

**Intel® Itanium® Architecture  
Software Developer's Manual  
Revision 2.3  
Volume 2: System Architecture**







# **Intel<sup>®</sup> Itanium<sup>®</sup> Architecture Software Developer's Manual**

**Volume 2: System Architecture**

---

**Revision 2.3**

***May 2010***

THIS DOCUMENT IS PROVIDED "AS IS" WITH NO WARRANTIES WHATSOEVER, INCLUDING ANY WARRANTY OF MERCHANTABILITY, FITNESS FOR ANY PARTICULAR PURPOSE, OR ANY WARRANTY OTHERWISE ARISING OUT OF ANY PROPOSAL, SPECIFICATION OR SAMPLE.

Information in this document is provided in connection with Intel® products. No license, express or implied, by estoppel or otherwise, to any intellectual property rights is granted by this document. Except as provided in Intel's Terms and Conditions of Sale for such products, Intel assumes no liability whatsoever, and Intel disclaims any express or implied warranty, relating to sale and/or use of Intel products including liability or warranties relating to fitness for a particular purpose, merchantability, or infringement of any patent, copyright or other intellectual property right. Intel products are not intended for use in medical, life saving, or life sustaining applications.

Intel may make changes to specifications and product descriptions at any time, without notice.

Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them.

Intel® processors based on the Itanium architecture may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Contact your local Intel sales office or your distributor to obtain the latest specifications and before placing your product order.

Copies of documents which have an order number and are referenced in this document, or other Intel literature, may be obtained by calling 1-800-548-4725, or by visiting Intel's website at <http://www.intel.com>.

Intel, Itanium, Pentium, VTune and MMX are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Copyright © 1999-2010, Intel Corporation

\*Other names and brands may be claimed as the property of others.

*Intel® Itanium® Architecture Software Developer's Manual, Rev. 2.3*

# Contents

---

## Part I: System Architecture Guide

<b>1</b>	<b>About this Manual</b>	<b>2:3</b>
1.1	Overview of <a href="#">Volume 1: Application Architecture</a>	2:3
1.1.1	Part 1: Application Architecture Guide	2:3
1.1.2	Part 2: Optimization Guide for the Intel® Itanium® Architecture	2:3
1.2	Overview of <a href="#">Volume 2: System Architecture</a>	2:4
1.2.1	Part 1: System Architecture Guide	2:4
1.2.2	Part 2: System Programmer's Guide	2:5
1.2.3	Appendices	2:6
1.3	Overview of <a href="#">Volume 3: Intel® Itanium® Instruction Set Reference</a>	2:6
1.4	Overview of <a href="#">Volume 4: IA-32 Instruction Set Reference</a>	2:6
1.5	Terminology	2:7
1.6	Related Documents	2:7
1.7	Revision History	2:8
<b>2</b>	<b>Intel® Itanium® System Environment</b>	<b>2:13</b>
2.1	Processor Boot Sequence	2:13
2.2	Intel® Itanium® System Environment Overview	2:14
<b>3</b>	<b>System State and Programming Model</b>	<b>2:17</b>
3.1	Privilege Levels	2:17
3.2	Serialization	2:17
3.2.1	Instruction Serialization	2:18
3.2.2	Data Serialization	2:18
3.2.3	Definition of In-flight Resources	2:19
3.3	System State	2:20
3.3.1	System State Overview	2:20
3.3.2	Processor Status Register (PSR)	2:23
3.3.3	Control Registers	2:29
3.3.4	Global Control Registers	2:31
3.3.5	Interrupt Control Registers	2:36
3.3.6	External Interrupt Control Registers	2:42
3.3.7	Banked General Registers	2:42
3.4	Processor Virtualization	2:44
<b>4</b>	<b>Addressing and Protection</b>	<b>2:45</b>
4.1	Virtual Addressing	2:45
4.1.1	Translation Lookaside Buffer (TLB)	2:47
4.1.2	Region Registers (RR)	2:58
4.1.3	Protection Keys	2:59
4.1.4	Translation Instructions	2:60
4.1.5	Virtual Hash Page Table (VHPT)	2:61
4.1.6	VHPT Hashing	2:65
4.1.7	VHPT Environment	2:67
4.1.8	Translation Searching	2:69
4.1.9	32-bit Virtual Addressing	2:71
4.1.10	Virtual Aliasing	2:72
4.2	Physical Addressing	2:73
4.3	Unimplemented Address Bits	2:73
4.3.1	Unimplemented Physical Address Bits	2:73
4.3.2	Unimplemented Virtual Address Bits	2:74
4.3.3	Instruction Behavior with Unimplemented Addresses	2:75



4.4	Memory Attributes	2:75
4.4.1	Virtual Addressing Memory Attributes	2:75
4.4.2	Physical Addressing Memory Attributes	2:76
4.4.3	Cacheability and Coherency Attribute	2:77
4.4.4	Cache Write Policy Attribute	2:78
4.4.5	Coalescing Attribute	2:78
4.4.6	Speculation Attributes	2:79
4.4.7	Sequentiality Attribute and Ordering	2:82
4.4.8	Not a Thing Attribute (NaTPage)	2:86
4.4.9	Effects of Memory Attributes on Memory Reference Instructions	2:86
4.4.10	Effects of Memory Attributes on Advanced/Check Loads	2:87
4.4.11	Memory Attribute Transition	2:88
4.5	Memory Datum Alignment and Atomicity	2:93
<b>5</b>	<b>Interruptions</b>	<b>2:95</b>
5.1	Interruption Definitions	2:95
5.2	Interruption Programming Model	2:97
5.3	Interruption Handling during Instruction Execution	2:98
5.4	PAL-based Interruption Handling	2:101
5.5	IVA-based Interruption Handling	2:101
5.5.1	Efficient Interruption Handling	2:102
5.5.2	Non-access Instructions and Interruptions	2:103
5.5.3	Single Stepping	2:104
5.5.4	Single Instruction Fault Suppression	2:104
5.5.5	Deferral of Speculative Load Faults	2:105
5.6	Interruption Priorities	2:108
5.6.1	IA-32 Interruption Priorities and Classes	2:111
5.7	IVA-based Interruption Vectors	2:113
5.8	Interrupts	2:114
5.8.1	Interrupt Vectors and Priorities	2:118
5.8.2	Interrupt Enabling and Masking	2:119
5.8.3	External Interrupt Control Registers	2:121
5.8.4	Processor Interrupt Block	2:127
5.8.5	Edge- and Level-sensitive Interrupts	2:131
<b>6</b>	<b>Register Stack Engine</b>	<b>2:133</b>
6.1	RSE and Backing Store Overview	2:133
6.2	RSE Internal State	2:135
6.3	Register Stack Partitions	2:136
6.4	RSE Operation	2:137
6.5	RSE Control	2:139
6.5.1	Register Stack Configuration Register	2:139
6.5.2	Register Stack NaT Collection Register	2:140
6.5.3	Backing Store Pointer Application Registers	2:141
6.5.4	RSE Control Instructions	2:142
6.5.5	Bad PFS used by Branch Return	2:143
6.6	RSE Interruptions	2:144
6.7	RSE Behavior on Interruptions	2:146
6.8	RSE Behavior with an Incomplete Register Frame	2:146
6.9	RSE and ALAT Interaction	2:146
6.10	Backing Store Coherence and Memory Ordering	2:147
6.11	RSE Backing Store Switches	2:147
6.11.1	Switch from Interrupted Context	2:148
6.11.2	Return to Interrupted Context	2:148
6.11.3	Synchronous Backing Store Switch	2:148
6.12	RSE Initialization	2:149
<b>7</b>	<b>Debugging and Performance Monitoring</b>	<b>2:151</b>
7.1	Debugging	2:151
7.1.1	Data and Instruction Breakpoint Registers	2:152

7.2	7.1.2	Debug Address Breakpoint Match Conditions . . . . .	2:154
		Performance Monitoring . . . . .	2:155
	7.2.1	Generic Performance Counter Registers . . . . .	2:156
	7.2.2	Performance Monitor Overflow Status Registers (PMC[0]..PMC[3]) . . . . .	2:160
	7.2.3	Performance Monitor Events . . . . .	2:162
	7.2.4	Implementation-independent Performance Monitor Code Sequences. . . . .	2:162
<b>8</b>		<b>Interruption Vector Descriptions . . . . .</b>	<b>2:165</b>
8.1		Interruption Vector Descriptions . . . . .	2:165
8.2		ISR Settings . . . . .	2:165
8.3		Interruption Vector Definition . . . . .	2:166
<b>9</b>		<b>IA-32 Interruption Vector Descriptions . . . . .</b>	<b>2:213</b>
9.1		IA-32 Trap Code . . . . .	2:213
9.2		IA-32 Interruption Vector Definitions. . . . .	2:213
<b>10</b>		<b>Itanium® Architecture-based Operating System Interaction Model with IA-32 Applications</b>	<b>2:239</b>
10.1		Instruction Set Transitions . . . . .	2:239
10.2		System Register Model. . . . .	2:239
10.3		IA-32 System Segment Registers. . . . .	2:241
	10.3.1	IA-32 Current Privilege Level . . . . .	2:243
	10.3.2	IA-32 System EFLAG Register. . . . .	2:243
	10.3.3	IA-32 System Registers . . . . .	2:246
10.4		Register Context Switch Guidelines for IA-32 Code . . . . .	2:252
	10.4.1	Entering IA-32 Processes. . . . .	2:253
	10.4.2	Exiting IA-32 Processes . . . . .	2:253
10.5		IA-32 Instruction Set Behavior Summary . . . . .	2:253
10.6		System Memory Model. . . . .	2:259
	10.6.1	Virtual Memory References . . . . .	2:260
	10.6.2	IA-32 Virtual Memory References . . . . .	2:261
	10.6.3	IA-32 TLB Forward Progress Requirements . . . . .	2:261
	10.6.4	Multiprocessor TLB Coherency . . . . .	2:262
	10.6.5	IA-32 Physical Memory References . . . . .	2:262
	10.6.6	Supervisor Accesses . . . . .	2:263
	10.6.7	Memory Alignment . . . . .	2:263
	10.6.8	Atomic Operations . . . . .	2:264
	10.6.9	Multiprocessor Instruction Cache Coherency. . . . .	2:264
	10.6.10	IA-32 Memory Ordering . . . . .	2:265
10.7		I/O Port Space Model. . . . .	2:267
	10.7.1	Virtual I/O Port Addressing . . . . .	2:268
	10.7.2	Physical I/O Port Addressing . . . . .	2:270
	10.7.3	IA-32 IN/OUT instructions. . . . .	2:271
	10.7.4	I/O Port Accesses by Loads and Stores. . . . .	2:272
10.8		Debug Model . . . . .	2:273
	10.8.1	Data Breakpoint Register Matching . . . . .	2:274
	10.8.2	Instruction Breakpoint Register Matching. . . . .	2:274
10.9		Interruption Model. . . . .	2:275
	10.9.1	Interruption Summary. . . . .	2:275
	10.9.2	IA-32 Numeric Exception Model. . . . .	2:277
10.10		Processor Bus Considerations for IA-32 Application Support . . . . .	2:277
	10.10.1	IA-32 Compatible Bus Transactions . . . . .	2:278
<b>11</b>		<b>Processor Abstraction Layer . . . . .</b>	<b>2:279</b>
11.1		Firmware Model . . . . .	2:279
	11.1.1	Processor Abstraction Layer (PAL) Overview . . . . .	2:280
	11.1.2	Firmware Entrypoints . . . . .	2:281
	11.1.3	PAL Entrypoints . . . . .	2:282
	11.1.4	SAL Entrypoints . . . . .	2:282



