Instructions	888
SSE2 64-Bit and 128-Bit SIMD Integer Instructions.....	888
SSE2 128-Bit SIMD Integer Instruction Extensions	889
Your Choice: Accuracy or Speed (DAZ).....	889
The Cache Line Flush Instruction.....	890
Fence Instructions	891
MFENCE Instruction.....	891
LFENCE Instruction	893
General	893
LFENCE Ordering Rules	894
SFENCE Instruction	894
Non-Temporal Store Instructions.....	894
Introduction	894
MOVNTDQ Instruction	894
MOVNTPD Instruction.....	895
MOVNTI Instruction.....	896
MASKMOVDQU Instruction.....	897
General	897
When a Mask of All Zeros Is Used.....	898
PAUSE Instruction.....	898
Thread Synchronization.....	898
The Problem.....	899
The Fix	900
When a Thread Is Idle	901
Spin-Lock Optimization.....	901
Branch Hints	901
SSE3 Instruction Set	902
Introduction	902
Improved x87 FP-to-Integer Conversion Instruction	902
The Problem.....	902
The Solution.....	903
New Complex Arithmetic Instructions.....	903

Improved Motion Estimation Performance	904
The Problem.....	904
The Solution.....	905
The Downside	905
Instructions to Improve Processing of a Vertex Database	906
MONITOR/MWAIT Instruction Pair	907
Background.....	907
Monitor Instruction	908
Mwait Instruction	908
Example Code Usage	909
The Wake Up Call.....	909
SSSE3, SSE 4.1, and 4.2.....	910

Part 3: IA-32e OS Kernel Environment

Chapter 23: IA-32e OS Environment

The Big Picture	914
Mode Switching Overview	916
Booting Into Protected Mode	916
Initial Switch from IA-32 to IA-32e Mode	917
CS D and L Bits Control IA-32e SubMode Selection	920
Old and New Applications Running Under a 64-bit OS	922
Things You Lose In IA-32e Mode (hint: not much).....	923
Old Applications Live in an Expanded Universe	923
Old Legacy Universe = 4GB or 64GB.....	923
IA-32e Universe Is At Least 16 Times Larger	924
Virtual Memory Addressing in IA-32e Mode	925
Virtual Address in Compatibility Mode.....	925
Virtual Address in 64-bit Mode	926
In Compatibility Mode, Segmentation Is Operative	926
In 64-bit Mode, Hardware-Enforced Flat Model.....	927
General.....	927
New Segment Selector Causes Descriptor Read	927
Segment Register Usage in 64-bit Mode	927
64-bit Mode: No Limit Checking = No Limits?	934
Table Limit Checks Are Performed.....	935
Stack Management.....	935
Stack Management in Compatibility Mode	935
Stack Management in 64-bit Mode.....	936
Push/Pop Size is 64-bits	936
Address Translation Replaces Limit Checking	936

Contents

Segment Override Prefixes Other Than FS/GS Are Ignored	937
Protection Provided by Paging	938
Segment Registers Preserved On Mode Switch	938
64-bit Instruction Pointer.....	938
Instruction Fetching.....	938
RIP-Relative Data Accesses	939
Changes To Kernel-Related Registers and Structures.....	939
Address Translation Mechanism.....	939
Basic Description.....	939
Top-Level Directory Placement	940
Detailed Description.....	940
GDT/LDT Descriptor Changes	940
GDT and GDTR Changes.....	947
GDT Descriptor Types	947
Executing LGDT in 64-bit Mode.....	950
Unaligned Accesses to GDT or LDT	951
LDT and LDTR Changes.....	952
LDT Descriptor Types.....	952
LDTR Contents in IA-32e Mode	953
Unaligned Accesses to LDT.....	953
IDT/IDTR and Interrupt/Exception Changes	955
IDT Descriptor Types.....	955
Interrupt/Trap Gate Operational Changes	958
General	958
Interrupt/Exception Stack Switch	959
Motivation for the IST.....	960
IRET Behavior	960
Executing LIDT in 64-bit Mode	963
All Accesses to IDT Are Properly Aligned	964
IA-32e Call Gate Operation	964
General	964
IA-32e Call Gate Detailed Operation.....	965
IA-32e Call Gate Stack Switch.....	966
TR and TSS Changes.....	968
Real World TSS Usage.....	968
Illegal For Jump or Call To Select a TSS Descriptor.....	973
Executing LTR in Compatibility Mode.....	973
Executing LTR in 64-bit Mode	974
Revised TSS Structure	974
TSS Usage.....	976
General	976
Call Gate Stack Switch	976

Interrupt/Exception Stack Switch	976
Register Set Expansion (in 64-bit Mode)	976
Scheduler's Software-Based Task Switching Mechanism.....	977
Switching to a 64-bit Task.....	977
Switching to a Legacy Task	979
General	979
Data Segment Register Initialization	980
CS and Instruction Pointer Initialization.....	980
The Switch.....	981

Chapter 24: IA-32e Address Translation

Theoretical Address Space Size	984
Limitation Imposed by Current Implementations	985
Four-Level Lookup Mechanism	985
Address Space Partitioning	985
The Address Translation.....	988
Initializing CR3	988
Step 1: PML4 Lookup.....	988
Step 2: PDPT Lookup	990
Step 3: Page Directory Lookup	992
Step 3a: PDE Points to a Page Table.....	994
Step 3b: PDE Points to a 2MB Physical Page.....	994
Step 4: Page Table Lookup	997
Page Protection Mechanisms in IA-32e Mode	999
Page Protection in Compatibility Mode	999
Page Protection in 64-bit Mode.....	999
Don't Forget the Execute Disable Feature!.....	1000
TLBs Are More Important Than Ever	1004
No 4MB Page Support.....	1005

Part 4: Compatibility Mode

Chapter 25: Compatibility Mode

Initial Entry to Compatibility Mode	1010
Switching Between Compatibility Mode and 64-bit Mode.....	1010
Differences Between IA-32 Mode and Compatibility Mode.....	1011
IA-32 Background	1011
Unsupported IA-32 Features.....	1011
Changes to the OS Environment.....	1011
Memory Addressing.....	1013
Segmentation	1013

Contents

FS/GS Segments.....	1014
Virtual Address	1014
Address Translation	1014
Register Set.....	1015
Visible Registers	1015
No Access to Additional or Extended Registers	1016
Control Register Accesses.....	1016
Debug Register Accesses.....	1016
Register Preservation Across Mode Switches.....	1016
Exception and Interrupt Handling.....	1017
OS Kernel Calls	1017
Call Gates	1017
Kernel Call Instruction Usage	1018
SysEnter Instruction	1018
SysCall Instruction.....	1018
Odds and Ends.....	1019
IRET Changes	1019
Segment Load Instructions.....	1019

Part 5: 64-bit Mode

Chapter 26: 64-bit Register Overview

Overview of 64-bit Register Set.....	1024
EFER (Extended Features Enable) Register.....	1025
Sixteen 64-bit Control Registers.....	1027
64-bit Rflags Register	1033
Sixteen 64-bit GPRs	1033
Kernel Data Structure Registers in 64-bit Mode	1036
SSE Register Set Expanded in 64-bit Mode	1037
Debug Breakpoint Registers.....	1038
Local APIC Register Set	1039
x87 FPU/MMX Register Set	1039
Architecturally-Defined MSRs.....	1039

Chapter 27: 64-bit Operands and Addressing

Helpful Background	1042
Switching to 64-bit Mode	1042
The Defaults.....	1042
The REX Prefix.....	1043
Problem 1: Addressing New Registers	1043
Problem 2: Using 16- and 64-bit Operands	1045

Solution: The Rex Prefix.....	1045
Making Room for REX	1046
REX Prefix Placement.....	1046
When You Need REX... ..	1047
...and when you don't.....	1048
Anatomy of a REX Prefix.....	1049
General	1049
The Width Bit	1052
The Register Bit	1053
Description	1053
An Example	1053
The Index and Base Bits	1055
Description	1055
An Example	1056
Addressing Registers Using REX[B] + Opcode[Reg]	1058
Addressing Registers Using REX[B] + ModRM[RM].....	1058
Byte-Register Addressing Limitations.....	1058
Sometimes, REX Fields Have No Effect.....	1059
Addressing Memory in 64-bit Mode	1059
64-bit Mode Uses a Hardware-Enforced Flat Model.....	1059
CS, DS, ES, and SS Segments Start at Virtual Address 0.....	1059
CS/DS/ES/SS Segment Override Prefixes Ignored.....	1060
FS and GS Segments Can Start at Non-Zero Base Addresses	1060
FS/GS Segment Override Prefixes Matter	1060
Default Virtual Address Size (and overriding it).....	1061
Actual Address Size Support: Theory vs. Practice.....	1062
Canonical Address.....	1063
General	1063
32- (and 16-) bit Addressing Limited to Lower 4GB	1064
32-bit Address Treatment in 64-bit Mode	1064
Address Treatment in Compatibility Mode	1064
Memory-based Operand Address Computation	1065
RIP-relative Data Addressing	1069
Near and Far Branch Addressing.....	1070
Immediate Data Values in 64-bit Mode.....	1073
Displacements in 64-bit Mode.....	1074

Chapter 28: 64-bit Odds and Ends

New Instructions	1076
General.....	1076
SwapGS Instruction	1076
The Problem.....	1076

Contents

The SwapGS Solution.....	1077
MOVSXD Instruction: Stretch It Out.....	1078
Enhanced Instructions.....	1078
Invalid Instructions	1079
Reassigned Instructions.....	1081
LAHF/SAHF Instruction Support	1081
Instructions That Default to a 64-bit Operand Size	1082
Stack Operations	1082
Near Branches.....	1083
Branching in 64-bit Mode	1083
Short/Near Branches Default to 64-bit Operand Size.....	1083
Unconditional Jumps in 64-bit Mode.....	1084
Calls/Ret/Iret in 64-bit Mode.....	1087
Instruction Forms in 64-bit Mode.....	1087
Example Call/Return Operations	1089
64-bit Near Call/Return	1089
32-bit Level 3 Code Calls 64-bit Level 3 Procedure	1089
32-bit Level 3 Code Calls 64-bit Level 2 Procedure	1091
Previous Example Plus Call to Level 0 Procedure.....	1093
Conditional Branches in 64-bit Mode	1093
NOP Instruction	1097
FXSAVE/FXRSTOR	1097
General.....	1097
Fast FxSave/Restore Feature (AMD-only).....	1097
The Nested Task Bit (Rflags[NT])	1101
SMM Save Area.....	1102
IA-32 Processor SM Save Area.....	1102
Intel 64 Processor SM Save Area.....	1106

Part 6: Mode Switching Detail

Chapter 29: Transitioning to Protected Mode

Real Mode Peculiarities That Affect the OS Boot Process.....	1114
Example OS Characteristics	1114
Flat Model With Paging	1115
Software-Based Task Switching.....	1115
Protected Mode Transition Primer	1116
GDT Must Be In Place Before Switch to Protected Mode	1116
No Interrupts or Exceptions During Mode Switch.....	1119
Creation of Protected Mode IDT.....	1120

Other Protected Mode Structures	1121
TSS	1121
Address Translation Mechanism	1122
Protected Mode Is a Prerequisite	1122
Identity Mapping	1122
Which Translation Mechanism?	1123
Optional Structure: LDT	1124
Enable A20 Gate	1125
Load Initial Code and Handlers Into Memory	1125
The Switch to Protected Mode	1125
Loading Segment Registers With GDT Descriptors	1125
Load TSS Descriptor Into TR	1127
Enable Interrupts	1127
Load Application Into Memory	1128
Create Task's Address Translation Tables	1128
Switching From OS Scheduler to First Task	1128
Example: Linux Startup	1128
1. Bootsect	1129
2. Setup	1129
3a. Startup_32 in boot/compressed/head.s	1131
3b. Startup_32 in kernel/head.s	1132

Chapter 30: Transitioning to IA-32e Mode

No Need to Linger in Protected Mode	1140
Entering Compatibility Mode	1140
Switch to 64-bit Mode	1142

Part 7: Other Topics

Chapter 31: Introduction to Virtualization Technology

Just an Introduction?	1148
Detailed Coverage of Virtualization	1148
The Intel Model	1149
OS: I Am the God of All Things!	1149
Virtualization Supervisor: Sure You Are (:<)	1150
Root versus Non-Root Mode	1150
Detecting VMX Capability	1151
Entering/Exiting VMX Mode	1152
Entering VMX Mode	1152
Exiting VMX Mode	1152
Virtualization Elements/Terminology	1153

Contents

Introduction to the VT Instructions	1154
Introduction to the VMCS Data Structure	1156
Preparing to Launch a Guest OS	1160
Launching a Guest OS	1161
Guest OS Suspension.....	1162
Handling Timeslice Expiration	1163
Handling a Sensitive Operation or a VMCALL	1163
Resuming a Guest OS	1163
Some Warnings Regarding VMCS Accesses	1165

Chapter 32: System Management Mode (SMM)

What Falls Under the Heading of System Management?	1168
The Genesis of SMM.....	1169
SMM Has Its Own Private Memory Space	1170
The Basic Elements of SMM.....	1170
A Very Simple Example Scenario	1171
How the Processor Knows the SM Memory Start Address	1171
Normal Operation, (Including Paging) Is Disabled.....	1172
The Organization of SM RAM	1172
General.....	1172
IA-32 Processor SM State Save Area	1173
Intel 64 Processor SM Save Area.....	1178
Protecting Access to SM Memory	1182
Entering SMM	1183
The SMI Interrupt Is Generated	1183
No Interruptions Please	1183
General	1183
Exceptions and Software Interrupts Permitted but Not Recommended.....	1184
Servicing Maskable Interrupts While in the Handler.....	1184
Single-Stepping through the SM Handler.....	1184
If Interrupts/Exceptions Permitted, Build an IDT.....	1185
SMM Uses Real Mode Address Formation.....	1185
NMI Handling While in SMM	1186
Default NMI Handling	1186
How to Re-Enable NMI Recognition in the SM Handler	1186
If an SMI Occurs within the NMI Handler	1187
Informing the Chipset SM Mode Has Been Entered	1188
General	1188
A Note Concerning Memory-Mapped IO Ports.....	1188
The Context Save.....	1188
General	1188
Although Saved, Some Register Images Are Forbidden Territory	1189

Contents

Special Actions Required on a Request for Power Down.....	1189
The Register Settings on Initiation of the SM Handler	1190
The SMM Revision ID	1191
The Body of the Handler.....	1192
Exiting SMM	1192
The Resume Instruction	1192
Informing the Chipset That SMM Has Been Exited.....	1193
The Auto Halt Restart Feature	1193
Executing the HLT Instruction in the SM Handler	1194
The IO Instruction Restart Feature	1195
Introduction.....	1195
An Example Scenario	1195
The Detail.....	1195
Back-to-Back SMIs During IO Instruction Restart	1196
Multiprocessor System Presents a Problem.....	1196
Caching from SM Memory.....	1198
Background.....	1198
The Physical Mapping of SM RAM Accesses	1199
FLUSH# and SMI#	1203
Setting Up the SMI Handler in SM Memory	1203
Relocating the SM RAM Base Address	1204
Description.....	1204
In an MP System, Each Processor Must Have a Separate State Save Area.....	1204
Accessing SM Memory Above the First MB	1205
SMM in an MP System	1205
SM Mode and Virtualization.....	1205

Chapter 33: Machine Check Architecture (MCA)

Why This Subject Is Included	1209
MCA = Hardware Error Logging Capability	1209
The MCA Elements.....	1210
The Machine Check Exception.....	1210
The MCA Register Set	1211
The Global Registers	1212
Introduction	1212
The Global Count and Present Register.....	1212
The Global Status Register.....	1213
The Global Control Register	1214
The Extended MC State MSRs	1214
The Composition of a Register Bank	1217
Overview	1217

Contents

The Bank Control Register.....	1217
General	1217
P6 and Core Processors.....	1218
The Bank Status Register.....	1218
General	1218
Error Valid Bit	1220
Overflow Bit	1220
Uncorrectable Error Bit	1221
Error Enabled Bit	1221
Miscellaneous Register Valid Bit	1221
Address Register Valid Bit	1221
Processor Context Corrupt Bit	1222
MCA Error Code and Model Specific Error Code	1222
Other Information.....	1222
The Bank Address Register	1222
The Bank Miscellaneous Register	1222
Control 2 Register	1222
The Error Code.....	1223
The Error Code Fields	1223
Simple MCA Error Codes.....	1223
Compound MCA Error Codes.....	1224
General	1224
Correction Report Filtering Bit	1225
Example External Interface Error Interpretation.....	1229
Cache Error Reporting.....	1232
Green/Yellow Cache Health Indicator.....	1232
Background.....	1232
TES (Threshold Error Status) Feature	1232
Interrupt On Soft Error Threshold Match	1233
Before CMCI, Soft Error Logging Required Periodic Scan	1233
CMCI Eliminates MC Register Scan	1234
Determining Processor's CMCI Support.....	1235
Determining a Bank's CMCI Support.....	1235
CMCI Interrupt Is Separate and Distinct From MC Exception.....	1235
CMC Interrupt May Affect Multiple Cores/Logical Processors	1235
CMC Interrupt Should Only Be Serviced Once	1235
MC Exception Is Generally Not Recoverable.....	1236
Machine Check and BINIT#	1237
Additional Error Logging Notes.....	1237
Error Buffering Capability	1237
Additional Information for Each Log Entry.....	1237

Chapter 34: The Local and IO APICs

APIC and the IA-32 Architecture	1240
Definition of IO and Local APICs	1241
Hardware Context Is Essential	1241
A Short History of the APIC's Evolution	1241
APIC Introduction	1241
Pentium Pro APIC Enhancements.....	1241
The Pentium II and Pentium III	1242
Pentium 4 APIC Enhancements: xAPIC	1242
The x2APIC Architecture.....	1243
Before the APIC	1244
MP Systems Need a Better Interrupt Distribution Mechanism.....	1246
Legacy Interrupt Delivery System Is Inefficient.....	1246
The APIC Interrupt Distribution Mechanism.....	1248
Introduction.....	1248
Message Types	1249
Inter-Processor Interrupt (IPI) Messages	1249
NMI, SMI and Init Messages.....	1249
Legacy Interrupt Message	1249
Message Transfer Mechanism Prior to the Pentium 4.....	1250
Message Transfer Mechanism Starting with the Pentium 4.....	1250
Message Transfer Mechanism in Intel QPI-based Systems	1251
Processors Reside in Clusters.....	1253
Each Core/Logical Processor Has a Dedicated Local APIC	1254
Introduction to the Message Addressing Modes	1254
Detecting Presence/Version/Capabilities of Local APIC	1255
Presence	1255
Version.....	1255
x2APIC Capability Verification	1256
Local APIC's Initial State.....	1256
Enabling/Disabling the Local APIC	1257
General.....	1257
Disabling Local APIC for Remainder of Power-Up Session.....	1257
Dynamically Enabling/Disabling Local APIC	1258
Mode Selection	1260
The Local APIC Register Set.....	1261
Register Access in xAPIC Mode: MMIO	1261
General	1261
Local and IO APIC xAPIC Register Areas Are Uncacheable	1261
xAPIC Register Access Alignment.....	1261
Register Access in x2APIC Mode: MSR.....	1262

Contents

Introduction to the Local APIC's Register Set	1262
Local APIC ID Assignments and Addressing	1277
ID Assignment in xAPIC Mode	1277
Introduction	1277
Cluster ID Assignment.....	1277
Physical/Logical Processor and Local APIC ID Assignment	1278
Example Xeon MP System: Hyper-Threading Disabled.....	1278
Example Xeon MP System: Hyper-Threading Enabled.....	1279
Dual Processor System: Hyper-Threading Enabled.....	1279
A Single-Processor System: Hyper-Threading Enabled	1280
xAPIC ID Register.....	1281
BIOS/OS Reassignment of xAPIC ID	1282
Logical xAPIC Address Assignment	1282
Maximum Number of xAPICs.....	1282
ID Assignment in x2APIC Mode	1283
Two Hardware-Assigned Local APIC IDs.....	1283
x2APIC ID (Physical Local APIC ID).....	1283
General	1283
Some Interesting Questions	1284
CUID Provides the Answers	1284
x2APIC ID and Physical Destination Mode.....	1285
Obtaining the x2APIC ID	1285
Logical x2APIC ID	1285
xAPIC Logical Addressing Background	1285
Logical x2APIC ID Formation	1286
Logical x2APIC ID Usage.....	1287
Local APIC Addressing.....	1287
Physical Addressing: Single Target.....	1287
Logical Addressing: Multiple Targets	1288
Introduction.....	1288
x2APIC Cluster Model.....	1288
Flat Model.....	1289
Flat Cluster Model.....	1290
Hierarchical Cluster Model.....	1291
Message Addressing Summary	1292
Lowest-Priority Delivery Mode	1294
General	1295
Warnings Related to Lowest-Priority Delivery Mode.....	1295
Chipset-Assisted Lowest-Priority Delivery	1296
Local APIC IDs Are Stored in the MP and ACPI Tables	1297
Accessing the Local APIC ID.....	1298

An Introduction to the Interrupt Sources.....	1298
Local Interrupts	1298
Remote Interrupt Sources	1299
Introduction to Interrupt Priority	1300
General.....	1300
Definition of a User-Defined Interrupt.....	1301
User-Defined Interrupt Priority	1302
Definition of Fixed Interrupts	1305
Masking User-Defined Interrupts	1305
Task and Processor Priority.....	1305
Introduction	1305
The Task Priority Register (TPR)	1306
The Processor Priority Register (PPR).....	1306
The User-Defined Interrupt Eligibility Test	1307
CR8 (Alternative TPR).....	1308
Interrupt Message Format	1309
IO/Local APICs Cooperate on Interrupt Handling.....	1313
The Purpose of the IO APIC.....	1313
Overview of Edge-Triggered Interrupt Handling	1316
Assumptions.....	1316
Description.....	1316
Overview of Level-Sensitive Interrupt Handling	1321
Assumptions.....	1321
Description.....	1322
Higher-Priority Fixed Interrupt Preempts Handler	1326
IO APIC Register Set	1331
IO APIC Register Set Base Address	1331
IO APIC Register Set Description.....	1332
IRQ Pin Assertion Register.....	1335
IO APIC EOI Register and Shared Interrupts.....	1336
Non-Shareable IRQ Lines.....	1336
Shareable IRQ Lines	1336
Linked List of Interrupt Handlers.....	1336
How It Works.....	1337
Broadcast Versus Directed EOI	1338
IO APIC ID Register	1340
IO APIC Version Register.....	1340
IO APIC Redirection Table (RT) Register Set	1341
IO APIC Interrupt Delivery Order Is Rotational.....	1344
Message Signaled Interrupts (MSI).....	1345
General.....	1345
Using the IO APIC as a Surrogate Message Sender.....	1346

Contents

Direct-Delivery of an MSI.....	1346
Memory Already Sync'd When Interrupt Handler Entered	1347
The Problem.....	1347
Old Solution.....	1347
How MSI Solves the Problem	1348
Interrupt Delivery from Legacy 8259a Interrupt Controller.....	1348
Virtual Wire Mode A.....	1348
Virtual Wire Mode B.....	1350
SW-Initiated Interrupt Message Transmission.....	1351
Introduction	1351
Sending a Message From the Local APIC	1353
ICR in xAPIC Mode.....	1353
ICR in x2APIC Mode.....	1353
x2APIC Mode's Self IPI Feature.....	1359
Locally Generated Interrupts.....	1360
Introduction	1360
The Local Vector Table.....	1360
The Pentium Family's LVT.....	1361
The P6 Family's LVT	1361
The Pentium 4 Family's LVT.....	1361
Core Processor's LVT	1361
LVT Register State After Reset, INIT, or Software Disable.....	1362
Local Interrupt 0 (LINT0).....	1362
The Mask Bit.....	1362
The Trigger Mode and the Input Pin Polarity	1362
The Delivery Mode	1363
The Vector Field	1364
The Remote IRR Bit	1365
The Delivery Status	1365
Local Interrupt 1 (LINT1).....	1365
The Local APIC Timer.....	1366
General	1366
The Divide Configuration Register	1367
One Shot Mode.....	1367
Periodic Mode	1367
The Performance Counter Overflow Interrupt.....	1368
The Thermal Sensor Interrupt.....	1370
Correctable Machine Check (CMC) Interrupt	1372
The Local APIC's Error Interrupt	1373
Local APIC Error LVT Register	1373
Error Status Register (ESR) Operation in xAPIC Mode	1373
Error Status Register Operation in x2APIC Mode	1373

Contents

The Spurious Interrupt Vector	1376
The Problem.....	1376
Solution.....	1376
Additional Spurious Vector Register Features	1377
Boot Strap Processor (BSP) Selection	1378
Introduction	1378
The BSP Selection Process.....	1379
Pre-QPI BSP Selection Process	1379
Intel QPI BSP Selection Process	1381
How the APs are Discovered and Configured	1381
AP Detection and Configuration	1382
Introduction.....	1382
BIOS AP Discovery Procedure.....	1382
Uni-Processor OS and the APs.....	1384
MP OS and the APs	1384
The FindAndInitAllCPUs Routine	1387
Glossary.....	1391
Index	1469

Contents

1-1	Processor, Core, Logical Processor	15
2-1	Execution Mode Diagram.....	22
2-2	Switching to IA-32e Mode	31
2-3	16-bit, 286-style CS Descriptor Format	35
2-4	32-bit Code Segment Descriptor Format	36
2-5	64-bit Code Segment Descriptor	37
2-6	Protected Mode and Compatibility Mode Consists of Two SubModes	39
3-1	Major Milestones in Evolution of Software Environment	42
3-2	16-bit, 286-style Code Segment Descriptor	51
3-3	32-bit, 386-style Code Segment Descriptor	55
4-1	IA-32 Register Set.....	65
5-1	8086 Register Set	80
5-2	8086 GPRs	81
5-3	8086 Flag Register	81
5-4	286 Register Set	82
5-5	286 Machine Status Word Register (MSW)	83
5-6	286 Flags Register	83
5-7	IA-32 Register Set.....	85
5-8	CR0.....	86
5-9	CR3.....	87
5-10	CR4.....	87
5-11	32-bit Eflags Register	88
5-12	IA-32 GPRs.....	91
5-13	MMX SIMD Solution Increases Throughput	103
7-1	General Instruction Format	163
7-2	8086 Opcode Map	169
7-3	Format of Instructions with Single Opcode Byte	171
7-4	Reg Select Field in Primary Opcode Byte.....	172
7-5	Instructions With 2 Opcode Bytes Use 2-level Lookup.....	174
7-6	Format of Instructions With 2 Opcode Bytes.....	175
7-7	Instructions With 3 Opcode Bytes Use 3-Level Lookup	178
7-8	Format of Instructions With 3 Opcode Bytes.....	179
7-9	The ModRM Byte.....	181
7-10	Micro-Maps (i.e., Groups) Associated with 1-byte Opcodes.....	182
7-11	Micro-Maps (i.e., Groups) Associated with 2-byte Opcodes.....	183
7-12	x87 FP Instructions Inhabit Opcode Mini-Maps	188
7-13	The Primary Opcode Byte	192
7-14	The Width, Sign-Extension and Direction Bits.....	192
7-15	Reg Select Field in Primary Opcode Byte.....	193
7-16	SIB Byte Usage.....	204
7-17	The Scale/Index/Base (SIB) Byte	205
8-1	Task/OS Relationship.....	235

Figures

8-2	Real Mode Register Set	237
8-3	Control Register 0 (CR0) in IA-32 Mode	238
8-4	Control Register 2 (CR2) in IA-32 Mode	243
8-5	Control Register 3 (CR3) in IA-32 Mode	244
8-6	Control Register 4 (CR4) in IA-32 Mode	244
8-7	XCR0 (also referred to XFEM).....	250
8-8	Eflags Register in IA-32 Mode	251
8-9	Stack Usage in C Function Call.....	258
8-10	General Purpose Registers (GPRs) in IA-32 Mode	259
8-11	Instruction Pointer Register	260
8-12	x87 FPU and MMX Registers	261
8-13	SSE Register Set in IA-32 Mode	262
8-14	Debug Register Set (available in all modes)	269
8-15	Local APIC Register Set	271
8-16	Segment Registers.....	289
8-17	Example Data Segment Access	289
8-18	IP and EIP Registers	294
8-19	Example Code Fetch in Real Mode	295
8-20	Push Operation in Real Mode.....	299
8-21	Stack Segment.....	299
8-22	Pop Operation in Real Mode.....	300
8-23	Example Data Segment Access	304
8-24	Example Usage of Segment Registers in Real Mode	306
8-25	A20 Gate	309
8-26	Real Mode Interrupt Table	318
8-27	Real Mode Event Handling	326
8-28	Return from Real Mode Handler to Interrupted Real Mode Application	327
8-29	Exception Handling in Real Mode	329
8-30	Real Mode Stack on Entry to Handler	329
8-31	Maskable Interrupt Handling in Real Mode.....	336
9-1	CR0.....	341
9-2	The x87 FPU Register Set.....	343
9-3	The Double Extended Precision (DEP) FP Numeric Format.....	345
9-4	64-bit DP FP Numeric Format.....	345
9-5	32-bit SP FP Numeric Format.....	345
9-6	The FPU's FCW Register	352
9-7	The FPU's FSW Register	354
9-8	The FPU's FTW Register	355
9-9	The x87 FPU's Opcode Register.....	357
9-10	IBM PC-AT FP Error Reporting Mechanism	359
9-11	Ignoring FP Errors	360
11-1	Segment Descriptor Selection	373

13-1	Segment Register Contents in Real Mode	385
13-2	Relationship of a Segment Register and GDT, GDTR, LDT, and LDTR.....	388
13-3	Segment Register's Visible and Invisible Elements	389
13-4	The Global Descriptor Table (GDT)	394
13-5	The GDT and the LDTs	395
13-6	LDT Structure	397
13-7	Format of an LDT Descriptor (must be in the GDT).....	398
13-8	Local Descriptor Table Register (LDTR)	398
13-9	General Format of a Segment Descriptor	399
13-10	16-bit, 286-Style Code Segment Descriptor	401
13-11	32-bit Code Segment Descriptor	402
13-12	32-bit Data Segment Descriptor Format	406
13-13	The Flat Memory Model	411
13-14	Creating a Flat Memory Model.....	411
14-1	Segment Register	418
14-2	32-bit Code Segment Descriptor Format	421
14-3	16-bit, 286-style CS Descriptor Format	422
14-4	Sample Code Segment Descriptor.....	426
14-5	Example Value in CS Register	426
14-6	Privilege Check on Far Call or Jmp.....	441
14-7	Example Scenario.....	453
14-8	32-bit Call Gate Descriptor Format	456
14-9	16-bit Segment Selector in Far Call Selects LDT Entry 12.....	457
14-10	Example Call Gate Descriptor.....	458
14-11	Call Gate in LDT Entry 12 Contains This CS Selector	459
14-12	Descriptor for 32-bit CS Containing Called Procedure	460
14-13	Bird's Eye View of Example Far Call through a Call Gate	461
14-14	Task State Segment (TSS) Format.....	464
14-15	Automatic Privilege Check and Stack Build (assumes called procedure resides in a 32-bit code segment).....	465
14-16	Calling 16-bit Procedure From 32-bit Code Using Far Call With Operand Size Override Prefix.....	467
14-17	Calling 16-bit Procedure From 32-bit Code Using a 16-bit Call Gate	469
14-18	Far Return From 16-bit Procedure to 32-bit Caller	470
14-19	16-bit, 286-Compliant Call Gate Descriptor	471
14-20	Calling 32-bit Procedure From 16-bit Code Using Far Call With Operand Size Override Prefix.....	472
14-21	32-bit Call Gate Descriptor Format	473
14-22	Calling 32-bit Procedure From 16-bit Code Using a 32-bit Call Gate	474
14-23	Far Return From 32-bit Procedure to 16-bit Caller	475
15-1	Data Segment Descriptor Pre-Load Privilege Check.....	482
15-2	Example Value in DS Register	484

Figures

15-3	Example Data Segment Descriptor	484
15-4	Expand Up Stack Approaching a Full Condition	486
15-5	Example Value in SS Register	487
15-6	Example Stack Segment Descriptor	487
15-7	Copying to a Larger Stack Renders Stored Pointers Incorrect.....	488
15-8	Enlarging the Stack by Lowering Stack Base Renders Stored Pointers Incorrect.....	489
15-9	Enlarging Stack by Increasing Limit Won't Head Off a Stack Overflow	489
15-10	Expand-Down Stack Approaching Full	491
15-11	Decreasing Limit Lowers the Stack's Artificial Floor	492
16-1	Life Without Address Translation.....	500
16-2	Life With Address Translation.....	501
16-3	Address Translation Redirects Access 1.....	502
16-4	Address Translation Redirects Access 2.....	503
16-5	Address Translation Redirects Access 3.....	504
16-6	Address Translation Redirects Access 4.....	505
16-7	Address Translation Redirects Access 5.....	506
16-8	First Generation 4GB Virtual Address Space Partitioning	507
16-9	Example Virtual Buffer	509
16-10	Paging Redirects DOS Accesses to a Discrete 1MB Area.....	513
16-11	Example Virtual Address	514
16-12	Accesses to Virtual Page 34 in Virtual 4MB Region Number 4	518
16-13	Accesses to Virtual Page 35 in Virtual 4MB Region Number 4	519
16-14	Accesses to Virtual Page 36 in Virtual 4MB Region Number 4	520
16-15	Accesses to Virtual Page 37 in Virtual 4MB Region Number 4	522
16-16	Code and Data TLBs.....	524
16-17	CR4[PGE] Enables/Disables the Global Page Feature.....	527
16-18	32-bit Page Table Entry (PTE)	527
16-19	Control Register 3 (CR3)	528
16-20	Control Register 0 (CR0)	529
16-21	32-bit Page Directory Entry (PDE) Pointing to a Page Table	530
16-22	32-bit Page Table Entry (PTE)	532
16-23	Effective Read/Write Permission Determination.....	534
16-24	User/Supervisor Permission Determination.....	535
16-25	Page Table Not in Memory (1-of-3).....	539
16-26	Page Table Not in Memory (2-of-3).....	540
16-27	Page Table Not in Memory (3-of-3).....	541
16-28	Page Fault Register (CR2).....	541
16-29	PDE or PTE when Page Table (or page) not Present in Memory	542
16-30	Page Not in Memory (1-of-3)	543
16-31	Page Not in Memory (2-of-3)	544
16-32	Page Not in Memory (3-of-3)	545

16-33	Page Fault Error Code Format	548
16-34	PDE Pointing to 4MB Page	551
16-35	4MB Page Address Translation	552
16-36	CR4[PAE] Enables/Disables PAE-36 Mode Feature	554
16-37	First-Generation Virtual Address Space Partitioning	555
16-38	Second-Generation (PAE-36) Virtual Address Space Partitioning.....	557
16-39	PAE-36 Mode Uses 3-Level Lookup	558
16-40	CR3 Format With PAE-36 Mode Enabled	559
16-41	Step 1: IA-32 PAE Mode: PDPTE Selection	561
16-42	IA-32 PAE Mode: PDPT Entry (PDPTE) Format.....	561
16-43	Step 2: IA-32 PAE Mode: PDE Selection.....	562
16-44	IA-32 PAE Mode: PDE Pointing to a 4KB Page Table	563
16-45	IA-32 PAE Mode: 2MB Physical Page Selected	564
16-46	IA-32 PAE Mode: PDE Pointing to a 2MB Physical Page	564
16-47	IA-32 PAE Mode: PTE Selection	565
16-48	IA-32 PAE Mode: 4KB Page Location Selection	566
16-49	IA-32 PAE Mode: PTE Pointing to a 4KB Physical Page.....	566
16-50	Read/Write Permission Determination	567
16-51	User/Supervisor Permission Determination.....	568
16-52	CR0.....	569
16-53	Write-Protection (Intel approach)	570
16-54	Write-Protection (AMD approach)	571
16-55	CR4.....	576
16-56	PDE Points to a 4MB Page Below or Above the 4GB Address Boundary	578
16-57	Stack Usage in C Function Call.....	581
16-58	The Exploit.....	582
16-59	Execute Disable Feature Enable Is in the EFER Register	583
16-60	PDE Pointing to a 4KB Page Table	584
16-61	PDE Pointing to a 2MB Physical Page	584
16-62	PTE Pointing to a 4KB Physical Page.....	585
16-63	IA32_CR_PAT MSR.....	588
17-1	MTRRCAP Register	603
17-2	MTRRDefType Register	604
17-3	First MB of Memory Space	607
17-4	Format of Variable-Range MTRRPhys Register Pair.....	608
17-5	Handling of Posted Writes by Memory Type.....	620
17-6	Four Empty WCBs.....	621
17-7	4-byte Write to WC Memory	621
17-8	1-byte Write to WC Memory	622
17-9	1-byte Write to WC Memory	623
17-10	Sixteen 4-byte Writes to WC Memory	624
17-11	8-byte Write to WC Memory	625