11.1.5	OS Entrypoints	2:283
11.1.6	Firmware Address Space	2:283
11.2	PAL Power On/Reset	2:289
11.2.1	PALE_RESET	2:289
11.2.2	PALE_RESET Exit State	2:289
11.2.3	PAL Self-test Control Word	2:295
11.3	Machine Checks	2:296
11.3.1	PALE_CHECK	2:296
11.3.2	PALE_CHECK Exit State	2:297
11.3.3	Returning to the Interrupted Process	2:305
11.4	PAL Initialization Events	2:306
11.4.1	PALE_INIT	2:306
11.4.2	PALE_INIT Exit State	2:306
11.5	Platform Management Interrupt (PMI)	2:310
11.5.1	PMI Overview	2:310
11.5.2	PALE_PMI Exit State	2:312
11.5.3	Resume from the PMI Handler	2:313
11.6	Power Management	2:313
11.6.1	Power/Performance States (P-states)	2:315
11.7	PAL Virtualization Support	2:324
11.7.1	Virtual Processor Descriptor (VPD)	2:325
11.7.2	Interrupt Handling in a Virtual Environment	2:331
11.7.3	PAL Intercepts in Virtual Environment	2:332
11.7.4	Virtualization Optimizations	2:335
11.8	PAL Glossary	2:350
11.9	PAL Code Memory Accesses and Restrictions	2:352
11.10	PAL Procedures	2:353
11.10.1	PAL Procedure Summary	2:354
11.10.2	PAL Calling Conventions	2:358
11.10.3	PAL Procedure Specifications	2:365
11.11	PAL Virtualization Services	2:486
11.11.1	PAL Virtualization Service Invocation Convention	2:486
11.11.2	PAL Virtualization Service Specifications	2:488

## Part II: System Programmer's Guide

<b>1</b>	<b>About the System Programmer's Guide</b>	<b>2:503</b>
1.1	Overview of the System Programmer's Guide	2:503
1.2	Related Documents	2:505
<b>2</b>	<b>MP Coherence and Synchronization</b>	<b>2:507</b>
2.1	An Overview of Intel® Itanium® Memory Access Instructions	2:507
2.1.1	Memory Ordering of Cacheable Memory References	2:507
2.1.2	Loads and Stores	2:508
2.1.3	Semaphores	2:508
2.1.4	Memory Fences	2:510
2.2	Memory Ordering in the Intel® Itanium® Architecture	2:510
2.2.1	Memory Ordering Executions	2:511
2.2.2	Memory Attributes	2:524
2.2.3	Understanding Other Ordering Models: Sequential Consistency and IA-32	2:525
2.3	Where the Intel® Itanium® Architecture Requires Explicit Synchronization	2:526
2.4	Synchronization Code Examples	2:526
2.4.1	Spin Lock	2:527
2.4.2	Simple Barrier Synchronization	2:528
2.4.3	Dekker's Algorithm	2:529
2.4.4	Lamport's Algorithm	2:530
2.5	Updating Code Images	2:531
2.5.1	Self-modifying Code	2:532
2.5.2	Cross-modifying Code	2:533

	2.5.3	Programmed I/O .....	2:534
	2.5.4	DMA .....	2:536
2.6		References .....	2:536
<b>3</b>		<b>Interruptions and Serialization .....</b>	<b>2:537</b>
3.1		Terminology .....	2:537
3.2		Interruption Vector Table .....	2:538
3.3		Interruption Handlers .....	2:539
	3.3.1	Execution Environment .....	2:539
	3.3.2	Interruption Register State .....	2:540
	3.3.3	Resource Serialization of Interrupted State .....	2:542
	3.3.4	Resource Serialization upon rfi. ....	2:543
3.4		Interruption Handling .....	2:543
	3.4.1	Lightweight Interruptions .....	2:543
	3.4.2	Heavyweight Interruptions .....	2:544
	3.4.3	Nested Interruptions .....	2:546
<b>4</b>		<b>Context Management .....</b>	<b>2:549</b>
4.1		Preserving Register State across Procedure Calls .....	2:549
	4.1.1	Preserving General Registers .....	2:550
	4.1.2	Preserving Floating-point Registers .....	2:551
4.2		Preserving Register State in the OS .....	2:551
	4.2.1	Preservation of Stacked Registers in the OS .....	2:552
	4.2.2	Preservation of Floating-point State in the OS .....	2:553
4.3		Preserving ALAT Coherency .....	2:554
4.4		System Calls .....	2:555
	4.4.1	epc/Demoting Branch Return .....	2:555
	4.4.2	break/rfi .....	2:556
	4.4.3	NaT Checking for NaTs in System Calls .....	2:556
4.5		Context Switching .....	2:557
	4.5.1	User-level Context Switching .....	2:557
	4.5.2	Context Switching in an Operating System Kernel .....	2:558
<b>5</b>		<b>Memory Management .....</b>	<b>2:561</b>
5.1		Address Space Model .....	2:561
	5.1.1	Regions .....	2:561
	5.1.2	Protection Keys .....	2:564
5.2		Translation Lookaside Buffers (TLBs) .....	2:565
	5.2.1	Translation Registers (TRs) .....	2:566
	5.2.2	Translation Caches (TCs) .....	2:567
5.3		Virtual Hash Page Table .....	2:571
	5.3.1	Short Format .....	2:572
	5.3.2	Long Format .....	2:573
	5.3.3	VHPT Updates .....	2:573
5.4		TLB Miss Handlers .....	2:573
	5.4.1	Data/Instruction TLB Miss Vectors .....	2:573
	5.4.2	VHPT Translation Vector .....	2:575
	5.4.3	Alternate Data/Instruction TLB Miss Vectors .....	2:576
	5.4.4	Data Nested TLB Vector .....	2:576
	5.4.5	Dirty Bit Vector .....	2:577
	5.4.6	Data/Instruction Access Bit Vector .....	2:577
	5.4.7	Page Not Present Vector .....	2:577
	5.4.8	Data/Instruction Access Rights Vector .....	2:577
5.5		Subpaging .....	2:577
<b>6</b>		<b>Runtime Support for Control and Data Speculation .....</b>	<b>2:579</b>
6.1		Exception Deferral of Control Speculative Loads .....	2:579
	6.1.1	Hardware-only Deferral .....	2:580
	6.1.2	Combined Hardware/Software Deferral .....	2:580
	6.1.3	Software-only Deferral .....	2:580



6.2	Speculation Recovery Code Requirements	2:580
6.3	Speculation Related Exception Handlers	2:581
6.3.1	Unaligned Handler	2:581
<b>7</b>	<b>Instruction Emulation and Other Fault Handlers</b>	<b>2:583</b>
7.1	Unaligned Reference Handler	2:583
7.2	Unsupported Data Reference Handler	2:584
7.3	Illegal Dependency Fault	2:584
7.4	Long Branch	2:585
<b>8</b>	<b>Floating-point System Software</b>	<b>2:587</b>
8.1	Floating-point Exceptions in the Intel® Itanium® Architecture	2:587
8.1.1	Software Assistance Exceptions (Faults and Traps)	2:587
8.1.2	The IEEE Floating-point Exception Filter	2:590
8.2	IA-32 Floating-point Exceptions	2:593
<b>9</b>	<b>IA-32 Application Support</b>	<b>2:595</b>
9.1	Transitioning between Intel® Itanium® and IA-32 Instruction Sets	2:596
9.1.1	IA-32 Code Execution Environments	2:596
9.1.2	br.ia	2:596
9.1.3	JMPE	2:597
9.1.4	Procedure Calls between Intel® Itanium® and IA-32 Instruction Sets	2:597
9.2	IA-32 Architecture Handlers	2:599
9.3	Debugging IA-32 and Itanium® Architecture-based Code	2:600
9.3.1	Instruction Breakpoints	2:600
9.3.2	Data Breakpoints	2:600
9.3.3	Single Step Traps	2:601
9.3.4	Taken Branch Traps	2:601
<b>10</b>	<b>External Interrupt Architecture</b>	<b>2:603</b>
10.1	External Interrupt Basics	2:603
10.2	Configuration of External Interrupt Vectors	2:604
10.3	External Interrupt Masking	2:604
10.3.1	PSR.i	2:604
10.3.2	IVR Reads and EOI Writes	2:605
10.3.3	Task Priority Register (TPR)	2:605
10.3.4	External Task Priority Register (XTPR)	2:605
10.4	External Interrupt Delivery	2:606
10.5	Interrupt Control Register Usage Examples	2:607
10.5.1	Notation	2:608
10.5.2	TPR and XPTR Usage Example	2:608
10.5.3	EOI Usage Example	2:609
10.5.4	IRR Usage Example	2:609
10.5.5	Interval Timer Usage Example	2:609
10.5.6	Resource Utilization Counter Usage Example	2:611
10.5.7	Local Redirection Example	2:611
10.5.8	Inter-processor Interrupts Layout and Example	2:612
10.5.9	INTA Example	2:612
<b>11</b>	<b>I/O Architecture</b>	<b>2:615</b>
11.1	Memory Acceptance Fence (mf.a)	2:615
11.2	I/O Port Space	2:616
<b>12</b>	<b>Performance Monitoring Support</b>	<b>2:619</b>
12.1	Architected Performance Monitoring Mechanisms	2:619
12.2	Operating System Support	2:620
<b>13</b>	<b>Firmware Overview</b>	<b>2:623</b>

13.1	Processor Boot Flow Overview	2:623
13.1.1	Firmware Boot Flow	2:623
13.1.2	Operating System Boot Steps	2:625
13.2	Runtime Procedure Calls	2:628
13.2.1	PAL Procedure Calls	2:628
13.2.2	SAL Procedure Calls	2:630
13.2.3	UEFI Procedure Calls	2:630
13.2.4	ACPI Control Methods	2:631
13.2.5	Physical and Virtual Addressing Mode Considerations	2:631
13.3	Event Handling in Firmware	2:632
13.3.1	Machine Check Abort (MCA) Flows	2:632
13.3.2	INIT Flows	2:635
13.3.3	PMI Flows	2:637
13.3.4	P-state Feedback Mechanism Flow Diagram	2:637
<b>A</b>	<b>Code Examples</b>	<b>2:639</b>
A.1	OS Boot Flow Sample Code	2:639
<b>Index</b>		<b>2:643</b>

## Figures

### Part I: System Architecture Guide

2-1	System Environment Boot Flow	2:13
2-2	Intel® Itanium® System Environment	2:14
3-1	System Register Model	2:22
3-2	Processor Status Register (PSR)	2:23
3-3	Default Control Register (DCR – CR0)	2:31
3-4	Interval Time Counter (ITC – AR44)	2:32
3-5	Interval Timer Match Register (ITM – CR1)	2:32
3-6	Interval Timer Offset Register (ITO – CR4)	2:34
3-7	Interrupt Vector Address (IVA – CR2)	2:35
3-8	Page Table Address (PTA – CR8)	2:35
3-9	Interrupt Status Register (ISR – CR17)	2:36
3-10	Interrupt Instruction Bundle Pointer (IIP – CR19)	2:38
3-11	Interrupt Faulting Address (IFA – CR20)	2:39
3-12	Interrupt TLB Insertion Register (ITIR)	2:39
3-13	Interrupt Instruction Previous Address (IIPA – CR22)	2:40
3-14	Interrupt Function State (IFS – CR23)	2:41
3-15	Interrupt Immediate (IIM – CR24)	2:41
3-16	Interrupt Hash Address (IHA – CR25)	2:41
3-17	Interrupt Instruction Bundle Registers (IIB0-1, – CR26, 27)	2:42
3-18	Banked General Registers	2:43
4-1	Virtual Address Spaces	2:46
4-2	Conceptual Virtual Address Translation for References	2:47
4-3	TLB Organization	2:47
4-4	Conceptual Virtual Address Searching for Inserts and Purges	2:51
4-5	Translation Insertion Format	2:54
4-6	Translation Insertion Format – Not Present	2:56
4-7	Region Register Format	2:58
4-8	Protection Key Register Format	2:59
4-9	Virtual Hash Page Table (VHPT)	2:62
4-10	VHPT Short Format	2:63



4-11	VHPT Not-present Short Format . . . . .	2:64
4-12	VHPT Long Format. . . . .	2:64
4-13	VHPT Not-present Long Format. . . . .	2:65
4-14	Region-based VHPT Short-format Index Function. . . . .	2:66
4-15	VHPT Long-format Hash Function . . . . .	2:66
4-16	TLB/VHPT Search . . . . .	2:70
4-17	32-bit Address Generation using addp4. . . . .	2:72
4-18	Physical Address Bit Fields . . . . .	2:73
4-19	Virtual Address Bit Fields . . . . .	2:74
4-20	Physical Addressing Memory . . . . .	2:76
4-21	Addressing Memory Attributes . . . . .	2:77
5-1	Interrupt Classification . . . . .	2:97
5-2	Interrupt Processing. . . . .	2:99
5-3	Interrupt Architecture Overview . . . . .	2:115
5-4	PAL-based Interrupt States . . . . .	2:117
5-5	External Interrupt States. . . . .	2:118
5-6	Local ID (LID – CR64) . . . . .	2:122
5-7	External Interrupt Vector Register (IVR – CR65) . . . . .	2:123
5-8	Task Priority Register (TPR – CR66) . . . . .	2:124
5-9	End of External Interrupt Register (EOI – CR67) . . . . .	2:124
5-10	External Interrupt Request Register (IRR0-3 – CR68, 69, 70, 71) . . . . .	2:125
5-11	Interval Timer Vector (ITV – CR72) . . . . .	2:125
5-12	Performance Monitor Vector (PMV – CR73) . . . . .	2:126
5-13	Corrected Machine Check Vector (CMCV – CR74) . . . . .	2:126
5-14	Local Redirection Register (LRR – CR80,81) . . . . .	2:127
5-15	Processor Interrupt Block Memory Layout . . . . .	2:128
5-16	Address Format for Inter-processor Interrupt Messages . . . . .	2:129
5-17	Data Format for Inter-processor Interrupt Messages . . . . .	2:129
6-1	Relationship Between Physical Registers and Backing Store . . . . .	2:134
6-2	Backing Store Memory Format. . . . .	2:134
6-3	Four Partitions of the Register Stack . . . . .	2:137
7-1	Data Breakpoint Registers (DBR). . . . .	2:152
7-2	Instruction Breakpoint Registers (IBR) . . . . .	2:152
7-3	Performance Monitor Register Set . . . . .	2:156
7-4	Generic Performance Counter Data Registers (PMD[4]..PMD[p]) . . . . .	2:157
7-5	Generic Performance Counter Configuration Register (PMC[4]..PMC[p]) . . . . .	2:157
7-6	Performance Monitor Overflow Status Registers (PMC[0]..PMC[3]) . . . . .	2:161
7-7	Performance Monitor Interrupt Service Routine (Implementation Independent) . . . . .	2:163
7-8	Performance Monitor Overflow Context Switch Routine . . . . .	2:164
9-1	IA-32 Trap Code . . . . .	2:213
9-2	IA-32 Trap Code . . . . .	2:213
9-3	IA-32 Intercept Code . . . . .	2:234
10-1	IA-32 System Segment Register Descriptor Format (LDT, GDT, TSS) . . . . .	2:241
10-2	IA-32 EFLAG Register . . . . .	2:243
10-3	Control Flag Register (CFLG, AR27) . . . . .	2:246
10-4	Virtual Memory Addressing . . . . .	2:260
10-5	Physical Memory Addressing . . . . .	2:262
10-1	I/O Port Space Model . . . . .	2:268
10-2	I/O Port Space Addressing . . . . .	2:269
11-1	Firmware Model . . . . .	2:280
11-2	Firmware Entrypoints Logical Model . . . . .	2:281

11-3	Firmware Address Space	2:284
11-4	Firmware Address Space with Processor-specific PAL_A Components	2:285
11-5	Firmware Interface Table	2:287
11-6	Firmware Interface Table Entry	2:288
11-7	SALE_ENTRY State Parameter	2:291
11-8	Geographically Significant Processor Identifier	2:293
11-9	Self Test State Parameter	2:293
11-10	Self-test Control Word	2:295
11-11	Processor State Parameter	2:299
11-1	Processor Min-state Save Area Layout	2:303
11-2	Processor State Saved in Min-state Save Area	2:304
11-3	NaT Bits for Saved GRs	2:305
11-4	SALE_ENTRY State Parameter	2:305
11-5	Processor State Parameter	2:308
11-6	SALE_ENTRY State Parameter	2:310
11-7	PMI Entrypoints	2:311
11-8	Power States	2:314
11-9	Power and Performance Characteristics for P-states	2:316
11-10	Example of a P-state Transition Policy	2:317
11-11	Computation of <i>performance_index</i>	2:321
11-12	Interaction of P-states with HALT State	2:324
11-13	Virtualization Acceleration Control ( <i>vac</i> )	2:329
11-14	Virtualization Disable Control ( <i>vdc</i> )	2:330
11-15	PAL Virtualization Intercept Handoff Opcode (GR25)	2:335
11-1	operation Parameter Layout	2:371
11-2	config_info_1 Return Value	2:374
11-3	config_info_2 Return Value	2:375
11-4	config_info_1 Return Value	2:378
11-5	config_info_2 Return Value	2:378
11-6	config_info_3 Return Value	2:379
11-7	<i>cache_protection</i> Fields	2:379
11-8	Layout of line_id Return Value	2:380
11-9	Layout of proc_n_cache_info1 Return Value	2:383
11-10	Layout of proc_n_cache_info2 Return Value	2:383
11-11	Layout of line_id Return Value	2:385
11-12	Return values	2:393
11-13	I/O Size and Type Information Layout	2:399
11-14	Layout of power_buffer Return Value	2:401
11-15	Layout of log_overview Return Value	2:405
11-16	Layout of proc_n_log_info1 Return Value	2:405
11-17	Layout of proc_n_log_info2 Return Value	2:406
11-18	Pending Return Parameter	2:407
11-19	<i>level_index</i> Layout	2:411
11-20	cache_check Layout	2:414
11-21	tlb_check Layout	2:415
11-22	bus_check Layout	2:417
11-23	reg_file_check Layout	2:418
11-24	uarch_check Layout	2:420
11-25	<i>err_type_info</i>	2:421
11-26	resources Return Value	2:423
11-27	err_struct_info – Cache	2:424



11-28	capabilities vector for cache	2:425
11-29	Buffer pointed to by <i>err_data_buffer</i> – Cache	2:426
11-30	<i>err_struct_info</i> – TLB	2:427
11-31	capabilities vector for TLB	2:428
11-32	Buffer pointed to by <i>err_data_buffer</i> – TLB	2:428
11-33	<i>err_struct_info</i> – Register File	2:428
11-34	capabilities Vector for Register File	2:430
11-35	Buffer pointed to by <i>err_data_buffer</i> – Register File	2:430
11-36	<i>err_struct_info</i> – Bus/Processor Interconnect	2:431
11-37	capabilities vector for Bus/Processor Interconnect	2:431
11-38	Layout of <i>hw_track</i> Return Value	2:432
11-39	Layout of <i>attrib</i> Return Value	2:437
11-40	Layout of <i>pm_info</i> Return Value	2:440
11-41	Layout of <i>pstate_buffer</i> Entry	2:451
11-42	Layout of <i>dd_info</i> Parameter	2:452
11-43	Layout of <i>hints</i> Return Value	2:455
11-44	Layout of <i>test_info</i> Argument	2:462
11-45	Layout of <i>test_param</i> Argument	2:463
11-46	Layout of <i>min_pal_ver</i> and <i>current_pal_ver</i> Return Values	2:465
11-47	Layout of <i>tc_info</i> Return Value	2:466
11-48	Layout of <i>vm_info_1</i> Return Value	2:468
11-49	Layout of <i>vm_info_2</i> Return Value	2:469
11-50	Layout of <i>TR_valid</i> Return Value	2:470

## Part II: System Programmer's Guide

2-1	Intel® Itanium® Ordering Semantics	2:512
2-2	Interaction of Ordering and Accesses to Sequential Locations	2:524
2-3	Why a Fence During Context Switches is Required in the Intel® Itanium® Architecture	2:526
2-4	Spin Lock Code	2:527
2-5	Sense-reversing Barrier Synchronization Code	2:528
2-6	Dekker's Algorithm in a 2-way System	2:530
2-7	Lamport's Algorithm	2:531
2-8	Updating a Code Image on the Local Processor	2:532
2-9	Supporting Cross-modifying Code without Explicit Serialization	2:533
2-10	Updating a Code Image on a Remote Processor	2:535
5-1	Self-mapped Page Table	2:572
5-2	Subpaging	2:578
8-1	Overview of Floating-point Exception Handling in the Intel® Itanium® Architecture	2:589
13-1	Firmware Model	2:624
13-2	Control Flow of Boot Process in a Multiprocessor Configuration	2:626
13-3	Correctable Machine Check Code Flow	2:633
13-4	Uncorrectable Machine Check Code Flow	2:633
13-5	INIT Flow	2:636
13-6	Flowchart Showing P-state Feedback Policy	2:638

## Tables

### Part I: System Architecture Guide

3-1	Processor Status Register Instructions	2:23
-----	--	------

3-2	Processor Status Register Fields . . . . .	2:24
3-3	Control Registers . . . . .	2:29
3-4	Control Register Instructions . . . . .	2:30
3-5	Default Control Register Fields . . . . .	2:31
3-6	Page Table Address Fields . . . . .	2:35
3-7	Interrupt Status Register Fields . . . . .	2:37
3-8	ITIR Fields . . . . .	2:39
3-9	Interrupt Function State Fields . . . . .	2:41
3-10	Virtualized Instructions. . . . .	2:44
4-1	Purge Behavior of TLB Inserts and Purges . . . . .	2:52
4-2	Purge behavior of VHPT Inserts. . . . .	2:53
4-3	Translation Interface Fields . . . . .	2:54
4-4	Page Access Rights . . . . .	2:56
4-5	Architected Page Sizes . . . . .	2:58
4-6	Region Register Fields . . . . .	2:58
4-7	Protection Register Fields . . . . .	2:59
4-8	Translation Instructions . . . . .	2:60
4-9	VHPT Long-format Fields . . . . .	2:64
4-10	TLB and VHPT Search Faults . . . . .	2:70
4-11	Virtual Addressing Memory Attribute Encodings . . . . .	2:76
4-12	Physical Addressing Memory Attribute Encodings . . . . .	2:77
4-13	Permitted Speculation . . . . .	2:80
4-14	Register Return Values on Non-faulting Advanced/Speculative Loads . . . . .	2:80
4-15	Ordering Semantics and Instructions . . . . .	2:83
4-16	Ordering Semantics . . . . .	2:84
4-17	ALAT Behavior on Non-faulting Advanced/Check Loads . . . . .	2:88
5-1	ISR Settings for Non-access Instructions . . . . .	2:104
5-2	Programming Models . . . . .	2:105
5-3	Exception Qualification . . . . .	2:106
5-4	Qualified Exception Deferral . . . . .	2:107
5-5	Spontaneous Deferral . . . . .	2:107
5-6	Interrupt Priorities. . . . .	2:109
5-7	Interrupt Vector Table (IVT) . . . . .	2:113
5-8	Interrupt Priorities, Enabling, and Masking . . . . .	2:119
5-9	External Interrupt Control Registers . . . . .	2:122
5-10	Local ID Fields. . . . .	2:122
5-11	Task Priority Register Fields. . . . .	2:124
5-12	Interval Timer Vector Fields . . . . .	2:125
5-13	Performance Monitor Vector Fields . . . . .	2:126
5-14	Corrected Machine Check Vector Fields . . . . .	2:126
5-15	Local Redirection Register Fields . . . . .	2:127
5-16	Address Fields for Inter-processor Interrupt Messages . . . . .	2:129
5-17	Data Fields for Inter-processor Interrupt Messages . . . . .	2:129
6-1	RSE Internal State . . . . .	2:135
6-2	RSE Operation Instructions and State Modification . . . . .	2:138
6-3	RSE Modes (RSC.mode) . . . . .	2:139
6-4	Backing Store Pointer Application Registers . . . . .	2:142
6-5	RSE Control Instructions . . . . .	2:143
6-6	RSE Interruption Summary . . . . .	2:145
7-1	Debug Breakpoint Register Fields (DBR/IBR). . . . .	2:153
7-2	Debug Instructions. . . . .	2:153