Figures

17-12	1-byte Write to WC Memory.....	626
17-13	Example of Write Collapsing in WC Memory.....	627
18-1	Task A's TDS, LDT and TSS.....	634
18-2	CR3 Points to Task A's Address Translation Tables	635
18-3	The GDT and the GDTR Register.....	636
18-4	The LDTR and Task A's LDT Register	637
18-5	The TR Register and Task A's TSS	638
18-6	Task A's Suspension.....	640
18-7	Scheduler's Resumption	641
18-8	The 32-bit TSS Descriptor Format	649
18-9	The Task Gate Format	651
18-10	The IDT (Interrupt Descriptor Table)	651
18-11	The Task Register.....	654
18-12	32-bit Task State Segment (TSS) Format.....	656
18-13	CR3 Format When Using First-Generation Address Translation.....	658
18-14	CR3 Format When Using Second-Generation Address Translation.....	659
18-15	Eflags Register.....	660
18-16	Task Switch Flowchart (1-of-3)	668
18-17	Task Switch Flowchart (2-of-3)	669
18-18	Task Switch Flowchart (3-of-3)	670
18-19	Task A Running	672
18-20	Task A Calls Task B	673
18-21	Task B Calls Task C	674
18-22	Task C Executes IRET.....	675
18-23	Task B Executes IRET.....	676
19-1	Real Mode Interrupt Handling.....	683
19-2	Return From Real Mode Handler To Interrupted Real Mode Application	684
19-3	Protected Mode Interrupt Descriptor Table (IDT)	687
19-4	Interrupt Descriptor Table Register (IDTR)	688
19-5	32-bit Interrupt Gate Descriptor Format	693
19-6	32-bit Trap Gate Format.....	694
19-7	Interrupt/Trap Gate Operation (1-of-2)	695
19-8	Interrupt/Trap Gate Operation (2-of-2)	696
19-9	Task Gate Format.....	697
19-10	Event Detected	700
19-11	Interrupted Program and Handler at Same Privilege Level	701
19-12	Handler More-Privileged than Interrupted Program	702
19-13	VM86 Mode Program Interrupted	703
19-14	Interrupt Causes Task Switch.....	704
19-15	Same Privilege Level and No Error Code	706
19-16	Same Privilege Level with Error Code	707
19-17	Privilege Level Switch without Error Code	707

19-18	Privilege Level Switch with Error Code	707
19-19	32-bit Task State Segment (TSS) Format	708
19-20	Return From Protected Mode Handler To Interrupted Protected Mode Program	710
19-21	Return to Interrupted Task (Interrupt/Exception Caused a Task Switch)	711
19-22	Return From a VM86 Mode Handler (i.e., a Real Mode Handler) to an Interrupted VM86 Mode Program	711
19-23	Return From Protected Mode Handler to Interrupted VM86 Mode Program	712
19-24	Hardware Device Interrupt Assignments in a PC-Compatible Platform	719
19-25	Standard Error Code Format	739
19-26	Page Fault Error Code Format	739
19-27	Page Fault Exception Error Code Format	768
20-1	Task State Segment (TSS)	786
20-2	Eflags Register	787
20-3	Real Mode IDT	791
20-4	DOS Task's Perception of the 1st MB of Memory Space	793
20-5	Paging Mechanism Used to Redirect DOS Task Memory Accesses	795
20-6	Using CLI/STI Instructions to Disable/Enable Interrupt Recognition	798
20-7	DOS Tasks Use INT Instructions	800
20-8	Solving the IO Problem: When VM = 1, IOPL is don't care	802
20-9	Protected Mode IDT	808
20-10	VMM Handling of CLI Instruction	814
20-11	Real Mode Task Is Interrupted	815
20-12	When the Real Mode Task's Timeslice Expires	815
20-13	CR4	819
20-14	Efficient Handling of the CLI/STI Instructions	823
20-15	Interrupt Received After CLI Execution	824
20-16	Privilege Level 0 Stack After Interrupt/Exception in VM86 Mode	824
20-17	VMM Passes Control to a Real Mode Handler (1-of-2)	830
20-18	VMM Passes Control to a Real Mode Handler (2-of-2)	831
21-1	MMX Register Set	836
21-2	Example Operation on Dual Frame Buffers	840
21-3	MMX SIMD Solution Increase Throughput	841
21-4	Dealing with Unpacked Data	842
21-5	Conditional Branches Can Severely Decrease Performance	844
21-6	Example MMX Operation (1-of-4)	846
21-7	Example MMX Operation (2-of-4)	846
21-8	Example MMX Operation (3-of-4)	847
21-9	Example MMX Operation (4-of-4)	847
22-1	The SSE Register Set	859
22-2	SSE Data Types	859
22-3	The MXCSR Register	860

Figures

22-4	Example SSE SIMD FP Operation on Packed 32-bit SP FP Numbers	864
22-5	Example SSE SIMD Scalar Operation	864
22-6	The MOVNTPS Instruction	871
22-7	The MOVNTQ Instruction	872
22-8	The MASKMOVQ Instruction	873
22-9	Stores to WC Memory Are Posted in the WCBs	875
22-10	Stores to WB Memory Create Modified Cache Lines	876
22-11	SFENCE Blocks the Logical Processor from Executing Downstream Stores...	876
22-12	Logical Processor Can Execute Downstream Stores After Buffers Are Flushed.....	877
22-13	x87 FP/MMX/SSE Register Save Data Structure	883
22-14	OSFXSR and OSXMMEXCPT Bits Added to CR4	885
22-15	SSE2 XMM Data types	887
22-16	64-bit DP FP Numeric Format.....	888
22-17	The MXCSR Register	890
22-18	The EBX Register After Executing a CPUID Request Type 1.....	891
22-19	The MOVNTDQ Instruction	895
22-20	The MOVNTPD Instruction	896
22-21	The MOVNTI Instruction	897
22-22	The MASKMOVDQU Instruction	898
22-23	Example Horizontal FP Math Operation	907
23-1	64-bit OS Environment.....	916
23-2	Switching to IA-32e Mode	919
23-3	CS Descriptor Interpretation in 64-bit Mode	922
23-4	DS/ES/SS Segment Descriptor Interpretation in 64-bit Mode	933
23-5	FS/GS Segment Descriptor Interpretation in 64-bit Mode	934
23-6	Instruction Pointer Register	939
23-7	IA-32e Mode Call Gate Descriptor	945
23-8	LDT Descriptor in IA-32e Mode	946
23-9	TSS Descriptor in IA-32e Mode	947
23-10	GDTR Contents After Loading in Compatibility Mode.....	950
23-11	GDTR Contents After Loading in 64-bit Mode	951
23-12	GDT and LDT Can Contain Mix of 8- and 16-byte Descriptors.....	952
23-13	Legacy 8-byte LDT Descriptor (see Figure 23-8 on page 946 for IA-32e Version)	954
23-14	LDTR in IA-32e Mode	955
23-15	Interrupt Gate Descriptor in IA-32e Mode.....	957
23-16	Trap Gate Descriptor in IA-32e Mode	958
23-17	Interrupt/Exception Flow in IA-32e Mode (1 of 2)	961
23-18	Interrupt/Exception Flow in IA-32e Mode (2 of 2)	962
23-19	Handler's Stack After Pushes.....	962
23-20	IDTR Contents After Loading in Compatibility Mode	963

23-21	IDTR Contents After Loading in 64-bit Mode	964
23-22	IA-32e Call Gate Operation	966
23-23	IA-32e Call Gate	968
23-24	Relationship of TSS, GDT and TSS Descriptor	970
23-25	IA-32e TSS Data Structure	971
23-26	IA-32e Mode TSS Descriptor	972
23-27	TR in IA-32e Mode.....	973
23-28	Lower 8-bytes of a 16-byte IA-32e TSS Descriptor.....	974
23-29	Legacy IA-32 TSS Data Structure	975
23-30	Stack Contents When IRET Executed to Start Privilege Level 3, 64-bit Task...	979
23-31	Stack Contents When IRET Executed to Start Legacy Task	981
24-1	IA-32e 3rd Generation Address Translation Mechanism	987
24-2	IA-32e Address Translation Step 1.....	989
24-3	IA-32e Mode: PML4 Entry (PML4E) Format	990
24-4	IA-32e Address Translation Step 2.....	991
24-5	IA-32e Mode: PDPT Entry (PDPTE) Format.....	992
24-6	IA-32e Address Translation Step 3.....	993
24-7	IA-32e Mode: PD Entry (PDE) Format (points to Page Table).....	994
24-8	IA-32e Mode: Page PD Entry (PDE) Format (points to 2MB page).....	995
24-9	2MB Physical Page Selected	996
24-10	PT Entry (PTE) Format.....	997
24-11	IA-32e Address Translation Step 4.....	998
24-12	Read/Write Permission Determination in 64-bit Mode.....	1001
24-13	User/Supervisor Permission Determination in 64-bit Mode.....	1002
24-14	Write-Protection in 64-bit Mode (Intel)	1003
24-15	Write-Protection in 64-bit Mode (AMD)	1004
25-1	Register Set in Compatibility Mode	1015
26-1	Intel 64 Register Set	1025
26-2	EFER (Extended Features Enable) Register	1026
26-3	CR0 in IA-32 Mode and Compatibility Mode	1030
26-4	CR0 in 64-bit Mode	1030
26-5	CR2 in IA-32e Mode	1030
26-6	CR3 in IA-32e Mode	1031
26-7	In IA-32e Mode, CR3 Points to Top-Level Paging Directory	1031
26-8	CR4 in IA-32 Mode and Compatibility Mode	1032
26-9	CR4 in 64-bit Mode	1032
26-10	CR8 (Task Priority Register).....	1033
26-11	Rflags Register (only in 64-bit Mode)	1033
26-12	The A, B, C and D Registers (in 64-bit Mode).....	1034
26-13	The BP, SI, DI and SP Registers (in 64-bit Mode).....	1034
26-14	Registers R8 - R15 (only in 64-bit Mode)	1035
26-15	Result Storage in 64-bit Mode.....	1035

Figures

26-16	Upper 32-bits of the First Eight GPRs Do Not Survive Mode Change	1036
26-17	SSE Register Set in IA-32 and Compatibility Modes	1037
26-18	SSE Register Set in 64-bit Mode	1037
26-19	Debug Register Set (all modes in an Intel 64 Processor)	1038
27-1	The ModRM Byte's Operand 1 and 2 Fields Are Each 3-bits Wide	1044
27-2	The Register Select Field in the Primary Opcode Byte Is 3-bits Wide	1044
27-3	The Scale/Index/Base (SIB) Byte	1045
27-4	Placement of REX Prefix	1047
27-5	REX Prefix Format	1050
27-6	The REX Prefix	1051
27-7	The ModRM Byte.....	1054
27-8	64-bit Code Segment Descriptor.....	1062
27-9	Example of Canonical Address Formation	1065
27-10	ModRM Byte.....	1068
28-1	FXSAVE Structure Legacy Format	1098
28-2	FXSAVE Structure in 64-bit Mode with REX[W] = 1.....	1099
28-3	EFER Register	1100
28-4	FXSAVE Format When EFER[FFXSR] = 1 in 64-bit Mode at Privilege Level 0.....	1101
28-5	Processor's SM RAM Memory Map.....	1103
29-1	Initial Flat Model GDT	1117
29-2	32-bit CS Descriptor.....	1118
29-3	32-bit DS Descriptor	1119
29-4	Protected Mode IDT	1121
29-5	CR0.....	1123
31-1	VMX Mode Transitions.....	1151
31-2	IA32_VMX_Basic MSR.....	1152
32-1	SM RAM Memory Map	1174
32-2	MTRRCAP Register.....	1183
32-3	SM Range Register Pair.....	1183
32-4	The SMM Revision ID	1192
32-5	The Auto Halt Restart Field	1194
32-6	IO State Field in SM State Save Area	1197
32-7	SM RAM Example One	1200
32-8	SM RAM Example Two	1201
32-9	SM RAM Example Three	1201
32-10	SM RAM Example Four.....	1202
32-11	SM RAM Example Five.....	1202
33-1	Machine Check Exception Enable/Disable Bit (CR4[MCE]).....	1210
33-2	MCA Register Set.....	1211
33-3	MC Global Count and Present Register	1213
33-4	MC Global Status Register.....	1214

33-5	MCi Control Register	1218
33-6	MCi Status Register Detail.....	1220
33-7	The MCi_CTL2 Register.....	1234
34-1	An External Hardware Interrupt Delivered to the Processor's INTR Pin.....	1245
34-2	Legacy PC-AT Compatible Interrupt Controllers.....	1246
34-3	Legacy Interrupt Delivery Mechanism Inefficient.....	1248
34-4	The APIC Bus (Pentium and P6)	1249
34-5	The Pentium 4 Eliminated the APIC Bus	1250
34-6	Intel QPI System With Single Physical Processor	1251
34-7	Intel QPI System With Dual Physical Processors.....	1252
34-8	Intel QPI System With Four Physical Processors.....	1253
34-9	The Big Picture	1254
34-10	Local APIC Version Register.....	1256
34-11	IA32_APIC_BASE MSR	1258
34-12	The Spurious Vector Register.....	1259
34-13	Local APIC Register Set	1264
34-14	Cluster ID Assignment.....	1278
34-15	Assignment of Agent ID and Local APIC ID (Xeon MP System).....	1280
34-16	Assignment of Agent ID and Local APIC ID (Dual-Processor System)	1281
34-17	The xAPIC ID Register.....	1281
34-18	EBX Contents After a CPUID Request Type 1.....	1282
34-19	Logical Destination Register (LDR) in xAPIC Mode.....	1286
34-20	Logical Address in LDR Is Formed From x2APIC ID	1287
34-21	x2APIC Logical Address Compare	1289
34-22	xAPIC Logical Address Compare When Using Flat or Hierarchical Cluster Model.....	1291
34-23	Logical Destination Register (LDR) in xAPIC Mode.....	1292
34-24	Destination Format Register (DFR) in xAPIC Mode	1292
34-25	Local Interrupt Sources.....	1299
34-26	Local APIC Error Status Register	1302
34-27	Task Priority Register (TPR).....	1306
34-28	Processor Priority Register (PPR)	1308
34-29	CR8.....	1308
34-30	IO APIC and Pentium/P6 Processors Communicate Via 3-Wire APIC Bus .	1314
34-31	IO APIC and Pentium 4 & Later Processors Communicate Via FSB	1315
34-32	IO APIC & Latest Processors Communicate Via Intel QPI.....	1315
34-33	No User-Defined (Fixed) Interrupts Pending.....	1317
34-34	RT Entry Format	1318
34-35	Fixed Interrupt 20h (32d) Received.....	1319
34-36	Fixed Interrupt 20h (32d) Forwarded to Logical Processor and Handler Starts Execution.....	1320
34-37	Handler 20h (32d) Issues EOI to Local APIC.....	1321

Figures

34-38	Fixed Interrupt 20h (32d) Received.....	1323
34-39	Fixed Interrupt 20h (32d) Forwarded to Logical Processor and Handler Starts Execution	1324
34-40	Handler 20h (32d) Issues EOI to Local APIC.....	1325
34-41	Fixed Interrupt 20h/32d Received and Registered	1326
34-42	Highest-Priority, So Shifted to ISR and Delivered to Logical Processor	1327
34-43	Fixed Interrupt FBh/251d Received and Registered.....	1328
34-44	Handler 20h/32d Interrupted by Fixed Interrupt FBh/251d	1329
34-45	FBh Handler Issues EOI to Local APIC	1330
34-46	20h Handler Issues EOI to Local APIC	1331
34-47	The IRQ Pin Assertion Register	1335
34-48	The IO APIC's EOI Register	1338
34-49	Shareable IRQ Line	1339
34-50	Local APIC Version Register	1339
34-51	The Spurious Vector Register.....	1340
34-52	The IO APIC's Version Register	1340
34-53	RT Entry Format	1342
34-54	Virtual Wire Mode A.....	1349
34-55	Virtual Wire Mode B	1351
34-56	xAPIC Interrupt Command Register (ICR)	1354
34-57	x2APIC Interrupt Command Register (ICR)	1354
34-58	x2APIC Self IPI MSR Register.....	1359
34-59	Local Interrupt Sources.....	1360
34-60	LVT LINT0 or LINT1 Register	1366
34-61	Local APIC Timer Register Set.....	1368
34-62	Performance Counter Overflow, Thermal Sensor & CMC LVT Registers	1370
34-63	The IA32 Thermal Status MSR.....	1371
34-64	The IA32 Thermal Interrupt MSR	1372
34-65	LVT Error Register.....	1374
34-66	Local APIC Error Status Register (ESR) in xAPIC Mode.....	1374
34-67	Local APIC Error Status Register (ESR) in x2APIC Mode.....	1375
34-68	Local APIC's Spurious Vector Register	1378
34-69	Pentium 4 BSP Selection Process	1380
34-70	The IA32_APIC_BASE MSR.....	1381
34-71	BIOS's AP Discovery Procedure Part 1	1385
34-72	BIOS's AP Discovery Procedure Part 2	1386
34-73	The Local APIC's Spurious Vector Register	1387
34-74	AP Setup Program Part 1	1389
34-75	AP Setup Program Part 2	1390

1	Trademarks.....	6
1-1	x86 Software Architectures.....	13
1-2	Physical, Virtual and Linear Memory Address Space	17
1-3	Some Other Useful Terms	18
2-1	Basic Execution Modes.....	23
2-2	IA-32 SubModes.....	25
2-3	IA-32e SubModes.....	28
2-4	Terminology	32
3-1	Major Evolutionary Developments.....	43
3-2	8086 Real Mode Characteristics	46
3-3	286 16-bit Protected Mode Characteristics (& short-comings)	48
3-4	386 32-bit Protected Mode Characteristics	52
3-5	Intel Processor Microarchitecture Families.....	56
3-6	x86 Family Members (as of February 2009)	57
4-1	Logical Processor State After Removal of Reset.....	66
6-1	386 Instruction Set.....	111
6-2	Current-Day Instruction Set.....	118
7-1	Effective Operand Size in 16- or 32-bit Mode (without prefix).....	157
7-2	Effective Operand Size in 16- or 32-bit Mode (with prefix).....	158
7-3	Determination of Effective Operand Size in 64-bit Mode.....	159
7-4	Information Related to an Instruction	161
7-5	Instruction Elements.....	164
7-6	0F 00 Micro-Map	183
7-7	0F 01 Micro-Map	184
7-8	0F 18 Micro-Map	184
7-9	0F 71 Micro-Map	184
7-10	0F 72 Micro-Map	185
7-11	0F 73 Micro-Map	185
7-12	0F AE Micro-Map	186
7-13	0F B9 Micro-Map Is Reserved for Future Use	186
7-14	0F BA Micro-Map	186
7-15	0F C7 Micro-Map	187
7-16	Fields in Primary Opcode or ModRM Byte of Some Instructions	189
7-17	tttn Condition Code Definition Field.....	193
7-18	Examples of Register Operand Specification.....	195
7-19	Register Is Selected Based on Reg Value (plus Opcode's W bit if present).....	198
7-20	Summary of Memory Addressing Modes.....	199
7-21	ModRM Interpretation When Effective Address Size = 16-bits.....	201
7-22	ModRM Interpretation When Effective Address Size = 32-bits.....	202
7-23	Effective Address = Base + (Index * Scale Factor).....	205
7-24	Branch Forms in 32-bit Mode.....	207

Tables

7-25	Example Instruction Execution in 32-bit Mode (with and without Operand Size Override prefix)	212
7-26	Example Instruction Execution in 16-bit Mode (with and without override) .	214
7-27	Format(s) Associated With Instruction Sets.....	222
8-1	Basic Differences Between 8086 Operation and Real Mode	229
8-2	Expanded/Enhanced Real Mode Instructions	230
8-3	CR0 Bit Assignment.....	238
8-4	CR0 FPU Control Bits.....	241
8-5	CR4 Bit Assignment.....	245
8-6	XCR0 Bit Assignment.....	250
8-7	Eflags Register Bit Assignment.....	252
8-8	Definition of DR7 Bits Fields.....	264
8-9	Interpretation of the R/W Field In DR7	267
8-10	Interpretation of the LEN Field In DR7	268
8-11	Debug Status Register Bits.....	268
8-12	List of Currently-Defined Architectural MSRs.....	273
8-13	Physical Memory Address Formation in Real Mode	290
8-14	Segment Register Usage in Real Mode	291
8-15	Reset Code Byte Values	313
8-16	IDT Entry Assignments.....	318
8-17	Software Exception Categories	330
9-1	CR0 x87 FPU-related Bit Fields.....	341
9-2	Representation of Special Values	348
9-3	The First Example	349
9-4	The Second Example	350
9-5	FCW Register Fields	350
9-6	FSW Register Fields.....	352
12-1	Protection Mechanisms	379
13-1	Logical Processor Actions When a Segment Register Is Loaded	389
13-2	Descriptor Table Types	391
13-3	Data/Stack Segment Types (C/D = 0).....	404
13-4	Code Segment Types (C/D = 1)	405
13-5	Types of System Segments	407
14-1	Actions That Cause a Switch to a New CS.....	417
14-2	Code Segment Descriptor Format	419
14-3	Short/Near Jump Forms.....	428
14-4	Condition Code Encoding in Least-Significant Primary Opcode Byte.....	429
14-5	Conditional Branch Forms.....	430
14-6	Loop Instruction Forms	432
14-7	Far Jump Forms.....	436
14-8	Near Call/Return Forms	445
14-9	Far Call/Return Forms	448

14-10	32-bit Call Gate Descriptor Elements.....	455
14-11	Example Call Gate Descriptor Elements (see Figure 14-10)	457
14-12	Example Code Segment Descriptor	459
14-13	Far Return Forms (in Protected Mode).....	477
16-1	Paging Evolution	495
16-2	Page Directory Caching Policy	528
16-3	Layout of PDE Pointing to a Second-Level Page Table.....	530
16-4	Layout of PTE.....	532
16-5	Effect of PDE/PTE U/S Bit Settings	536
16-6	Effect of R/W Bit Settings.....	537
16-7	Page Fault Exception Error Code Status Bit Interpretation.....	548
16-8	PCD and PWT Bit Settings	559
16-9	32-bit Windows PAE Support.....	572
16-10	Address Translation Table Caching Policy	586
16-11	Bit Field Selection in the IA32_CR_PAT MSR	588
16-12	Memory Types That Can Be Encoded in a IA32_CR_PAT MSR Entry	589
16-13	Default Memory Types in the IA32_CR_PAT MSR Entries	590
16-14	Effective Memory Type Determination.....	591
16-15	Pre-PAT Interpretation of the PCD and PWT Bits.....	597
17-1	The Fixed Range MTRRs	606
17-2	Memory Type Determination Using MTRRs	616
18-1	Hardware-based Task Switching's Key Elements	643
18-2	Events that Cause a Task Switch	646
18-3	Additional Information Related to Task Switch Triggers.....	665
19-1	Introduction to the IDT Gate Types	689
19-2	Elements of Interrupt Gate Descriptor	692
19-3	Interrupt/Exception Handler State Save Cases	705
19-4	The Interrupt Return (IRET) Instruction	709
19-5	PC-Compatible IRQ Assignment.....	715
19-6	Results of CLI/STI Execution	721
19-7	Software Exception Types	729
19-8	Exception Categories.....	729
19-9	Exceptions that Return Error Codes	736
19-10	Interrupt/Exception Priority	739
19-11	Debug Exception Conditions and Exception Class	744
19-12	Interrupt and Exception Classes.....	754
19-13	Exception Combinations Resulting in a Double Fault Exception.....	755
19-14	Invalid TSS Conditions	757
19-15	FPU Handling of Masked Errors.....	771
19-16	Alignment Requirements by Data Type.....	775
19-17	SIMD FP Error Priorities.....	781
20-1	Hardware Interrupt Types	809

Tables

20-2	INT nn Handling Methods in VM86 Mode	826
21-1	Data Range Limits for Saturation	843
21-2	MMX Instruction Set Summary, Part 1	848
21-3	MMX Instruction Set Summary, Part 2	849
22-1	Evolution of SIMD Model	852
22-2	The MXCSR Register Bit Field Definitions	861
22-3	Prefetch Instruction Behavior	867
22-4	Processors Actions on a Store to Cacheable Memory	869
22-5	CR4[OSFXSR] Bit	884
22-6	SSE Instructions	886
23-1	Target CS's D and L Bits Control Mode Selection	921
23-2	Virtual Address Calculation in Compatibility Mode	925
23-3	Virtual Address Calculation in 64-bit Mode	926
23-4	Segment Register Operation in 64-bit Mode	929
23-5	Instruction Pointer Usage	938
23-6	IA-32e GDT/LDT Descriptor Changes	941
23-7	GDT/LDT Descriptor Types in Protected Mode and IA-32e Mode	943
23-8	GDT Descriptor Types in IA-32e Mode	948
23-9	LDT Descriptor Types in IA-32e Mode	953
23-10	IDT Descriptor Types in IA-32e Mode	956
23-11	IA-32e Call Gate Changes	964
23-12	Call Gate Stack Switch Legacy/IA-32e Mode Differences	967
23-13	TSS Usage in IA-32e Mode	976
25-1	OS Environment Changes in IA-32e Mode	1012
25-2	IRET Characteristics in IA-32e Mode	1019
26-1	EFER Register Bit Assignment	1026
26-2	Control Registers in Compatibility and 64-bit Modes	1028
27-1	Instructions That Don't Require the REX Prefix	1048
27-2	Effective Operand Size in 64-bit Mode	1052
27-3	REX[R] + ModRM[Reg], or REX[B] + Opcode[Reg], or REX[B] + ModRM[RM] = 4-bit Register Select Field	1054
27-4	64-bit Effective Address = Base + (Index * Scale Factor)	1056
27-5	Theory vs. Practice	1063
27-6	Summary of Memory Addressing Modes Available in 64-bit Mode	1066
27-7	Mod + RM Interpretation When Effective Address Size = 64-bits	1068
27-8	Branch Forms in 64-bit Mode	1071
28-1	Instructions Enhanced in 64-bit Mode With REX[W] = 1	1078
28-2	Instructions Which Are Invalid In 64-bit Mode	1080
28-3	Instructions Reassigned In 64-bit Mode	1081
28-4	Instructions That Reference RSP (64-bit Mode) & Default to 64-bit Operand Size	1082
28-5	Unconditional Branches in 64-bit Mode	1084

28-6	Calls, RET, and IRET in 64-bit Mode	1087
28-7	Conditional Branches in 64-bit Mode	1093
28-8	IA-32 Processor SMRAM State Save Area (shown from top down)	1104
28-9	Intel 64 Processor SMRAM State Save Area	1107
29-1	Segment Register Reload Procedure	1126
29-2	Startup GDT Content	1130
29-3	Startup_32 GDT Content	1133
31-1	VMX-related Instructions	1154
31-2	VMCS Regions	1157
32-1	IA-32 SM RAM State Save Area (shown from the top down)	1174
32-2	Intel 64 Processor SMRAM State Save Area	1178
32-3	Register Set Values After Entering SMM	1190
32-4	Processor Actions on RSM Are Defined by the Bit Setting	1194
33-1	IA-32 Processor's Extended MC State MSRs	1215
33-2	Intel 64 Processor's Extended MC State MSRs	1216
33-3	Simple MCA Error Codes	1223
33-4	Compound Error Codes Forms	1225
33-5	Transaction Type Sub-Field (TT)	1226
33-6	Memory Hierarchy Sub-Field (LL)	1226
33-7	Request Sub-Field (RRRR)	1226
33-8	Definition of the Bus and Interconnect-related PP, T, and II Fields	1227
33-9	Definition of Integrated Memory Controller MMM and CCCC Fields	1228
33-10	MCi_STATUS Breakdown for FSB-related Errors	1229
33-11	Cache ECC Green/Yellow Light Indicator	1233
34-1	Legacy Interrupt Delivery Inefficiencies	1247
34-2	Local APIC Operational Modes	1260
34-3	A Brief Description of the Local APIC Registers	1265
34-4	Quad Xeon MP System with Hyper-Threading Disabled	1279
34-5	Quad Xeon MP System with Hyper-Threading Enabled	1279
34-6	Dual Processor System with Hyper-Threading Enabled	1280
34-7	Addressing Summary Table	1292
34-8	Interrupt Priorities	1301
34-9	User-Defined Interrupt Priority Scheme	1303
34-10	Interrupt Message Address Format	1309
34-11	Interrupt Message Data Format	1311
34-12	The IO APIC Register Set	1332
34-13	RT Register Format	1343
34-14	ICR Bit Assignment	1355
34-15	State of LVT Registers	1362
34-16	Permissible Delivery Modes for the LINT Inputs	1363
34-17	Local APIC Error Status Register	1375

About This Book

Is This the Book for You?

If you're looking for a book designed specifically for those who need to come up to speed on the 32-/64-bit x86 Instruction Set Architecture (ISA) as quickly and painlessly as possible, then consider this book.

On the other hand, the author fully realizes that a certain segment of the technical population rejects and is, indeed, deeply offended by any attempt to present arcane technical material in a learning-friendly manner. Having been exposed to the occasional criticism from such individuals, if you fall into this category I can only forewarn that this book is not for you. Do not waste your money or your time.

A Moving Target

The reader should keep in mind that MindShare books often deal with rapidly evolving technologies. This being the case, it should be recognized that this book is a *snapshot* of the state of the x86 programming environment at the time that the book was completed (November, 2009).

x86 Instruction Set Architecture (ISA)

Throughout this book, the term *ISA* (*Instruction Set Architecture*) refers to the current execution environment defined by the x86 ISA specification:

- Any reference to the term *IA-32 ISA* refers to the facilities visible to the programmer when the processor is operating in 32-bit mode (referred to by Intel as IA-32 Mode and by AMD as Legacy Mode).
- Any reference to the term *Intel 64 ISA* refers to the facilities visible to the programmer when the processor is operating in 64-bit Mode (referred to by Intel as Intel 64 Mode and by AMD as Long Mode).

x86 Instruction Set Architecture

Glossary of Terms

A comprehensive glossary may be found on page 1391.

32-/64-bit x86 Instruction Set Architecture Specification

As of this writing (February, 2009), the ISA specification is embodied in the *Intel 64 and IA-32 Architectures Software Developer's Manual* which currently consist of the following five volumes:

- *Basic Architecture*; order number 253665.
- *Instruction Set Reference A-M*; order number 253666.
- *Instruction Set Reference N-Z*; order number 253667.
- *System Programming Guide, Part 1*; order number 253668.
- *System Programming Guide, Part 2*; order number 253669.

Alternatively, the specification is also embodied in the equivalent manuals available from AMD.

While the specification does define the register and instruction sets, interrupt and software exception handling, and standard processor facilities such as memory address generation and translation, the processor modes of operation, multitasking and protection mechanisms, etc., it does *not* specify processor-specific features such as the following:

- Whether or not a processor design includes caches and, if so, the number of, size of, and architecture of the caches.
- Whether or not a processor design includes one or more TLBs (Translation Lookaside Buffers) and, if so, the number of, size of, and architecture of the TLBs.
- The type of interface that connects the processor to the system.
- The number and types of instruction execution units.
- The implementation-specific aspects of a processor's microarchitecture.
- Various other performance enhancement features (branch prediction mechanisms, etc.).

The Specification Is the Final Word

This book represents the author's interpretation of the Intel x86 ISA specification. When in doubt, the Intel specification is the final word.

Book Organization

This book is organized in seven parts:

“Part 1: Introduction”, intended as a back-drop to the detailed discussions that follow, consists of the following chapters:

- Chapter 1, entitled "Basic Terms and Concepts," on page 11.
- Chapter 2, entitled "Mode/SubMode Introduction," on page 21.
- Chapter 3, entitled "A (very) Brief History," on page 41.
- Chapter 4, entitled "State After Reset," on page 63.

“Part 2: IA-32 Mode” provides a detailed description of two IA-32 Mode sub-modes—Real Mode and Protected Mode—and consists of the following chapters:

- Chapter 5, entitled "Intro to the IA-32 Ecosystem," on page 79.
- Chapter 6, entitled "Instruction Set Expansion," on page 109.
- Chapter 7, entitled "32-bit Machine Language Instruction Format," on page 155.
- Chapter 8, entitled "Real Mode (8086 Emulation)," on page 227.
- Chapter 9, entitled "Legacy x87 FP Support," on page 339.
- Chapter 10, entitled "Introduction to Multitasking," on page 361.
- Chapter 11, entitled "Multitasking-Related Issues," on page 367.
- Chapter 12, entitled "Summary of the Protection Mechanisms," on page 377.
- Chapter 13, entitled "Protected Mode Memory Addressing," on page 383.
- Chapter 14, entitled "Code, Calls and Privilege Checks," on page 415.
- Chapter 15, entitled "Data and Stack Segments," on page 479.
- Chapter 16, entitled "IA-32 Address Translation Mechanisms," on page 493.
- Chapter 17, entitled "Memory Type Configuration," on page 599.
- Chapter 18, entitled "Task Switching," on page 629.
- Chapter 19, entitled "Protected Mode Interrupts and Exceptions," on page 681.
- Chapter 20, entitled "Virtual 8086 Mode," on page 783.
- Chapter 21, entitled "The MMX Facilities," on page 835.
- Chapter 22, entitled "The SSE Facilities," on page 851.

“Part 3: IA-32e OS Kernel Environment” provides a detailed description of the IA-32e OS kernel environment and consists of the following chapters:

- Chapter 23, entitled "IA-32e OS Environment," on page 913.
- Chapter 24, entitled "IA-32e Address Translation," on page 983.

x86 Instruction Set Architecture

“Part 4: Compatibility Mode” provides a detailed description of the Compatibility submode of IA-32e Mode and consist of the following chapter:

- Chapter 25, entitled "Compatibility Mode," on page 1009.

“Part 5: 64-bit Mode” provides a detailed description of the 64-bit submode of IA-32e Mode and consists of the following chapters:

- Chapter 26, entitled "64-bit Register Overview," on page 1023.
- Chapter 27, entitled "64-bit Operands and Addressing," on page 1041.
- Chapter 28, entitled "64-bit Odds and Ends," on page 1075.

“Part 6: Mode Switching Detail” provides a detailed description of:

- Switching from Real Mode to Protected Mode. This topic is covered in Chapter 29, entitled "Transitioning to Protected Mode," on page 1113.
- Switching from Protected Mode to IA-32e Mode. This topic is covered in Chapter 30, entitled "Transitioning to IA-32e Mode," on page 1139.

“Part 7: Other Topics” provides detailed descriptions of the following topics:

- Chapter 31, entitled "Introduction to Virtualization Technology," on page 1147.
- Chapter 32, entitled "System Management Mode (SMM)," on page 1167.
- Chapter 33, entitled "Machine Check Architecture (MCA)," on page 1207.
- Chapter 34, entitled "The Local and IO APICs," on page 1239.