7-3	Generic Performance Counter Data Register Fields . . . . .	2:157
7-4	Generic Performance Counter Configuration Register Fields (PMC[4]..PMC[p]) . . . . .	2:157
7-5	Reading Performance Monitor Data Registers . . . . .	2:158
7-6	Performance Monitor Instructions . . . . .	2:159
7-7	Performance Monitor Overflow Register Fields (PMC[0]...PMC[3]) . . . . .	2:161
8-1	Writing of Interruption Resources by Vector . . . . .	2:166
8-2	ISR Values on Interruption . . . . .	2:168
8-3	ISR.code Fields on Intel® Itanium® Traps . . . . .	2:170
8-4	Interruption Vectors Sorted Alphabetically . . . . .	2:171
9-1	Intercept Code Definition . . . . .	2:234
9-2	Segment Prefix Override Encodings . . . . .	2:234
9-3	Gate Intercept Trap Code Identifier . . . . .	2:235
9-4	System Flag Intercept Instruction Trap Code Instruction Identifier . . . . .	2:236
10-1	IA-32 System Register Mapping . . . . .	2:240
10-2	IA-32 System Segment Register Fields (LDT, GDT, TSS) . . . . .	2:242
10-3	IA-32 EFLAG Field Definition . . . . .	2:244
10-4	IA-32 Control Register Field Definition . . . . .	2:247
10-5	IA-32 Instruction Summary . . . . .	2:254
10-6	Instruction Cache Coherency Rules . . . . .	2:265
10-7	IA-32 Load/Store Sequentiality and Ordering . . . . .	2:265
10-8	IA-32 Interruption Vector Summary . . . . .	2:275
10-9	IA-32 Interruption Summary . . . . .	2:275
11-1	FIT Entry Types . . . . .	2:288
11-2	GR38 Reset Layout . . . . .	2:290
11-3	function Field Values . . . . .	2:291
11-4	status Field Values . . . . .	2:292
11-5	Geographically Significant Processor Identifier Fields . . . . .	2:293
11-6	state Field Values . . . . .	2:294
11-7	Processor State Parameter Fields . . . . .	2:299
11-8	Software Recovery Bits in Processor State Parameter . . . . .	2:301
11-9	PSP Bit Settings for Unconsumed Data-poisoning Events on MCA . . . . .	2:302
11-10	NaT Bits for Saved GRs . . . . .	2:305
11-11	function Field Values . . . . .	2:305
11-12	Processor State Parameter Fields . . . . .	2:308
11-13	function Field Values . . . . .	2:310
11-14	PMI Events and Priorities . . . . .	2:311
11-15	PMI Message Vector Assignments . . . . .	2:311
11-16	Virtual Processor Descriptor (VPD) . . . . .	2:326
11-17	Virtual Processor Descriptor (VPD) – VPSR . . . . .	2:328
11-18	Virtual Processor Descriptor (VPD) – VCR[0-127] . . . . .	2:329
11-19	Virtualization Acceleration Control (vac) Fields . . . . .	2:329
11-20	Virtualization Disable Control (vdc) Fields . . . . .	2:330
11-21	IVA Settings after PAL Virtualization-related Procedures and Services . . . . .	2:332
11-22	PAL Virtualization Intercept Handoff Cause (GR24) . . . . .	2:334
11-23	Global Virtualization Optimizations Summary . . . . .	2:336
11-24	Synchronization Requirements for Virtualization Opcode Optimization . . . . .	2:336
11-25	Behavior of Guest MOV-from-AR.ITC Instruction in Virtual Environment . . . . .	2:337
11-26	Virtualization Accelerations Summary . . . . .	2:338
11-27	Detection of Virtual External Interrupts . . . . .	2:339
11-28	Synchronization Requirements for Virtual External Interrupt Optimization . . . . .	2:339
11-29	Interruptions when Virtual External Interrupt Optimization is Enabled . . . . .	2:340

11-30	Synchronization Requirements for Interruption Control Register Read Optimization . . . . .	2:340
11-31	Interruptions when Interruption Control Register Read Optimization is Enabled . . . . .	2:341
11-32	Synchronization Requirements for Interruption Control Register Write Optimization . . . . .	2:341
11-33	Interruptions when Interruption Control Register Write Optimization is Enabled . . . . .	2:341
11-34	Synchronization Requirements for MOV-from-PSR Optimization . . . . .	2:342
11-35	Interruptions when MOV-from-PSR Optimization is Enabled . . . . .	2:342
11-36	Synchronization Requirements for MOV-from-CPUID Optimization. . . . .	2:343
11-37	Interruptions when MOV-from-CPUID Optimization is Enabled . . . . .	2:343
11-38	Synchronization Requirements for Cover Optimization . . . . .	2:343
11-39	Interruptions when Cover Optimization is Enabled . . . . .	2:343
11-40	Synchronization Requirements for Bank Switch Optimization. . . . .	2:344
11-41	Interruptions when Bank Switch Optimization is Enabled . . . . .	2:344
11-42	Impact of clearing VCPUID bits with the <i>a_tf</i> optimization. . . . .	2:345
11-43	Synchronization Requirements for Test Feature Optimization . . . . .	2:345
11-44	Synchronization Requirements for Interrupt Collection and User Mask Optimization . . . . .	2:346
11-45	Interruptions when Interrupt Collection and User Mask Optimization is Enabled . . . . .	2:346
11-46	Virtualization Disables Summary . . . . .	2:346
11-47	Supported Virtualization Optimization Combinations . . . . .	2:349
11-48	PAL Procedure Index Assignment. . . . .	2:354
11-49	PAL Cache and Memory Procedures . . . . .	2:354
11-50	PAL Processor Identification, Features, and Configuration Procedures. . . . .	2:355
11-51	PAL Machine Check Handling Procedures . . . . .	2:356
11-52	PAL Power Information and Management Procedures . . . . .	2:356
11-53	PAL Processor Self Test Procedures . . . . .	2:357
11-54	PAL Support Procedures . . . . .	2:357
11-55	PAL Virtualization Support Procedures . . . . .	2:357
11-56	State Requirements for PSR . . . . .	2:359
11-57	Definition of Terms. . . . .	2:360
11-58	System Register Conventions . . . . .	2:361
11-59	General Registers – Static Calling Convention . . . . .	2:362
11-60	General Registers – Stacked Calling Conventions . . . . .	2:362
11-61	Application Register Conventions . . . . .	2:363
11-62	Processor Brand Information Requested . . . . .	2:366
11-63	Processor Bus Features . . . . .	2:368
11-64	<i>cache_type</i> Encoding . . . . .	2:370
11-65	Cache Line State when <i>inv</i> = 0 . . . . .	2:371
11-66	Cache Line State when <i>inv</i> = 1 . . . . .	2:372
11-67	Cache Memory Attributes . . . . .	2:374
11-68	Cache Store Hints . . . . .	2:375
11-69	Cache Load Hints . . . . .	2:375
11-70	PAL_CACHE_INIT level Argument Values . . . . .	2:376
11-71	PAL_CACHE_INIT restrict Argument Values . . . . .	2:376
11-72	<i>method</i> Values. . . . .	2:379
11-73	<i>t_d</i> Values . . . . .	2:379
11-75	<i>part</i> Input Values and corresponding <i>data</i> Return Values. . . . .	2:381
11-76	<i>mesi</i> Return Values . . . . .	2:381
11-74	<i>part</i> Input Values. . . . .	2:381
11-77	<i>part</i> Input Values. . . . .	2:386
11-78	<i>mesi</i> Return Values . . . . .	2:386
11-79	Interpretation of <i>data</i> Input Field. . . . .	2:386
11-80	Hardware policies returned in <i>cur_policy</i> . . . . .	2:395



11-81	PAL_GET_PSTATE <i>type</i> Argument . . . . .	2:397
11-82	I/O Detail Pointer Description . . . . .	2:399
11-83	I/O Type Definition . . . . .	2:400
11-84	I/O Size Definition . . . . .	2:400
11-85	Pending Return Parameter Fields . . . . .	2:407
11-86	<i>info_index</i> Values . . . . .	2:411
11-87	<i>level_index</i> Fields . . . . .	2:412
11-88	<i>err_type_index</i> Values . . . . .	2:412
11-89	<i>error_info</i> Return Format when <i>info_index</i> = 2 and <i>err_type_index</i> = 0 . . . . .	2:413
11-90	<i>cache_check</i> Fields . . . . .	2:414
11-91	<i>tlb_check</i> Fields . . . . .	2:415
11-92	<i>bus_check</i> Fields . . . . .	2:417
11-93	<i>reg_file_check</i> Fields . . . . .	2:418
11-94	<i>uarch_check</i> Fields . . . . .	2:420
11-95	<i>err_type_info</i> . . . . .	2:422
11-96	<i>resources</i> Return Value . . . . .	2:424
11-97	<i>err_struct_info</i> – Cache . . . . .	2:424
11-98	<i>capabilities</i> vector for cache . . . . .	2:425
11-99	Buffer pointed to by <i>err_data_buffer</i> – Cache . . . . .	2:426
11-100	<i>err_struct_info</i> – TLB . . . . .	2:427
11-101	<i>capabilities</i> vector for TLB . . . . .	2:428
11-102	Buffer pointed to by <i>err_data_buffer</i> – TLB . . . . .	2:428
11-103	<i>err_struct_info</i> – Register File . . . . .	2:429
11-104	<i>capabilities</i> Vector for Register File . . . . .	2:430
11-105	Buffer pointed to by <i>err_data_buffer</i> – Register File . . . . .	2:430
11-106	<i>err_struct_info</i> – Bus/Processor Interconnect . . . . .	2:431
11-107	<i>capabilities</i> vector for Bus/Processor Interconnect . . . . .	2:431
11-108	<i>hw_check</i> Fields . . . . .	2:432
11-109	<i>control_word</i> Layout . . . . .	2:438
11-110	<i>pm_info</i> Fields . . . . .	2:440
11-111	<i>pm_buffer</i> Layout . . . . .	2:440
11-112	Processor Features . . . . .	2:447
11-113	Values for <i>ddt</i> Field . . . . .	2:452
11-114	<i>info_request</i> Return Value . . . . .	2:454
11-115	RSE Hints Implemented . . . . .	2:455
11-116	Processor Hardware Sharing Policies . . . . .	2:456
11-117	<i>notify_platform</i> Layout . . . . .	2:460
11-118	<i>vp_env_info</i> – Virtual Environment Information Parameter . . . . .	2:473
11-119	<i>config_options</i> – Global Configuration Options . . . . .	2:479
11-120	PAL Virtualization Services . . . . .	2:486
11-121	State Requirements for PSR for PAL Virtualization Services . . . . .	2:487
11-122	Virtual Processor Settings in Architectural Resources for PAL_VPS_RESUME_NORMAL and PAL_VPS_RESUME_HANDLER2:489	
11-123	Processor Status Register Settings for Virtual Processor Execution . . . . .	2:490
11-124	<i>vhpi</i> – Virtual Highest Priority Pending Interrupt . . . . .	2:495