Topics Outside the Scope of This Book

The CPUID Instruction

The CPUID instruction is referred to numerous times in this book. For a detailed description of its usage, refer to the Intel publication entitled *Intel 64 and IA-32 Architectures Software Developer's Manual, Volume 2A, Instruction Set Reference A-M*.

Detailed Description of Hyper-Threading

The Intel Hyper-Threading facility is not covered in this book because it is not part of the x86 ISA. For a detailed description of this facility, refer to Chapter 39 in the MindShare book entitled *The Unabridged Pentium 4*.

Detailed Description of Performance Monitoring

The Intel Performance Monitoring facility is not covered in this book because it is not part of the x86 ISA. For a detailed description of this facility, refer to the section entitled *The Performance Monitoring Facility* in Chapter 56 of the MindShare book entitled *The Unabridged Pentium 4*.

Documentation Conventions

The conventions used in this book for numeric values are defined below:

- Hexadecimal Notation. All hex numbers are followed by an “h.” Examples:
 - 9A4Eh
 - 0100h
- Binary Notation. All binary numbers are followed by a “b.” Examples:
 - 0001 0101b
 - 01b
- Decimal Notation. Numbers without any suffix are decimal. When required for clarity, decimal numbers may be followed by a *d*. Examples:
 - 16
 - 255
 - 256d
 - 128d

Other commonly used designations are defined below:

- *lsb* refers to the least-significant bit.
- *LSB* refers to the least-significant byte.
- *msb* refers to the most-significant bit.
- *MSB* refers to the most-significant byte.
- Bit Fields. In many cases, bit fields are documented in the following manner: CR0[15:8] refers to Control Register 0 bits 8 - 15.
- Notations such as *CSDesc[BaseAddress]* are interpreted as the *Base Address field in the Code Segment Descriptor*.

Trademarks

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Those trademark designations known to MindShare Press are listed in Table 1-1 on page 6.

x86 Instruction Set Architecture

Table 1-1: Trademarks

Trademarked Terms	Trademark Owner
AMD, AMD64, Opteron	AMD
Atom, Core, Core 2, Core 2 Duo, Core 2 Quad, Core 2 Solo, Core i7, Core Solo, Hyper-Threading, Intel, Itanium, MMX, NetBurst, Pentium, QPI or QuickPath Interconnect, SpeedStep, SSE, VTune, Xeon.	Intel
Apple, OS X	Apple Computer
FrameMaker	Adobe Systems
IBM, PC-AT, PS/2	IBM
Linux	Linus Torvalds
PCI, PCI Express, PCIe, PCI-X	PCI SIG
SIMD	?
Unix	The Open Group, SCO, ? (I'll leave this one to the lawyers; your guess is as good as mine)
Word	Microsoft

Visit Our Web Site

Our web site (www.mindshare.com) provides detailed descriptions of all of our products including:

- Books and ebooks.
- On-site classes.
- Public, open-enrollment classes.
- Virtual, instructor-led classes.
- Self-paced eLearning modules.
- Technical papers and the MindShare newsletter.

We Want Your Feedback

MindShare values your comments, questions and suggestions. You can contact us via mail, phone, fax or email.

Phone: (719) 487-1417 and in the U.S. (800) 633-1440

Fax: (719) 487-1434

Email: support@mindshare.com

Mailing Address:

MindShare, Inc.
4285 Slash Pine Drive
Colorado Springs, CO 80908

1 Basic Terms and Concepts

This Chapter

This chapter provides a basic definition of the Instruction Set Architecture (ISA), differentiates between the IA-32 and Intel 64 processor architectures, and defines some other important terms and concepts.

The Next Chapter

The next chapter introduces the execution modes and submodes as well as mode switching basics.

ISA Definition

Wikipedia Definition: The Instruction Set Architecture, or ISA, is defined as that part of the processor architecture related to programming, including the native data types, instructions, registers, addressing modes, memory architecture, interrupt and exception handling, and external IO.

This Book Focuses on the Common Intel/AMD ISA

With the exception of some small deviations and differences in terminology, the Intel and AMD x86 processors share a common ISA. This book focuses on their shared attributes and does not cover those areas where the two companies have chosen widely divergent, non-x86 ISA-compliant, solutions.

For Simplicity, Intel Terminology Is Used Throughout

Rather than confusing matters by using both Intel and AMD terminology throughout the book, the author has chosen to use only Intel terminology.

Some Terms in This Chapter May Be New To the Reader

Someone new to the x86 software environment will almost certainly encounter some unfamiliar terms in this chapter. Don't let it disturb you. Every term and concept *will* be described in detail at the appropriate place in the book. The important things to take away from this chapter are the broader concepts.

Two x86 ISA Architectures

All Intel x86-family processors introduced since the advent of the 386 can be divided into two categories (see Table 1-1 on page 13):

- Those that cannot execute 64-bit code—defined by Intel as **IA-32 Architecture** processors,
- and those that can—defined by Intel as **Intel 64 Architecture** processors.

This distinction is an important one but is not always referred to correctly—even in the vendor's own documentation. As an example, in section 3.2.1 of the *Intel 64 and IA-32 Architectures Software Developer's Manual Volume 1: Basic Architecture* manual, it states:

“A task or program running in 64-bit mode on an IA-32 processor can address linear address space of up to 2^{64} bytes (subject to the canonical addressing requirement described in Section 3.3.7.1) and physical address space of up to 2^{40} bytes.”

There is no such thing as 64-bit Mode on an IA-32 processor. Consistent use of terms is critical to a clear understanding of any subject. For someone learning the fundamentals of the x86 programming environment, misleading statements such as this can lead to monumental confusion.

Chapter 1: Basic Terms and Concepts

Table 1-1: x86 Software Architectures

Processor Family	Description
x86 Software Architectures All Intel x86-family processors introduced since the advent of the 386 can be divided into two categories:	
IA-32 Processor	Implements only the Intel IA-32 Architecture which supports the execution of 16- and 32-bit x86 code.
Intel 64 Processor	Implements the Intel 64 Architecture, a superset of the IA-32 Architecture: <ul style="list-style-type: none">• When operating in IA-32 Mode, the processor supports the execution of 16- and 32-bit x86 code.• When operating in IA-32e Mode, the processor supports the execution of 16-, 32- and 64-bit x86 code.

Processors, Cores and Logical Processors

For many, many years, life was simple: a physical processor package contained a single *core*: i.e., the engine that fetched machine language instructions (i.e., a program) from memory, decoded them, dispatched them to the appropriate execution units and then committed their results to the core's register set. This required:

- A single register set.
- A single set of execution units.
- A set of facilities to handle things like:
 - Virtual-to-physical memory address translation.
 - Interrupts and exceptions.
 - Protection.
 - etc.

The advent of multi-core processors and Hyper-Threading (more in a moment) has inevitably led to a confusion of terminology. As an example, consider the case where a dual core processor contains two cores each of which represents a stand-alone fetch, decode, dispatch, execution engine. Each implements its own register set, instruction fetcher, decoders, dispatcher, and execution units. So, in this scenario, the term *processor* really refers to a package containing two cores, each of which represents a separate processing engine. In all likelihood, though, the two cores may share some resources (typically, one or more caches).

Refer to Figure 1-1 on page 15. To further muddy the waters, a core may implement Hyper-Threading capability, in which case, from a programmer's perspective, a single core would implement two or more independent execution engines (referred to as logical processors):

- Each of which implements its own register set and dedicated resources. This includes a dedicated Local APIC (see "APIC" on page 19) to handle interrupt and exception events for its associated logical processor.
- All of which, invisible to software, may share some resources.

As if that's not confusing enough, if the physical processor's Hyper-Threading capability is disabled, then the second logical processor in each core (referred to as the *secondary logical processor*; the first is referred to as the *primary logical processor*) is disabled and each core functions as a single logical processor.

To sum it up, a physical processor contains one or more cores and, if it implements Hyper-Threading *and* it has been enabled, each core appears to software as two or more separate processors (i.e., logical processors). During a given period of time, all of the logical processors could be executing separate program threads.

Fundamental Processing Engine: Logical Processor

Rather than sprinkling hundreds of references to processors, cores and logical processors throughout the remainder of the book, the fundamental processing engine will heretofore be referred to as a *logical processor* (unless, of course, I am specifically discussing the physical processor package or a core, rather than a logical processor).

2 *Mode/SubMode Introduction*

The Previous Chapter

The previous chapter provided a basic definition of the Instruction Set Architecture (ISA), differentiated between the IA-32 and Intel 64 processor architectures, and defined some other important terms and concepts.

This Chapter

This chapter introduces the execution modes and submodes and mode switching basics.

The Next Chapter

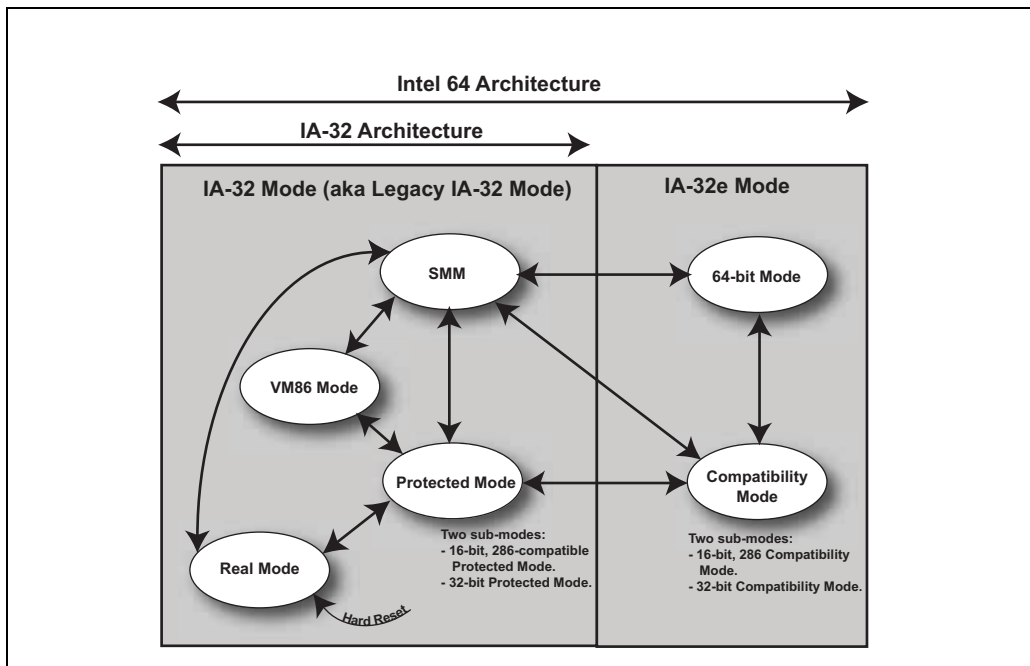
The next chapter introduces the evolution of the x86 ISA, as well as the basic operational characteristics of 8086 Real Mode, 286 Protected Mode, and 386 Protected Mode. It also introduces the Intel microarchitecture families including a product introduction timeline.

Basic Execution Modes

Figure 2-1 on page 22 illustrates the execution modes supported on processors based on the IA-32 architecture versus those based on the Intel 64 architecture. Table 2-1 on page 23 provides an elementary description of the two basic execution modes—IA-32 Mode and IA-32e Mode.

x86 Instruction Set Architecture

Figure 2-1: Execution Mode Diagram



Chapter 2: Mode/SubMode Introduction

Table 2-1: Basic Execution Modes

Mode	Description
IA-32 Mode (also referred to as Legacy IA-32 Mode)	<ul style="list-style-type: none">• IA-32 Architecture processors are always in IA-32 Mode which consists of the following SubModes:<ul style="list-style-type: none">– Real Mode.– System Management Mode (SMM).– Protected Mode.– VM86 Mode.• At a given moment in time, an Intel 64 Architecture processor is operating in either:<ul style="list-style-type: none">– IA-32 Mode, or– IA-32e (IA-32 Extended) Mode. <p>“IA-32 SubModes” on page 25 describes the IA-32 execution SubModes.</p> <p>Problems associated with IA-32 Mode:</p> <p>Some of the problems associated with IA-32 Mode are:</p> <ul style="list-style-type: none">• The instruction set syntax uses a 3-bit field to specify a source or destination register. As a result, there are only eight addressable General Purpose Registers (GPRs), Control registers, Debug registers, or XMM registers.• The maximum width of each GPR is 32-bits limiting the amount of data each can hold.• Virtual memory address space available for each application is limited to 4GB by the 32-bit width of the linear (i.e., virtual) address.• The virtual-to-physical memory address translation mechanism limits the maximum addressable physical memory address space to 64GB.• The 32-bit Extended Instruction Pointer (EIP) register limits each application’s code space to 4GB.• The x86 family’s segmented memory model is complex and difficult to work with. Virtually all of today’s OSs utilize a Flat Memory Model that effectively disables the segmented memory model.• The hardware-assisted task switching mechanism defined by the IA-32 ISA is slow and cumbersome.• IA-32 Mode permits virus code to be loaded into a stack or data segment from which it can then be executed.• Lacks the ability to address code-local data by specifying an address relative to the current EIP value.

x86 Instruction Set Architecture

Table 2-1: Basic Execution Modes (Continued)

Mode	Description
IA-32e Mode	<p>IA-32 Extended Mode is comprised of two submodes:</p> <ul style="list-style-type: none">– 64-bit Mode.– Compatibility Mode. <ul style="list-style-type: none">• Must be enabled by a 64-bit capable OS.• Provides an environment for the execution of 64-bit applications, as well as existing 32- and 16-bit Protected Mode applications.• Doesn't support the execution of VM86 Mode applications (i.e., MS-DOS applications).• Provides a fast transition between a 32-bit environment (Compatibility Mode) and a 64-bit environment (64-bit Mode).• Implements the Intel 64 extensions (formerly known as x86-64 or EM64T). <p>"IA-32e SubModes" on page 28 describes the IA-32e execution SubModes.</p> <p>Some benefits associated with IA-32e Mode:</p> <p>The following are some of the benefits realized when the logical processor is executing in IA-32e Mode:</p> <ul style="list-style-type: none">• Backward compatible with the IA-32 code environment. Intel's earlier attempt at a 64-bit architecture (Itanium) is not.• Expands the size of the virtual memory address space from 2^{32} (4GB) to 2^{64} (16EB; EB = exabytes).• Expands the size of the physical memory address space to 2^{52} (4PB; PB = petabytes).• The larger number of data registers permits a greater number of data variables to be accessed/manipulated rapidly:<ul style="list-style-type: none">– Yields faster data set accessibility.– Widening and increasing the number of registers diminishes the number of accesses to memory and translates into improved performance.– The degree of improvement in kernel code efficiency depends on a kernel rewrite to manage memory better and to utilize 64-bit (rather than 32-bit) data variables.– The degree of improvement in application code efficiency depends on utilization of 64-bit data variables and, for large scale applications, capitalizing on the enlarged virtual address space.

3 *A (very) Brief History*

The Previous Chapter

The previous chapter introduced the execution modes and submodes and mode switching basics.

This Chapter

This chapter introduces the evolution of the x86 ISA, as well as the basic operational characteristics of 8086 Real Mode, 286 Protected Mode, and 386 Protected Mode. It also introduces the Intel microarchitecture families including a product introduction timeline.

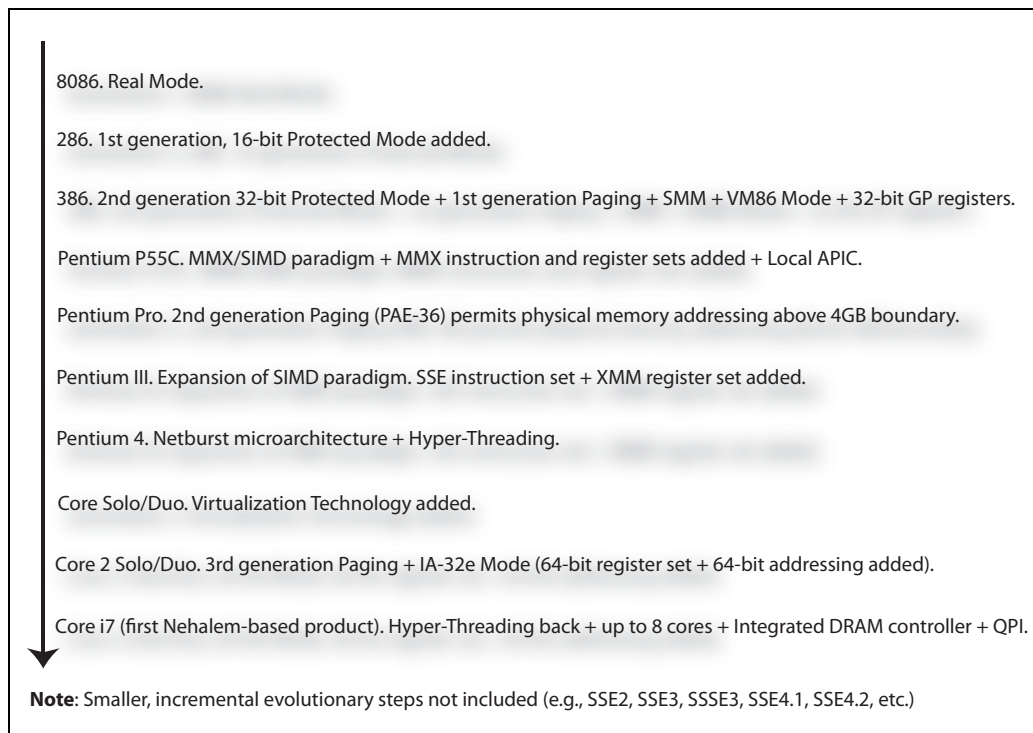
The Next Chapter

The next chapter defines the state of a logical processor immediately after the removal of reset and introduces the concept of a soft reset (also referred to as an INIT). It also describes the initial code fetches performed by the BootStrap Processor as well as the methodology utilized by software to discover and configure all of the logical processors in the system.

Major Evolutionary Developments

Through the years, the x86 software architecture has steadily evolved with the introduction of new x86 processors. Some changes were small, evolutionary ones while others made significant additions to the architecture. Figure 3-1 on page 42 illustrates (and Table 3-1 on page 43 describes) those that, in the author's opinion, fall into the latter category.

Figure 3-1: Major Milestones in Evolution of Software Environment



Chapter 3: A (very) Brief History

Table 3-1: Major Evolutionary Developments

Introduced In	Major Enhancements to the x86 Instruction Set Architecture
8086	<p>At the root of the x86 family tree lies the 8086, a processor with the following characteristics:</p> <ul style="list-style-type: none">• Modes: Real Mode.• Addressable Memory: 1MB.• Data Transfer Width: 16-bits.• Programming Model: 16-bit.
286	<p>Next, we have the 286, a processor with the following characteristics:</p> <ul style="list-style-type: none">• Modes: Real Mode and first generation 16-bit Protected Mode.• Missing three critical capabilities:<ul style="list-style-type: none">– It did not implement a virtual-to-physical address translation facility (i.e., the Paging mechanism).– Could not address more than 16MB of physical memory.– It did not implement Virtual 8086 (VM86) Mode, an exemption that effectively crippled the processor's ability to run ill-behaved DOS applications under a multi-tasking OS.• Addressable Physical Memory: 16MB.• Data Transfer Width: 16-bits.• Programming Model: 16-bit (16-bit GP registers and 16-bit addressing).

x86 Instruction Set Architecture

Table 3-1: Major Evolutionary Developments (Continued)

Introduced In	Major Enhancements to the x86 Instruction Set Architecture
386	<p>The introduction of the 386 contributed the following major architectural changes:</p> <ul style="list-style-type: none">• 2nd generation 32-bit Protected Mode. Permits 32-bit addressing (rather than the 16-bit addressing supported by the 286's first generation Protected Mode).• 1st generation virtual-to-physical address translation mechanism (i.e., Paging). Permitted a 32-bit virtual memory address to be translated into a 32-bit physical memory address.• System Management Mode (SMM) first appeared in the 386SX processor. If the platform logic detects a platform-specific issue (e.g., a thermal zone is warming up), the chipset generates an SM Interrupt (SMI) to the processor which interrupts the currently-running program, saves the processor's register set and executes the SMM handler. The handler checks chipset status to determine the nature of the problem, handles the problem (e.g., by turning on a fan) and then restores the processor's register set and resumes execution of the interrupted program.• VM86 Mode. This mechanism permits the processor hardware to monitor the execution of ill-behaved DOS applications on an instruction-by-instruction basis. If an instruction that could destabilize the multi-tasking OS is detected, the DOS program is interrupted and a special program, the VMM (Virtual Machine Monitor), is executed to determine the nature of the problem and fix it. The OS then resumes execution of the DOS application.• Addressable Physical Memory: 4GB.• Data Transfer Width: 32-bits.• Programming Model: 32-bit (32-bit GP registers and 32-bit addressing).

4 *State After Reset*

The Previous Chapter

The previous chapter introduced the evolution of the x86 ISA, as well as the basic operational characteristics of 8086 Real Mode, 286 Protected Mode, and 386 Protected Mode. It also introduced the Intel microarchitecture families including a product introduction timeline.

This Chapter

This chapter defines the state of a logical processor immediately after the removal of reset and introduces the concept of a soft reset (also referred to as an INIT). It also describes the initial code fetches performed by the BootStrap Processor and introduces the methodology utilized by software to discover and configure all of the logical processors in the system.

The Next Chapter

The next chapter provides a very basic introduction to the various facilities that support the IA-32 computing environment. These facilities include:

- Pre-386 Register Sets (this section is provided for historical background).
- IA-32 Register Set Overview.
- Control Registers.
- Status/Control Register (Eflags).
- Instruction Fetch Facilities.
- General Purpose Data Registers.
- Defining Memory Regions/Characteristics.
- Interrupt/Exception Facilities.
- Kernel Facilities.
- Address Translation Facilities.
- Legacy FP Facilities.
- MMX Facilities.

x86 Instruction Set Architecture

- SSE Facilities.
- Model-Specific Registers.
- Debug Facilities.
- Automatic Task Switching Mechanism.

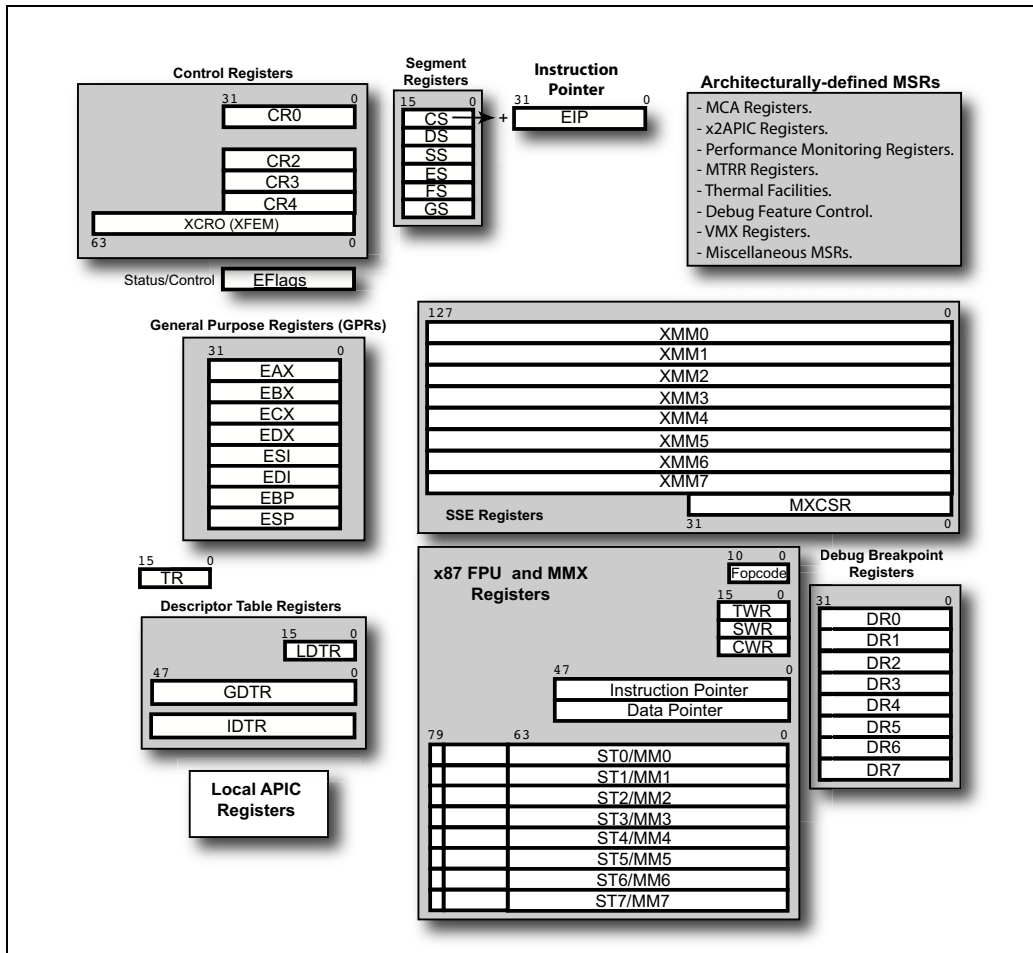
State After Reset

Table 4-1 on page 66 defines the state of the logical processor's registers (the IA-32 register set is shown in Figure 4-1 on page 65) and resources immediately after the removal of reset. To summarize:

- The logical processor is in Real Mode (Protected Mode and Paging are disabled).
- Its caches are empty and caching is disabled.
- All of the feature bits in CR4 are cleared disabling most of the new features introduced after the advent of the 386.
- Recognition of external hardware interrupts is disabled.
- No instructions have been fetched from memory.
- The x87 FPU is disabled.
- All x87 FPU and SSE exceptions are disabled.
- The Machine Check and Alignment Check exceptions are disabled.
- The first instruction will be fetched from location FFFFFFF0h.

Chapter 4: State After Reset

Figure 4-1: IA-32 Register Set



x86 Instruction Set Architecture

Table 4-1: Logical Processor State After Removal of Reset

Register or Resource	Effect(s)
Instruction Pipeline	No instructions have been fetched from memory yet.
ROB	Since no instructions have been fetched from memory yet to be translated into micro-ops, the Reorder Buffer is empty and the instruction dispatch logic is idle.
BTB Cache	The Branch Target Buffer maintains history on branch execution (i.e., whether branches were taken or not taken) is empty.
CR0 register	<p>Contains 60000010h after reset:</p> <ul style="list-style-type: none">• CR0[PE] = 0, disabling Protected Mode. The logical processor is therefore in Real Mode.• CR0[PG] = 0, disabling Paging (virtual-to-physical address translation services).• CR0[CD&NW] = 11b, disabling the caching logic.• CR0[EM&MP] = 11b, indicating the x87 FPU isn't present and the logical processor should therefore permit software to emulate it (by executing integer-only code).• CR0[TS] = 0, indicating a task switch has not occurred.• CR0[ET] = 1, indicating that the integrated x87 FPU is compatible with the 387 FPU.• CR0[NE] = 0, disabling the logical processor's ability to generate an exception 16 if an x87 FP exception occurs. Instead, the logical processor signals the event using the DOS-compatible method (i.e., by asserting the processor's FERR# output which causes the assertion of IRQ13 to the 8259A PIC (Programmable Interrupt Controller).• CR0[WP] = 0. This has no effect at this time because Paging is disabled. When paging is enabled, setting this bit to a one prevents privilege level 0 software from writing to read-only pages.• CR0[AM] = 0. Clearing the Alignment Mask bit disables the logical processor's ability to generate the Alignment Check exception when a mis-aligned multi-byte memory access is detected.

5 *Intro to the IA-32 Ecosystem*

The Previous Chapter

The previous chapter defined the state of a logical processor immediately after the removal of reset and introduced the concept of a soft reset (also referred to as an INIT). It also described the initial code fetches performed by the BootStrap Processor and introduced the methodology utilized by software to discover and configure all of the logical processors in the system.

This Chapter

This chapter provides a very basic introduction to the various facilities that support the IA-32 computing environment. These facilities include:

- Pre-386 Register Sets (this section is provided for historical background).
- IA-32 Register Set Overview.
- Control Registers.
- Status/Control Register (Eflags).
- Instruction Fetch Facilities.
- General Purpose Data Registers.
- Defining Memory Regions/Characteristics.
- Interrupt/Exception Facilities.
- Kernel Facilities.
- Address Translation Facilities.
- Legacy FP Facilities.
- MMX Facilities.
- SSE Facilities.
- Model-Specific Registers.
- Debug Facilities.
- Automatic Task Switching Mechanism.

The Next Chapter

The next chapter highlights the dramatic expansion of the x86 instruction set since the advent of the 386 by listing both the 386 instruction set as well as the current-day instruction set.

The Pre-386 Register Sets

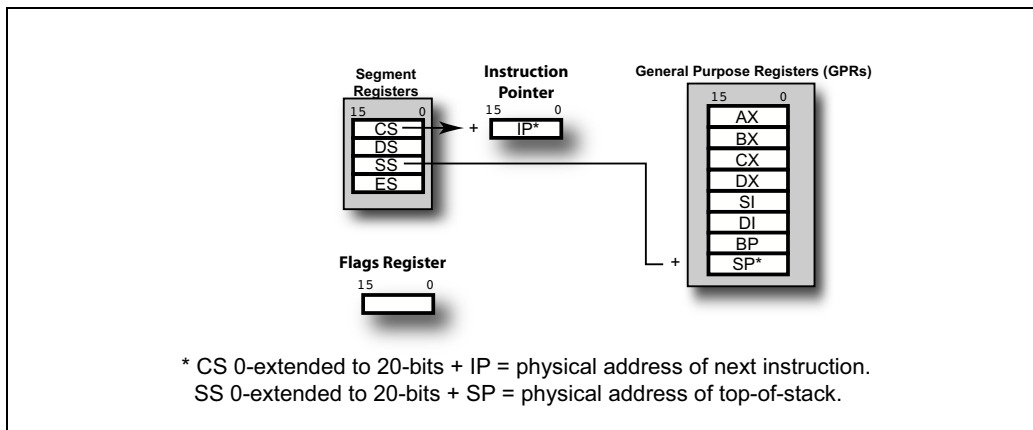
This section provides a little background regarding the baseline register set that was a precursor to the expanded register set found in today's x86 processors. Basic descriptions of these registers are included in this chapter.

8086 Register Set

The 8086 register set (see Figure 5-1 on page 80) consisted of the following registers:

- Eight General Purpose Registers (GPRs). See Figure 5-2 on page 81.
- Flags register. See Figure 5-3 on page 81.
- Four Segment registers.
- Instruction Pointer (IP) register.

Figure 5-1: 8086 Register Set



Chapter 5: Intro to the IA-32 Ecosystem

Figure 5-2: 8086 GPRs

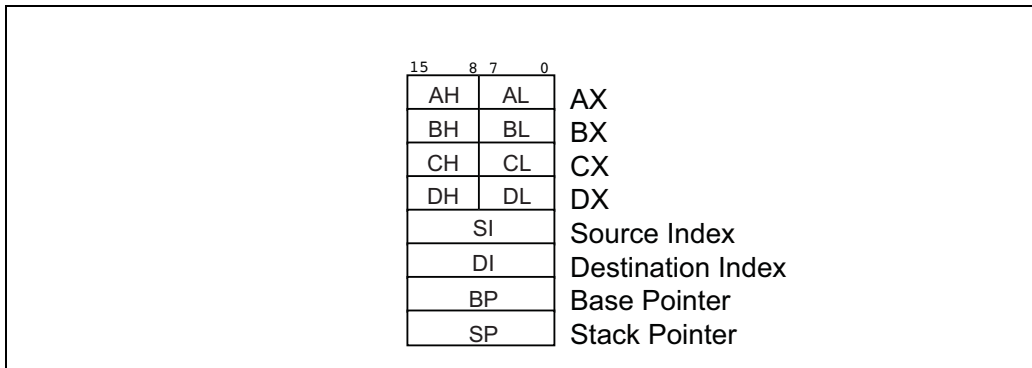
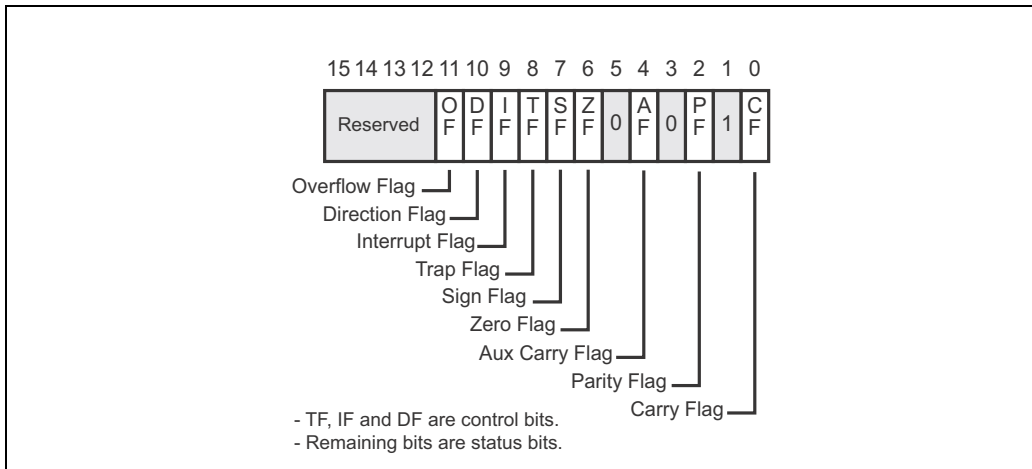


Figure 5-3: 8086 Flag Register



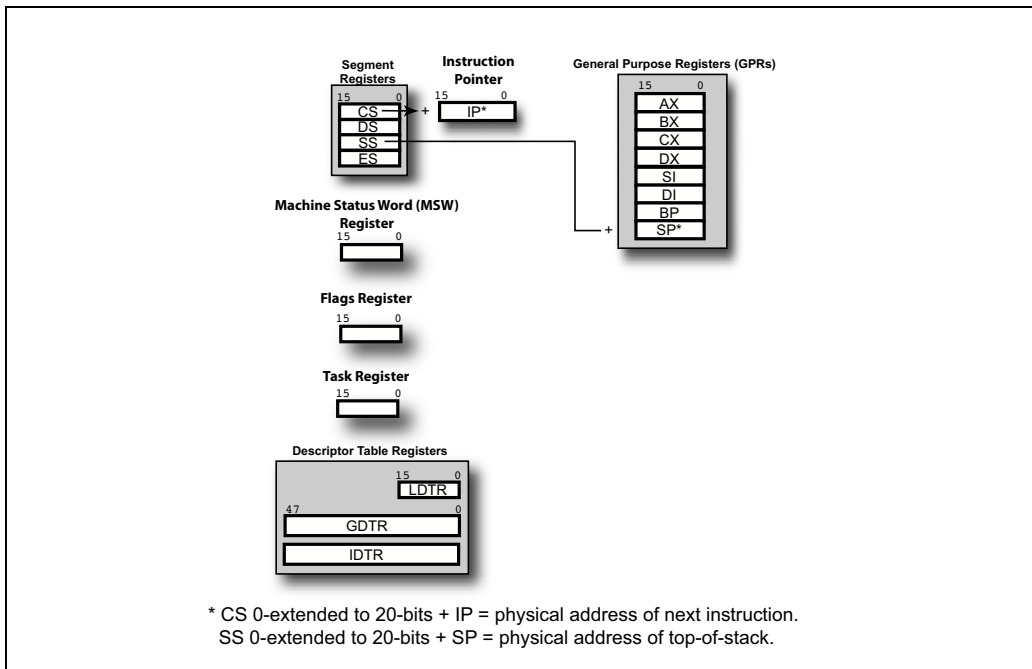
286 Register Set

The 286 added the following registers (see Figure 5-4 on page 82):

- Machine Status Word (MSW) register. See Figure 5-5 on page 83.
- Kernel-related registers:
 - Task Register (TR).
 - Interrupt Descriptor Table Register (IDTR).
 - Global Descriptor Table Register (GDTR).
 - Local Descriptor Table Register (LDTR).