## Part II: System Programmer's Guide

2-1	Intel® Itanium® Architecture Provides a Relaxed Ordering Model . . . . .	2:512
2-2	Acquire and Release Semantics Order Intel® Itanium® Memory Operations . . . . .	2:513
2-3	Loads May Pass Stores to Different Locations . . . . .	2:514
2-4	Loads May Not Pass Stores in the Presence of a Memory Fence . . . . .	2:514

2-5	Dependencies Do Not Establish MP Ordering (1)	2:515
2-6	Memory Ordering and Data Dependency	2:516
2-7	Memory Ordering and Data Dependency Through a Predicate Register	2:516
2-8	Memory Ordering and Data and Control Dependencies	2:517
2-9	Memory Ordering and Control Dependency	2:517
2-10	Store Buffers May Satisfy Loads if the Stored Data is Not Yet Globally Visible	2:518
2-11	Preventing Store Buffers from Satisfying Local Loads	2:519
2-12	Bypassing to a Semaphore Operation	2:521
2-13	Bypassing from a Semaphore Operation	2:521
2-14	Enforcing the Same Visibility Order to All Observers in a Coherence Domain	2:522
2-15	Intel® Itanium® Architecture Obeys Causality	2:523
2-16	Potential Pipeline Behaviors of the Branch at x from <a href="#">Figure 2-9</a>	2:534
3-1	Interrupt Handler Execution Environment (PSR and RSE.CFLE Settings)	2:540
4-1	Preserving Intel® Itanium® General and Floating-point Registers	2:549
4-2	Register State Preservation at Different Points in the OS	2:552
5-1	Comparison of VHPT Formats	2:572
6-1	Speculation Recovery Code Requirements	2:581
9-1	IA-32 Vectors that need Itanium® Architecture-based OS Support	2:599

§



# ***Part I: System Architecture Guide***





The Intel® Itanium® architecture is a unique combination of innovative features such as explicit parallelism, predication, speculation and more. The architecture is designed to be highly scalable to fill the ever increasing performance requirements of various server and workstation market segments. The Itanium architecture features a revolutionary 64-bit instruction set architecture (ISA) which applies a new processor architecture technology called EPIC, or Explicitly Parallel Instruction Computing. A key feature of the Itanium architecture is IA-32 instruction set compatibility.

The *Intel® Itanium® Architecture Software Developer's Manual* provides a comprehensive description of the programming environment, resources, and instruction set visible to both the application and system programmer. In addition, it also describes how programmers can take advantage of the features of the Itanium architecture to help them optimize code.

## 1.1 Overview of Volume 1: Application Architecture

This volume defines the Itanium application architecture, including application level resources, programming environment, and the IA-32 application interface. This volume also describes optimization techniques used to generate high performance software.

### 1.1.1 Part 1: Application Architecture Guide

Chapter 1, "About this Manual" provides an overview of all volumes in the *Intel® Itanium® Architecture Software Developer's Manual*.

Chapter 2, "Introduction to the Intel® Itanium® Architecture" provides an overview of the architecture.

Chapter 3, "Execution Environment" describes the Itanium register set used by applications and the memory organization models.

Chapter 4, "Application Programming Model" gives an overview of the behavior of Itanium application instructions (grouped into related functions).

Chapter 5, "Floating-point Programming Model" describes the Itanium floating-point architecture (including integer multiply).

Chapter 6, "IA-32 Application Execution Model in an Intel® Itanium® System Environment" describes the operation of IA-32 instructions within the Itanium System Environment from the perspective of an application programmer.

### 1.1.2 Part 2: Optimization Guide for the Intel® Itanium® Architecture

Chapter 1, "About the Optimization Guide" gives an overview of the optimization guide.

Chapter 2, “Introduction to Programming for the Intel® Itanium® Architecture” provides an overview of the application programming environment for the Itanium architecture.

Chapter 3, “Memory Reference” discusses features and optimizations related to control and data speculation.

Chapter 4, “Predication, Control Flow, and Instruction Stream” describes optimization features related to predication, control flow, and branch hints.

Chapter 5, “Software Pipelining and Loop Support” provides a detailed discussion on optimizing loops through use of software pipelining.

Chapter 6, “Floating-point Applications” discusses current performance limitations in floating-point applications and features that address these limitations.

## 1.2 Overview of Volume 2: System Architecture

This volume defines the Itanium system architecture, including system level resources and programming state, interrupt model, and processor firmware interface. This volume also provides a useful system programmer's guide for writing high performance system software.

### 1.2.1 Part 1: System Architecture Guide

Chapter 1, “About this Manual” provides an overview of all volumes in the *Intel® Itanium® Architecture Software Developer's Manual*.

Chapter 2, “Intel® Itanium® System Environment” introduces the environment designed to support execution of Itanium architecture-based operating systems running IA-32 or Itanium architecture-based applications.

Chapter 3, “System State and Programming Model” describes the Itanium architectural state which is visible only to an operating system.

Chapter 4, “Addressing and Protection” defines the resources available to the operating system for virtual to physical address translation, virtual aliasing, physical addressing, and memory ordering.

Chapter 5, “Interruptions” describes all interruptions that can be generated by a processor based on the Itanium architecture.

Chapter 6, “Register Stack Engine” describes the architectural mechanism which automatically saves and restores the stacked subset (GR32 – GR 127) of the general register file.

Chapter 7, “Debugging and Performance Monitoring” is an overview of the performance monitoring and debugging resources that are available in the Itanium architecture.

Chapter 8, “Interruption Vector Descriptions” lists all interruption vectors.

[Chapter 9, “IA-32 Interruption Vector Descriptions”](#) lists IA-32 exceptions, interrupts and intercepts that can occur during IA-32 instruction set execution in the Itanium System Environment.

[Chapter 10, “Itanium® Architecture-based Operating System Interaction Model with IA-32 Applications”](#) defines the operation of IA-32 instructions within the Itanium System Environment from the perspective of an Itanium architecture-based operating system.

[Chapter 11, “Processor Abstraction Layer”](#) describes the firmware layer which abstracts processor implementation-dependent features.

## **1.2.2 Part 2: System Programmer’s Guide**

[Chapter 1, “About the System Programmer’s Guide”](#) gives an introduction to the second section of the system architecture guide.

[Chapter 2, “MP Coherence and Synchronization”](#) describes multiprocessing synchronization primitives and the Itanium memory ordering model.

[Chapter 3, “Interruptions and Serialization”](#) describes how the processor serializes execution around interruptions and what state is preserved and made available to low-level system code when interruptions are taken.

[Chapter 4, “Context Management”](#) describes how operating systems need to preserve Itanium register contents and state. This chapter also describes system architecture mechanisms that allow an operating system to reduce the number of registers that need to be spilled/filled on interruptions, system calls, and context switches.

[Chapter 5, “Memory Management”](#) introduces various memory management strategies.

[Chapter 6, “Runtime Support for Control and Data Speculation”](#) describes the operating system support that is required for control and data speculation.

[Chapter 7, “Instruction Emulation and Other Fault Handlers”](#) describes a variety of instruction emulation handlers that Itanium architecture-based operating systems are expected to support.

[Chapter 8, “Floating-point System Software”](#) discusses how processors based on the Itanium architecture handle floating-point numeric exceptions and how the software stack provides complete IEEE-754 compliance.

[Chapter 9, “IA-32 Application Support”](#) describes the support an Itanium architecture-based operating system needs to provide to host IA-32 applications.

[Chapter 10, “External Interrupt Architecture”](#) describes the external interrupt architecture with a focus on how external asynchronous interrupt handling can be controlled by software.

[Chapter 11, “I/O Architecture”](#) describes the I/O architecture with a focus on platform issues and support for the existing IA-32 I/O port space.

Chapter 12, “Performance Monitoring Support” describes the performance monitor architecture with a focus on what kind of support is needed from Itanium architecture-based operating systems.

Chapter 13, “Firmware Overview” introduces the firmware model, and how various firmware layers (PAL, SAL, UEFI, ACPI) work together to enable processor and system initialization, and operating system boot.

### 1.2.3 Appendices

Appendix A, “Code Examples” provides OS boot flow sample code.

## 1.3 Overview of Volume 3: Intel® Itanium® Instruction Set Reference

This volume is a comprehensive reference to the Itanium instruction set, including instruction format/encoding.

Chapter 1, “About this Manual” provides an overview of all volumes in the *Intel® Itanium® Architecture Software Developer’s Manual*.

Chapter 2, “Instruction Reference” provides a detailed description of all Itanium instructions, organized in alphabetical order by assembly language mnemonic.

Chapter 3, “Pseudo-Code Functions” provides a table of pseudo-code functions which are used to define the behavior of the Itanium instructions.

Chapter 4, “Instruction Formats” describes the encoding and instruction format instructions.

Chapter 5, “Resource and Dependency Semantics” summarizes the dependency rules that are applicable when generating code for processors based on the Itanium architecture.

## 1.4 Overview of Volume 4: IA-32 Instruction Set Reference

This volume is a comprehensive reference to the IA-32 instruction set, including instruction format/encoding.

Chapter 1, “About this Manual” provides an overview of all volumes in the *Intel® Itanium® Architecture Software Developer’s Manual*.

Chapter 2, “Base IA-32 Instruction Reference” provides a detailed description of all base IA-32 instructions, organized in alphabetical order by assembly language mnemonic.



Chapter 3, “IA-32 Intel® MMX™ Technology Instruction Reference” provides a detailed description of all IA-32 Intel® MMX™ technology instructions designed to increase performance of multimedia intensive applications. Organized in alphabetical order by assembly language mnemonic.

Chapter 4, “IA-32 SSE Instruction Reference” provides a detailed description of all IA-32 SSE instructions designed to increase performance of multimedia intensive applications, and is organized in alphabetical order by assembly language mnemonic.

## 1.5 Terminology

The following definitions are for terms related to the Itanium architecture and will be used throughout this document:

**Instruction Set Architecture (ISA)** – Defines application and system level resources. These resources include instructions and registers.

**Itanium Architecture** – The new ISA with 64-bit instruction capabilities, new performance-enhancing features, and support for the IA-32 instruction set.

**IA-32 Architecture** – The 32-bit and 16-bit Intel architecture as described in the *Intel® 64 and IA-32 Architectures Software Developer’s Manual*.

**Itanium System Environment** – The operating system environment that supports the execution of both IA-32 and Itanium architecture-based code.

**Itanium Architecture-based Firmware** – The Processor Abstraction Layer (PAL) and System Abstraction Layer (SAL).

**Processor Abstraction Layer (PAL)** – The firmware layer which abstracts processor features that are implementation dependent.

**System Abstraction Layer (SAL)** – The firmware layer which abstracts system features that are implementation dependent.

## 1.6 Related Documents

The following documents can be downloaded at the Intel’s Developer Site at <http://developer.intel.com>:

- **Dual-Core Update to the Intel® Itanium® 2 Processor Reference Manual for Software Development and Optimization**– Document number 308065 provides model-specific information about the dual-core Itanium processors.
- **Intel® Itanium® 2 Processor Reference Manual for Software Development and Optimization** – This document (Document number 251110) describes model-specific architectural features incorporated into the Intel® Itanium® 2 processor, the second processor based on the Itanium architecture.
- **Intel® Itanium® Processor Reference Manual for Software Development** – This document (Document number 245320) describes model-specific architectural features incorporated into the Intel® Itanium® processor, the first processor based on the Itanium architecture.

- **Intel® 64 and IA-32 Architectures Software Developer’s Manual** – This set of manuals describes the Intel 32-bit architecture. They are available from the Intel Literature Department by calling 1-800-548-4725 and requesting Document Numbers 243190, 243191 and 243192.
- **Intel® Itanium® Software Conventions and Runtime Architecture Guide** – This document (Document number 245358) defines general information necessary to compile, link, and execute a program on an Itanium architecture-based operating system.
- **Intel® Itanium® Processor Family System Abstraction Layer Specification** – This document (Document number 245359) specifies requirements to develop platform firmware for Itanium architecture-based systems.

The following document can be downloaded at the Unified EFI Forum website at <http://www.uefi.org>:

- **Unified Extensible Firmware Interface Specification** – This document defines a new model for the interface between operating systems and platform firmware.

## 1.7 Revision History

Date of Revision	Revision Number	Description
March 2010	2.3	<p>Added information about illegal virtualization optimization combinations and IIPA requirements.</p> <p>Added Resource Utilization Counter and PAL_VP_INFO.</p> <p>PAL_VP_INIT and VPD.vpr changes.</p> <p>New PAL_VPS_RESUME_HANDLER parameter to indicate RSE Current Frame Load Enable setting at the target instruction.</p> <p>PAL_VP_INIT_ENV implementation-specific configuration option.</p> <p>Minimum Virtual address increased to 54 bits.</p> <p>New PAL_MC_ERROR_INFO health indicator.</p> <p>New PAL_MC_ERROR_INJECT implementation-specific bit fields.</p> <p>MOV-to_SR.L reserved field checking.</p> <p>Added virtual machine disable.</p> <p>Added variable frequency mode additions to ACPI P-state description.</p> <p>Removed <i>pal_proc_vector</i> argument from PAL_VP_SAVE and PAL_VP_RESTORE.</p> <p>Added PAL_PROC_SET_FEATURES data speculation disable.</p> <p>Added Interruption Instruction Bundle registers.</p> <p>Min-state save area size change.</p> <p>PAL_MC_DYNAMIC_STATE changes.</p> <p>PAL_PROC_SET_FEATURES data poisoning promotion changes.</p> <p>ACPI P-state clarifications.</p> <p>Synchronization requirements for virtualization opcode optimization.</p> <p>New priority hint and multi-threading hint recommendations.</p>

Date of Revision	Revision Number	Description
August 2005	2.2	<p>Allow register fields in CR.LID register to be read-only and CR.LID checking on interruption messages by processors optional. See Vol 2, Part I, Ch 5 “Interruptions” and Section 11.2.2 PALE_RESET Exit State for details.</p> <p>Relaxed reserved and ignored fields checkings in IA-32 application registers in Vol 1 Ch 6 and Vol 2, Part I, Ch 10.</p> <p>Introduced visibility constraints between stores and local purges to ensure TLB consistency for UP VHPT update and local purge scenarios. See Vol 2, Part I, Ch 4 and description of <code>ptc.l</code> instruction in Vol 3 for details.</p> <p>Architecture extensions for processor Power/Performance states (P-states). See Vol 2 PAL Chapter for details.</p> <p>Introduced Unimplemented Instruction Address fault.</p> <p>Relaxed ordering constraints for VHPT walks. See Vol 2, Part I, Ch 4 and 5 for details.</p> <p>Architecture extensions for processor virtualization.</p> <p>All instructions which must be last in an instruction group results in undefined behavior when this rule is violated.</p> <p>Added architectural sequence that guarantees increasing ITC and PMD values on successive reads.</p> <p>Addition of PAL_BRAND_INFO, PAL_GET_HW_POLICY, PAL_MC_ERROR_INJECT, PAL_MEMORY_BUFFER, PAL_SET_HW_POLICY and PAL_SHUTDOWN procedures.</p> <p>Allows IPI-redirection feature to be optional.</p> <p>Undefined behavior for 1-byte accesses to the non-architected regions in the IPI block.</p> <p>Modified insertion behavior for TR overlaps. See Vol 2, Part I, Ch 4 for details.</p> <p>“Bus parking” feature is now optional for PAL_BUS_GET_FEATURES.</p> <p>Introduced low-power synchronization primitive using <code>hint</code> instruction.</p> <p>FR32-127 is now preserved in PAL calling convention.</p> <p>New return value from PAL_VM_SUMMARY procedure to indicate the number of multiple concurrent outstanding TLB purges.</p> <p>Performance Monitor Data (PMD) registers are no longer sign-extended.</p> <p>New memory attribute transition sequence for memory on-line delete. See Vol 2, Part I, Ch 4 for details.</p> <p>Added 'shared error' (se) bit to the Processor State Parameter (PSP) in PAL_MC_ERROR_INFO procedure.</p> <p>Clarified PMU interrupts as edge-triggered.</p> <p>Modified 'proc_number' parameter in PAL_LOGICAL_TO_PHYSICAL procedure.</p> <p>Modified <code>pal_copy_info</code> alignment requirements.</p> <p>New bit in PAL_PROC_GET_FEATURES for variable P-state performance.</p> <p>Clarified descriptions for <code>check_target_register</code> and <code>check_target_register_sof</code>.</p> <p>Various fixes in dependency tables in Vol 3 Ch 5.</p> <p>Clarified effect of sending IPIs to non-existent processor in Vol 2, Part I, Ch 5.</p> <p>Clarified instruction serialization requirements for interruptions in Vol 2, Part II, Ch 3.</p> <p>Updated performance monitor context switch routine in Vol 2, Part I, Ch 7.</p>

Date of Revision	Revision Number	Description
August 2002	2.1	<p>Added Predicate Behavior of <code>alloc</code> Instruction Clarification (Section 4.1.2, Part I, Volume 1; Section 2.2, Part I, Volume 3).</p> <p>Added New <code>fc.i</code> Instruction (Section 4.4.6.1, and 4.4.6.2, Part I, Volume 1; Section 4.3.3, 4.4.1, 4.4.5, 4.4.6, 4.4.7, 5.5.2, and 7.1.2, Part I, Volume 2; Section 2.5, 2.5.1, 2.5.2, 2.5.3, and 4.5.2.1, Part II, Volume 2; Section 2.2, 3, 4.1, 4.4.6.5, and 4.4.10.10, Part I, Volume 3).</p> <p>Added Interval Time Counter (ITC) Fault Clarification (Section 3.3.2, Part I, Volume 2).</p> <p>Added Interruption Control Registers Clarification (Section 3.3.5, Part I, Volume 2).</p> <p>Added Spontaneous NaT Generation on Speculative Load (<code>ld.s</code>) (Section 5.5.5 and 11.9, Part I, Volume 2; Section 2.2 and 3, Part I, Volume 3).</p> <p>Added Performance Counter Standardization (Sections 7.2.3 and 11.6, Part I, Volume 2).</p> <p>Added Freeze Bit Functionality in Context Switching and Interrupt Generation Clarification (Sections 7.2.1, 7.2.2, 7.2.4.1, and 7.2.4.2, Part I, Volume 2)</p> <p>Added <code>IA_32_Exception</code> (Debug) IIPA Description Change (Section 9.2, Part I, Volume 2).</p> <p>Added capability for Allowing Multiple <code>PAL_A_SPEC</code> and <code>PAL_B</code> Entries in the Firmware Interface Table (Section 11.1.6, Part I, Volume 2).</p> <p>Added BR1 to Min-state Save Area (Sections 11.3.2.3 and 11.3.3, Part I, Volume 2).</p> <p>Added Fault Handling Semantics for <code>lfetch.fault</code> Instruction (Section 2.2, Part I, Volume 3).</p>
December 2001	2.0	<p>Volume 1:</p> <p>Faults in <code>ld.c</code> that hits ALAT clarification (Section 4.4.5.3.1).</p> <p>IA-32 related changes (Section 6.2.5.4, Section 6.2.3, Section 6.2.4, Section 6.2.5.3).</p> <p>Load instructions change (Section 4.4.1).</p>

Date of Revision	Revision Number	Description
		<p>Volume 2:</p> <p>Class pr-writers-int clarification (Table A-5).</p> <p>PAL_MC_DRAIN clarification (Section 4.4.6.1).</p> <p>VHPT walk and forward progress change (Section 4.1.1.2).</p> <p>IA-32 IBR/DBR match clarification (Section 7.1.1).</p> <p>ISR figure changes (pp. 8-5, 8-26, 8-33 and 8-36).</p> <p>PAL_CACHE_FLUSH return argument change – added new status return argument (Section 11.8.3).</p> <p>PAL self-test Control and PAL_A procedure requirement change – added new arguments, figures, requirements (Section 11.2).</p> <p>PAL_CACHE_FLUSH clarifications (Chapter 11).</p> <p>Non-speculative reference clarification (Section 4.4.6).</p> <p>RID and Preferred Page Size usage clarification (Section 4.1).</p> <p>VHPT read atomicity clarification (Section 4.1).</p> <p>IIP and WC flush clarification (Section 4.4.5).</p> <p>Revised RSE and PMC typographical errors (Section 6.4).</p> <p>Revised DV table (Section A.4).</p> <p>Memory attribute transitions – added new requirements (Section 4.4).</p> <p>MCA for WC/UC aliasing change (Section 4.4.1).</p> <p>Bus lock deprecation – changed behavior of DCR 'lc' bit (Section 3.3.4.1, Section 10.6.8, Section 11.8.3).</p> <p>PAL_PROC_GET/SET_FEATURES changes – extend calls to allow implementation-specific feature control (Section 11.8.3).</p> <p>Split PAL_A architecture changes (Section 11.1.6).</p> <p>Simple barrier synchronization clarification (Section 13.4.2).</p> <p>Limited speculation clarification – added hardware-generated speculative references (Section 4.4.6).</p> <p>PAL memory accesses and restrictions clarification (Section 11.9).</p> <p>PSP validity on INITs from PAL_MC_ERROR_INFO clarification (Section 11.8.3).</p> <p>Speculation attributes clarification (Section 4.4.6).</p> <p>PAL_A FIT entry, PAL_VM_TR_READ, PSP, PAL_VERSION clarifications (Sections 11.8.3 and 11.3.2.1).</p> <p>TLB searching clarifications (Section 4.1).</p> <p>IA-32 related changes (Section 10.3, Section 10.3.2, Section 10.3.2, Section 10.3.3.1, Section 10.10.1).</p> <p>IPSR.ri and ISR.ei changes (Table 3-2, Section 3.3.5.1, Section 3.3.5.2, Section 5.5, Section 8.3, and Section 2.2).</p>
		<p>Volume 3:</p> <p>IA-32 CPUID clarification (p. 5-71).</p> <p>Revised figures for extract, deposit, and alloc instructions (Section 2.2).</p> <p>RCPSS, RCPSS, RSQRTPS, and RSQRTSS clarification (Section 7.12).</p> <p>IA-32 related changes (Section 5.3).</p> <p>tak, tpa change (Section 2.2).</p>
July 2000	1.1	<p>Volume 1:</p> <p>Processor Serial Number feature removed (Chapter 3).</p> <p>Clarification on exceptions to instruction dependency (Section 3.4.3).</p>

Date of Revision	Revision Number	Description
		<p>Volume 2:</p> <p>Clarifications regarding “reserved” fields in ITIR (Chapter 3).</p> <p>Instruction and Data translation must be enabled for executing IA-32 instructions (Chapters 3,4 and 10).</p> <p>FCR/FDR mappings, and clarification to the value of PSR.ri after an RFI (Chapters 3 and 4).</p> <p>Clarification regarding ordering data dependency.</p> <p>Out-of-order IPI delivery is now allowed (Chapters 4 and 5).</p> <p>Content of EFLAG field changed in IIM (p. 9-24).</p> <p>PAL_CHECK and PAL_INIT calls – exit state changes (Chapter 11).</p> <p>PAL_CHECK processor state parameter changes (Chapter 11).</p> <p>PAL_BUS_GET/SET_FEATURES calls – added two new bits (Chapter 11).</p> <p>PAL_MC_ERROR_INFO call – Changes made to enhance and simplify the call to provide more information regarding machine check (Chapter 11).</p> <p>PAL_ENTER_IA_32_Env call changes – entry parameter represents the entry order; SAL needs to initialize all the IA-32 registers properly before making this call (Chapter 11).</p> <p>PAL_CACHE_FLUSH – added a new cache_type argument (Chapter 11).</p> <p>PAL_SHUTDOWN – removed from list of PAL calls (Chapter 11).</p> <p>Clarified memory ordering changes (Chapter 13).</p> <p>Clarification in dependence violation table (Appendix A).</p>
		<p>Volume 3:</p> <p>fmix instruction page figures corrected (Chapter 2).</p> <p>Clarification of “reserved” fields in ITIR (Chapters 2 and 3).</p> <p>Modified conditions for alloc/loadrs/flushrs instruction placement in bundle/ instruction group (Chapters 2 and 4).</p> <p>IA-32 JMPE instruction page typo fix (p. 5-238).</p> <p>Processor Serial Number feature removed (Chapter 5).</p>
January 2000	1.0	Initial release of document.