It also added two additional bit fields (Nested Task and IO Privilege Level) to the Flags register (see Figure 5-6 on page 83).

Figure 5-4: 286 Register Set



6 *Instruction Set Expansion*

The Previous Chapter

The previous chapter provided a very basic introduction to the various facilities that support the IA-32 computing environment. These facilities include:

- Pre-386 Register Sets (this section is provided for historical background).
- IA-32 Register Set Overview.
- Control Registers.
- Status/Control Register (Eflags).
- Instruction Fetch Facilities.
- General Purpose Data Registers.
- Defining Memory Regions/Characteristics.
- Interrupt/Exception Facilities.
- Kernel Facilities.
- Address Translation Facilities.
- Legacy FP Facilities.
- MMX Facilities.
- SSE Facilities.
- Model-Specific Registers.
- Debug Facilities.
- Automatic Task Switching Mechanism.

This Chapter

This chapter illustrates the expansion of the x86 instruction set since the advent of the 386 by listing both the 386 instruction set as well as the current-day instruction set.

The Next Chapter

The next chapter provides a detailed explanation of the structure of an IA-32 instruction and covers the following topics:

- Effective Operand Size.
- Instruction Composition.
- Instruction Format Basics.
- Opcode (Instruction Identification).
 - In the Beginning.
 - 1-byte Opcodes.
 - 2-byte Opcodes Use 2-Level Lookup.
 - 3-byte Opcodes Use 3-Level Lookup.
 - Opcode Micro-Maps (Groups).
 - x87 FP Opcodes Inhabit Opcode Mini-Maps.
 - Special Opcode Fields.
- Operand Identification.
 - Specifying Registers as Operands.
 - Addressing a Memory-Based Operand.
 - Specifying an Immediate Value As an Operand.
- Instruction Prefixes.
 - Operand Size Override Prefix (66h).
 - Address Size Override Prefix (67h).
 - Lock Prefix.
 - Repeat Prefixes.
 - Segment Override Prefix.
 - Branch Hint Prefix.
- Summary of Instruction Set Formats.

Why a Comprehensive Instruction Set Listing Isn't Included

Since the Intel and AMD x86 instruction set reference guides already do a fine job fulfilling this role, this chapter does not provide a comprehensive description of each instruction in the x86 instruction set. Rather, it is intended as an introduction to the instruction set. To lend historical perspective, it begins with a listing of the entire 386 instruction set (all 128 of them) and then continues with a listing of the current instruction set (well over 400 instructions as of March 2009) sorted by category. It should be stressed that the instruction set is constantly evolving, maintaining backward-compatibility even as successive generations of x86 processors continue to add new instructions—sometimes in

Chapter 6: Instruction Set Expansion

small numbers; at other times, with substantial additions to the instruction repertoire (e.g., MMX and SSE).

386 Instruction Set

To lend historical perspective, Table 6-1 on page 111 lists the 386 processor's instruction set organized by category.

Table 6-1: 386 Instruction Set

Instruction	Description
Data Transfer—General Purpose	
MOV	Move operand
PUSH	Push operand onto stack
POP	Pop operand off stack
PUSHA	Push all registers on stack
POPA	Pop all registers off stack
XCHG	Exchange Operand, Register
XLAT	Translate
Data Transfer—Conversion	
MOVZX	Move byte or Word, DW, with zero extension
MOVSX	Move byte or Word, DW, sign extended
CBW	Convert byte to Word, or Word to DW
CWD	Convert Word to DW
CWDE	Convert Word to DW extended
CDQ	Convert DW to QW
Data Transfer—Input/Output	
IN	Input operand from I/O space

x86 Instruction Set Architecture

Table 6-1: 386 Instruction Set (Continued)

Instruction	Description
OUT	Output operand to I/O space
Data Transfer—Address Object	
LEA	Load effective address
LDS	Load pointer into D segment register
LES	Load pointer into E segment register
LFS	Load pointer into F segment register
LGS	Load pointer into G segment register
LSS	Load pointer into S (Stack) segment register
Data Transfer—Flag Manipulation	
LAHF	Load A register from Flags
SAHF	Store A register in Flags
PUSHF	Push flags onto stack
POPF	Pop flags off stack
PUSHFD	Push Eflags onto stack
POPFD	Pop Eflags off stack
CLC	Clear Carry Flag
CLD	Clear Direction Flag
CMC	Complement Carry Flag
STC	Set Carry Flag
STD	Set Direction Flag
Arithmetic Instructions - Addition	
ADD	Add operands
ADC	Add with carry

7

32-bit Machine Language Instruction Format

The Previous Chapter

The previous chapter illustrated the expansion of the x86 instruction set since the advent of the 386 by listing both the 386 instruction set as well as the current-day instruction set.

This Chapter

This chapter provides a detailed explanation of the structure of an IA-32 instruction and covers the following topics:

- Effective Operand Size.
- Instruction Composition.
- Instruction Format Basics.
- Opcode (Instruction Identification).
 - In the Beginning.
 - 1-byte Opcodes.
 - 2-byte Opcodes Use 2-Level Lookup.
 - 3-byte Opcodes Use 3-Level Lookup.
 - Opcode Micro-Maps (Groups).
 - x87 FP Opcodes Inhabit Opcode Mini-Maps.
 - Special Opcode Fields.
- Operand Identification.
 - Specifying Registers as Operands.
 - Addressing a Memory-Based Operand.
 - Specifying an Immediate Value As an Operand.
- Instruction Prefixes.
 - Operand Size Override Prefix (66h).
 - Address Size Override Prefix (67h).
 - Lock Prefix.
 - Repeat Prefixes.

x86 Instruction Set Architecture

- Segment Override Prefix.
- Branch Hint Prefix.
- Summary of Instruction Set Formats.

The Next Chapter

The next chapter provides a detailed description of Real Mode operation and covers the following topics:

- 8086 Emulation.
- Unused Facilities.
- Real Mode OS Environment.
- Running Real Mode Applications Under a Protected Mode OS.
- Real Mode Applications Aren't Supported in IA-32e Mode.
- Real Mode Register Set.
- IO Space versus Memory Space.
- IO and Memory-Mapped IO Operations.
- Operand Size Selection.
- Address Size Selection.
- Real Mode Memory Addressing.
- Real Mode Interrupt/Exception Handling.
- Summary of Real Mode Limitations.

64-bit Machine Language Instruction Format

As its name implies, the current chapter provides a detailed description of the 32-bit machine language instruction format. The 64-bit extensions to the machine language instruction format are covered in:

- “64-bit Operands and Addressing” on page 1041.
- “64-bit Odds and Ends” on page 1075.

A Complex Instruction Set with Roots in the Past

As mentioned earlier in “IA Instructions vs. Micro-ops” on page 15, the x86 machine language instruction set is quite complex. Depending on the type of instruction, the number of operands it specifies, and the operand types (memory- and/or register-based), a single instruction may consist of anywhere between one and fifteen bytes. Beginning with the advent of the Pentium Pro processor, all x86 processors incorporate a translator that converts each IA-32

Chapter 7: 32-bit Machine Language Instruction Format

machine language instruction into a series of one or more simple, fixed-length micro-ops which are then executed by the logical processor.

Effective Operand Size

Introduction

In order to limit the number of opcodes, the same opcode is used for an instruction whether it operates on an 8-, 16- or 32-bit operand. As an example:

```
mov ax,dx
mov eax,edx
```

both use the same basic opcode. This naturally brings up a question: how, then, does the logical processor determine which registers are being referenced? The answer is simple and is described in the next two sections.

Operand Size in 16- and 32-bit Code Segments

Assuming that an instruction is not prefaced by a Operand Size Override prefix byte (66h), the logical processor behaves as outlined in Table 7-1 on page 157. Adding the prefix byte before the instruction's first opcode byte alters its behavior as defined in Table 7-2 on page 158.

Table 7-1: Effective Operand Size in 16- or 32-bit Mode (without prefix)

State of D-bit in active CS Descriptor	Effective Operand Size and Instruction Behavior
0: 16-bit, 286 CS descriptor.	16-bits. Instruction operates on 16-bits (a word) in a 16-bit register (e.g., AX) and either of the following: <ul style="list-style-type: none">• Two sequential memory locations.• Another 16-bit register.
1: 32-bit, 386 CS descriptor.	32-bits. Instruction operates on 32-bits (a dword) in a 32-bit register (e.g., EAX) and either of the following: <ul style="list-style-type: none">• Four sequential memory locations.• Another 32-bit register.

x86 Instruction Set Architecture

Table 7-1: Effective Operand Size in 16- or 32-bit Mode (without prefix) (Continued)

State of D-bit in active CS Descriptor	Effective Operand Size and Instruction Behavior
For some instructions, the opcode contains a width (W) bit: <ul style="list-style-type: none">W = 0. The operand size is 8- rather than 16- or 32-bits.W = 1. The operand size is either 16- or 32-bits (based on the state of CSDesc[D] and the presence or absence of the Operand Size Override prefix).	

Table 7-2: Effective Operand Size in 16- or 32-bit Mode (with prefix)

State of D-bit in active CS Descriptor	Effective Operand Size and Instruction Behavior
0: 16-bit, 286 CS descriptor.	Inclusion of the Operand Size Override prefix before the instruction flips the effective operand size from 16- to 32-bits. Instruction operates on 32-bits (a dword) in a 32-bit register (e.g., EAX) and either of the following: <ul style="list-style-type: none">Four sequential memory locations.Another 32-bit register.
1: 32-bit, 386 CS descriptor.	Inclusion of the Operand Size Override prefix before the instruction flips the effective operand size from 32- to 16-bits. Instruction operates on 16-bits (a word) in a 16-bit register (e.g., AX) and either of the following: <ul style="list-style-type: none">Two sequential memory locations.Another 16-bit register.
For some instructions, the opcode contains a width (W) bit: <ul style="list-style-type: none">W = 0. The operand size is 8- rather than 16- or 32-bits.W = 1. The operand size is either 16- or 32-bits (based on the state of CSDesc[D] and the presence or absence of the Operand Size Override prefix).	

Operand Size in 64-bit Code Segments

The default data operand size when executing code from a 64-bit code segment (code segment descriptor's L bit = 1) is 32-bits and its default address size is 64-bits. In other words, unless instructed otherwise, the logical processor assumes

8 *Real Mode* (*8086 Emulation*)

The Previous Chapter

The previous chapter provided a detailed explanation of the structure of an IA-32 instruction and covered the following topics:

- Effective Operand Size.
- Instruction Composition.
- Instruction Format Basics.
- Opcode (Instruction Identification).
 - In the Beginning.
 - 1-byte Opcodes.
 - 2-byte Opcodes Use 2-Level Lookup.
 - 3-byte Opcodes Use 3-Level Lookup.
 - Opcode Micro-Maps (Groups).
 - x87 FP Opcodes Inhabit Opcode Mini-Maps.
 - Special Opcode Fields.
- Operand Identification.
 - Specifying Registers as Operands.
 - Addressing a Memory-Based Operand.
 - Specifying an Immediate Value As an Operand.
- Instruction Prefixes.
 - Operand Size Override Prefix (66h).
 - Address Size Override Prefix (67h).
 - Lock Prefix.
 - Repeat Prefixes.
 - Segment Override Prefix.
 - Branch Hint Prefix.
- Summary of Instruction Set Formats.

This Chapter

This chapter provides a detailed description of Real Mode operation and covers the following topics:

- 8086 Emulation.
- Unused Facilities.
- Real Mode OS Environment.
- Running Real Mode Applications Under a Protected Mode OS.
- Real Mode Applications Aren't Supported in IA-32e Mode.
- Real Mode Register Set.
- IO Space versus Memory Space.
- IO and Memory-Mapped IO Operations.
- Operand Size Selection.
- Address Size Selection.
- Real Mode Memory Addressing.
- Real Mode Interrupt/Exception Handling.
- Summary of Real Mode Limitations.

The Next Chapter

The next chapter provides a detailed description of the x87 FPU and covers the following topics:

- A Little History.
- x87 FP Instruction Format.
- FPU-Related CR0 Bit Fields.
- x87 FPU Register Set.
 - The FP Data Registers.
 - x87 FPU's Native Data Operand Format.
 - 32-bit SP FP Numeric Format.
 - DP FP Number Representation.
 - FCW Register.
 - FSW Register.
 - FTW Register.
 - Instruction Pointer Register.
 - Data Pointer Register.
 - Fopcode Register.
- FP Error Reporting.
 - Precise Error Reporting.
 - Imprecise (Deferred) Error Reporting.
 - Why Deferred Error Reporting Is Used.
 - The WAIT/FWAIT Instruction.
 - CR0[NE].
 - Ignoring FP Errors.

Chapter 8: Real Mode (8086 Emulation)

8086 Emulation

Real Mode was introduced with the advent of the 8086/8088 processors, and, due to the huge success of the IBM PC and the proliferation of software written for the Real Mode environment, Intel could ill-afford to leave it behind. Immediately after the removal of reset, all subsequent x86 processors emulate the operation of the 8086 by initiating operation in Real Mode. There are, of course, some differences (see Table 8-1 on page 229).

Table 8-1: Basic Differences Between 8086 Operation and Real Mode

Difference	Basic Description
Speed of execution	In the IBM PC, the processor clock ran at 4.77MHz. Today's processors execute Real Mode code hundreds of times faster.
Cache boost	Today's processors enjoy a substantial performance boost due to on-chip caches.
Accessing extended memory	While the 8086/8088 processors were strictly limited to a 1MB address space due to an address bus width of 20-bits, current-day x86 processors can, even in Real Mode, access significantly more memory: <ul style="list-style-type: none">• See "Accessing Extended Memory in Real Mode" on page 307.• See "Big Real Mode" on page 310.
Operand and address size	Operand Size. While the default data operand size in Real Mode is 16-bits, 32-bit operands can be specified by prefacing an instruction with the Operand Size Override prefix. The logical processor can then access 32-bit operands in memory as well as the 32-bit GPR registers: EAX, EBX, ECX, EDX, ESP, EBP, ESI, and EDI. The 8086/8088 GPR registers were only 16-bits wide. Address Size. While the default address size for memory-based operands in Real Mode is 16-bits, a 32-bit address can be specified by prefacing an instruction with the Address Size Override prefix.

x86 Instruction Set Architecture

Table 8-1: Basic Differences Between 8086 Operation and Real Mode (Continued)

Difference	Basic Description
Number of memory data segments	Using the ES, FS and GS Segment Override prefixes, the programmer can access memory-based data operands in the E, F and G data segments. The 8086/8088 only implemented the DS data segment.
Integrated x87 FPU	Today's logical processors incorporate an integrated x87 FPU. The 8087 FPU was implemented as an optional, external companion device to the 8086/8088. Upon detection of a FP instruction, the processor had to forward the instruction to the x87 FPU by performing a series of IO write transactions on its external interface. Very slow, indeed.
Debug register set	The address breakpoint facility is available through the Debug register set (not implemented on the 8086/8088).
Additional instructions available	The instructions listed in Table 8-2 on page 230, although not available on the 8086/8088, can be used in Real Mode.

Table 8-2: Expanded/Enhanced Real Mode Instructions

Instructions Available (that weren't present in the 8086/8088)
The MMX instruction set as well as the MMX data registers (MM0 - MM7).
The SSE1, SSE2, SSE3, SSSE3, SSE4.1 and SSE4.2 instruction sets as well as the SSE register set and the SSE FP exception (exception 19).
MOV instructions that operate on the Control and Debug registers.
Load segment register instructions: LSS, LFS, and LGS.
Generalized multiply and multiply immediate data instructions.
Shift and rotate by immediate counts.
PUSHA, PUSHAD, POPA and POPAD, and PUSH immediate data stack instructions.
MOVSX and MOVZX Move with sign extension instructions.

9 *Legacy x87 FP Support*

The Previous Chapter

The previous chapter provided a detailed description of Real Mode operation and covered the following topics:

- 8086 Emulation.
- Unused Facilities.
- Real Mode OS Environment.
- Running Real Mode Applications Under a Protected Mode OS.
- Real Mode Applications Aren't Supported in IA-32e Mode.
- Real Mode Register Set.
- IO Space versus Memory Space.
- IO and Memory-Mapped IO Operations.
- Operand Size Selection.
- Address Size Selection.
- Real Mode Memory Addressing.
- Real Mode Interrupt/Exception Handling.
- Summary of Real Mode Limitations.

This Chapter

This chapter provides a detailed description of the x87 FPU and covers the following topics:

- A Little History.
- x87 FP Instruction Format.
- FPU-Related CR0 Bit Fields.
- x87 FPU Register Set.
 - The FP Data Registers.
 - x87 FPU's Native Data Operand Format.

x86 Instruction Set Architecture

- 32-bit SP FP Numeric Format.
- DP FP Number Representation.
- FCW Register.
- FSW Register.
- FTW Register.
- Instruction Pointer Register.
- Data Pointer Register.
- Fopcode Register.
- FP Error Reporting.
 - Precise Error Reporting.
 - Imprecise (Deferred) Error Reporting.
 - Why Deferred Error Reporting Is Used.
 - The WAIT/FWAIT Instruction.
 - CR0[NE].
 - Ignoring FP Errors.

The Next Chapter

The next chapter provides an introduction to the concept of multitasking and covers the following topics:

- Concept.
- An Example—Timeslicing.
- Another Example—Awaiting an Event.
 - 1. Task Issues Call to OS for Disk Read.
 - 2. Device Driver Initiates Disk Read.
 - 3. OS Suspends Task.
 - 4. OS Makes Entry in Event Queue.
 - 5. OS Starts or Resumes Another Task.
 - 6. Disk-Generated Interrupt Causes Jump to OS.
 - 7. Interrupted Task Suspended.
 - 8. Task Queue Checked.
 - 9. OS Resumes Task.

A Little History

Prior to the advent of the 486DX processor, x86 processors did not include an on-die FPU. In order to perform floating-point (FP) math operations the end user had to add an external x87 FPU chip to the system which the processor treated as a specialized IO device. When the processor encountered a FP instruction in the currently-executing program, it would perform a series of one or more IO write transactions on its external bus to forward the instruction to the off-chip FPU for execution. Obviously, this was very inefficient.

Chapter 9: Legacy x87 FP Support

The 486DX was the first x86 processor to integrate the x87 FPU (all subsequent x86 processors include it). The sections that follow provide a description of the FPU's register set and the format in which FP numbers are represented.

x87 FP Instruction Format

This topic is covered in “x87 FP Opcodes Inhibit Opcode Mini-Maps” on page 187.

FPU-Related CR0 Bit Fields

Refer to Figure 9-1 on page 341, Table 9-1 on page 341, and Table 8-4 on page 241 for a description of the CR0 bit fields related to the x87 FPU.

Figure 9-1: CR0

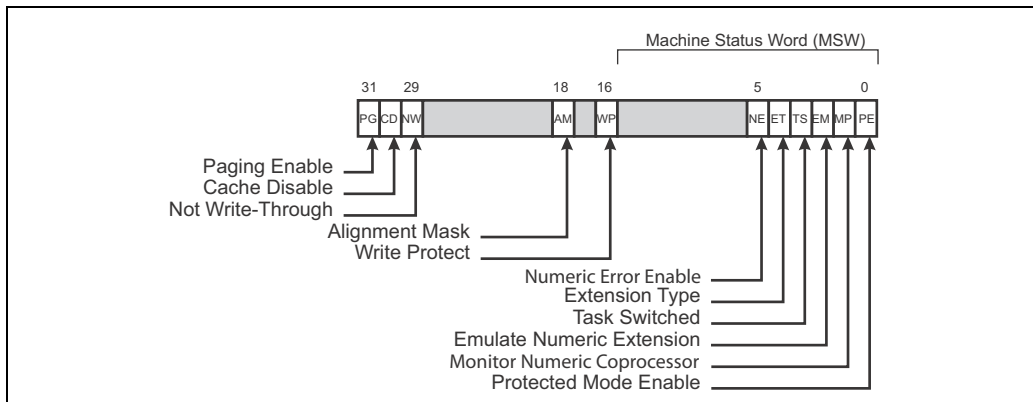


Table 9-1: CR0 x87 FPU-related Bit Fields

Bit(s)	Name	Description
2:1	EM, MP	The OS uses these x87 FPU-related bits to indicate whether the logical processor is: <ul style="list-style-type: none">• Running DOS.• Running a multitasking OS.• Neither of the above (the FPU isn't present or is disabled); in this case, a software exception is generated when the logical processor detects an x87 instruction and software emulates the x87 FPU. See Table 8-4 on page 241.

x86 Instruction Set Architecture

Table 9-1: CR0 x87 FPU-related Bit Fields (Continued)

Bit(s)	Name	Description
3	TS	x87 Task Switch status bit. The x87 FPU registers are not saved on an automatic task switch. When a hardware-based task switch occurs: <ul style="list-style-type: none">• CR0[TS] is set to one.• Most but not all registers are automatically saved in the TSS (Task State Segment) data structure associated with the currently-running task.• When an attempt is made to execute an x87 FP instruction or an MMX instruction while CR0[TS] = 1, a DNA (Device Not Available) Exception 7 is generated.• The DNA exception handler executes FSAVE or FXSAVE to save the x87 and, possibly, the SSE registers in the TSS of the task that last used the x87. This information is saved in an OS-designated area of the previous task's TSS.• The DNA exception handler then clears CR0[TS] and executes an IRET instruction to return to the x87 instruction in the current task that caused the DNA exception. It now executes successfully.
4	ET	x87 FPU type. In processors prior to the 486DX, CR0[ET] (ET = Extension Type) was a read/write bit used by software to indicate the type of numeric coprocessor installed on the system board (287 or 387 FPU-compatible). Since the advent of the 486DX this bit is hardwired to one indicating that the logical processor incorporates a 387-style FPU.
5	NE	Numeric Exception. Controls whether FP errors are reported using the DOS-compatible method (IRQ13; see “DOS-Compatible FP Error Reporting” on page 359) or by generating an exception 16. The OS kernel sets NE = 1 if it incorporates an x87 FP exception handler. Any x87 FPU error then causes the logical processor to generate an internal exception 16 (rather than using the DOS-compatible method and asserting its external FERR# signal).

10 Introduction to Multitasking

The Previous Chapter

The previous chapter provided a detailed description of the x87 FPU and covered the following topics:

- A Little History.
- x87 FP Instruction Format.
- FPU-Related CR0 Bit Fields.
- x87 FPU Register Set.
 - The FP Data Registers.
 - x87 FPU's Native Data Operand Format.
 - 32-bit SP FP Numeric Format.
 - DP FP Number Representation.
 - FCW Register.
 - FSW Register.
 - FTW Register.
 - Instruction Pointer Register.
 - Data Pointer Register.
 - Fopcode Register.
- FP Error Reporting.
 - Precise Error Reporting.
 - Imprecise (Deferred) Error Reporting.
 - Why Deferred Error Reporting Is Used.
 - The WAIT/FWAIT Instruction.
 - CR0[NE].
 - Ignoring FP Errors.

This Chapter

This chapter provides an introduction to the concept of multitasking and covers the following topics:

- Concept.
- An Example—Timeslicing.
- Another Example—Awaiting an Event.
 - 1. Task Issues Call to OS for Disk Read.
 - 2. Device Driver Initiates Disk Read.
 - 3. OS Suspends Task.
 - 4. OS Makes Entry in Event Queue.
 - 5. OS Starts or Resumes Another Task.
 - 6. Disk-Generated Interrupt Causes Jump to OS.
 - 7. Interrupted Task Suspended.
 - 8. Task Queue Checked.
 - 9. OS Resumes Task.

The Next Chapter

The next chapter introduces the concepts of hardware-based task switching, global and local memory, privilege checking, read/write protection, IO port protection, interrupt masking, and BIOS call interception. The following topics are covered:

- Hardware-based Task Switching Is Slow!
- Private (Local) and Global Memory.
- Preventing Unauthorized Use of OS Code.
- With Privilege Comes Access.
 - Program Privilege Level.
 - The CPL.
 - Calling One of Your Equals.
 - Calling a Procedure to Act as Your Surrogate.
 - Data Segment Protection.
 - Data Segment Privilege Level.
 - Read-Only Data Areas.
- Some Code Segments Contain Data, Others Don't.
- IO Port Anarchy.
- No Interrupts, Please!
- BIOS Calls.

Concept

A multitasking OS does not run multiple programs (i.e., tasks) simultaneously. In reality, it loads a task into memory, permits it to run for a while and then suspends it. The program is suspended by creating a snapshot, or image, containing the contents of all or many of the logical processor's registers in memory (frequently referred to as the processor context). If the OS *were* using the hardware-based task switching mechanism included in the IA-32 architecture, (most modern OSs do not) the logical processor would automatically store the image in a special data structure in memory referred to as a Task State Segment (TSS) by performing an automated series of memory writes. In other words, the state of the logical processor at the point of suspension would be saved in memory. In reality, most modern OSs save the register contents (under software control) in an OS-specific data structure (rather than the TSS) which the author will refer to as the Task Data Structure. Some of the TSS's functionality is, in fact, used, however (this will be covered later).

Having effectively saved a snapshot that indicates the point of suspension and the logical processor's state at that time, the logical processor could then initiate another task by loading it into memory and jumping to its entry point, or it could resume a previously-suspended task by reloading its register set from that task's Task Data Structure. Based on OS-specific criteria, the OS could at some point decide to suspend this task as well. As before, the state of the logical processor would be saved in memory (in the Task Data Structure) as a snapshot of the task's state at its point of suspension.

When the OS decides to resume a previously-suspended task, the logical processor's registers would be restored from the Task Data Structure under software control by performing a series of memory reads. The logical processor would then use the address pointer stored in the CS:EIP register pair to fetch the next instruction, thereby resuming program execution at the point where it had been suspended earlier.

The circumstances under which an OS decides to suspend a task is specific to the OS. It may simply use *timeslicing* wherein each task is permitted to execute for a fixed amount of time (e.g., 10ms). At the end of that period of time, the currently executing task is suspended and the next task in the queue is started or resumed. The OS might be designed to suspend the currently executing program when it requests something that is not immediately available (e.g., when it attempts an access to a page of information that resides on a mass storage device and is currently not in memory). It starts or resumes another task and, when the event previously requested by the now-suspended task occurs, typi-

cally signaled by an interrupt, the current task is interrupted and suspended and the previously suspended task resumed. This is commonly referred to as *preemptive multitasking*.

An Example—Timeslicing

Prior to starting or resuming execution of a task:

1. The OS *task scheduler* would initialize a hardware timer (typically, the timer incorporated within the Local APIC associated with the logical processor) to interrupt program execution after a defined period of time (e.g., 10ms).
2. The scheduler then causes the logical processor to initiate or resume execution of the task.
3. The logical processor proceeds to fetch and execute the instructions comprising the task for 10ms.
4. When the hardware timer expires it generates an interrupt, causing the logical processor to suspend execution of the currently executing task and to switch back to the OS's task scheduler.
5. The scheduler then determines which task to initiate or resume next.

Another Example—Awaiting an Event

1. Task Issues Call to OS for Disk Read

The application program calls the OS requesting that a block of data be read from a disk drive into memory. The OS then forwards the request to the disk driver, and, having done so, suspends the task and either starts or resumes another one.

Rather than awaiting the completion of the disk read, the OS scheduler would better utilize the machine's resources by suspending the task that originated the request and transferring control to another program so work can be accomplished while the disk operation is in progress.

2. Device Driver Initiates Disk Read

The driver issues a call to malloc (the OS's memory allocation manager) requesting the allocation of a memory buffer to hold the requested data. After

11 Multitasking-Related Issues

The Previous Chapter

The previous chapter provided an introduction to the concept of multitasking and covered the following topics:

- Concept.
- An Example—Timeslicing.
- Another Example—Awaiting an Event.
 - 1. Task Issues Call to OS for Disk Read.
 - 2. Device Driver Initiates Disk Read.
 - 3. OS Suspends Task.
 - 4. OS Makes Entry in Event Queue.
 - 5. OS Starts or Resumes Another Task.
 - 6. Disk-Generated Interrupt Causes Jump to OS.
 - 7. Interrupted Task Suspended.
 - 8. Task Queue Checked.
 - 9. OS Resumes Task.

This Chapter

This chapter introduces the concept of hardware-based task switching, global and local memory, privilege checking, read/write protection, IO port protection, interrupt masking, and BIOS call interception. The following topics are covered:

- Hardware-based Task Switching Is Slow!
- Private (Local) and Global Memory.
- Preventing Unauthorized Use of OS Code.

- With Privilege Comes Access.
 - Program Privilege Level.
 - The CPL.
 - Calling One of Your Equals.
 - Calling a Procedure to Act as Your Surrogate.
 - Data Segment Protection.
 - Data Segment Privilege Level.
 - Read-Only Data Areas.
- Some Code Segments Contain Data, Others Don't.
- IO Port Anarchy.
- No Interrupts, Please!
- BIOS Calls.

The Next Chapter

The next brief chapter summarizes various situations that can destabilize a multitasking OS environment and the x86 protection mechanisms that exist to address each of them.

Hardware-based Task Switching Is Slow!

The multitasking OS loads multiple tasks into different areas of memory and permits each to run for a slice of time. As described in the previous chapter, it permits a task to run for its assigned timeslice, suspends it, permits another task to run for a timeslice, suspends it, etc. If the OS is executing on a fast processor with fast access to memory, this task switching can be accomplished so quickly that all of the tasks *appear* to be executing simultaneously.

While the logical processor is executing a task, the OS kernel and all of the other dormant tasks are resident in memory. When each of the tasks (and the OS kernel's scheduler) were suspended earlier in time, the logical processor created a snapshot of its register image in memory at the moment of task suspension. In the IA-32 environment, the typical OS sets up a separate Task Data Structure for each task. If the OS designers had chosen to utilize the x86 processor's hardware-based task switching mechanism, the processor would automatically save its register set in and restore it from a task's TSS when suspending or resuming a task. *In fact, though, due to the inefficiency of this mechanism, no modern, mainstream OSs use the hardware-based task switch mechanism.* Rather, the OS task scheduler performs the register set save and restore in software using a task-specific data structure the author refers to as the Task Data Structure.

While x86 processors support the hardware-based mechanism in IA-32 Mode to ensure that any software that does use it will function correctly, *the hardware mechanism is not supported in IA-32e Mode.*

Private (Local) and Global Memory

The currently executing application is typically only aware of two entities—itsself and the OS that manages it—and is unaware of the existence of any other tasks that, although partially or fully present in memory, are currently suspended. The currently executing application should only be permitted to access its own, private memory and, perhaps, one or more areas of memory that the OS has designated as globally-accessible by multiple applications to permit data and/or code sharing. If it were permitted to perform memory writes anywhere in memory, it is entirely probable that it will corrupt the code, stack or data areas of programs that are in memory but currently suspended. Consider what would happen when the OS resumes execution of a task that had been corrupted while in suspension. Its program and/or data would have been corrupted, causing it to behave unpredictably when it resumes execution. The OS must protect suspended tasks (including itself!) from the currently executing task. If it doesn't, multitasking will not work reliably.

When an application is loaded into memory, the OS memory manager designates certain areas of memory for its use:

- Some areas of memory are designated as private (i.e., local) to the application. These segments could be defined by entries (segment descriptors) in the application's Local Descriptor Table (LDT; see Figure 11-1 on page 373).
- The OS may also designate one or more areas of memory that are globally accessible by multiple applications (thereby permitting the sharing of data or code). These segments could be defined by entries in the Global Descriptor Table (GDT; Figure 11-1 on page 373).

In addition to defining the accessibility of memory areas using the logical processor's segmentation mechanism, the OS memory manager can also accomplish this using the virtual-to-physical address translation mechanism (i.e., Paging).

Preventing Unauthorized Use of OS Code

The OS maintains the integrity of the system. It manages all shared resources and decides what task will run next and for how long. It should be fairly obvious that the person in charge must have more authority (i.e., greater privileges) than the other tasks. It would be ill-conceived to permit normal tasks to access certain logical processor control registers, OS-related tables in memory, etc.

This form of protection can be accomplished in two ways: assignment of privilege levels to programs and assignment of ownership to areas of memory. IA-32 processors utilize both methods. There are four privilege levels:

- **Level zero.** Greatest amount of privilege. Assigned to the heart, or kernel, of the OS. It handles the task queues, memory management, etc.
- **Level one.** Typically assigned to OS services that provide services to the application programs and device drivers.
- **Level two.** Typically assigned to device drivers that the OS uses to communicate with peripheral devices.
- **Level three.** Least-privileged. Assigned to application programs.

The application program operates at the lowest privilege level (3) because its actions must be restricted. The OS kernel has the highest privilege level (0) so that it can accomplish its job of managing every aspect of the system. The integrity of the system would be compromised if an application program could call highly-privileged parts of the OS code to accomplish things it shouldn't be able to do. This implies that the logical processor must have some way of comparing the privilege level of the calling program to that of the program being called. To gain entry into the called program, the calling program's privilege level (CPL, or Current Privilege Level) must equal or exceed the privilege level of the program it is calling. IA-32 processors incorporate this feature.

With Privilege Comes Access

Privilege level 0 code has access to all of the logical processor's facilities: it can execute any instruction, access any register, and access all memory data segments. Privilege level 3 code (i.e., application code), on the other hand, only has access to a subset of the instruction set and register set.

Program Privilege Level

The CPL

Refer to Figure 11-1 on page 373. When a far jump or far call is performed, the new value loaded into the 16-bit CS register selects a segment descriptor in either the GDT or the LDT. The 2-bit DPL (Descriptor Privilege Level) field in the selected code segment descriptor becomes the logical processor's CPL (Current Privilege Level).

There are multiple cases where the OS wishes to restrict access to the code that resides in a code segment and they are introduced in the sections that follow.

12 Summary of the Protection Mechanisms

The Previous Chapter

The previous chapter introduced the concept of hardware-based task switching, global and local memory, privilege checking, read/write protection, IO port protection, interrupt masking, and BIOS call interception. The following topics were covered:

- Hardware-based Task Switching Is Slow!
- Private (Local) and Global Memory.
- Preventing Unauthorized Use of OS Code.
- With Privilege Comes Access.
 - Program Privilege Level.
 - The CPL.
 - Calling One of Your Equals.
 - Calling a Procedure to Act as Your Surrogate.
 - Data Segment Protection.
 - Data Segment Privilege Level.
 - Read-Only Data Areas.
- Some Code Segments Contain Data, Others Don't.
- IO Port Anarchy.
- No Interrupts, Please!
- BIOS Calls.

This Chapter

This chapter summarizes various situations that can destabilize a multitasking OS environment and the x86 protection mechanisms that exist to address each of them.

The Next Chapter

The next chapter introduces segment register usage in Protected Mode and the roles of segment descriptors, the GDT, the LDTs, the IDT, and the general segment descriptor format. It also introduces the concept of the flat memory model. The following topics are covered:

- Real Mode Segment Limitations.
- An Important Reminder: Segment Base + Offset = Virtual Address.
- Descriptor Contains Detailed Segment Description.
- Segment Register—Selects Descriptor Table and Entry.
- The Descriptor Tables.
- General Segment Descriptor Format.
- Goodbye to Segmentation.

Protection-Related Mechanisms

This chapter is not intended as a detailed discussion of the various protection mechanisms available in the x86 architecture. Rather, they've been collected here in one place for ease of reference and as an introduction to the various topics related to protection.

Some of the protection mechanisms were introduced in the previous chapter (privilege level assignment, local and global memory segments, data segment write protection, preventing data accesses to code segments, and guarding access to IO ports). All of the protection mechanisms are listed in Table 12-1 on page 379.

Chapter 12: Summary of the Protection Mechanisms

Table 12-1: Protection Mechanisms

Condition	Mechanism
Attempted access to a segment by a program with insufficient privilege	<p>Some examples:</p> <ul style="list-style-type: none">• Currently-running program attempts to access a data segment with a higher privilege level (i.e., the data segment's descriptor[DPL] value is numerically less than the CPL of the program).• Currently-running program executes a far jump or a far call to a procedure in a code segment with a higher privilege level (e.g., a privilege level 3 program attempts to jump to or call a procedure in a privilege level 0 code segment). <p>This would result in an exception.</p>
Attempted write to a read-only data segment	<p>The data segment descriptor's W bit = 0 indicating it is a read-only data segment. This would result in an exception.</p>
Attempted data read from a code segment	<p>The code segment descriptor's R bit = 0 indicating it's an execute-only code segment. This would result in an exception (if R = 1, then the code segment contains read-only data as well as code).</p>
Out-of-range access to a segment	<p>The offset address specified exceeds the segment size specified in the target segment's descriptor. This would result in an exception.</p>
Attempted page access by an under-privileged program	<p>A program with a privilege level of 3 attempted to access a page whose PTE[U/S] bit = 0. This would result in an exception. Note: U/S stands for User/Supervisor.</p>

x86 Instruction Set Architecture

Table 12-1: Protection Mechanisms (Continued)

Condition	Mechanism
Attempted write to a read-only page	<p>There are two possibilities:</p> <ul style="list-style-type: none"> • A privilege level 3 program attempted to write to a page whose PTE[R/W] bit = 0 marking it as a read-only page. • A program with supervisor privileges (i.e., it has a privilege level of 0, 1 or 2) attempted to write to a user page (its PTE[U/S] bit = 1) with PTE[R/W] bit = 0 marking it as a read-only user page <i>and</i> CR0[WP] = 1 (indicating supervisor programs are not permitted to write into read-only user pages). <p>This would result in an exception.</p>
Access to an absent page	<p>A virtual memory address selected a PTE with the Page Present bit (bit 0) = 0 indicating that the target physical page isn't currently in memory. This would result in an Page Fault exception.</p>
Attempted direct access to an IO or memory-mapped IO port	<ul style="list-style-type: none"> • Access to an IO port: <ul style="list-style-type: none"> – By a Protected Mode task (other than a VM86 task). Any attempt by a program (other than a VM86 task) with a privilege level numerically greater than the Eflags[IOPL] threshold to execute the IN, INS, OUT and OUTS instructions will trigger a General Protection exception. – By a VM86 task. When an IN, INS, OUT or OUTS instruction is executed, the logical processor uses the 16-bit IO port address to index into the IO Permission bit map in the task's TSS data structure. The state of the selected bit determines whether the IO instruction is executed or a General Protection exception is generated. • Access to a memory-mapped IO port. This form of protection can be provided by the virtual-to-physical address translation mechanism (i.e., Paging). The virtual addresses of memory-mapped IO ports could be grouped into a page and the virtual page address could select a PTE that indicates the page isn't present in memory. Any attempted access within the page would then result in a Page Fault exception.

13 Protected Mode Memory Addressing

The Previous Chapter

The previous chapter summarized various situations that can destabilize a multitasking OS environment and the x86 protection mechanisms that exist to address each of them.

This Chapter

This chapter introduces segment register usage in Protected Mode and the roles of segment descriptors, the GDT, the LDTs, the IDT, and the general segment descriptor format. It also introduces the concept of the flat memory model. The following topics are covered:

- Real Mode Segment Limitations.
- An Important Reminder: Segment Base + Offset = Virtual Address.
- Descriptor Contains Detailed Segment Description.
- Segment Register—Selects Descriptor Table and Entry.
- The Descriptor Tables.
- General Segment Descriptor Format.
- Goodbye to Segmentation.

The Next Chapter

The next chapter provides a detailed description of code segments (both Conforming and Non-Conforming), privilege checking, and Call Gates. The following topics are covered:

- Selecting the Active Code Segment.
- CS Descriptor.
- Accessing the Code Segment.
- Short/Near Jumps.
- Unconditional Far Jumps.
- Privilege Checking.
- Jumping from a Higher-to-Lesser Privileged Program.
- Direct Procedure Calls.
- Indirect Procedure Far Call Through a Call Gate.
- Automatic Stack Switch.
- Far Call From 32-bit CS to 16-bit CS.
- Far Call From 16-bit CS to 32-bit CS.
- Far Returns.

Real Mode Segment Limitations

Figure 13-1 on page 385 illustrates the contents of a segment register while operating in Real Mode; i.e., the upper 16 bits of the segment's 20-bit base address (aligned on a 16-byte address boundary) in the first megabyte of memory space. The logical processor automatically appends four bits of zero to the lower end to form the base address. As an example, if the programmer moves the value 1010h into the DS register

```
mov ax, 1010
mov ds, ax
```

this would set the start address of the data segment to 10100h.

As stated earlier in the book, when in Protected Mode the OS memory manager must be able to define a number of segment properties in addition to its base address (and this is not possible in a 16-bit register).

In Real Mode, a segment has the following characteristics:

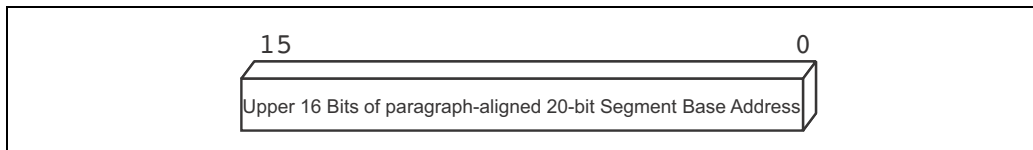
- Its base address must be in the **first megabyte of memory space**. In order to have the maximum flexibility in memory allocation while operating in Protected Mode, the OS must be able to define a program's segments as residing anywhere within physical memory (above or below the 1MB address boundary).
- The **segment length is fixed at 64KB**. Unless they are incredibly small, programs and the data they manipulate virtually always occupy more than 64KB of memory space, but each segment has a fixed length of 64KB in Real

Chapter 13: Protected Mode Memory Addressing

Mode. If the OS only requires a very small segment for a program's code, data or stack area, the only size available is still fixed at 64KB. This can waste memory space (albeit, not very much). If the code comprising a particular program is larger than 64KB, the programmer must set up and jump back and forth between multiple code segments. The data that a program acts upon may also occupy multiple data segments. This is a very inefficient memory organization model and one that forces the programmer to think in a very fragmented manner. It's one of the major things programmers dislike about Real Mode segmentation.

- The **segment can be read or written by any program**. In Real Mode, a segment can be accessed by any program. This is an invitation for one program to inadvertently trash another's code, data or stack area. In addition, any program can call procedures within any other program. There is no concept of restricting access to certain programs.

Figure 13-1: Segment Register Contents in Real Mode



An Important Reminder: Segment Base + Offset = Virtual Address

While this chapter (along with the two chapters that follow) provides a detailed description of segmentation, keep in mind that the address produced by adding an offset to a segment base address may *not*, in fact, be the memory address that is used to access physical memory. If the virtual-to-physical address translation mechanism (i.e., paging) is enabled, it is treated as a *virtual* (also referred to as *linear*) memory address that is subsequently submitted to the virtual-to-physical address translation logic. Upon receipt of a virtual memory address, the Paging logic uses it to perform a lookup in special address translation tables created in memory by the OS kernel. The Page Table Entry (PTE) selected by the virtual address contains the information used to translate the virtual address into a physical memory address.

The two chapters immediately following this one provide detailed discussions of code, data and stack segments, after which the subsequent chapter provides a detailed description of the address translation mechanism.

Descriptor Contains Detailed Segment Description

In a Protected Mode environment, the OS programmer must be able to specify the following characteristics of each segment:

- The base address anywhere in a 4GB virtual address range.
- Segment length (anywhere from one byte to 4GB in length).
- How the segment may be accessed:
 - A read-only data segment.
 - An execute-only code segment (contains only code and no data).
 - A code segment that also contains read-only data.
 - A read/writable data segment (contains both code and read-only data).
- The minimum privilege level a program must have in order to access the segment.
- Whether it's a code or data segment, or a special segment used only by the OS kernel and the logical processor.
- Whether the segment of information is currently present in memory or not.

In Protected Mode, it requires eight bytes (64-bits) of information to describe all of these characteristics. The Protected Mode OS memory manager must create an eight byte descriptor for each memory segment to be used by each program (including those used by the OS itself). Obviously, it would consume a great deal of processor real estate to keep the descriptors of all segments in use by all tasks on the processor chip itself. For this reason, the descriptors are stored in special tables in memory. The next section provides a description of these descriptor tables.

Segment Register—Selects Descriptor Table and Entry

When a programmer wishes to gain access to an area of memory, the respective segment register (the CS, SS, or one of the data segment registers: DS, ES, FS, or GS) must be loaded with a 16-bit value that identifies the area of memory. In Real Mode, the value loaded into the segment register represents the upper 16 bits of the 20-bit start address of the segment in memory. In Protected Mode, the value loaded into a segment register is referred to as the segment selector, illustrated in the upper part (i.e., the segment register's visible part) of Figure 13-3 on page 389:

- **RPL field.** The Requester Privilege Level (RPL) field is described in “RPL Definition” on page 439 and “RPL Definition” on page 439.

14 Code, Calls and Privilege Checks

The Previous Chapter

The previous chapter introduced segment register usage in Protected Mode and the role of segment descriptors, the GDT, the LDTs, the IDT, and the general segment descriptor format. It also introduced the concept of the flat memory model. The following topics were covered:

- Real Mode Segment Limitations.
- An Important Reminder: Segment Base + Offset = Virtual Address.
- Descriptor Contains Detailed Segment Description.
- Segment Register—Selects Descriptor Table and Entry.
- The Descriptor Tables.
- General Segment Descriptor Format.
- Goodbye to Segmentation.

This Chapter

This chapter provides a detailed description of code segments (both Conforming and Non-Conforming), privilege checking, and Call Gates. The following topics are covered:

- Selecting the Active Code Segment.
- CS Descriptor.
- Accessing the Code Segment.
- Short/Near Jumps.
- Unconditional Far Jumps.
- Privilege Checking.
- Jumping from a Higher-to-Lesser Privileged Program.
- Direct Procedure Calls.

- Indirect Procedure Far Call Through a Call Gate.
- Automatic Stack Switch.
- Far Call From 32-bit CS to 16-bit CS.
- Far Call From 16-bit CS to 32-bit CS.
- Far Returns.

The Next Chapter

The next chapter provides a detailed description of Data and Stack segments (including Expand-Up and Expand-Down Stacks) and privilege checking. The following topics are covered:

- The Data Segments.
 - General.
 - Two-Step Permission Check.
 - An Example.
- Selecting and Accessing a Stack Segment.
 - Introduction.
 - Expand-Up Stack.
 - Expand-Down Stack.
 - The Problem.
 - Expand-Down Stack Description.
 - An Example.
 - Another Example.

Abbreviation Alert

In many cases in this chapter, the abbreviation *CS* is substituted for *code segment*.

Selecting the Active Code Segment

To execute code from a specific area of memory, the programmer must tell the logical processor what code segment the instructions are to be fetched from. This is accomplished by loading a 16-bit value (a selector) into the Code Segment (CS) register. In Real Mode, this value represents the upper 16-bits of the 20-bit zero-extended segment base address. In Protected Mode, the value loaded into the 16-bit visible portion of the CS register (see Figure 14-1 on page 418) selects an entry in either the GDT or the LDT. The selected 8-byte CS descriptor is automatically read from memory and stored in the invisible part of the CS register.

Chapter 14: Code, Calls and Privilege Checks

Any of the actions listed in Table 14-1 on page 417 loads a new selector into CS and causes the logical processor to begin fetching instructions from a new code segment in memory.

Table 14-1: Actions That Cause a Switch to a New CS

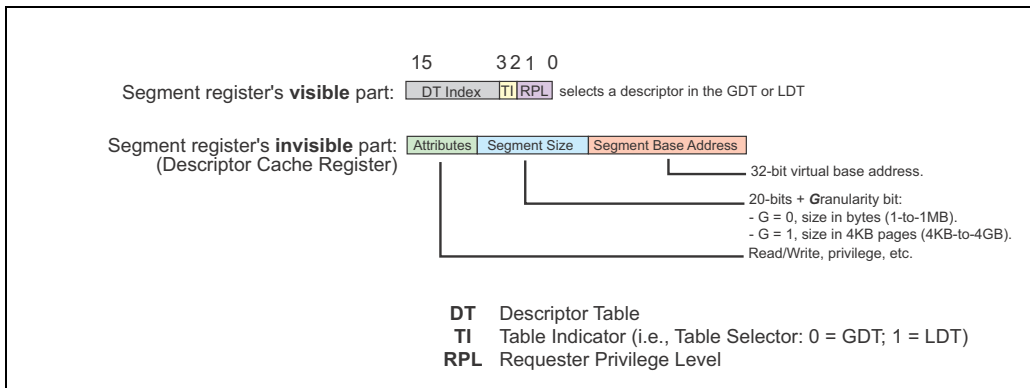
Action	Description
Execution of a far jump instruction	Loads the CS/Instruction Pointer register pair with new values.
Execution of a far call instruction	
A hardware interrupt or a software exception	In response, the logical processor loads new values into the CS/Instruction Pointer register pair from the IDT entry selected by the interrupt or exception vector.
Execution of a software interrupt instruction	<ul style="list-style-type: none">• INT nn. In response, the logical processor loads new values into the CS/Instruction Pointer register pair from the IDT entry selected by the instruction's 8-bit operand.• INT3. In response, the logical processor loads new values into the CS/Instruction Pointer register pair from IDT entry 3.• INTO. In response, the logical processor loads new values into the CS/Instruction Pointer register pair from IDT entry 4.• BOUND. In response, the logical processor loads new values into the CS/Instruction Pointer register pair from IDT entry 5.
Initiation of a new task or resumption of a previously-suspended task	During a task switch by the x86 processor's hardware-based task switching mechanism (which most modern OSs do not use), the logical processor loads most of its registers, including CS:EIP, with values from the TSS associated with the task being started or resumed.
Execution of a far RET instruction	The far return address (CS and Instruction Pointer) is popped from the stack and loaded into the CS/Instruction Pointer register pair.

x86 Instruction Set Architecture

Table 14-1: Actions That Cause a Switch to a New CS (Continued)

Action	Description
Execution of an Interrupt Return instruction (IRET)	The far return address is popped from the stack and loaded into the CS/Instruction Pointer register pair.

Figure 14-1: Segment Register



CS Descriptor

CS Descriptor Selector

The value loaded into the visible part of CS (Figure 14-1 on page 418) identifies:

- The descriptor table that contains the code segment descriptor:
 - TI = 0 selects the GDT.
 - TI = 1 selects the LDT.
- The entry in the specified descriptor table. The DT (Descriptor Table) Index field selects one of 8192d entries in the selected table.
- The privilege level of the program that created the selector in the CS register. This is referred to as the Requester Privilege Level (RPL).

15 Data and Stack Segments

The Previous Chapter

The previous chapter provided a detailed description of code segments (both Conforming and Non-Conforming), privilege checking, and Call Gates. The following topics were covered:

- Selecting the Active Code Segment.
- CS Descriptor.
- Accessing the Code Segment.
- Short/Near Jumps.
- Unconditional Far Jumps.
- Privilege Checking.
- Jumping from a Higher-to-Lesser Privileged Program.
- Direct Procedure Calls.
- Indirect Procedure Far Call Through a Call Gate.
- Automatic Stack Switch.
- Far Call From 32-bit CS to 16-bit CS.
- Far Call From 16-bit CS to 32-bit CS.
- Far Returns.

This Chapter

This chapter provides a detailed description of Data and Stack segments (including Expand-Up and Expand-Down Stacks) and privilege checking when accessing data or stack segments. The following topics are covered:

- The Data Segments.
 - General.
 - Two-Step Permission Check.

- An Example.
- Selecting and Accessing a Stack Segment.
 - Introduction.
 - Expand-Up Stack.
 - Expand-Down Stack.
 - The Problem.
 - Expand-Down Stack Description.
 - An Example.
 - Another Example.

The Next Chapter

The next chapter covers the following topics:

- Summarizes the evolution of the virtual-to-physical address translation facilities on the x86 processors and provides a backgrounder on memory and disk management.
- The concept of virtual memory is introduced as well as the advantages of address translation.
- The first and second generation virtual-to-physical address translation mechanisms are described in detail.
- The role of the Translation Lookaside Buffer (TLB) is described, as well as the Global Page feature and TLB maintenance.
- Page Directory Entries (PDEs) and Page Table Entries (PTEs) are described in detail.
- Page access permission.
- Missing page or Page Table.
- Page access history.
- 4MB pages.
- PSE-36 Mode.
- Execute Disable feature.
- Page caching rules.
- Page write protection.

A Note Regarding Stack Segments

While the stack segment is, in reality, nothing more than a read/writable data segment, it is treated separately in this chapter because it is used differently than the typical data segment.

Data Segments

General

x86 processors introduced after the 286 implement four data segment registers (as opposed to just one, DS, in the 286): DS, ES, FS and GS. They permit software to identify up to four separate data segments (in memory) that can be accessed by the currently executing program.

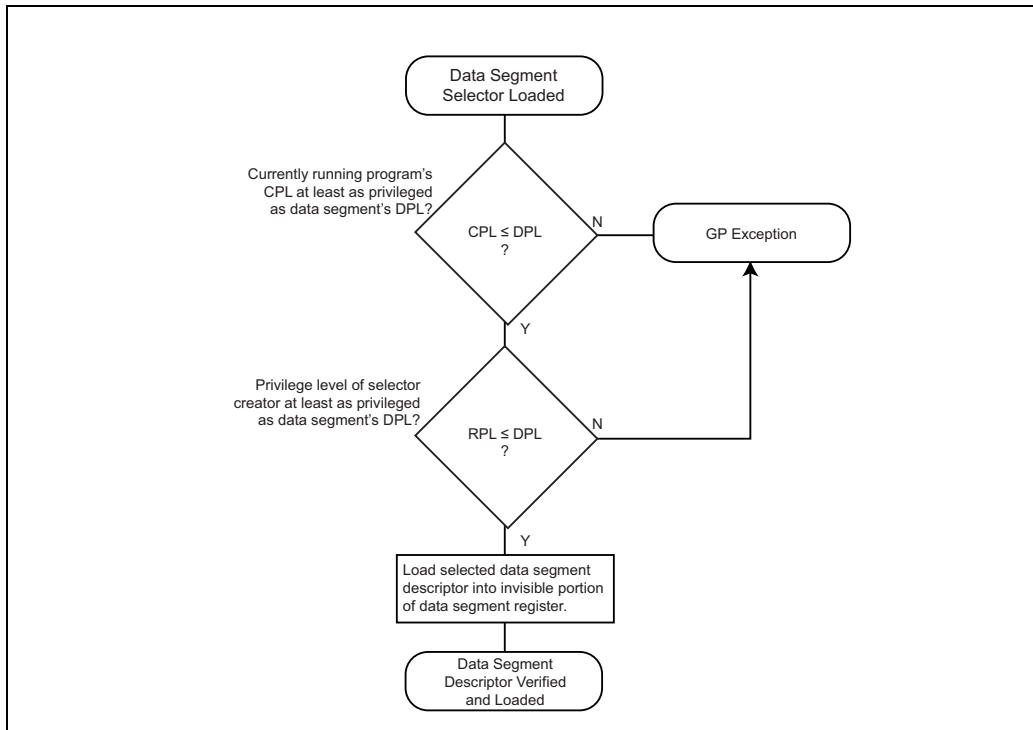
To access data within any of the four data segments, the programmer must first load a 16-bit descriptor selector into the respective data segment register. In Real Mode, the value in a data segment register specifies the upper 16-bits of the 20-bit zero-extended memory start address of the data segment. In Protected Mode, the value selects a data segment descriptor in either the GDT or LDT. Figure 15-3 on page 484 illustrates the format of a 32-bit data segment descriptor (in a 286-style 16-bit data segment descriptor, bytes 6 and 7 are reserved).

Two-Step Permission Check

In order to successfully access one or more locations in a data segment, two permission checks must be passed:

1. **Descriptor pre-load privilege check.** Refer to Figure 15-1 on page 482. The currently running program must have sufficient privilege to select the target data segment descriptor in the GDT or LDT. Assuming it does, the selected data segment descriptor is loaded into the invisible portion of the respective data segment register.
2. **Access type/limit checks.** Before any subsequent access is permitted within a data segment, the logical processor must verify that the access type is permitted (e.g., that a write is permitted) and must also verify that the specified location (i.e., offset) falls within the bounds of the targeted data segment.

Figure 15-1: Data Segment Descriptor Pre-Load Privilege Check



An Example

Consider this example (assumes code is fetched from a Non-Conforming CS with a DPL of 2):

```
mov ax, 4f36    ;load ds register
mov ds, ax      ;
mov al, [0100]  ;read 1 byte from data segment into al
mov [2100], al  ;write 1 byte to data segment from al
```

The value 4F36h in the DS register is interpreted by the logical processor as indicated in Figure 15-2 on page 484. The logical processor accesses LDT entry 2534 to obtain the data segment descriptor. The selector's RPL = 2 indicating that a privilege level 2 program created the selector value. Figure 15-3 on page 484

16 IA-32 Address Translation Mechanisms

The Previous Chapter

The previous chapter provided a detailed description of Data and Stack segments (including Expand-Up and Expand-Down Stacks) and privilege checking. The following topics were covered:

- The Data Segments.
 - General.
 - Two-Step Permission Check.
 - An Example.
- Selecting and Accessing a Stack Segment.
 - Introduction.
 - Expand-Up Stack.
 - Expand-Down Stack.
 - The Problem.
 - Expand-Down Stack Description.
 - An Example.
 - Another Example.

This Chapter

This chapter covers the following topics:

- Summarizes the evolution of the virtual-to-physical address translation facilities on the x86 processors and provides a backgrounder on memory and disk management.

- The concept of virtual memory is introduced as well as the advantages of address translation.
- The first and second generation virtual-to-physical address translation mechanisms are described in detail.
- The role of the Translation Lookaside Buffer (TLB) is described, as well as the Global Page feature and TLB maintenance.
- Page Directory Entries (PDEs) and Page Table Entries (PTEs) are described in detail.
- Page access permission.
- Missing page or Page Table.
- Page access history.
- 4MB pages.
- PSE-36 Mode.
- Execute Disable feature.
- Page caching rules.
- Page write protection.

The Next Chapter

The next chapter describes the operational characteristics of various types of memory targets (UC, WC, WP, WT, and WB) and the role of the Memory Type and Range Registers (MTRRs). It defines the concept of speculatively executed loads and describes issues related to the logical processor's Posted Memory Write Buffer (PMWB) and Write-Combining Buffers (WCBs).

Three Generations

Over the years, the x86 address translation mechanism has experienced three major evolutionary changes (as well as a number of smaller, incremental changes). Consequently, the author has divided the discussion into three major sections:

- **1st-generation paging.** The address translation mechanism was first introduced in the x86 product line with the advent of the 386 processor. This mechanism (including some minor enhancements added in the 486 and Pentium) is what the author refers to as the *first-generation paging* mechanism. It should be noted, however, that page address translation was actually first introduced in mainframe computers many years earlier.
- **2nd-generation paging.** The next major evolutionary jump, PAE-36 Mode, was first implemented in the Pentium Pro processor.
- **3rd-generation paging.** Part of the Intel 64 architecture.

The first and second generations are covered in this chapter. The third generation is covered in "IA-32e Address Translation" on page 983.

Demand Mode Paging Evolution

Since the advent of the 386 processor, a number of enhancements have been made to the Paging mechanism. Table 16-1 on page 495 tracks the evolutionary changes that appeared in successive generations of the x86 processor family.

Table 16-1: Paging Evolution

Processor	Enhancement	Described in
386	-	First-generation address translation. Virtual-to-physical address translation was first introduced to the x86 product family with the advent of the 386 processor. Using this mechanism, a 2-level lookup is used to translate a 32-bit virtual address into a 32-bit physical memory address.
486	Write Protect feature	A complete description of this minor enhancement can be found in “Example Usage: Unix Copy-on-Write Strategy” on page 569.
	Caching Rules	Minor enhancements: <ul style="list-style-type: none">• CR3[PCD] and CR3[PWT] were added. A complete description can be found in “Translation Table Caching Rules” on page 585.• PCD and PWT bits were added to each Page Directory Entry (PDE). Refer to “Defining a Page’s Caching Rules” on page 585.• PCD and PWT bits were added to each Page Table Entry (PTE). Refer to “Defining a Page’s Caching Rules” on page 585.
Pentium	4MB Pages	The Page Size Extension (PSE) feature was added. This minor enhancement was first implemented in the Pentium and was migrated into the later versions of the 486. A complete description can be found in “4MB Pages” on page 550.
	Global pages	A minor enhancement. A complete description can be found in “Global Pages” on page 526.

Table 16-1: Paging Evolution (Continued)

Processor	Enhancement	Described in
Pentium Pro	PAE-36 Mode (2nd generation paging)	Second-generation address translation. Using this mechanism, a 3-level lookup is used to translate a 32-bit virtual address into a 36-bit physical memory address. A complete description can be found in “Second-Generation Paging” on page 553.
Pentium II	PSE-36 Mode	This was a minor enhancement added in the Pentium II Xeon (the very first Xeon processor). A complete description can be found in “PSE-36 Mode Background” on page 575.
	PAT feature	Page Attribute Table feature (a minor enhancement). A complete description can be found in “PAT Feature (Page Attribute Table)” on page 587.
Pentium 4	Intel 64	Third-generation address translation. The Intel 64 architecture introduced the third-generation address translation mechanism. Using this mechanism, a 4-level lookup can translate a 48-bit virtual address into a physical memory address up to 48-bits in width (<i>note</i> : current implementations translate a 48-bit virtual address into a 40-, 41-, or 48-bit physical address).

Background

Memory and Disk: Block-Oriented Devices

Mass storage devices are block-oriented devices. Information is stored on a disk as a series of fixed-length blocks of information and the OS manages disks in that manner. From the perspective of the OS kernel, memory is also managed as a series of fixed-length blocks—referred to as pages—of storage.

Definition of a Page

The OS kernel’s memory manager (frequently referred to as the *malloc*, or memory allocation, facility) manages memory as a series of pages of information, each of a uniform size, each starting on an address boundary divisible by its

17 Memory Type Configuration

The Previous Chapter

The previous chapter covered the following topics:

- Evolution of demand mode paging on the x86 processors. Backgrounder on memory and disk management.
- Virtual memory concept and advantages of address translation.
- First and second generation virtual-to-physical address translation mechanisms.
- Role of the Translation Lookaside Buffer (TLB). Global Page feature and TLB maintenance.
- Page Directory Entries (PDEs) and Page Table Entries (PTEs).
- Page access permission.
- Missing page or Page Table.
- Page access history.
- 4MB pages.
- PSE-36 Mode.
- Execute Disable feature.
- Page caching rules.
- Page write protection.

This Chapter

This chapter describes the operational characteristics of various types of memory targets (UC, WC, WP, WT, and WB) and the role of the Memory Type and Range Registers (MTRRs). It defines the concept of speculatively executed loads and describes issues related to the logical processor's Posted Memory Write Buffer (PMWB) and Write-Combining Buffers (WCBs).

The Next Chapter

The next chapter contrasts hardware- versus software-based tasking switching and provides a conceptual overview of task switching as well as a detailed description of the hardware-based task switching mechanism. The following topics are covered:

- Hardware- vs. Software-Based Task Switching
- A Condensed Conceptual Overview
- A More Comprehensive Overview
- Hardware-Based Task Switching
 - It's Slow
 - Why Didn't OSs Use It?
 - Why Wasn't It Improved?
 - Why Does It Still Exist?
 - Introduction to the Key Elements
 - The Trigger Events
 - The Descriptors
 - The Task Register
 - TSS Data Structure Format
 - Comprehensive Task Switch Description
 - Calling Another Task
 - Task Switching and Address Translation
 - Switch from Higher-Privilege Code to Lower

Characteristics of Memory Targets

Introduction

When the logical processor must perform a memory access, it is important that it understand the operational characteristics of the target device in order to ensure proper operation. If it does not, the manner in which the memory access is accomplished may result in improper operation of the device or of the program.

Example Problem: Caching from MMIO

As an example, assume that an area of memory is populated with a series of memory-mapped IO (MMIO) registers associated with one or more devices.

Chapter 17: Memory Type Configuration

Now assume that the program performs a 4-byte memory read to obtain the status of a device from its 32-bit, device-specific status register. If the logical processor were to assume that the region of memory being accessed is cacheable, it would perform a lookup in its caches and, in the event of a cache miss, would initiate a memory read to obtain not only the four requested locations, but would in fact read from all locations that encompass the line within which the desired four locations reside. This could result in a serious problem. The contents of *all* of the memory-mapped IO ports within that line of memory space would be read and cached in the processor. If the program subsequently issued a request to access any of those locations, it would result in a cache hit and:

- **If it's a read:** the requested data is supplied from the cache, *not* from the actual IO device that implements that memory-mapped IO port. This means that the data or status obtained would not represent the current, up-to-date contents of the location read. This desynchronization between a device driver and its related device can result in erroneous operation.
- **If it's a write:** the line in the cache is updated but, if the memory area is designated as WB (cacheable Write-Back) memory, the data is not written to memory. The actual memory-mapped IO device therefore does not receive the write.

Early Processors Implemented Primitive Mechanism

The example just described is but one case wherein the logical processor's lack of knowledge regarding the rules of conduct it must follow within a given memory area can result in spurious operation. In a very limited sense, the 486 and Pentium processors possessed a mechanism that permitted the OS kernel to define the characteristics of a region of memory. Each PTE (Page Table Entry) contained two bits, PCD and PWT, that permitted the OS to define a 4KB memory page as cacheable Write Through (WT), cacheable Write Back (WB), or uncacheable (UC) memory. This solution was insufficient for two reasons:

- The OS typically is not platform-specific and therefore doesn't necessarily know the characteristics of the various devices that populate memory space. The BIOS, on the other hand, *is* platform-specific but it is the OS and not the BIOS that sets up and maintains the Page Tables in memory.
- There are many different types of devices and some require different processor operation than that defined using the PTE's PCD and PWT bits (the WB, WT and UC memory types). Be advised that the later addition of the PAT feature [see "PAT Feature (Page Attribute Table)" on page 587] permitted the OS to assign any memory type to a page).

Solution/Problem: Chipset Memory Type Registers

When a program executing on the 486 or the Pentium had to initiate a memory access, the processor's internal hardware consulted the PTE[PCD] and PTE[PWT] bits to determine the rules of conduct to follow within the addressed memory page. If the memory access necessitated the performance of a transaction on the FSB (Front-Side Bus), during the memory transaction the processor transmitted the state of the PCD and PWT bits on its PCD and PWT output pins. Using the memory address output by the processor, the chipset would consult a chipset design-specific register set to determine the rules of conduct to be followed within the addressed memory area. If there was a disagreement between the OS-defined rules (as output on PCD and PWT) and the chipset's rules (as defined by the contents of its register set), the chipset would defer to the more conservative memory type (i.e., the less aggressive of the two memory types). As an example, if the processor initiated a cache line read on the FSB and the chipset said it was UC (uncacheable) memory while the processor said it was WB (cacheable Write Back) memory, the chipset would inform the processor that the entire line would *not* be returned (as the processor requested), but rather just the requested data item that caused a cache miss would be returned.

The chipset's register set was programmed by the BIOS at startup time. The problem with this approach is that the chipset's register set was implemented in a chipset-specific manner outside the scope of any industry standard specification. There would therefore have to be a separate version of the BIOS to cover all of the possible chipset types that would be used on system boards incorporating the BIOS.

Solution: Memory Type Register Set

With the advent of the Pentium Pro, Intel migrated the memory type configuration register set that had historically resided in the chipset into the processor itself. This register set is referred to as the Memory Type and Range Registers (MTRRs). While the MTRRs were, in fact, implemented identically in all members of the P6 and Pentium 4 processor families, they were *not* part of the x86 ISA specification and therefore not guaranteed to be implemented identically (or, for that matter, at all) in any given processor model. With the advent of the Pentium 4, however, the MTRRs were officially defined as part of the x86 ISA (and the register names are preceded by *IA32*). They are implemented as MSRs and are accessed using the RDMSR and WRMSR instructions.

18 *Task Switching*

The Previous Chapter

The previous chapter described the operational characteristics of various types of memory targets (UC, WC, WP, WT, and WB) and the role of the Memory Type and Range Registers (MTRRs). It defined the concept of speculatively executed loads and described issues related to the logical processor's Posted Memory Write Buffer (PMWB) and Write-Combining Buffers (WCBs).

This Chapter

This chapter contrasts hardware- versus software-based tasking switching and provides a conceptual overview of task switching before providing a detailed description of the hardware-based task switching mechanism. The following topics are covered:

- Hardware- vs. Software-Based Task Switching
- A Condensed Conceptual Overview
- A More Comprehensive Overview
- Hardware-Based Task Switching
 - It's Slow
 - Why Didn't OSs Use It?
 - Why Wasn't It Improved?
 - Why Does It Still Exist?
 - Introduction to the Key Elements
 - The Trigger Events
 - The Descriptors
 - The Task Register
 - TSS Data Structure Format
 - Comprehensive Task Switch Description
 - Calling Another Task
 - Task Switching and Address Translation
 - Switch from Higher-Privilege Code to Lower

The Next Chapter

The next chapter provides a detailed description of interrupt and exception handling in Protected Mode. This includes detailed coverage of:

- The IDT.
- Interrupt and Trap Gate operation.
- Task Gate operation.
- Interrupt and exception event categories.
- State save (and stack selection).
- The IRET instruction.
- Maskable hardware interrupts.
- Non-Maskable Interrupt (NMI).
- Machine Check exception.
- SM interrupt (SMI).
- Software interrupt instructions.
- Software exceptions.
- Interrupt/exception priority.

Hardware- vs. Software-Based Task Switching

The 386 introduced a number of well-received features, but for many OS vendors, its hardware-based task switching mechanism was not considered one of them. On the one hand, it permits the automation of the OS scheduler's job of switching from one task to another after the current task's timeslice has expired. On the other hand, the hardware mechanism's indulgence in an excessive validity checks renders it ponderously slow.

As a result, major OS vendors chose not to utilize the hardware mechanism and instead implemented task switching solely under the control of software (the exact implementation is OS design-specific). In tacit recognition of this reality, the hardware switching mechanism is disabled in IA-32e Mode (any attempted use of it is, in fact, considered illegal and results in an exception). The hardware mechanism is, however, supported in IA-32 Mode and is, for completeness, described in detail in this chapter. The two sections in this chapter entitled:

- “A Condensed Conceptual Overview” on page 631,
- and “A More Comprehensive Overview” on page 631

provide an introduction to task switching applicable to both the software- and hardware-based mechanisms. A description of software-based task switching can be found in “Scheduler's Software-Based Task Switching Mechanism” on page 977.

A Condensed Conceptual Overview

The task switching concept is simple:

1. The OS scheduler selects the next task to run.
2. It initializes the logical processor's register set with the appropriate startup values.
3. It triggers a hardware timer (the timeslice timer; typically the Local APIC timer is used) configured to run for the timeslice assigned to that task (e.g., 10ms).
4. The logical processor starts executing the task and continues to do so until either:
 - The timer expires and generates an interrupt.
 - The task requires something that will take a while to complete (e.g., it issues a request to the OS to load some information from disk to memory). In this case, the OS scheduler will suspend the task (more on this later).
5. Assuming the task's timeslice has expired (i.e., the timer generates an interrupt), the event interrupts the execution of the task and returns control back to the OS kernel (specifically, to the task scheduler).
6. The scheduler suspends the task by recording the state of the logical processor's register set in a special data structure the scheduler has associated with that task.
7. It then selects the next task to start or resume and goes back to step 2.

A More Comprehensive Overview

It should be stressed that the details of task switching are OS design-specific. This discussion is conceptual (and general) in nature and applies to both software- and hardware-based task switching.