§



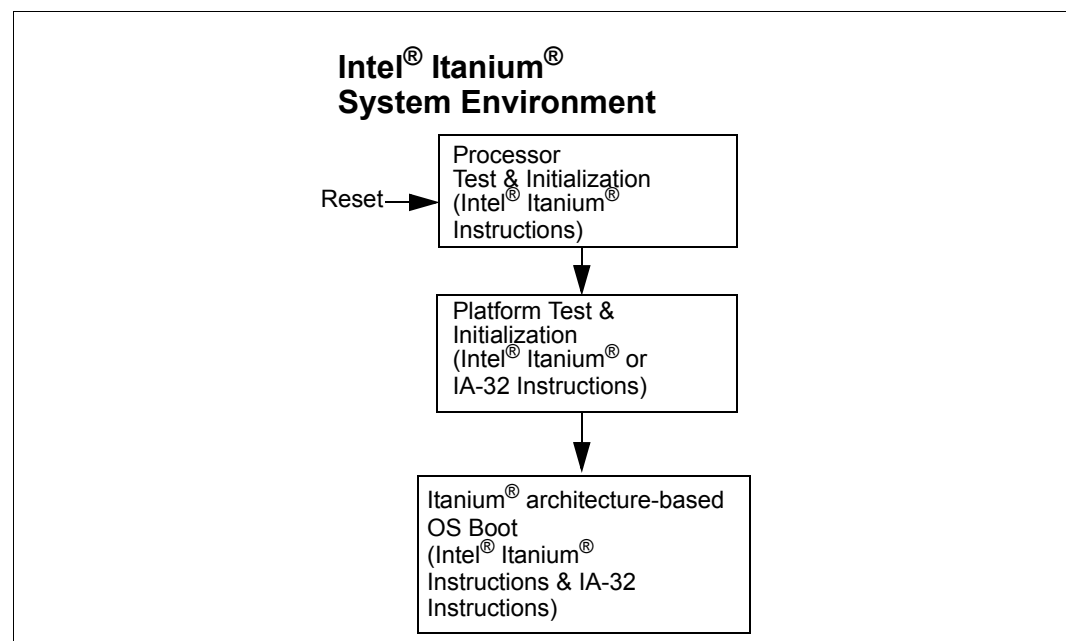
As described in [Section 2.1, “Operating Environments” on page 1:13](#), the Itanium System Environment supports Itanium architecture-based operating systems. The architectural model also supports a mixture of IA-32 and Itanium architecture-based application code within an Itanium architecture-based operating system.

The system environment determines the set of processor system resources seen by the operating system. These resources include: virtual memory management, physical memory attributes, external interrupt mechanisms, exception and interrupt delivery, machine check architectures, debug, performance monitoring, control registers, and the set of privileged instructions.

## 2.1 Processor Boot Sequence

[Figure 2-1](#) shows the defined boot sequence. Unlike IA-32 processors, which power up in 32-bit Real Mode, processors in the Itanium processor family power up in the Itanium System Environment running Itanium architecture-based code. Processor initialization, testing, memory, and platform initialization/testing are performed by processor firmware. Mechanisms are provided to execute Real Mode IA-32 boot BIOSs and device drivers during the boot sequence.

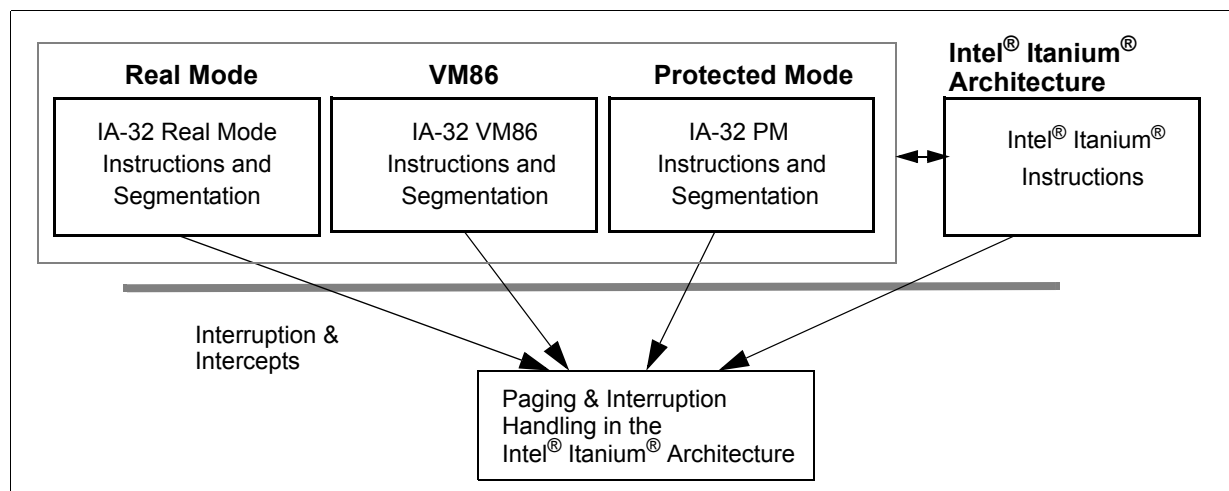
**Figure 2-1. System Environment Boot Flow**



## 2.2 Intel® Itanium® System Environment Overview

The Itanium System Environment is designed to support execution of Itanium architecture-based operating systems running IA-32 or Itanium architecture-based applications. IA-32 applications can interact with Itanium architecture-based operating systems, applications and libraries within this environment. Both IA-32 application level code and Itanium instructions can be executed by the operating system and user level software. The entire machine state, including the IA-32 general registers and floating-point registers, segment selectors and descriptors is accessible to Itanium architecture-based code. As shown in Figure 2-2, all major IA-32 operating modes are fully supported.

Figure 2-2. Intel® Itanium® System Environment



In the Itanium system environment, Itanium architecture operating system resources supersede all IA-32 system resources. Specifically, the IA-32 defined set of control, test, debug, machine check registers, privilege instructions, and virtual paging algorithms are replaced by the Itanium architecture system resources. When IA-32 code is running on an Itanium architecture-based operating system, the processor directly executes all performance critical but non-sensitive IA-32 application level instructions. Accesses to sensitive system resources (interrupt flags, control registers, TLBs, etc.) are intercepted into the Itanium architecture-based operating system. Using this set of intervention hooks, an Itanium architecture-based operating system can emulate or virtualize an IA-32 system resource for an IA-32 application, OS, or device driver.

The Itanium system architecture features are presented in the following chapters:

- Chapter 3, "System State and Programming Model" describes system resources.
- Chapter 4, "Addressing and Protection" describes the virtual memory architecture.
- Chapter 5, "Interruptions" defines the interrupt and exception architecture.
- Chapter 6, "Register Stack Engine" describes the register stack engine.
- Chapter 7, "Debugging and Performance Monitoring" describes debug and performance monitoring hooks.
- Chapter 8, "Interrupt Vector Descriptions" describes interruption handler entry points.

Additional support for IA-32 applications in the Itanium system environment is defined by chapters:

- [Chapter 9](#) describes IA-32 interruption handler entry points.
- [Chapter 10, "Itanium® Architecture-based Operating System Interaction Model with IA-32 Applications"](#) describes how IA-32 applications interact with Itanium architecture-based operating systems.

§



This chapter describes the architectural state visible only to an operating system and defines system state programming models. It covers the functional descriptions of all the system state registers, descriptions of individual fields in each register, and their serialization requirements. The virtual and physical memory management details are described in [Chapter 4, “Addressing and Protection.”](#) Interruptions are described in [Chapter 5, “Interruptions.”](#)

**Note:** Unless otherwise noted, references to “interruption” in this chapter refer to IVA-based interruptions. See [“Interruption Definitions” on page 2:95.](#)

## 3.1 Privilege Levels

Four privilege levels, numbered from 0 to 3, are provided to control access to system instructions, system registers and system memory areas. Level 0 is the most privileged and level 3 the least privileged. Application instructions and registers can be accessed at any privilege level. System instructions and registers defined in this chapter can only be accessed at privilege level 0; otherwise, a Privilege Operation fault is raised. The processor maintains a Current Privilege Level (CPL) in the `cpl` field of the Processor Status Register (PSR). CPL can only be modified by controlled entry and exit points managed by the operating system. Virtual memory protection mechanisms control memory accesses based on the Privilege Level (PL) of the virtual page and the CPL.

## 3.2 Serialization

For all application and system level resources, apart from the control register file, the processor ensures values written to a register are observed by instructions in subsequent instruction groups. This is termed **data dependency**. For example, writes to general registers, floating-point and application registers are observed by subsequent reads of the same register. (See [“Control Registers” on page 2:29](#) for control register serialization requirements.) For modifications of application level resources with side effects, the side effects are ensured by the processor to be observed by subsequent instruction groups. This is termed **implicit serialization**. Application registers (ARs), with the exception of the Interval Time Counter, the User Mask, when modified by `sum`, `rum`, and `mov` to `psr.um`, and the Current Frame Marker (CFM), are implicitly serialized. PMD registers have special serialization requirements as described in [“Generic Performance Counter Registers” on page 2:156](#). All other application-level resources (GRs, FRs, PRs, BRs, IP, CPUID) have no side effects and so need not be serialized.

To avoid serialization overhead in privileged operating system code, system register resources are not implicitly serialized. The processor does not ensure modification of registers with side effects are observed by subsequent instruction groups. For system register resources other than control registers, the processor ensures data dependencies are honored (reads see the results of prior writes to the same register). See Section 3.3.3, “Control Registers” and [Table 3-3 on page 2:29](#) for control register

serialization requirements. This approach simplifies hardware and allows for more efficient software operations. For example, during a low level context switch where there is no immediate use of loaded system registers, these registers can be loaded without any serialization overhead. To ensure side effects are observed before a dependent instruction is fetched or executed, two serialization operations are provided: **instruction serialization** and **data serialization**.

### 3.2.1 Instruction Serialization

**Instruction serialization** ensures that modifications to processor resources are observed before subsequent instruction group fetches are re-initiated. Software must use an instruction serialization operation before any instruction group that is dependent upon the modified system resource. Resource side effects may be observed at any point before the explicit serialization operation.