The Scheduler and the Task Queue

Some of the critical components involved in task switching are:

- **Task Scheduler.** The kernel's task scheduler is responsible for managing the task switching environment.
- **Task Queue.** Maintained by the scheduler, the **task queue** is used to keep track of:

- The currently-running tasks.
- Any pending events associated with those tasks. A task may have been suspended earlier after issuing a request to the OS. As an example, while a task was running, it may have issued a request to the OS for a block of information to be read from disk and placed in memory. Since this would take quite a while to complete, pending completion of the request the OS would suspend the task and start or resume another one. The scheduler would create an entry in the event queue associating the pending disk controller completion interrupt with the resumption of the previously-suspended task.
- **Timer.** The Local APIC's programmable timer is used by the task scheduler to assign the amount of time a task is permitted to execute on the logical processor before it is interrupted and control is returned to the scheduler.

Setting Up a Task

In preparation for running a task, the scheduler must:

- Load the application's startup code and data into memory. The remainder of the application remains on disk and will only be read into memory on demand (i.e., if it's required).
- Set up a series of kernel tables in memory that are necessary to support a task.
- Set up the kernel registers that tell the logical processor the location and size of these tables.

The sections that follow provide additional information about the associated tables and registers.

The Task Data Structure

Refer to Figure 18-1 on page 634. The scheduler must create a register save/restore data structure (let's call it the Task Data Structure, or TDS) in memory for each task that will be run:

- When a task is initially set up, the TDS fields will be initialized with the values to be loaded into the logical processor's register set when the task is first started.
- When starting or resuming the execution of a task, many of the logical processor's registers will be loaded from this data structure.
- When suspending a task, the registers will be saved in the task's data structure.

19 *Protected Mode Interrupts and Exceptions*

The Previous Chapter

The previous chapter contrasted hardware- versus software-based tasking switching and provided a conceptual overview of task switching. It then provided a detailed description of the hardware-based task switching mechanism. The following topics were covered:

- Hardware- vs. Software-Based Task Switching
- A Condensed Conceptual Overview
- A More Comprehensive Overview
- Hardware-Based Task Switching
 - It's Slow
 - Why Didn't OSs Use It?
 - Why Wasn't It Improved?
 - Why Does It Still Exist?
 - Introduction to the Key Elements
 - The Trigger Events
 - The Descriptors
 - The Task Register
 - TSS Data Structure Format
 - Comprehensive Task Switch Description
 - Calling Another Task
 - Task Switching and Address Translation
 - Switch from Higher-Privilege Code to Lower

This Chapter

This chapter provides a detailed description of interrupt and exception handling in Protected Mode. This includes detailed coverage of:

- The IDT.
- Interrupt and Trap Gate operation.

x86 Instruction Set Architecture

- Task Gate operation.
- Interrupt and exception event categories.
- State save (and stack selection).
- The IRET instruction.
- Maskable hardware interrupts.
- Non-Maskable Interrupt (NMI).
- Machine Check exception.
- SM interrupt (SMI).
- Software interrupt instructions.
- Software exceptions.
- Interrupt/exception priority.

A detailed description of the Local and IO APICs can be found in “The Local and IO APICs” on page 1239.

The Next Chapter

The next chapter provides a detailed description of VM86 Mode (also known as Virtual 8086 Mode). This includes the following topics:

- Switching Between Protected Mode and VM86 Mode.
- Real Mode Application’s World View.
- Sensitive Instructions.
- Handling Direct IO.
- Handling Exceptions.
- Hardware Interrupt Handling in VM86 Mode
- Software Interrupt Instruction Handling
- Halt Instruction in VM86 Mode
- Protected Mode Virtual Interrupt Feature
- Registers Accessible in Real/VM86 Mode
- Instructions Usable in Real/VM86 Mode

Handler vs. ISR

The program executed to service a hardware interrupt or a software exception is commonly referred to as either a handler or an Interrupt Service Routine (ISR). For consistency and brevity’s sake, the author has elected to use the term *handler*.

Chapter 19: Protected Mode Interrupts and Exceptions

Real Mode Interrupt/Exception Handling

Real Mode handling of hardware and software interrupts as well as software exceptions was covered earlier in “Real Mode Interrupt/Exception Handling” on page 316. The following figures provide an overview of Real Mode event handling:

- Refer to Figure 19-1 on page 683.
- Refer to Figure 19-2 on page 684.

The remainder of this chapter focuses on interrupt and exception handling in Protected Mode.

Figure 19-1: Real Mode Interrupt Handling

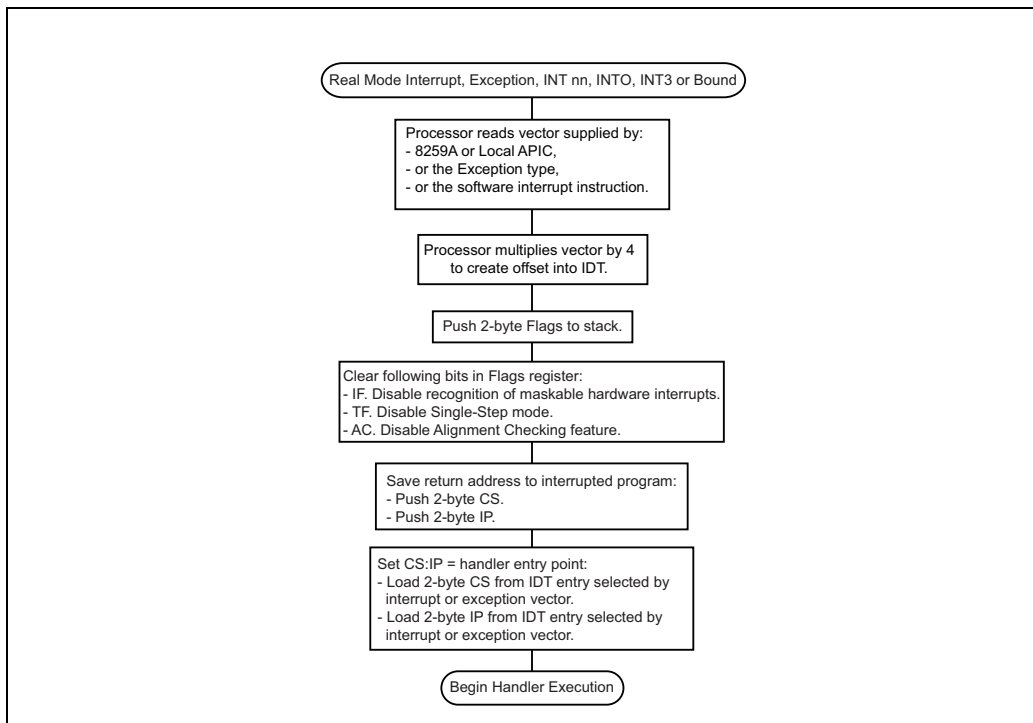
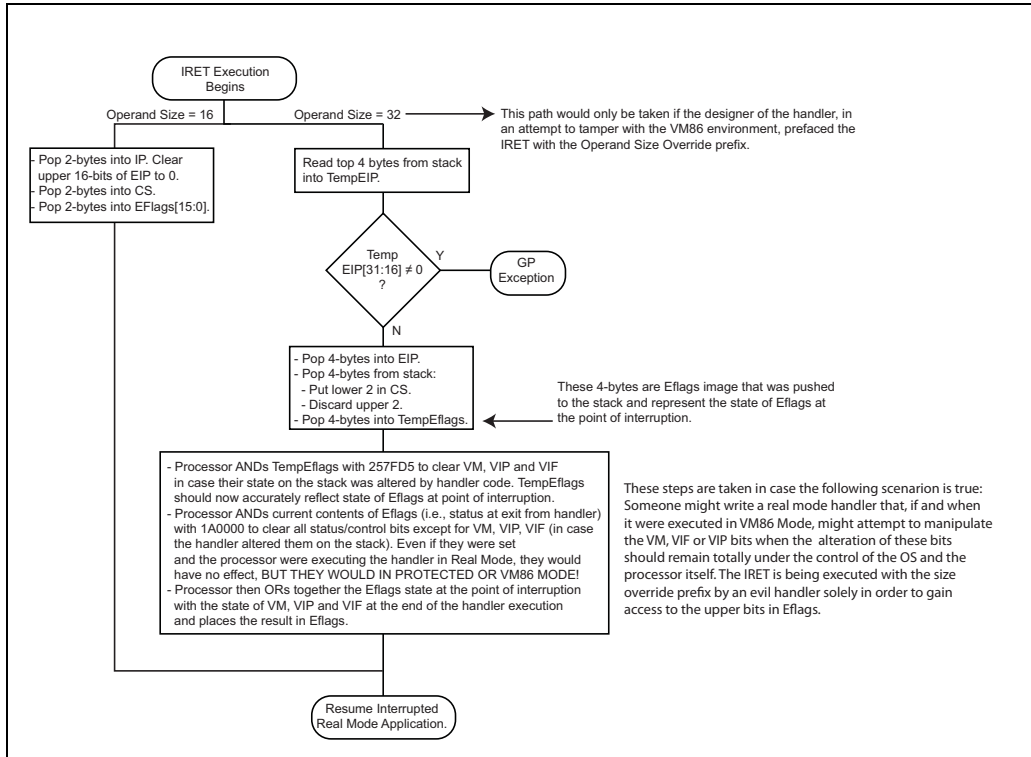


Figure 19-2: Return From Real Mode Handler To Interrupted Real Mode Application



The IDT

General

In Real Mode, the OS permits the logical processor to execute a single program at a time (i.e., multitasking is not supported). The BIOS, OS services, interrupt and exception handlers exist solely to support the program that is executing. This being the case, there is no need to restrict access to these services when the program executes a software interrupt instruction. The Protected Mode environment, on the other hand, was specifically designed to support a multitasking OS and must therefore provide protection from code being called by an entity with insufficient privilege. If the currently-executing program attempts to

20 *Virtual 8086 Mode*

The Previous Chapter

The previous chapter provided a detailed description of interrupt and exception handling in Protected Mode. This included detailed coverage of:

- The IDT.
- Interrupt and Trap Gate operation.
- Task Gate operation.
- Interrupt and exception event categories.
- State save (and stack selection).
- The IRET instruction.
- Maskable hardware interrupts.
- Non-Maskable Interrupt (NMI).
- Machine Check exception.
- SM interrupt (SMI).
- Software interrupt instructions.
- Software exceptions.
- Interrupt/exception priority.

This Chapter

This chapter provides a detailed description of VM86 Mode (also known as Virtual 8086 Mode). This includes the following topics:

- Switching Between Protected Mode and VM86 Mode.
- Real Mode Application's World View.
- Sensitive Instructions.
- Handling Direct IO.
- Handling Exceptions.
- Hardware Interrupt Handling in VM86 Mode
- Software Interrupt Instruction Handling
- Halt Instruction in VM86 Mode

- Protected Mode Virtual Interrupt Feature
- Registers Accessible in Real/VM86 Mode
- Instructions Usable in Real/VM86 Mode

The Next Chapter

The next chapter introduces the MMX register set and the original MMX instruction set. The SIMD programming model is introduced, how to deal with unpacked data as well as math underflows and overflows, and the elimination of conditional branches. Handling a task switch is described and the instruction set syntax is introduced.

A Special Note

The terms DOS task, VM86 task, and Real Mode task may be used interchangeably in this chapter (the vast majority of VM86 tasks are DOS tasks and, as such, intended to run in Real Mode). It should not be construed, however, that only DOS tasks are VM86 candidates. Any Real Mode application executed by a multitasking OS must be run under VM86 Mode.

Real Mode Applications Are Dangerous

The chapter entitled “Multitasking-Related Issues” on page 367 introduced some of the ways in which a Real Mode application might prove disruptive in a multitasking environment:

- Access memory belonging to currently-suspended programs.
- Communicate directly with IO ports (and thereby alter the state of device adapters).
- Call OS kernel code (including procedures it may not be allowed to access).
- Execute the CLI or STI instruction to disable or enable recognition of maskable hardware interrupts. The PUSHF and POPF instructions can also be used to change the state of the Eflags[IF] bit).
- Utilize a software interrupt instruction to call the BIOS or the Real Mode OS. A Real Mode application assumes it’s running under a Real Mode OS rather than a multitasking, Protected Mode OS. Consequently, all OS calls initiated by the application should be intercepted and passed to the host OS (or another program that substitutes for the DOS OS).

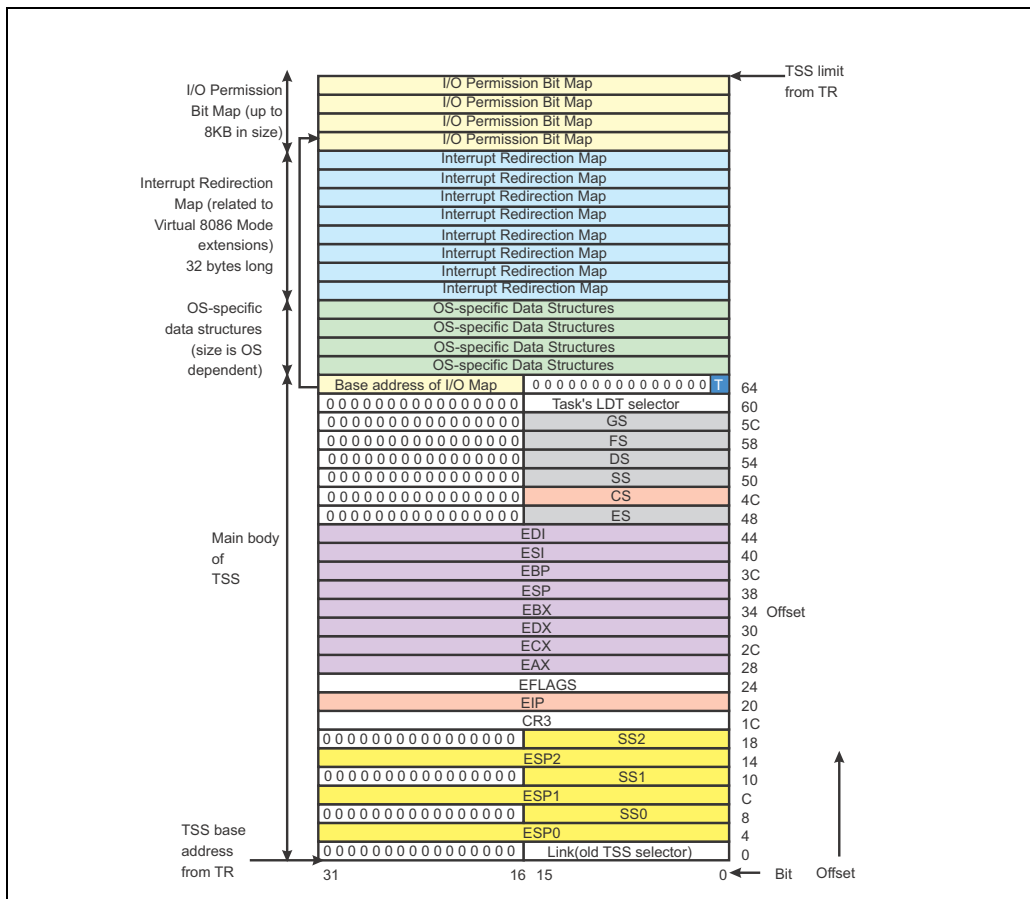
Solution: a Watchdog

When the scheduler switches to a Real Mode application, it sets the Eflags[VM] bit to one (see Figure 20-2 on page 787). This activates a logical processor mechanism (the VM86 logic) that monitors the behavior of the application on an instruction-by-instruction basis. Any operation that might prove destabilizing to the overall multitasking environment (referred to as a *sensitive* operation) is intercepted and an exception is generated to inform Virtual Machine Monitor (VMM) handler. There are a number of elements associated with this mechanism:

1. **Monitor logic.** The VM86 hardware detects the attempted execution of the sensitive instructions. First introduced in the 386, this mechanism was improved with the Pentium's addition of the VM86 Extensions (VME) feature.
2. **GP exception.** A GP exception is generated when a sensitive operation is detected.
3. **Monitor program.** If the GP exception handler determines it was invoked by the VM86 logic, it calls a special privilege level 0 procedure referred to as the VMM (Virtual Machine Monitor) to handle the event.
4. **TSS.** The scheduler creates a TSS data structure (see Figure 20-1 on page 786) for each Real Mode application. Several TSS elements are specifically-associated with VM86 Mode:
 - IO permission bitmap.
 - Interrupt redirection bitmap.
 - The VM bit in the Eflags register field.
 - The IOPL field in the Eflags register field.
5. **VM86 Extensions.** An OS may or may not activate the VM86 Mode Extensions by setting CR4[VME] = 1. If it is enabled, the following elements come into play:
 - Eflags[VIF] and Eflags[VIP] bits.
 - Interrupt redirection bitmap consultation.
6. **IOPL threshold.** The threshold value in Eflags[IOPL].

x86 Instruction Set Architecture

Figure 20-1: Task State Segment (TSS)



21 The MMX Facilities

The Previous Chapter

The previous chapter provided a detailed description of VM86 Mode (also known as Virtual 8086 Mode). This included the following topics:

- Switching Between Protected Mode and VM86 Mode.
- Real Mode Application's World View.
- Sensitive Instructions.
- Handling Direct IO.
- Handling Exceptions.
- Hardware Interrupt Handling in VM86 Mode
- Software Interrupt Instruction Handling
- Halt Instruction in VM86 Mode
- Protected Mode Virtual Interrupt Feature
- Registers Accessible in Real/VM86 Mode
- Instructions Usable in Real/VM86 Mode

This Chapter

This chapter introduces the MMX register set and the original MMX instruction set. The SIMD programming model is introduced, how to deal with unpacked data as well as math underflows and overflows, and the elimination of conditional branches. Handling a task switch is described and the instruction set syntax is introduced.

The Next Chapter

The next chapter describes the SSE, SSE2 and SSE 3 instruction sets and summarizes the SSSE3, SSE4.1 and SSE4.2 instruction sets.

Introduction

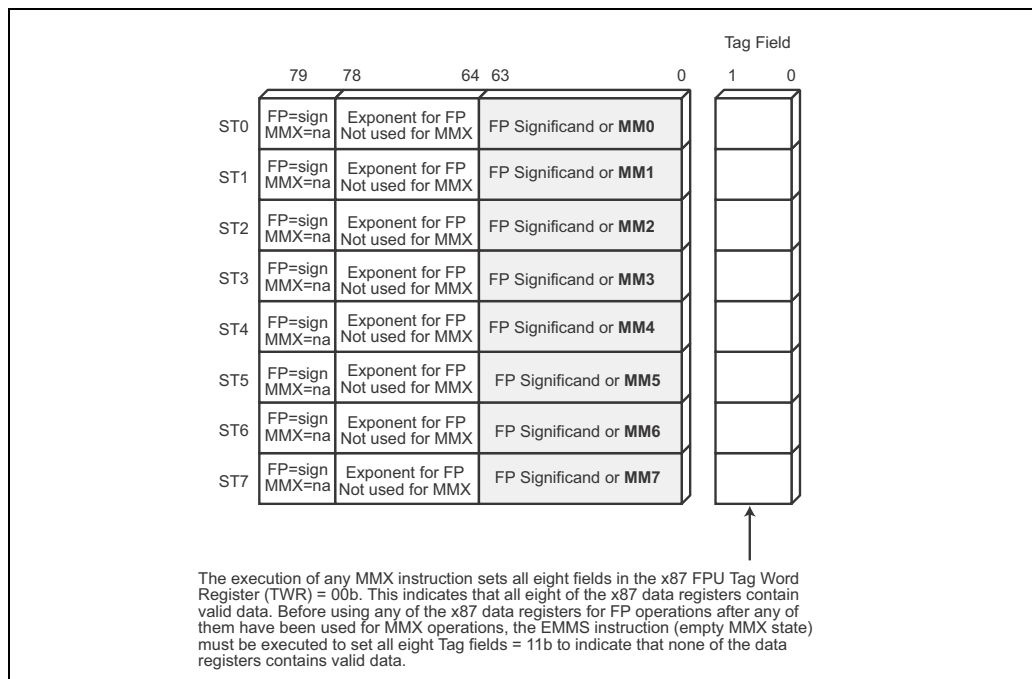
The MMX instruction set was first introduced in the P55C version of the Pentium and consisted of 47 new instructions. In addition, there are eight MMX data registers (MM0 - MM7; see Figure 21-1 on page 836). As shown in the illustration, the lower 64-bits of the x87 FPU data registers perform double-duty:

- They are used as MMX data registers when MMX code is executed.
- They are used as x87 FPU data registers when x87 FPU code is executed.

Over the years, the core concept introduced with the advent of MMX—instructions capable of simultaneously operating on multiple data items packed into wide registers—has continued to expand as Intel introduced the SSE (Streaming SIMD Extensions, where SIMD stands for **S**ingle **I**nstruction operating on **M**ultiple **D**ata items), SSE2, SSE3, SSSE3 (Supplemental SSE3), SSE4.1 and SSE4.2 instruction sets.

This chapter is not intended as an in-depth look at the initial MMX instruction set. Rather, it provides an overview of the basic concepts introduced with the advent of the MMX instruction set.

Figure 21-1: MMX Register Set



Detecting MMX Capability

Whether or not a processor supports MMX is detected by executing a CPUID request type 1 in response to which the processor capabilities bit mask is returned in the EDX register (bit 23 = 1 indicates the processor supports MMX).

The Basic Problem

Assumptions

Refer to Figure 21-2 on page 840. As an example, assume that there are two video frame buffers in memory (it should not be assumed, however, that MMX is only intended for processing video data) and that the current video mode has the following characteristics:

- Each location in the two buffers represents the color of one pixel. The first buffer location corresponds to the first pixel on the left end of the first line of pixels on the screen, the second buffer location corresponds to the second pixel on the left end of the first line of pixels on the screen, etc.
 - A single location contains 8-bits (one byte), so a pixel can be any one of 256 possible colors (as represented by the values 00h - FFh).
 - The video controller is currently operating at a resolution of 1024 x 786, so each of the two video frame buffers consists of 786,432 locations.
-

The Operation

Now assume that the programmer wants to:

1. Read the byte from the first location of one buffer,
2. Read the first location of the other buffer,
3. Add the two bytes together, and
4. Store the result back into the first location of the second buffer.

Repeat the operation for every pixel in the two frame buffers.

Example: Processing One Pixel Per Iteration

This could be accomplished in the following manner:

1. Read a byte (a pixel) from buffer one into a 1-byte register (e.g., AL).
2. Read the corresponding byte from buffer two into another 1-byte register (e.g., BL).
3. Add AL and BL together and store the result in the respective location in buffer two.
4. Since the add may result in the generation of a carry, the programmer has to decide whether to discard the carry or to factor it into the result. If the possibility of a carry must be dealt with, the programmer must include a conditional branch after the add that will either:
 - jump to the code that handles the carry,
 - or loop back to process the next pixel from the two buffers.

As indicated in the illustration, this would result in 786,432 x 2 memory reads and 786,432 memory writes. This code would generate a tremendous number of memory accesses which may or may not hit on the logical processor's internal caches. Any misses would result in memory transactions being performed on the processor's external interface. This would degrade performance in two ways:

- In a multiprocessor system wherein the processors share the same external interface, the interface bandwidth available to the other processor(s) could be substantially impacted.
- Since the external interface typically operates at a substantially slower rate of speed than the logical processor, the memory accesses would be time consuming.

Example: Processing Four Pixels Per Iteration

The number of memory accesses could be reduced by reading four bytes at a time from each buffer:

1. Read four bytes from one buffer into a 32-bit GPR register (e.g., EAX).
2. Read the corresponding four bytes from the other buffer into another 32-bit GPR register (e.g., EBX).
3. Add EAX and EBX together and store the result in one of the buffers.

There is a problem inherent in such a simplistic approach. The four bytes read

22 *The SSE Facilities*

The Previous Chapter

The previous chapter introduced the MMX register set and the original MMX instruction set. The SIMD programming model was introduced, how to deal with unpacked data as well as math underflows and overflows, and the elimination of conditional branches. Handling a task switch was described and the instruction set syntax was introduced.

This Chapter

This chapter describes the SSE, SSE2 and SSE 3 instruction sets and summarizes the SSSE3, SSE4.1 and SSE4.2 instruction sets. It also completes the discussion of the IA-32 programming environment.

The Next Chapter

The next chapter provides a detailed description of the IA-32e OS environment. The following topics are covered:

- Mode Switching Overview.
- Virtual Memory Addressing in IA-32e Mode.
- In 64-bit Mode, Hardware-Enforced Flat Model.
- 64-bit Instruction Pointer.
- Instruction Fetching.
- RIP-Relative Data Accesses.
- Changes To Kernel-Related Registers and Structures.
- Address Translation Mechanism.
- GDT/LDT Descriptor Changes.
- GDT and GDTR Changes.
- LDT and LDTR Changes.
- IDT/IDTR and Interrupt/Exception Changes.
- Interrupt/Trap Gate Operational Changes.

x86 Instruction Set Architecture

- IRET Behavior.
- IA-32e Call Gate Operation.
- TR and TSS Changes.
- Register Set Expansion (in 64-bit Mode).
- Scheduler's Software-Based Task Switching Mechanism.

Chapter Objectives

This chapter is not intended to provide a detailed description of each instruction in the SSE instruction sets. That role is already more than adequately fulfilled by the Intel and AMD instruction set reference manuals. Rather, the intention here is two-fold:

- To provide a fundamental understanding of the SSE architecture and how it works.
- To provide additional descriptions of some of the more odd or interesting instructions in the SSE instruction sets.

The SSE facilities are described in the same order in which they were introduced: SSE, SSE2, SSE3, SSSE3, SSE4.1, and SSE 4.2.

SSE: MMX on Steroids

As shown in Table 22-1 on page 852, application performance enhancement has been steadily addressed over the years by the expanding role of the SIMD programming model. It would be incorrect, however, to describe the SSE facilities solely as an expansion of MMX's SIMD programming model. As will be demonstrated in this chapter, while many of the SSE instructions do, in fact, expand on the SIMD programming model, many other non-SIMD instructions were added to address application-specific performance issues.

Table 22-1: Evolution of SIMD Model

Instruction Set	Introduced in	Description
MMX	Pentium P55C	47 new instructions. As described in "The MMX Facilities" on page 835, the SIMD concept was first introduced with the advent of MMX: <ul style="list-style-type: none">• Introduction of the SIMD model.• Eight 64-bit registers (MM0:MM7) available for SIMD operations on packed data.

Chapter 22: The SSE Facilities

Table 22-1: Evolution of SIMD Model (Continued)

Instruction Set	Introduced in	Description
SSE	Pentium III	<p>70 new instructions. The SIMD model was expanded with the introduction of the SSE (Streaming SIMD Extensions) instruction set and register set:</p> <ul style="list-style-type: none">• Eight dedicated 128-bit data registers (XMM0 - XMM7) available for SIMD operations on packed data.• The ability to perform SIMD packed and scalar FP operations on 32-bit DP FP numerical values.• MXCSR Control/Status register added to control SSE FP operations:<ul style="list-style-type: none">– SSE FP exception masking and status.– Enable/disable DAZ (Denormals-As-Zero) performance enhancement mode.– Enable/disable FTZ (Flush-to-Zero) performance enhancement mode. <p>When SSE was originally introduced, it was under the name Internet SSE (the word <i>Internet</i> was appended to just about everything during those crazy 1990s).</p>
SSE2	130nm Pentium 4	<p>144 new instructions.</p> <ul style="list-style-type: none">• Added the ability to perform both scalar and packed FP operations on 64-bit DP FP numbers.• The programmer can pack two, 64-bit DP FP numbers in each of two 128-bit XMM registers and then perform a packed FP operation on them (or between two numbers packed in an XMM register and two in memory).• MMX instructions enhanced to perform operations on data items packed in the XMM registers (prior to this, MMX instruction could only operate on data in MMX registers).• The CLFLUSH, MFENCE, LFENCE and new streaming store (commonly referred to as non-temporal store) instructions were added.• The PAUSE instruction was added to enhance performance when Hyper-Threading is enabled.

x86 Instruction Set Architecture

Table 22-1: Evolution of SIMD Model (Continued)

Instruction Set	Introduced in	Description
SSE3	90nm Pentium 4	<p>13 new instructions.</p> <ul style="list-style-type: none">• One x87 FPU instruction (FISTTP) that improves x87 FP-to-integer conversion.• One SIMD integer instruction providing a specialized 128-bit unaligned data load.• Nine new SIMD FP instructions:<ul style="list-style-type: none">– 3 instructions that enhance performance of Load/Move/Duplicate operations.– 2 instructions that perform simultaneous add/subtract operations on SP FP numbers packed into a pair of XMM registers.– 4 instructions that perform horizontal rather than vertical add and subtract operations on packed FP numbers.• Two thread-synchronization instructions (MONITOR and MWAIT) that provide a more elegant solution than the PAUSE instruction (added in SSE2) in applications employing Hyper-Threading.• They can use GPRs rather than MMX or SSE registers.

23 *IA-32e OS Environment*

The Previous Chapter

The previous chapter described the SSE, SSE2 and SSE 3 instruction sets and summarized the SSSE3, SSE4.1 and SSE4.2 instruction sets.

This Chapter

This chapter provides a detailed description of the IA-32e OS environment. The following topics are covered:

- Mode Switching Overview.
- Virtual Memory Addressing in IA-32e Mode.
- In 64-bit Mode, Hardware-Enforced Flat Model.
- 64-bit Instruction Pointer.
- Instruction Fetching.
- RIP-Relative Data Accesses.
- Changes To Kernel-Related Registers and Structures.
- Address Translation Mechanism.
- GDT/LDT Descriptor Changes.
- GDT and GDTR Changes.
- LDT and LDTR Changes.
- IDT/IDTR and Interrupt/Exception Changes.
- Interrupt/Trap Gate Operational Changes.
- IRET Behavior.
- IA-32e Call Gate Operation.
- TR and TSS Changes.
- Register Set Expansion (in 64-bit Mode).
- Scheduler's Software-Based Task Switching Mechanism.

The Next Chapter

The next chapter provides a detailed description of the third generation address translation mechanism utilized in IA-32e Mode.

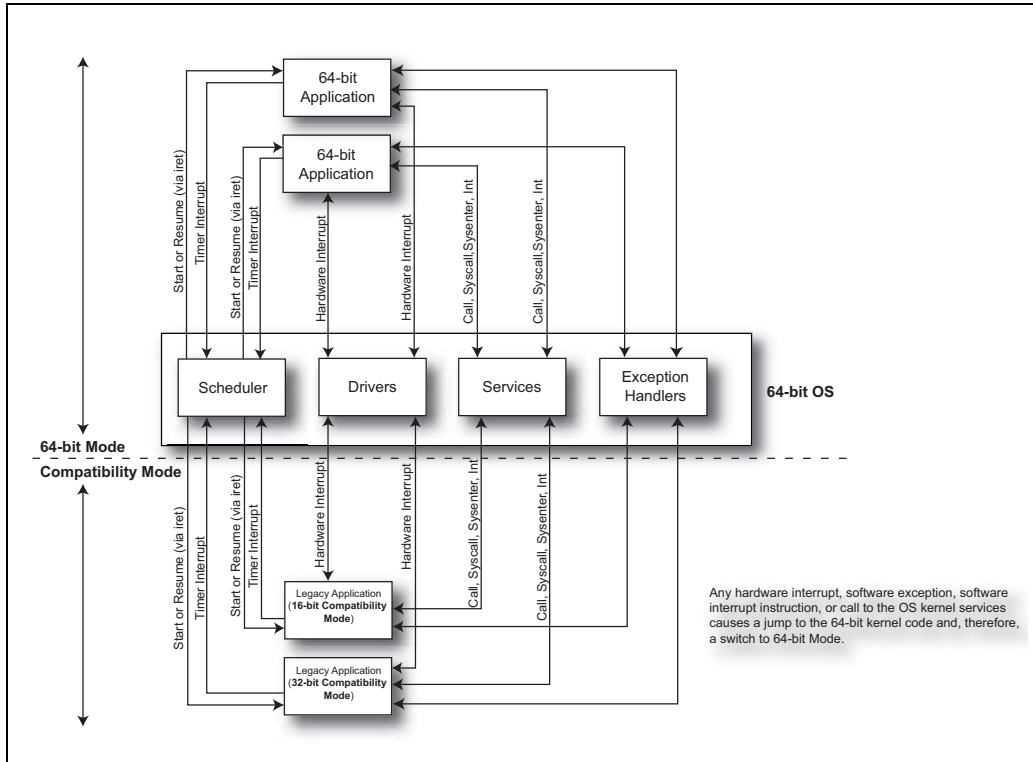
The Big Picture

Refer to Figure 23-1 on page 916. Ideally, all of the OS components are implemented as 64-bit code (i.e., they reside in 64-bit code segments and have full access to all of the logical processor's privileged, 64-bit facilities). Among other components, this would include:

- The **task scheduler**:
 1. Before starting or resuming a task, the scheduler would trigger the Local APIC timer.
 2. It then causes the logical processor to jump to the task:
 - If the task resides in a 64-bit code segment (i.e., the L bit in the CS descriptor = 1), the logical processor remains in 64-bit Mode.
 - If the task resides in a 16-bit legacy code segment (i.e., the CS descriptor's L and D bits both = 0), this causes the logical processor to automatically switch into 16-bit Compatibility Mode.
 - If the task resides in a 32-bit legacy code segment (the CS descriptor's L bit = 0 and D bit = 1), this causes the logical processor to automatically switch into 32-bit Compatibility Mode.
 3. In the background, while the logical processor executes the task, the timer continues to decrement.
 4. On timer expiration, the timer interrupt causes the logical processor to perform a far jump back to the scheduler. Since the CS descriptor selected by the far jump selects a CS descriptor wherein the L bit = 1, the logical processor switches back to 64-bit Mode (if it was in Compatibility Mode because the interrupted task was a legacy task).
- All **device drivers (including all hardware interrupt handlers)**. In IA-32e Mode, it is a rule that all interrupt handlers must reside in 64-bit code segments:
 - Upon detection of any hardware interrupt, software exception, or the attempted execution of a software interrupt instruction (INT nn, BOUND, INT3, or INTO), the logical processor would therefore reenter 64-bit Mode (if it was in Compatibility Mode because the interrupted task was a legacy task). Note: the INTO and BOUND instructions are illegal in 64-bit Mode.

- **OS services.** OS services are typically called in one of the following ways:
 - **Far Call through a Call Gate.** The execution of a far call instruction wherein the segment selector portion of the branch target address selects a Call Gate descriptor in the GDT or LDT. Since in IA-32e Mode it is a rule that the procedure pointed to by a Call Gate must reside in a 64-bit code segment, the far call (or a far jump for that matter) causes a switch to 64-bit Mode.
 - **Software Interrupt.** The execution of a software interrupt instruction will select either an Interrupt Gate or a Trap Gate descriptor in the IDT and, since it is a rule in IA-32e Mode that all IDT descriptors must point to handlers in 64-bit code segments, the interrupt causes a switch to 64-bit Mode.
 - **SYSCALL instruction.** Used to make a call to the OS services:
 - In Intel processors, the SYSCALL instruction can only be executed successfully by 64-bit applications. Otherwise it results in an Undefined Opcode exception.
 - In AMD processors, the SYSCALL instruction can be executed in any mode.
 - **SYSENTER instruction.** Used to make a call to the OS services:
 - In Intel processors, the SYSENTER instruction can be executed in any mode.
 - In AMD processors, the SYSENTER instruction can only be executed successfully in legacy Protected Mode. Otherwise it results in an Undefined Opcode exception.
- **Software exception handlers.** In IA-32e Mode, it is a rule that all software exception handlers must reside in 64-bit code segments. Upon detection of any software exception, the logical processor would therefore reenter 64-bit Mode (if it was in Compatibility Mode because the interrupted task was a legacy task).

Figure 23-1: 64-bit OS Environment



Mode Switching Overview

Booting Into Protected Mode

The basic boot sequence is as follows:

1. Immediately after system power-up, the reset signal remains asserted until the power supply voltages have achieved their required levels and stabilized.
2. The reset signal is deasserted to the processor. One of the logical processors is selected as the Bootstrap processor, begins operation in Real Mode and initiates code fetching from the boot ROM.

24 IA-32e Address Translation

The Previous Chapter

The previous chapter provided a detailed description of the IA-32e OS environment. The following topics were covered:

- Mode Switching Overview.
- Virtual Memory Addressing in IA-32e Mode.
- In 64-bit Mode, Hardware-Enforced Flat Model.
- 64-bit Instruction Pointer.
- Instruction Fetching.
- RIP-Relative Data Accesses.
- Changes To Kernel-Related Registers and Structures.
- Address Translation Mechanism.
- GDT/LDT Descriptor Changes.
- GDT and GDTR Changes.
- LDT and LDTR Changes.
- IDT/IDTR and Interrupt/Exception Changes.
- Interrupt/Trap Gate Operational Changes.
- IRET Behavior.
- IA-32e Call Gate Operation.
- TR and TSS Changes.
- Register Set Expansion (in 64-bit Mode).
- Scheduler's Software-Based Task Switching Mechanism.

This Chapter

This chapter provides a detailed description of the third generation address translation mechanism utilized in IA-32e Mode. This includes the following topics:

x86 Instruction Set Architecture

- Theoretical Address Space Size.
- Limitation Imposed by Current Implementation.
- Four-Level Lookup Mechanism.
 - Address Space Partitioning.
 - The Address Translation.
 - Initializing CR3.
 - Step 1: PML4 Lookup.
 - Step 2: PDPT Lookup.
 - Step 3: Page Directory Lookup.
 - Step 4: Page Table Lookup.
 - Page Protection Mechanisms in IA-32e Mode.
 - Page Protection in Compatibility Mode.
 - Page Protection in 64-bit Mode.
 - Don't Forget the Execute Disable Feature!
- TLBs Are More Important Than Ever.
- No 4MB Page Support.

The Next Chapter

The next chapter provides a detailed description of the Compatibility SubMode of IA-32e Mode. This includes the following topics:

- Initial Entry to Compatibility Mode.
- Switching Between Compatibility Mode and 64-bit Mode.
- Differences Between IA-32 Mode and Compatibility Mode.
- Memory Addressing.
- Register Set.
- Exception and Interrupt Handling.
- OS Kernel Calls.
- IRET Changes.
- Segment Load Instructions.

Theoretical Address Space Size

Theoretically, the 3rd generation address translation mechanism utilized in IA-32e Mode would support the translation of 64-bit virtual addresses to 52-bit physical addresses. This would provide the OS with the following virtual and physical memory space sizes:

- 2^{64} virtual addressing would permit the OS to assign virtual address ranges to applications within an 16EB (exabyte) virtual address space.

Chapter 24: IA-32e Address Translation

- 2^{52} physical addressing would permit the OS to map a 64-bit virtual address to any physical memory address in a 4PB (petabyte) physical memory address space.

Limitation Imposed by Current Implementations

Current implementations do not support the theoretical maximum virtual or physical address ranges, however:

- A 2^{48} (256TB—terabyte) virtual address space is currently supported.
- A 2^{40} (1TB) physical memory address space (and, in some high-end AMD products, 2^{48}) is currently supported.

In other words, in IA-32e Mode the 3rd generation address translation mechanism is presented with a 48-bit virtual address (sign-extended to 64-bits to form a 64-bit canonical address) which it translates into a 40-bit (or, in some high-end AMD products, a 41- or 48-bit) physical memory address.

Four-Level Lookup Mechanism

Address Space Partitioning

Refer to Figure 24-1 on page 987. In A-32e Mode, the partitioning of the 256TB virtual address space using a 48-bit address is viewed as follows:

- The overall 48-bit 256TB virtual space is divided into 512 blocks of 512GB each. Bits 47:39 identify the target 512GB block and selects the entry in the PML4 Directory associated with the addressed 512GB virtual address block (block 66 in the illustration). PML4 Entry 4 contains the physical base address of the Page Directory Pointer Table that catalogs the location of the 512 Page Directories (PDs) associated with the targeted 512GB block.
- Each 512GB block is sub-divided into 512 blocks of 1GB each. Bits 38:30 identify the target 1GB block and selects the entry in the Page Directory Pointer Table associated with the addressed 1GB virtual address block (block 97 in the illustration). Page Directory Pointer Table Entry 97 contains the physical base address of the Page Directory that catalogs the location of the 512 Page Tables associated with the targeted 1GB block.
- Each 1GB block is sub-divided into 512 blocks of 2MB each. Bits 29:21 identify the target 2MB block and selects the entry in the Page Directory associ-

ated with the addressed 2MB virtual address block (block 8 in the illustration). Page Directory Entry 8 contains either:

- The 4KB-aligned physical base address of the Page Table (PT) that catalogs the location of the 512 4KB pages in the targeted 2MB block;
- Or the physical base address of the targeted 2MB page in memory.
- Each 2MB block is sub-divided into 512 pages of 4KB each. Bits 20:12 identify the target 4KB page and selects the entry in the Page Table associated with the addressed 4KB virtual page (page 34 in the illustration). Page Table Entry 34 contains the physical base address of the target 4KB page.
- The lower 12-bits (11:0) identifies the target location within the page.

25 *Compatibility Mode*

The Previous Chapter

The previous chapter provided a detailed description of the third generation address translation mechanism utilized in IA-32e Mode. This included the following topics:

- Theoretical Address Space Size.
- Limitation Imposed by Current Implementation.
- Four-Level Lookup Mechanism.
 - Address Space Partitioning.
 - The Address Translation.
 - Initializing CR3.
 - Step 1: PML4 Lookup.
 - Step 2: PDPT Lookup.
 - Step 3: Page Directory Lookup.
 - Step 4: Page Table Lookup.
 - Page Protection Mechanisms in IA-32e Mode.
 - Page Protection in Compatibility Mode.
 - Page Protection in 64-bit Mode.
 - Don't Forget the Execute Disable Feature!
- TLBs Are More Important Than Ever.
- No 4MB Page Support.

This Chapter

This chapter provides a detailed description of the Compatibility SubMode of IA-32e Mode. This includes the following topics:

- Initial Entry to Compatibility Mode.
- Switching Between Compatibility Mode and 64-bit Mode.
- Differences Between IA-32 Mode and Compatibility Mode.
- Memory Addressing.
- Register Set.

x86 Instruction Set Architecture

- Exception and Interrupt Handling.
- OS Kernel Calls.
- IRET Changes.
- Segment Load Instructions.

The Next Chapter

The next chapter provides an overview of the following:

- 64-bit Register Set.
- EFER (Extended Features Enable) Register.
- Sixteen 64-bit Control Registers.
- 64-bit Rflags Register.
- Sixteen 64-bit GPRs.
- Kernel Data Structure Registers in 64-bit Mode.
- SSE Register Set Expanded in 64-bit Mode.
- Debug Breakpoint Registers.
- Local APIC Register Set.
- x87 FPU/MMX Register Set.
- Architecturally-Defined MSRs.

Initial Entry to Compatibility Mode

This subject was introduced in “Initial Switch from IA-32 to IA-32e Mode” on page 917. A detailed description may be found in “Transitioning to IA-32e Mode” on page 1139.

Switching Between Compatibility Mode and 64-bit Mode

Once the logical processor has entered into IA-32e Mode, its mode of operation is controlled by the state of the D (Default) and L (Long Mode) bits in the current code segment’s descriptor. Additional information about switching between Compatibility Mode and 64-bit Mode may be found in the following sections:

- “CS D and L Bits Control IA-32e SubMode Selection” on page 920.
- “Scheduler’s Software-Based Task Switching Mechanism” on page 977.

Differences Between IA-32 Mode and Compatibility Mode

IA-32 Background

Except for those differences cited in this chapter, Compatibility Mode works exactly like 16- and 32-bit Protected Mode. Detailed descriptions of Protected Mode operation may be found in “Part 2: IA-32 Mode”.

Unsupported IA-32 Features

The following IA-32 Mode features are not supported in IA-32e Mode (and are therefore not supported in Compatibility Mode):

- The hardware-based task switching mechanism. Due to this constraint, the following changes take effect:
 - Task Gates may not be selected in system tables (i.e., the GDT, LDTs and the IDT).
 - A far jump or a far call may not select:
 - A TSS descriptor in the GDT.
 - A Task Gate in the GDT or LDT.TSS descriptors are still used, however.
 - The TSS fields associated with automated task switching have been eliminated.
 - Execution of the IRET instruction when CR0[NT] = 1 does not cause a task switch.
- Real Mode.
- VM86 Mode.
- The 1st and 2nd generation address translation mechanisms.
- 4MB pages.

Changes to the OS Environment

As previously described in “IA-32e OS Environment” on page 913, the changes to the OS environment listed in Table 25-1 on page 1012 take effect in IA-32e Mode.

x86 Instruction Set Architecture

Table 25-1: OS Environment Changes in IA-32e Mode

Element	Description of Change
CS descriptor	A previously reserved bit, now defined as the L (Long Mode) bit, sets the logical processor's operating mode to 64-bit (L = 1) or Compatibility Mode (L = 0).
IDT	All entries in the IDT must consist of 16-byte Interrupt and Trap Gate descriptors that point to 64-bit handlers.
Call Gate	All Call Gate descriptors must be 16-byte descriptors that point to 64-bit OS services.
TSS descriptor	All TSS descriptors are 16-bytes long.
TSS	The TSS data structure format has been restructured to: <ul style="list-style-type: none">• Support the IST mechanism (refer to "Interrupt/Exception Stack Switch" on page 976).• Eliminate the register save/restore area used by the automated task switch mechanism (because it is not supported).• Eliminate the Interrupt Redirection bitmap (because VM86 Mode is not supported).• Eliminate the Link field used by the automated task switch mechanism.• Eliminate the debug Trap bit.
LDT descriptor	All LDT descriptors are 16-bytes long.
Address Translation	<ul style="list-style-type: none">• Only the 3rd generation translation mechanism is supported.• 4MB pages are not supported.
Virtual address	Although the virtual addresses generated by the legacy segmentation mechanism are 24- or 32-bits in length, they are zero-extended to form a 64-bit virtual address (in canonical form).
Physical address	Current implementations support a 2^{40} (1TB), 2^{41} , or 2^{48} physical address space.
CR3	CR3 is 64-bits wide enabling the top-level address translation table (the PML4) to be located anywhere in physical memory.

26 *64-bit Register Overview*

The Previous Chapter

The previous chapter provided a detailed description of the Compatibility Sub-Mode of IA-32e Mode. This included the following topics:

- Initial Entry to Compatibility Mode.
- Switching Between Compatibility Mode and 64-bit Mode.
- Differences Between IA-32 Mode and Compatibility Mode.
- Memory Addressing.
- Register Set.
- Exception and Interrupt Handling.
- OS Kernel Calls.
- IRET Changes.
- Segment Load Instructions.

This Chapter

This chapter provides an overview of the following:

- 64-bit Register Set.
- EFER (Extended Features Enable) Register.
- Sixteen 64-bit Control Registers.
- 64-bit Rflags Register.
- Sixteen 64-bit GPRs.
- Kernel Data Structure Registers in 64-bit Mode.
- SSE Register Set Expanded in 64-bit Mode.
- Debug Breakpoint Registers.
- Local APIC Register Set.
- x87 FPU/MMX Register Set.
- Architecturally-Defined MSRs.

The Next Chapter

The next chapter describes the following topics:

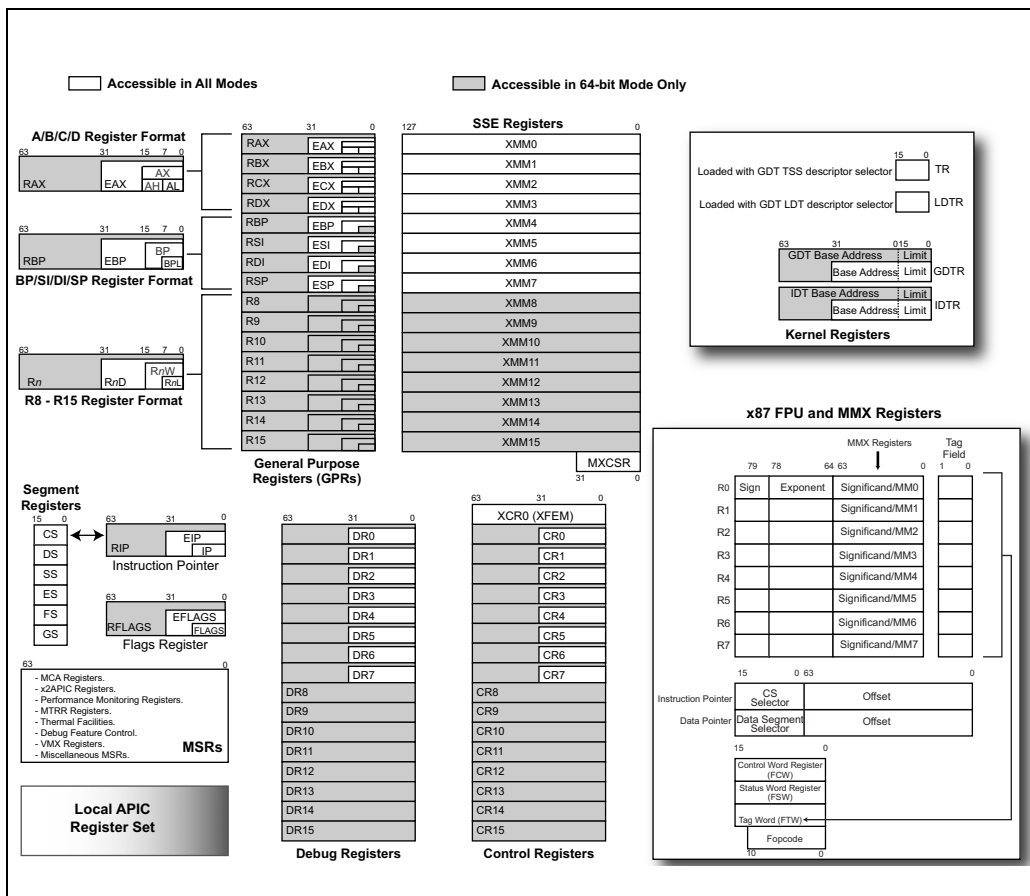
- Switching to 64-bit Mode.
- The Defaults.
- The REX Prefix.
- Addressing Memory in 64-bit Mode.
 - 64-bit Mode Uses a Hardware-Enforced Flat Model.
 - Default Virtual Address Size (and overriding it).
 - Actual Address Size Support: Theory vs. Practice.
 - Canonical Address.
 - Memory-based Operand Address Computation.
 - RIP-relative Data Addressing.
 - Near and Far Branch Addressing.
- Immediate Data Values in 64-bit Mode.
- Displacements in 64-bit Mode.

Overview of 64-bit Register Set

Figure 26-1 on page 1025 illustrates the registers that are visible to the programmer when the logical processor is operating in the 64-bit SubMode of IA-32e Mode. A description of the registers may be found in this chapter. A description of segment register usage in 64-bit Mode can be found in “Segment Register Usage in 64-bit Mode” on page 927.

Chapter 26: 64-bit Register Overview

Figure 26-1: Intel 64 Register Set



EFER (Extended Features Enable) Register

The EFER register (an MSR), pictured in Figure 26-2 on page 1026, plays a central role when switching a logical processor between legacy Protected Mode and IA-32e Mode. The bit critical to the switching process is the EFER[LME] bit (Table 26-1 on page 1026 describes the register's bit assignment).

x86 Instruction Set Architecture

Figure 26-2: EFER (Extended Features Enable) Register

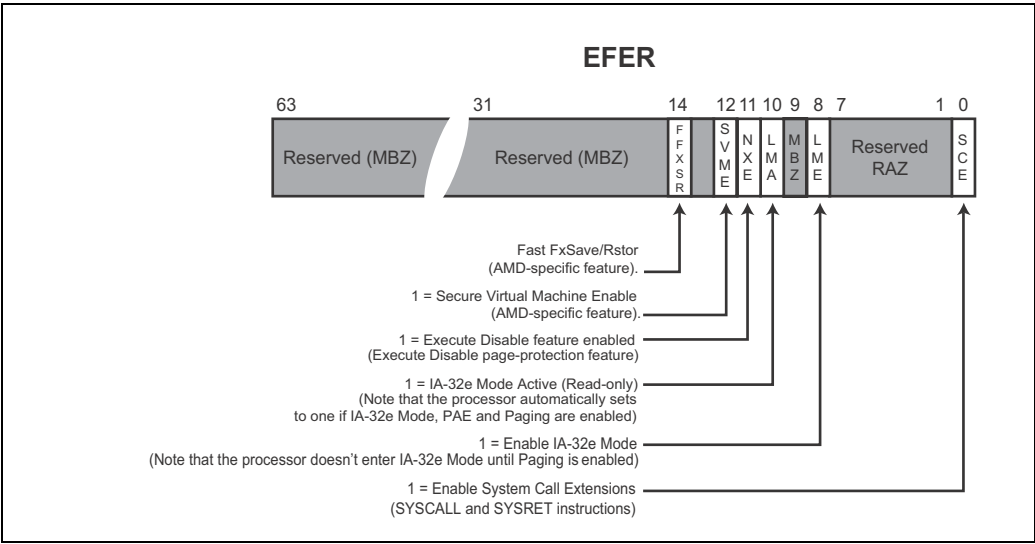


Table 26-1: EFER Register Bit Assignment

Bit	Description
SCE	System Call Enable. When set to one by the OS, enables the execution of the SYSCALL and SYSRET instructions which are used to make calls to the OS kernel. The OS sets this bit once it has set up the MSRs (STAR, LSTAR, CSTAR, SFMASK) used by these instructions: <ul style="list-style-type: none">• Due to low overhead, these instructions provide applications a way to perform OS kernel calls very quickly.• Accomplished using predefined call/return points. The logical processor skips many of the normal type and limit checks when changing segments (CS and SS).• The call entry points and return info are defined in a set of MSRs: STAR, LSTAR, CSTAR and SFMASK.• Refer to “SysCall Instruction” on page 1018 for more information.
LME	Enable IA-32e Mode. The OS kernel sets this bit before paging is enabled, but the logical processor doesn't actually enter IA-32e Mode until paging is subsequently turned on with physical address extensions enabled (CR0[PG] = 1 and CR4[PAE] = 1). As an interesting side-note, Intel uses AMD's acronym (Long Mode Enable, rather than IA-32e Mode Enable) for this bit.

27 *64-bit Operands and Addressing*

The Previous Chapter

The previous chapter provided an overview of the following:

- 64-bit Register Set.
- EFER (Extended Features Enable) Register.
- Sixteen 64-bit Control Registers.
- 64-bit Rflags Register.
- Sixteen 64-bit GPRs.
- Kernel Data Structure Registers in 64-bit Mode.
- SSE Register Set Expanded in 64-bit Mode.
- Debug Breakpoint Registers.
- Local APIC Register Set.
- x87 FPU/MMX Register Set.
- Architecturally-Defined MSRs.

This Chapter

This chapter covers the following topics:

- Switching to 64-bit Mode.
- The Defaults.
- The REX Prefix.
- Addressing Memory in 64-bit Mode.
 - 64-bit Mode Uses a Hardware-Enforced Flat Model.
 - Default Virtual Address Size (and overriding it).
 - Actual Address Size Support: Theory vs. Practice.
 - Canonical Address.
 - Memory-based Operand Address Computation.
 - RIP-relative Data Addressing.
 - Near and Far Branch Addressing.
- Immediate Data Values in 64-bit Mode.
- Displacements in 64-bit Mode.

The Next Chapter

The next chapter describes the following 64-bit related topics:

- New Instructions.
- Enhanced Instructions.
- Invalid Instructions.
- Reassigned Instructions.
- Instructions That Default to a 64-bit Operand Size.
- Branching in 64-bit Mode.
- NOP Instruction.
- FXSAVE/FXRSTOR.
- The Nested Task Bit (Rflags[NT]).
- SMM Save Area.

Helpful Background

An understanding of the 32-bit instruction format (see “32-bit Machine Language Instruction Format” on page 155) provides the background necessary for a complete understanding of the subject matter in this chapter.

Switching to 64-bit Mode

This subject was covered earlier in “Mode Switching Overview” on page 916.

The Defaults

Unless overridden by instruction prefixes, while the logical processor is executing code from a 64-bit code segment (i.e., the L bit = 1 in the code segment descriptor), its default assumptions are set as follows:

- The default operand size = 32-bits.
- The default address size = 64-bits.

It should be noted that some instructions have a default operand size of 64-bits without the use of the REX prefix.

The REX Prefix

Problem 1: Addressing New Registers

The IA-32 instruction set's ability to specify a register as an operand is limited by the following:

- As described in “Explicit Register Specification in ModRM Byte” on page 196 (see Figure 27-1 on page 1044; please note that this figure describes the instruction format in IA-32 Mode and Compatibility Mode, *not* in 64-bit Mode), the Operand 1 (i.e., the RM field) and Operand 2 (i.e., the Reg field) fields in the ModRM byte are each three bits wide.
- As described in “Explicit Register Specification in Opcode” on page 196 (see Figure 27-2 on page 1044), the register specification field found in the primary opcode byte of some instructions is a 3-bit field.

Obviously, the constraint imposed by a 3-bit register selection field limits the selection to 1 of 8 possible registers. In 64-bit Mode, however, the programmer has the ability to specify any of 16:

- GPRs.
- XMM registers.
- Control Registers.
- Debug Registers.

This obviously requires that the Reg, RM and the primary opcode byte's register specification fields be expanded from 3- to 4-bits wide in order to address the new registers when the logical processor is in 64-bit Mode.

Refer to Figure 27-3 on page 1045. In addition:

- The Base field in the SIB (Scale/Index/Base) byte used to specify the register containing the base address of a memory-based data structure is only a 3-bit field.
- The Index field in the SIB byte used to specify the register containing the location (i.e., the index) within the data structure is only a 3-bit field.

In order to specify any of the upper eight of the sixteen GPRs as the Index and Base registers, both of these bit fields must also be expanded from 3- to 4-bits.

x86 Instruction Set Architecture

Figure 27-1: The ModRM Byte's Operand 1 and 2 Fields Are Each 3-bits Wide

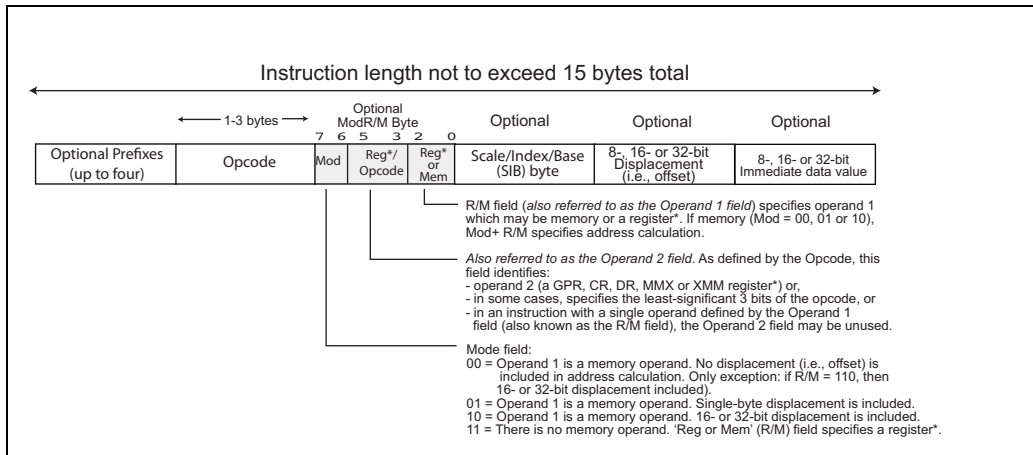
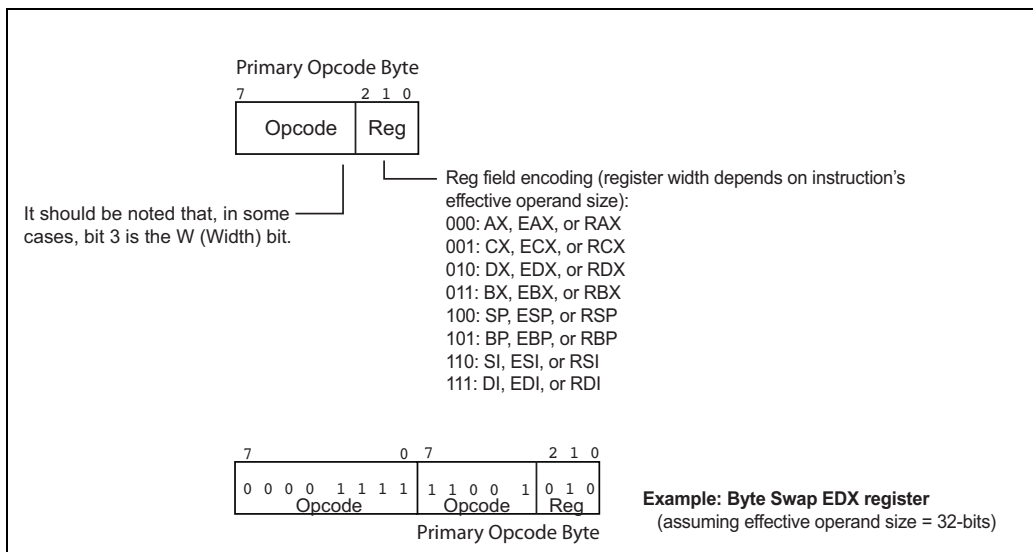


Figure 27-2: The Register Select Field in the Primary Opcode Byte Is 3-bits Wide



28 *64-bit Odds and Ends*

The Previous Chapter

The previous chapter described the following topics:

- Switching to 64-bit Mode.
- The Defaults.
- The REX Prefix.
- Addressing Memory in 64-bit Mode.
 - 64-bit Mode Uses a Hardware-Enforced Flat Model.
 - Default Virtual Address Size (and overriding it).
 - Actual Address Size Support: Theory vs. Practice.
 - Canonical Address.
 - Memory-based Operand Address Computation.
 - RIP-relative Data Addressing.
 - Near and Far Branch Addressing.
- Immediate Data Values in 64-bit Mode.
- Displacements in 64-bit Mode.

This Chapter

This chapter describes the following 64-bit related topics:

- New Instructions.
- Enhanced Instructions.
- Invalid Instructions.
- Reassigned Instructions.
- Instructions That Default to a 64-bit Operand Size.
- Branching in 64-bit Mode.
- NOP Instruction.
- FXSAVE/FXRSTOR.
- The Nested Task Bit (Rflags[NT]).
- SMM Save Area.

The Next Chapter

The next chapter describes the process of switching from Real Mode into Protected Mode. The following topics are covered:

- Real Mode Peculiarities That Affect the OS Boot Process.
- Typical OS Characteristics.
- Protected Mode Transition Primer.
- Example: Linux Startup.

New Instructions

General

Only two new instructions are defined when the logical processor is operating in 64-bit Mode:

- SwapGS.
- MOVSXD.

They are described in the next two sections.

SwapGS Instruction

The Problem

An application program may use the SysCall instruction to call the OS kernel services through a pre-defined entry point (the logical processor obtains the address from a special MSR register initialized by the OS). There is a problem, however. The kernel services must assume that the caller (i.e., the application program) expects the contents of the GPR registers to be preserved upon return from the kernel call. This being the case and considering that the kernel services will have to make use of one or more GPR registers in order to service the request, it will have to push the contents of one or more of the GPR registers to the stack. The problem lies in the following:

- With other system call mechanisms like Call Gates and software interrupts (INT nn), the switch to a new stack occurs automatically because there is a privilege level change occurring (unless the caller is a privilege level 0

entity). With SYSCALL, there is no automatic stack switch, so, without an instruction like this there is nowhere the called service can reliably save state information.

The SwapGS Solution

The SwapGS instruction may be used to solve this problem as shown in the following example. On entry to the kernel services, the following conditions are assumed to be true:

- The GS_Base MSR contains the 64-bit virtual base address of the caller's data area.
- The Kernel_GS_Base MSR contains the 64-bit virtual base address of a data structure within which the OS kernel stores critical information (e.g., a pointer to an empty stack area reserved for the kernel's use).

KernelServicesEntryPoint:

```
SwapGS          ;swap KernelGSBase MSR and GSBase MSR
mov gs:[SavedUserRSP], rsp;save caller's stack pointer
mov rsp, gs:[KernelStackPtr] ;set RSP = kernel stack ptr
push rax        ;save caller's GPR(s) to kernel stack
.
.               ;perform requested service
.
pop ---        ;restore caller's GPR(s) and stack pointer
mov rsp, gs:[SavedUserRSP];restore caller's stack pointer
SwapGS         ; restore caller's GSBase and KernelGSBase
ret
```

Although the stack problem doesn't exist for interrupt and exception handlers (see "Interrupt/Exception Stack Switch" on page 976), the SwapGS instruction can be used to quickly set the GS base address to point to the base address of a kernel-specific data structure that may contain information useful to the handler.

SwapGS has the following characteristics:

- SwapGS is a serializing instruction (see "Synchronizing Events" on page 618).
- The base address of the kernel-specific data structure is written to the KernelGSBase MSR (at MSR address C000_0102h) using the WRMSR instruction:
 - WRMSR may only be executed by privilege level 0 software.
 - The write will result in a GP exception if the address written to the register is not in canonical form.

x86 Instruction Set Architecture

- SwapGS uses a previously unused (and illegal) ModRM value accompanying the 2-byte opcode 0F01h (INVLPG; Invalidate Page Table Entry). Previously, only the memory forms (i.e., where the Mod field does not equal 11b) of this opcode were legal and the register forms (where the Mod field = 11b) were illegal. In 64-bit mode, when an 2-byte opcode of 0F01h is detected accompanied by a ModRM byte of 11 111 xxxb, the logical processor treats the xxxb bit field (i.e., the RM field) as an extension to the opcode which selects 1 of 8 instructions in a group of eight (see “Micro-Maps Associated with 2-byte Opcodes” on page 183). Currently, only RM = 000b is defined (as the SwapGS instruction) and the other seven values are currently undefined (and may be used to encode additional instructions in the future).

MOVSXD Instruction: Stretch It Out

In IA-32 and Compatibility Mode, the MOVSEX instruction—Move and Sign-Extend—sign-extends a byte or word operand to a full 32-bit dword. In 64-bit Mode, the ARPL (Adjust RPL field of segment selector) is reassigned as a new instruction, MOVSXD (Move Dword and Sign Extend to 64-bits). When used with the REX prefix, it sign-extends a 32-bit value to a full 64-bits.

Enhanced Instructions

Table 28-1 on page 1078 lists instructions that support a 64-bit operand size when prefaced by the REX prefix with REX[W] = 1.

Table 28-1: Instructions Enhanced in 64-bit Mode With REX[W] = 1

Mnemonic	Opcode (hex)	Description
CDQE	98	RAX = sign-extended EAX.
CMPSQ	A7	String compare operation. Compares quadword at address RSI with quadword at address RDI and sets the Rflags status flags accordingly.

29 *Transitioning to Protected Mode*

The Previous Chapter

The previous chapter described the following 64-bit related topics:

- New Instructions.
- Enhanced Instructions.
- Invalid Instructions.
- Reassigned Instructions.
- Instructions That Default to a 64-bit Operand Size.
- Branching in 64-bit Mode.
- NOP Instruction.
- FXSAVE/FXRSTOR.
- The Nested Task Bit (Rflags[NT]).
- SMM Save Area.

This Chapter

This chapter describes the process of switching from Real Mode into Protected Mode. The following topics are covered:

- Real Mode Peculiarities That Affect the OS Boot Process.
- Typical OS Characteristics.
- Protected Mode Transition Primer.
- Example: Linux Startup.

The Next Chapter

The next chapter describes the process of switching from Protected Mode into the Compatibility SubMode of IA-32e Mode. It then describes making the switch from Compatibility Mode into 64-bit Mode.

Real Mode Peculiarities That Affect the OS Boot Process

Immediately after the removal of reset, an x86 processor comes up in Real Mode and, to a goodly degree, emulates the 8086. Some of the logical processor's operational characteristics when operating in Real Mode are listed below:

- **It's a 16-bit world (unless overridden).** Unless overridden by prefacing an instruction with the Address Size and/or Operand Size Override prefixes (67h and 66h, respectively), the default segment offset address size and the default data operand size are both 16-bits.
- **An addressing anomaly:**
 - **Segment wrap-around** (see “Accessing Extended Memory in Real Mode” on page 307). When executing code on an 8086, specifying a segment base address of FFFF0h and any offset address between 0010h and FFFFh resulted in segment wraparound to the bottom of memory (specifically, to locations 00000h - 0FFEFh). The 8086 only had twenty address lines (19:0) and was therefore incapable of addressing memory above the 1MB address boundary.
 - **The HMA.** Every x86 processor since the advent of the 286, however, is capable of addressing memory above the 1MB boundary (referred to as extended memory). When operating in Real Mode, adding any offset value between 0010h - FFFFh to a segment base address of FFFF0h generates a carry bit on address bit 20 permitting software to address extended memory locations 100000h - 10FFEFh (the area of memory referred to as the HMA, or High Memory Area) without switching the logical processor into Protected Mode.
 - **A20 Gate.** Refer to “Accessing Extended Memory in Real Mode” on page 307. In order to boot the OS kernel into memory (the kernel will consume a large amount of system RAM), the logical processor must have to have the ability to address memory above the 1MB address boundary. If the A20 Gate is disabled, however, the A20 address line will always be 0 and, as a result, the logical processor will be unable to correctly address extended memory (i.e., memory above the 1MB address boundary).

Example OS Characteristics

This discussion makes the following assumptions:

- In preparation for booting the OS kernel into memory, the logical processor will be transitioned from Real Mode to Protected Mode. This is necessary

Chapter 29: Transitioning to Protected Mode

because today's highly-complex kernels are very large and will not fit in the 1MB of memory addressable in Real Mode.

- The Protected Mode memory model utilized will be a Flat Memory Model (see “IA-32 Flat Memory Model” on page 409) using the first generation virtual-to-physical address translation mechanism (i.e., paging).
- Virtually all of today's modern OSs utilize software-based task switching rather than the x86 processor's hardware-based tasking switching mechanism.
- The OS kernel and device drivers will execute at privilege level 0 while application programs will execute at level 3.

Flat Model With Paging

The discussion that follows assumes we will boot an OS that uses the Flat Memory Model and the virtual-to-physical address translation mechanism. This means that in the course of switching from Real Mode to Protected Mode, we will have to set up an appropriately formatted GDT as well as a set of virtual-to-physical address translation tables.

Software-Based Task Switching

Since the OS will not utilize the x86 hardware-based task switching mechanism, neither the GDT, the LDT (if the OS uses LDTs), nor the IDT will utilize Task Gate descriptors.

To give it full access to all of the logical processor's facilities, the OS code will execute at privilege level 0 while application programs will run at level 3. Since the logical processor is incapable of executing a jump or a call from a more-privileged to a less-privileged program, the OS task scheduler (which is privilege level 0 code) will have to use the software-based task switching mechanism described in “Scheduler's Software-Based Task Switching Mechanism” on page 977 to launch or resume an application program.

In order to avoid possible stack overflows, the OS kernel will utilize the logical processor's automatic stack switching ability (see “Automatic Stack Switch” on page 462) when making calls from one level to another. This being the case, we will have to define a TSS data structure (most modern OSs use one TSS for *all* tasks—see “Real World TSS Usage” on page 968; the TSS contains the pointers to the level 2, 1, and 0 stacks preallocated by the OS) and the TR (Task Register) is loaded with the GDT entry selector for the TSS descriptor.

Protected Mode Transition Primer

GDT Must Be In Place Before Switch to Protected Mode

At a minimum, at least a rudimentary GDT must be created in memory (see Figure 29-1 on page 1117) prior to switching to Protected Mode. The location of this initial GDT is dictated by the current state of the A20 gate (see “Accessing Extended Memory in Real Mode” on page 307):

- If the A20 Gate has not been enabled by software yet, the logical processor cannot access memory above the first MB while in Real Mode. The initial GDT must therefore be created in the first MB of memory space.
- If, on the other hand, the A20 Gate has been enabled by software, the logical processor is not solely restricted to the first MB of memory space but can also access extended memory locations 00100000h - 0010FFEFh (the HMA). In this case, the GDT may be placed in the HMA.

Keeping in mind that the desired OS memory configuration is the Flat Memory Model, the structure of the initial minimalist GDT is as follows (see Figure 29-1 on page 1117):

- **GDT size:** 24 bytes (8 bytes/entry x 3 entries).
- **Entry 0:** Entry 0 must be a null descriptor consisting of all zeros.
- **Entry 1 = CS Descriptor.** Refer to Figure 29-2 on page 1118. Entry 1 will be a code segment descriptor with the following characteristics:
 - Segment base address: 00000000h.
 - Segment size: 4GB.
 - Segment DPL: 0.
 - Other characteristics: Present bit = 1, S bit = 1, C bit = 0 for a Non-Conforming code segment, and R bit = 1 defining the CS as accessible for both code fetches and data reads.
- **Entry 2 = DS Descriptor.** Refer to Figure 29-3 on page 1119. Entry 2 will be a data segment descriptor with the following characteristics:
 - Segment base address: 00000000h.
 - Segment size: 4GB.
 - Segment DPL: 0.
 - Other characteristics: Present bit = 1, S bit = 1, R/W bit = 1 indicating it is a read/writable data segment, E bit = 0 indicating it can be used as an expand-up stack.

30 *Transitioning to IA-32e Mode*

The Previous Chapter

The previous chapter described the process of switching from Real Mode into Protected Mode. The following topics were covered:

- Real Mode Peculiarities That Affect the OS Boot Process.
- Typical OS Characteristics.
- Protected Mode Transition Primer.
- Example: Linux Startup.

This Chapter

This chapter describes the process of switching from Protected Mode into the Compatibility SubMode of IA-32e Mode. It then describes making the switch from Compatibility Mode into 64-bit Mode.

The Next Chapter

The next chapter provides a basic introduction to Virtualization Technology and covers the following topics:

- OS: I Am the God of All Things!
- Virtualization Supervisor: Sure You Are (:<)
- Root versus Non-Root Mode.
- Detecting VMX Capability.
- Entering/Exiting VMX Mode.
- Entering VMX Mode.
- Exiting VMX Mode.
- Virtualization Elements/Terminology.
- Introduction to the VT Instructions.

- Introduction to the VMCS Data Structure.
- Preparing to Launch a Guest OS.
- Launching a Guest OS.
- Guest OS Suspension.
- Resuming a Guest OS.
- Some Warnings Regarding VMCS Accesses.

No Need to Linger in Protected Mode

This chapter assumes that software will take the most efficient route possible from the removal of reset through Real Mode, Protected Mode, Compatibility Mode and, finally, to 64-bit Mode.

Entering Compatibility Mode

IA-32e Mode can only be entered by transitioning from legacy Protected Mode to Compatibility Mode. This transition is accomplished as follows:

1. Switch from Real Mode to legacy Protected Mode (this can be achieved using either 16- or 32-bit Protected Mode code). This topic was covered in “Transitioning to Protected Mode” on page 1113.
 - Note that there is no requirement that the address translation mechanism must be activated upon entering Protected Mode. Rather, the programmer may choose to immediately set up the 3rd generation address translation tables in preparation for the switch to the Compatibility SubMode of IA-32e Mode.

Note: This discussion assumes the logical processor is now fetching code from a 32- rather than a 16-bit code segment. The default operand and address sizes are therefore 32-bits.

DISABLE INTERRUPTS IN PREPARATION FOR SWITCH TO COMPATIBILITY MODE

2. Disable interrupts in preparation for switch from Protected Mode to IA-32e Mode:
 - Execute CLI to disable recognition of maskable hardware interrupts.
 - The programmer must ensure that the platform’s ability to deliver an NMI has been disabled. In a PC-compatible environment, this is accomplished by executing the following:
 - `mov al, 80`
 - `out 70h, al` ;performing an IO write to port 70h with bit 7 set to 1 will mask the platform’s ability to deliver an NMI to the logical processor.
 - The programmer also must ensure that no instructions generate software exceptions during the switch.

Chapter 30: Transitioning to IA-32e Mode

SET UP IA-32E COMPLIANT DATA STRUCTURES

3. Set up the 3rd generation address translation tables (see “IA-32e Address Translation” on page 983).
4. Point CR3 (see Figure 16-19 on page 528) to the top-level address translation table (i.e., the PML4 directory). Since CR3 is only 32-bits wide in Protected Mode, the PML4 directory’s physical base address must be in the lower 4GB.
5. Create an IA-32e compliant IDT containing 16-byte Interrupt Gates and Trap Gates (and no Task Gates).
6. Create an IA-32e compliant GDT containing (in addition to Protected-/Compatibility Mode-compliant data and stack segment descriptors):
 - An IA-32e compliant TSS descriptor (see Figure 23-9 on page 947).
 - An IA-32e compliant LDT descriptor (see Figure 23-8 on page 946).
 - IA-32e compliant Call Gate descriptors (see Figure 23-7 on page 945).
 - A 64-bit, privilege level 0 Non-Conforming code segment descriptor (see Figure 23-3 on page 922).
7. Create an IA-32e compliant LDT containing (in addition to Protected-/Compatibility Mode-compliant data and stack segment descriptors) IA-32e compliant Call Gate descriptors (see Figure 23-7 on page 945).
8. Create an IA-32e compliant TSS data structure (see Figure 23-9 on page 947).

EXECUTE 3-STEP PROCESS TO ENABLE IA-32E MODE

9. Enable 2nd generation address translation by setting CR4[PAE] to 1. This is the first required precondition for the transition to IA-32e Mode. Note that although the PAE feature is now enabled, address translation itself has not yet been activated (CR0[PG] is still 0).
10. Set EFER[LME] = 1 to enable IA-32e Mode. This is the second required precondition for the transition to IA-32e Mode. IA-32e Mode is not yet active, however.
11. Set CR0[PG] = 1 to activate paging. This is the third and final precondition. Since all three preconditions for IA-32e Mode activation have now been met, **IA-32e Mode is now activated**. The three prerequisites are:
 - CR4[PAE] = 1.
 - EFER[LME] = 1.
 - CR0[PG] = 1.
12. The L bit in the currently-active CS descriptor = 0, so the logical processor is not in 64-bit Mode. Rather, based on the state of the D bit in the selected CS descriptor, it is now in either the 16- or 32-bit Compatibility SubMode of IA-32e Mode:
 - D = 0. The logical processor is in 16-bit Compatibility Mode.
 - D = 1. The logical processor is in 32-bit Compatibility Mode.

A NOTE REGARDING IDENTITY ADDRESS MAPPING

Up until this moment, address translation was disabled. The memory address that the MOV CR0 instruction was fetched from was therefore treated as a physical rather than a virtual address. Address translation is now enabled, however, so the memory address used to fetch the next instruction (i.e., the one immediately following the MOV CR0 instruction which activated address translation) is treated as a virtual address and is therefore translated into a physical address. In order to fetch the instruction that immediately follows the MOV CR0 in physical memory, the address translation tables must translate this virtual address into the identical physical memory address (virtual = physical; referred to as identity address mapping).

AFTER SWITCH TO IA-32E MODE, LOAD SYSTEM REGISTERS

13. Execute the LIDT instruction to load the IDTR with the 16-bit size and 32-bit base address of the IA-32e compliant IDT.
14. Execute the LGDT instruction to load the GDTR with the 16-bit size and 32-bit base address of the IA-32e compliant GDT.
15. Execute the LLDT instruction to load the LDTR with the 16-bit GDT selector that points to the IA-32e compliant LDT descriptor in the GDT. In response, the logical processor loads the LDT descriptor into the invisible portion of the LDTR.
16. Execute the LTR instruction to load the TR with the 16-bit GDT selector that points to the 16-byte IA-32e compliant TSS descriptor in the GDT. In response, the logical processor loads the TSS descriptor into the invisible portion of the TR.

Switch to 64-bit Mode

1. Execute a far jump wherein the CS selector portion of the branch target address selects a 64-bit privilege level 0 Non-Conforming code segment descriptor in the GDT:
 - When the logical processor loads the CS descriptor into the invisible portion of the CS register, the one in the descriptor's L bit switches it into 64-bit Mode.
 - Segmentation is now disabled and the flat memory model is hardware-enforced: all segments (with the exception of the FS and GS data segments) have an assumed base address of 0 and a length of 2^{64} locations.
 - Since the far jump is executed while the logical processor is still in 32-bit Compatibility Mode, the offset portion of the branch target address is only 32-bits wide. The 64-bit RIP register is therefore loaded with the 32-bit address 0-extended to 64-bits. The target 64-bit code entry point must therefore reside in the lower 4GB of the 64-bit code segment.

31 Introduction to Virtualization Technology

The Previous Chapter

The previous chapter described the process of switching from Protected Mode into the Compatibility SubMode of IA-32e Mode. It then described making the switch from Compatibility Mode into 64-bit Mode.

This Chapter

This chapter provides a basic introduction to Virtualization Technology and covers the following topics:

- OS: I Am the God of All Things!
- Virtualization Supervisor: Sure You Are (:<)
- Root versus Non-Root Mode.
- Detecting VMX Capability.
- Entering/Exiting VMX Mode.
- Entering VMX Mode.
- Exiting VMX Mode.
- Virtualization Elements/Terminology.
- Introduction to the VT Instructions.
- Introduction to the VMCS Data Structure.
- Preparing to Launch a Guest OS.
- Launching a Guest OS.
- Guest OS Suspension.
- Resuming a Guest OS.
- Some Warnings Regarding VMCS Accesses.

The Next Chapter

The next chapter provides a detailed description of System Management Mode (SMM). It includes the following topics:

- What Falls Under the Heading of System Management?
- The Genesis of SMM
- SMM Has Its Own Private Memory Space
- The Basic Elements of SMM
- A Very Simple Example Scenario
- How the Processor Knows the SM Memory Start Address
- Normal Operation, (Including Paging) Is Disabled
- The Organization of SM RAM
- Entering SMM
- Exiting SMM
- Caching from SM Memory
- Setting Up the SMI Handler in SM Memory
- Relocating the SM RAM Base Address
- SMM in an MP System
- SM Mode and Virtualization

Just an Introduction?

Yes, rather than a detailed description of every aspect of virtualization, this chapter provides an introduction. Complete coverage of all aspects of virtualization would entail the addition of several hundred additional pages to an already oversized book. As such, it warrants treatment as a separate topic.

Detailed Coverage of Virtualization

Comprehensive coverage of all aspects of virtualization is available in the following MindShare class offerings:

- Comprehensive PC Virtualization:
 - Instructor-led class. Duration: 4 days.
 - Instructor-led internet class. Duration: 5 days.
- Fundamentals of PC Virtualization:
 - Instructor-led class. Duration: 1 day.
 - Instructor-led internet class. Duration: 1 day.
- Introduction to Virtualization Technology:
 - Self-paced E-Learning Module.

Chapter 31: Introduction to Virtualization Technology

- Introduction to PCI Express IO Virtualization:
 - Self-paced E-Learning Module.
- Comprehensive IO Virtualization:
 - Instructor-led class. Duration: 2 days.
 - Instructor-led internet class. Duration: 3 days.

Detailed information about MindShare's training classes and E-Learning modules may be found at www.mindshare.com.

The Intel Model

Although the basic concepts are the same, Intel and AMD have implemented vendor-specific approaches to virtualization. This chapter focuses on the Intel model.

OS: *I Am the God of All Things!*

A traditional OS (e.g., Windows XP, Windows 7, Mac OS X) has complete control of all of the logical processor's facilities:

- It executes at privilege level 0 and can therefore:
 - Access any register.
 - Control the logical processor's operational mode (i.e., whether it is in Real Mode, Protected Mode, IA-32e Mode, etc.).
 - Execute any instruction in the instruction set.
- It manages memory for all software (including itself).
- Under software control, it handles task switching among the various program's that are currently being executed.
- **Permission violation.** It manages all of the x86 protection mechanisms. If an application program attempts to touch something beyond its permission level (e.g., memory, an IO port, a Control Register, etc.), a software exception is generated which immediately returns control back to the OS kernel.
- **Action evaluation.** The kernel then evaluates the attempted action and determines how to handle it:
 - **Action permitted.** If the action would not prove detrimental to any other currently-suspended software entity, the OS may decide to permit it. In that case, the OS can execute the offending instruction itself (unlike the interrupted application program, it has sufficient privilege to do so) and then return to the interrupted program at the instruction immediately after the one that caused the exception.

- **Forbidden action.** If the attempted action is one that *would* prove detrimental to other currently-suspended software entities, the OS can handle it in either of two ways:
 - **Emulate attempted action.** The OS might choose to achieve the same goal by performing a set of actions that will not result in chaos for one or more other software entities that are currently-suspended.
 - **Abort the application.** If, in the OS's opinion, the attempted action cannot safely be permitted, it may choose to issue an alert message to the end user and then abort the errant application.

In a nutshell, the OS believes itself to be lord of all it surveys. While this *is* true under ordinary circumstances, it's *not* so when virtualization is enabled.

Virtualization Supervisor: *Sure You Are* (:<)

When the logical processor's Virtualization Technology (VT) feature (referred to as the *Virtual Mode Extensions*, or *VMX*) is enabled, the old gods (i.e., OSs) are subjugated to a new, all powerful God—the *Virtual Machine Monitor*, or *VMM* (otherwise referred to as the *hypervisor*). The hypervisor permits guest OSs to run under its guidance, allowing each to run either for a preallocated period of time (a *timeslice*) or until the guest OS attempts a sensitive operation that might prove harmful to another, currently-suspended guest OS or to the hypervisor itself (i.e., a *sensitive* operation).

Root versus Non-Root Mode

Refer to Figure 31-1 on page 1151. When the logical processor is executing the hypervisor code, it is said to be in *VMX Root Mode*. Conversely, it is in *VMX Non-Root Mode* when it is executing one of the guest OSs.

32 *System Management Mode (SMM)*

The Previous Chapter

The previous chapter provided a basic introduction to Virtualization Technology and covered the following topics:

- OS: I Am the God of All Things!
- Virtualization Supervisor: Sure You Are (:<)
- Root versus Non-Root Mode.
- Detecting VMX Capability.
- Entering/Exiting VMX Mode.
- Entering VMX Mode.
- Exiting VMX Mode.
- Virtualization Elements/Terminology.
- Introduction to the VT Instructions.
- Introduction to the VMCS Data Structure.
- Preparing to Launch a Guest OS.
- Launching a Guest OS.
- Guest OS Suspension.
- Resuming a Guest OS.
- Some Warnings Regarding VMCS Accesses.

This Chapter

This chapter provides a detailed description of System Management Mode (SMM). It includes the following topics:

- What Falls Under the Heading of System Management?
- The Genesis of SMM

x86 Instruction Set Architecture

- SMM Has Its Own Private Memory Space
- The Basic Elements of SMM
- A Very Simple Example Scenario
- How the Processor Knows the SM Memory Start Address
- Normal Operation, (Including Paging) Is Disabled
- The Organization of SM RAM
- Entering SMM
- Exiting SMM
- Caching from SM Memory
- Setting Up the SMI Handler in SM Memory
- Relocating the SM RAM Base Address
- SMM in an MP System
- SM Mode and Virtualization

The Next Chapter

The next chapter provides a detailed description of the Machine Check Architecture (MCA):

- The MCA Elements.
- The Global Registers.
- The Composition of a Register Bank.
- The Error Code.
- Cache Error Reporting.
- MC Exception Is Generally Not Recoverable.

What Falls Under the Heading of System Management?

The types of operations that typically fall under the heading of System Management are power management and management of the system's thermal environment (e.g., temperature monitoring in the platform's various thermal zones and fan control). It should be stressed, however, that system management is not necessarily limited to these specific areas.

The following are some example situations that would require action by the SM handler program:

- A laptop chipset implements a timer that tracks how long it's been since the hard drive was last accessed. If this timer should elapse, the chipset generates an SMI (System Management Interrupt) to the processor to invoke the SM handler program. In the handler, software checks a chipset-specific status register to determine the cause of the SMI (in this case, a prolonged ces-

Chapter 32: System Management Mode (SMM)

sation of accesses to the hard drive). In response, the SM handler issues a command to the hard disk controller to spin down the spindle motor (to save on energy consumption).

- A laptop chipset implements a timer that tracks how long it's been since the keyboard and/or mouse was used. If this timer should elapse, the chipset generates an SMI to the processor to invoke the SM handler program. In the handler, software checks a chipset-specific status register to determine the cause of the SMI (in this case, a prolonged cessation of user interaction). In response, the SM handler issues a command to the display controller to dim or turn off the display's backlighting (to save on energy consumption).
- In a server platform, the chipset or system board logic detects that a thermal sensor in a specific zone of the platform is experiencing a rise in temperature. It generates an SMI to the processor to invoke the SM handler program. In the handler, software checks a chipset-specific status register to determine the cause of the SMI (in this case, a potential overheat condition). In response, the SM handler issues a command to the system board's fan control logic to turn on an exhaust fan in that zone.

The Genesis of SMM

Intel first implemented SMM in the 386SL processor and it has not changed very much since then. While it was not present in the earlier 486 models, it was implemented in all of the later models of the 486 and in all subsequent x86 processors. SMM is entered by generating an SMI (System Management Interrupt) to the processor. Prior to the P54C version of the Pentium, the chipset could only deliver the interrupt to the processor by asserting the processor's SMI# input pin. Starting with the P54C (which was the first IA-32 processor to incorporate the Local APIC) and up to and including the Pentium III, the chipset could also deliver the interrupt to the processor by sending an SMI IPI (Inter-Processor Interrupt) message to the processor over the 3-wire APIC bus. With the advent of the Pentium 4, the 3-wire APIC bus was eliminated and IPIs (including the SMI IPI) are sent to and from a logical processor by performing a special memory write transaction on the processor's external interface.

With the advent of the P54C processor, SMM was enhanced to include the IO Instruction Restart feature (described in this chapter).

The base address of the area of memory assigned to System Management Mode (SMM) has a default value of 30000h. While it could be reprogrammed on the earlier IA-32 processors, the newly-assigned address had to be aligned on an address that was evenly divisible by 32K. Starting with the Pentium Pro, this constraint was eliminated.

SMM Has Its Own Private Memory Space

Prior to the generation of an SMI to the logical processor, the chipset directs all memory accesses generated by the logical processor to system RAM memory:

- When interrupted by an SMI, the logical processor signals to the chipset that all subsequent memory accesses generated by the logical processor are to be directed to a special, separate area of memory referred to as SM RAM.
- Upon concluding the execution of the SM handler program, the logical processor signals to the chipset that all subsequent memory accesses generated by the logical processor are to be directed to system RAM memory rather than SM RAM.

The platform vendor's implementation of SM RAM can be up to 4GB in size.

The Basic Elements of SMM

The following is a list of the basic elements associated with SMM:

- The processor's SMI# input.
- The APIC SMI IPI message.
- The chipset/system board logic responsible for monitoring conditions within the platform that might require an invocation of the SM handler program.
- Chipset's ability to assert SMI# to invoke the SMI handler.
- The chipset's ability to send an SMI IPI message to the logical processor to invoke the SMI handler.
- The Resume (RSM) instruction which must always be the last instruction executed in the SM handler.
- The SM RAM area.
- The logical processor's context state save/restore area (i.e., data structure) in SM memory.
 - 512-bytes for an IA-32 processor.
 - 1024-bytes for an Intel 64 processor.
- The SMI Acknowledge message was added to the message repertoire of the Special transaction.
- On processors that utilize the FSB external interface, the processor's SMMEM# output (also referred to as the EFX4# output).
- The chipset's ability to discern when the processor is addressing regular RAM memory versus SM RAM memory. On processors that utilize the FSB external interface, it does this by monitoring for the processor's issuance of the SMI Acknowledge message and whether or not the processor is asserting the SMMEM# signal during a processor-initiated memory transaction.

33 *Machine Check Architecture (MCA)*

The Previous Chapter

The previous chapter provided a detailed description of System Management Mode (SMM). It included the following topics:

- What Falls Under the Heading of System Management?
- The Genesis of SMM
- SMM Has Its Own Private Memory Space
- The Basic Elements of SMM
- A Very Simple Example Scenario
- How the Processor Knows the SM Memory Start Address
- Normal Operation, (Including Paging) Is Disabled
- The Organization of SM RAM
- Entering SMM
- Exiting SMM
- Caching from SM Memory
- Setting Up the SMI Handler in SM Memory
- Relocating the SM RAM Base Address
- SMM in an MP System
- SM Mode and Virtualization

This Chapter

This chapter provides a detailed description of the Machine Check Architecture (MCA):

- The MCA Elements.
- The Global Registers.

- The Composition of a Register Bank.
- The Error Code.
- Cache Error Reporting.
- MC Exception Is Generally Not Recoverable.

The Next Chapter

The next chapter provides a complete description of the Local and IO APICs. It includes:

- A Short History of the APIC's Evolution.
- Before the APIC.
- MP Systems Need a Better Interrupt Distribution Mechanism.
- Detecting Presence/Version/Capabilities of Local APIC.
- Local APIC's Initial State.
- Enabling/Disabling the Local APIC.
- Mode Selection.
- The Local APIC Register Set.
- Local APIC ID Assignments and Addressing.
 - ID Assignment in xAPIC Mode.
 - ID Assignment in x2APIC Mode.
 - Local APIC Addressing.
 - Lowest-Priority Delivery Mode.
- Local APIC IDs Are Stored in the MP and ACPI Tables.
- Accessing the Local APIC ID.
- An Introduction to the Interrupt Sources.
- Introduction to Interrupt Priority.
- Task and Processor Priority.
- IO/Local APICs Cooperate on Interrupt Handling.
- Message Signaled Interrupts (MSI).
- Interrupt Delivery from Legacy 8259a Interrupt Controller.
- SW-Initiated Interrupt Message Transmission.
- x2APIC Mode's Self IPI Feature.
- Locally Generated Interrupts.
 - The Local Vector Table.
 - Local Interrupt 0 (LINT0).
 - Local Interrupt 1 (LINT1).
 - The Local APIC Timer.
 - The Performance Counter Overflow Interrupt.
 - The Thermal Sensor Interrupt.
 - Correctable Machine Check (CMC) Interrupt.
 - The Local APIC's Error Interrupt.

Chapter 33: Machine Check Architecture (MCA)

- The Spurious Interrupt Vector.
- Boot Strap Processor (BSP) Selection.
- How the APs are Discovered and Configured.

Why This Subject Is Included

When first introduced in x86 processors, the Machine Check Architecture error logging facility was not architecturally defined. Rather, it was a processor-specific addition to the Pentium and Pentium Pro processors with no guarantee that it would be implemented in subsequent processors or, if it was, that it would be implemented in the same manner. It was only with the advent of the Pentium 4 that it was defined as part of the x86 software architecture.

MCA = Hardware Error Logging Capability

The Machine Check Architecture facility consists of a set of error logging registers and the Machine Check exception. During a power up session, a logical processor may experience one or more hardware-related errors internally or on its external interface. Such errors can be divided into two basic categories:

- Soft errors that are automatically corrected by the processor hardware.
- Hard errors that cannot be automatically corrected.

It is expected that the OS will start a daemon that runs in background and periodically examines the MCA registers to determine if any soft errors have been logged since the last time the registers were examined. If so, the application snapshots the errors in a non-volatile storage medium (e.g., on the hard drive or in flash memory) and then clears the errors from the register set (the registers are then available to record any additional errors that may occur in the future).

If a hard error is detected and the Machine Check exception has been enabled, the logical processor records the error in the register set and also generates a Machine Check exception [see “Machine Check Exception (18)” on page 778 for a detailed description] to report it. In the Machine Check exception handler, the error recorded in the register set is read, recorded in the non-volatile storage medium, and possibly displayed on the console. Generally speaking, software cannot recover from most hard errors.

The MCA Elements

The Machine Check Architecture first appeared in the Pentium in rudimentary form (consisting of only two registers), but was greatly expanded with the advent of the Pentium Pro.

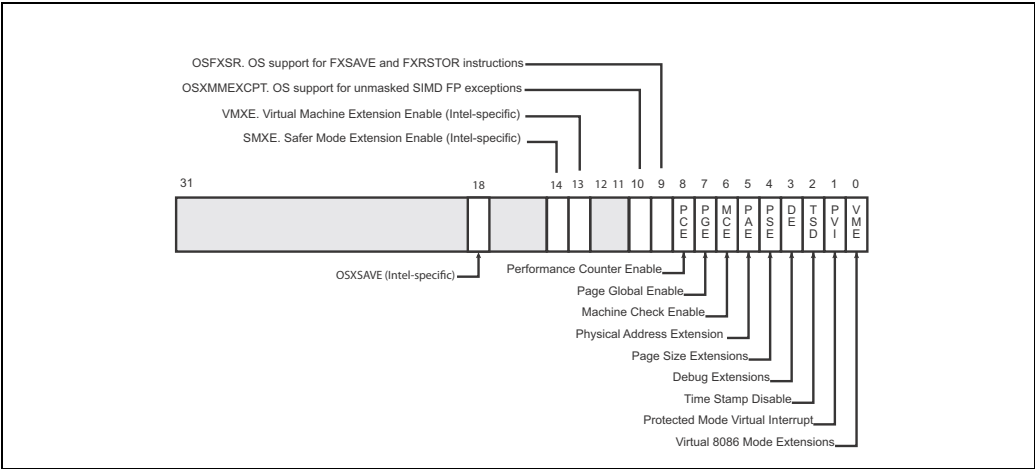
The MCA actually consists of two capabilities (detected by performing a CPUID request type 1 and checking the returned EDX capabilities bit mask): the ability to generate the Machine Check exception and the presence of the MCA register set.

The Machine Check Exception

Although it is optional whether or not a logical processor possesses the ability to generate a Machine Check exception when an uncorrectable hardware error has been detected, all current-day processors support this capability. Whether or not a processor supports this ability is indicated by executing a CPUID request type 1 and checking the EDX[MCE] bit. If it supports the generation of the Machine Check exception, this capability is enabled by setting CR4[MCE] = 1 (see Figure 33-1 on page 1210).

On the P6 processors, the Pentium 4 (and its Xeon and Celeron derivatives), and the Pentium M (and, to the author's knowledge, current-day processors), the Machine Check exception is not recoverable. The interrupted program cannot be safely resumed.

Figure 33-1: Machine Check Exception Enable/Disable Bit (CR4[MCE])



34 *The Local and IO APICs*

The Previous Chapter

The previous chapter provided a detailed description of the Machine Check Architecture:

- The MCA Elements.
- The Global Registers.
- The Composition of a Register Bank.
- The Error Code.
- Cache Error Reporting.
- MC Exception Is Generally Not Recoverable.

This Chapter

This chapter provides a complete description of the Local and IO APICs. It includes:

- A Short History of the APIC's Evolution.
- Before the APIC.
- MP Systems Need a Better Interrupt Distribution Mechanism.
- Detecting Presence/Version/Capabilities of Local APIC.
- Local APIC's Initial State.
- Enabling/Disabling the Local APIC.
- Mode Selection.
- The Local APIC Register Set.
- Local APIC ID Assignments and Addressing.
 - ID Assignment in xAPIC Mode.
 - ID Assignment in x2APIC Mode.
 - Local APIC Addressing.
 - Lowest-Priority Delivery Mode.
- Local APIC IDs Are Stored in the MP and ACPI Tables.

x86 Instruction Set Architecture

- Accessing the Local APIC ID.
- An Introduction to the Interrupt Sources.
- Introduction to Interrupt Priority.
- Task and Processor Priority.
- IO/Local APICs Cooperate on Interrupt Handling.
- Message Signaled Interrupts (MSI).
- Interrupt Delivery from Legacy 8259a Interrupt Controller.
- SW-Initiated Interrupt Message Transmission.
- x2APIC Mode's Self IPI Feature.
- Locally Generated Interrupts.
 - The Local Vector Table.
 - Local Interrupt 0 (LINT0).
 - Local Interrupt 1 (LINT1).
 - The Local APIC Timer.
 - The Performance Counter Overflow Interrupt.
 - The Thermal Sensor Interrupt.
 - Correctable Machine Check (CMC) Interrupt.
 - The Local APIC's Error Interrupt.
- The Spurious Interrupt Vector.
- Boot Strap Processor (BSP) Selection.
- How the APs are Discovered and Configured.

APIC and the IA-32 Architecture

The APIC's (Advanced Programmable Interrupt Controller's) role has been central to the x86 platform's inter-device communication scheme for many years. In addition to providing a communication channel between device adapters and their respective drivers, the APICs allow program threads running on different logical processors to communicate with each other by passing IPI (Inter-Processor Interrupt) messages to each other. Until the advent of the x2APIC architecture, however, the Local APIC's register set and operational characteristics were considered design-specific and outside the scope of the IA-32 architecture. With the introduction of the x2APIC architecture (in the Core i7 processor), the Local APIC's register set is accessible as a set of architecturally-defined MSRs. Even prior to the introduction of the x2APIC feature, however, the Local and IO APICs played a central role in system design. For that reason, this chapter provides a detailed description of the Local and IO APIC operation in both x2APIC Mode as well as the earlier APIC modes of operation.

Definition of IO and Local APICs

Basic definitions of the Local and IO APICs may be found in “APIC” on page 19.

Hardware Context Is Essential

In order to adequately explain the genesis and functionality of the Local and IO APIC modules, it is necessary to have some understanding of their role within the hardware platform. Towards this end, this chapter, where applicable, describes the hardware ecosystem within which the Local and IO APICs fulfill such an important role.

A Short History of the APIC's Evolution

The following is a *very* short history of the APIC.

APIC Introduction

The APIC was first introduced as a separate, stand-alone chip, the 82489DX. The Local APIC was first introduced in the P54C version of the Pentium.

Pentium Pro APIC Enhancements

The Pentium Pro implemented the following improvements to the Local APIC:

- **Performance Counter Interrupt.** The APIC can be enabled to generate an interrupt if a Performance Counter generates an overflow when incremented. The Performance Counter Overflow entry (see “The Performance Counter Overflow Interrupt” on page 1368) was added to the Local APIC's Local Vector Table (LVT register set) to support this feature.
- **APIC Base MSR added.** While the memory address range associated with the Local APIC's register set was hardwired on the Pentium (the base address of the 4KB range was hardwired at FEE00000h), the APIC_BASE MSR added (see Figure 34-11 on page 1258) in the Pentium Pro permits the programmer to specify (in APIC_BASE[35:12]) the register set's base address starting on any 4KB-aligned address range in the logical processor's physical memory address space.

- **Software Enable/Disable added.** The Pentium processor's Local APIC could not be enabled or disabled under software control. This could only be accomplished via hardware when the processor sampled an input on the trailing-edge of reset. Starting with the Pentium Pro, the APIC_BASE[EN] bit can be used by software for this purpose.
- **BSP bit added.** The APIC_BASE[BSP] bit is a read-only bit that *remembers* whether the logical processor was selected as the Boot Strap Processor or as an Application Processor (AP) at startup time. See “Boot Strap Processor (BSP) Selection” on page 1378 for more information on the BSP selection process.
- **APIC access propagation deleted.** When software executing on a Pentium performed a load or a store targeting the processor's Local APIC register set, the memory access was also propagated out onto the processor's external interface. This was eliminated with the advent of the P6 processor family.
- **Illegal Register Address error bit added** to the Local APIC's Error Status register (see Figure 34-26 on page 1302).
- **Remote Register Read capability was eliminated.**
- **SMI delivery added.** The ability to deliver an SMI to a logical processor via an Inter-Processor Interrupt (IPI) message was added.

The Pentium II and Pentium III

No enhancements were made to the APIC architecture in the Pentium II and Pentium III processors.

Pentium 4 APIC Enhancements: xAPIC

To differentiate the revised APIC architecture introduced with the advent of the Pentium 4 from the old APIC architecture, it is referred to as the xAPIC architecture. For the remainder of this chapter, the earlier APIC architecture is referred to simply as the *APIC architecture*, or the *legacy APIC architecture*. The xAPIC architecture introduced the following improvements:

- **Thermal Sensor interrupt added** (see “The Thermal Sensor Interrupt” on page 1370).
- **APIC ID Register enhanced.** In the P6 and Pentium, the APIC ID field was 4-bits, and encodings 0h - Eh could be used to uniquely identify 15 different processors connected to the APIC bus. In the Pentium 4, the xAPIC spec extended the local APIC ID field to 8 bits which can be used to identify up to 255 logical processors in the system.

Live Training Courses on the x86 Architecture



Comprehensive x86 Architecture (32/64-bit)

This lecture-based course describes the entire x86 architecture; everything from Real Mode addressing to 64-bit Mode interrupt handling. This course focuses on the architectural elements of x86, like segmentation, interrupt/exception handling, paging, etc. After completing this course, you will have a much deeper understanding of the x86 architecture.

Topics covered in this course:

- Instruction Set
- Register Sets
- Operating Modes
- Segmentation
- Task Management
- Interrupts and Exceptions
- Paging
- Memory Types
- Virtualization



System Programming for the x86 Architecture

This course teaches the x86 architecture (both 32-bit and 64-bit) through a mix of lectures and hands-on programming labs. All topics are explained in lecture format first and then the students are given programming labs in assembly to reinforce the concepts and to get hands-on experience working with x86 processors at a very low level. This course focuses mainly on the behavior of legacy Protected Mode, Compatibility Mode and 64-bit Mode as these are the modes most commonly used in modern operating systems.

The lab exercises range from printing to the screen using the flat memory model in legacy Protected Mode to setting up an interrupt driven, multitasking 64-bit Mode environment, with paging turned on.

You Will Learn:

- x86 programming basics like an overview of the instruction set, register set and operating modes
- The behavior of segmentation, how it was originally intended to be used and how it is actually used by operating systems today
- How to setup system calls using multiple methods (and benefits / side-effects of each)
- How to setup interrupt service routines for both software and hardware interrupts and implement a rudimentary scheduler
- How to implement paging in both the 32-bit environments as well as the 64-bit environments including using various page sizes

eLearning Course on the x86 Architecture



Intro to 32/64-bit x86 Architecture

MindShare's Intro to x86 eLearning course provides a great overview of the x86 instruction set architecture. It describes the concepts and implementation of the major architectural elements of x86 (paging, interrupt handling, protection, register set, address spaces, operating modes, virtualization, etc.).

You Will Learn:

- The different groupings of x86 instructions and basic instruction formats
- The different address spaces available in x86 and how each can be accessed
- The registers defined in x86 and how they have grown over time
- All the operating modes that x86-based processors can run in and characteristics of each
- The concepts of paging and its base implementation in x86
- How interrupts and exceptions are handled inside the core
- What virtualization is and about the hardware extensions that companies like Intel are adding their processors

What's Included:

- Unlimited access to the x86 eLearning modules for 90 days
- PDF of the course slides (yours to keep, does not expire)
- Q&A capability with the instructor

Benefits of MindShare eLearning:

- **Cost Effective** - Get the same information delivered in a live MindShare class at a fraction of the cost
- **Available 24/7** - MindShare eLearning courses are available when and where you need them
- **Learn at Your Pace** - MindShare eLearning courses are self-paced, so you can proceed when you're ready
- **Access to the Instructor** - Ask questions to the MindShare instructor that taught the course