Modification of the following system resources (if the modification affects instruction fetching) require instruction serialization: RR, PKR, ITR, ITC, IBR, PMC, PMD, PSR bits as defined in [“Processor Status Register \(PSR\)” on page 2:23](#) and Control Registers as defined in [“Control Registers” on page 2:29](#).

The instructions Return from Interruption (`rfi`) and Instruction Serialize (`srlz.i`) perform explicit instruction serialization.

An interruption performs an implicit instruction serialization operation, so the first instruction group in the interruption handler will observe the serialized state.

Instruction Serialization Example:

```
mov ibr[reg]= reg    // move to instruction debug register
;;                 // end of instruction group
srlz.i              // ensure subsequent instruction fetches observe
                   // modification
;;                 // end of instruction group
inst                // dependent instruction
```

**Note:** The serializing instruction, the instruction to be serialized, and any operations dependent on the serialization must be in three separate instruction groups.

### 3.2.2 Data Serialization

**Data serialization** ensures that modifications to processor resources affecting both execution and data memory accesses are observed. Software must issue a data serialize operation prior to the instruction dependent upon the modified resource. Data serialization can be issued within the same instruction group as the dependent instruction. Resource side effects may be observed at any point before the explicit serialization operation.

Modification of the following system resources require data serialization: RR, PKR, RUC, DTR, DTC, DBR, PMC, PMD, PSR bits as defined in [“Processor Status Register \(PSR\)” on page 2:23](#) and Control Registers as defined in [“Control Registers” on page 2:29](#).

The control registers are different from the general registers and other registers. Most control registers require an explicit data serialization between the writing of a control register and the reading of that same control register. (See [Table 3-3 on page 2:29](#) for serialization requirements for specific control registers.)

The Data Serialize (`srlz.d`) instruction performs explicit data serialization. Instruction serialization operations (`rfi`, `srlz.i`, and interruptions) also perform a data serialization operation.

Data Serialization Example:

```
mov rr[reg] = reg    //move into region register
;;                 //end of instruction group
srlz.d              //serialize region register modification
ld                 //perform a dependent load
```

The serializing instruction and the instruction to be serialized (the one writing the resource) must be in two different instruction groups. Operations dependent on the serialization and the serialization can be in the same instruction group, but the `srlz` instruction must be before the dependent instruction slot.

### 3.2.3 Definition of In-flight Resources

When the value of a resource that requires an explicit instruction or data serialization is changed by one or more writers, that resource is said to be **in-flight** until the required serialization is performed. There can be multiple in-flight values if multiple writers have occurred since the last serialization.

An instruction that reads an in-flight resource will see one of the in-flight values or the state prior to any of the unserialized writers. However, whether such a reader sees the original or one of the in-flight values is not predictable.

For a reader of an in-flight resource, this definition includes (but is not limited to) the following possible outcomes:

- The reader of an in-flight resource may see the most-recently-serialized value or any of the in-flight values each time it is executed – seeing the value from a particular writer one time does not guarantee that the same writer’s value will be seen by that reader the next time.
- Multiple readers of an in-flight resource may see different values – each may see the most-recently-serialized value or any of the in-flight values, independent of what other readers may see.
- If a single execution of an instruction reads an in-flight resource more than once during its execution, each read may see a different value.

Thus, the only way to guarantee that the latest value is seen by a reader is to perform the required serialization.

## 3.3 System State

The architecture provides a rich set of system register resources for process control, interruptions handling, protection, debugging, and performance monitoring. This section gives an overview of these resources.

### 3.3.1 System State Overview

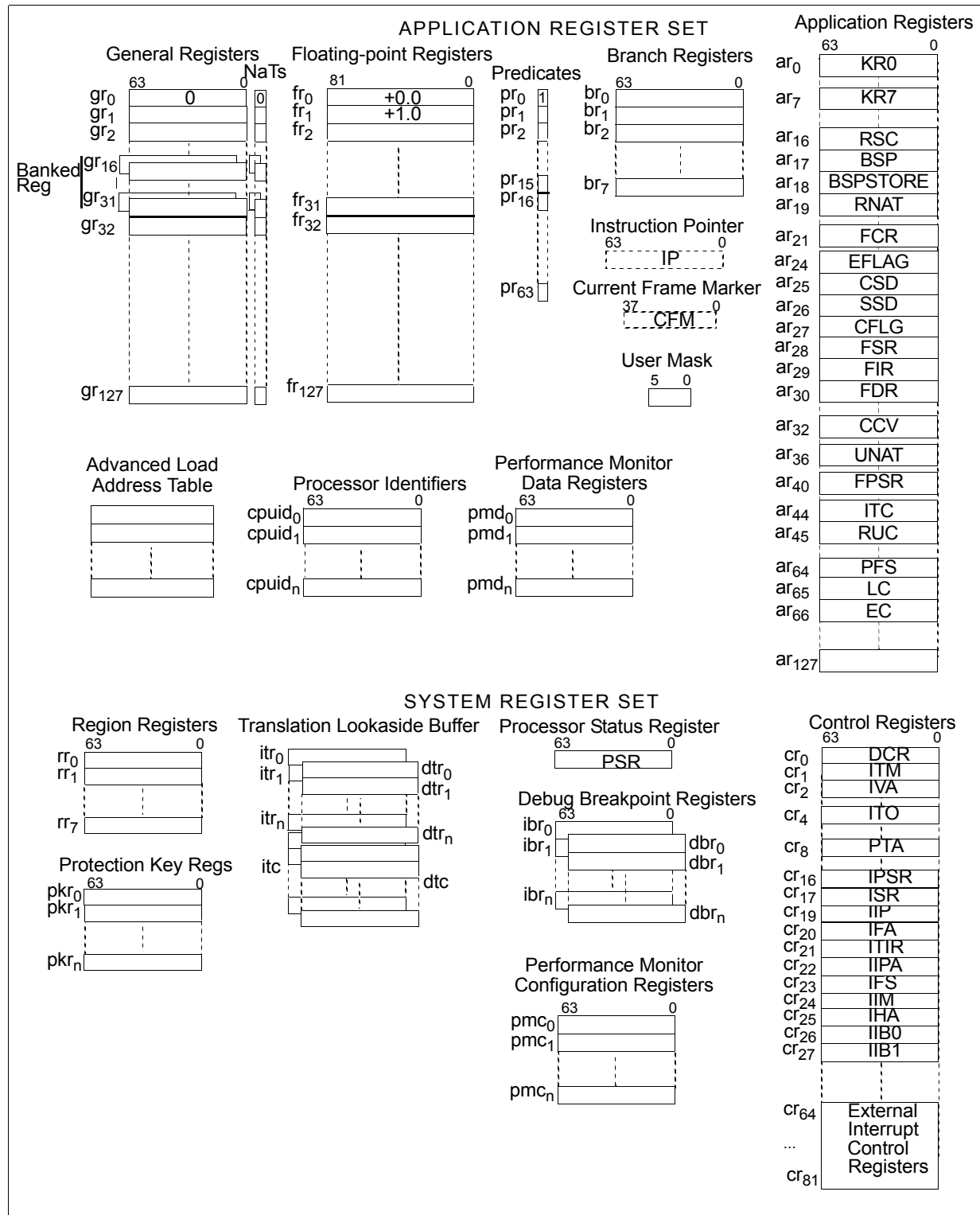
Figure 3-1 shows the set of all defined privileged system register resources. Application state as defined in “[Application Register State](#)” on page 1:23 is also accessible.

- **Processor Status Register (PSR)** – 64-bit register that maintains control information for the currently running process. See “[Processor Status Register \(PSR\)](#)” on page 2:23 for complete details.
- **Control Registers (CR)** – This register name space contains several 64-bit registers that capture the state of the processor on an interruption, enable system-wide features, and specify global processor parameters for interruptions and memory management. See “[Control Registers](#)” on page 2:29 for complete information.
- **Interrupt Registers** – These registers provide the capability of masking external interrupts, reading external interrupt vector numbers, programming vector numbers for internal processor asynchronous events and external interrupt sources. For complete information, see “[Interrupts](#)” on page 2:114.
- **Interval Timer Facilities** – A 64-bit interval timer is provided for privileged and non-privileged use and as a time base for performance measurements. Timing facilities are defined in detail in “[Interval Time Counter and Match Register \(ITC – AR44 and ITM – CR1\)](#)” on page 2:32.
- **Resource Utilization Facility** – A 64-bit resource utilization counter is provided for privileged and non-privileged use. This counts the number of Interval Timer cycles consumed by this logical processor. See [Section 3.1.8.11, “Resource Utilization Counter \(RUC – AR 45\)”](#) on page 1:31.
- **Debug Breakpoint Registers (DBR/IBR)** – 64-bit Data and 64-bit Instruction Breakpoint Register pairs (DBR, IBR) can be programmed to fault on reference to a range of virtual and physical addresses generated by either Itanium or IA-32 instructions. See “[Debugging](#)” on page 2:151 for details. The minimum number of DBR register pairs and IBR register pairs is 4 in any implementation. On some implementations, a hardware debugger may use two or more of these register pairs for its own use; see “[Data and Instruction Breakpoint Registers](#)” on page 2:152 for details.
- **Performance Monitor Configuration/Data Registers (PMC/PMD)** – Multiple performance monitors can be programmed to measure a wide range of user, operating system, or processor performance values. Performance monitors can be programmed to measure performance values from either IA-32 or Itanium instructions. Performance monitors are defined in “[Performance Monitoring](#)” on page 2:155. The minimum number of generic PMC/PMD register pairs in any implementation is 4.
- **Banked General Registers** – A set of 16 banked 64-bit general purpose registers, GR 16-GR 31, are available as temporary storage and register context when operating in low level interruption code. See “[Banked General Registers](#)” on page 2:42 for complete details.



- **Region Registers (RR)** – Eight 64-bit region registers specify the identifiers and preferred page sizes for multiple virtual address spaces. Refer to [“Region Registers \(RR\)” on page 2:58](#) for complete information.
- **Protection Key Registers (PKR)** – At least sixteen 64-bit protection key registers contain protection keys and read, write, execute permissions for virtual memory protection domains. Please see the processor-specific documentation for further information on the number of Protection Key Registers implemented on the Itanium processor. Refer to [“Protection Keys” on page 2:59](#) for details.
- **Translation Lookaside Buffer (TLB)** – Holds recently used virtual to physical address mappings. The TLB is divided into Instruction (ITLB), Data (DTLB), Translation Registers (TR) and Translation Cache (TC) sections. See [“Translation Lookaside Buffer \(TLB\)” on page 2:47](#) for complete details. Translation Registers are software managed portions of the TLB and the Translation Cache section of the TLB is directly managed by the processor.

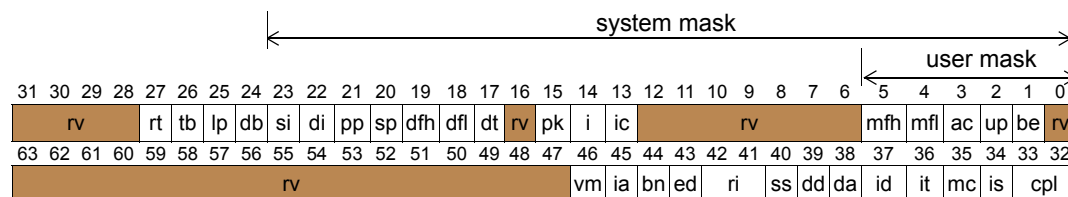
**Figure 3-1. System Register Model**



### 3.3.2 Processor Status Register (PSR)

The PSR maintains the current execution environment. The PSR is divided into four overlapping sections (See Figure 3-2): user mask bits (PSR{5:0}), system mask bits (PSR{23:0}), the lower half (PSR{31:0}), and the entire PSR (PSR{63:0}). PSR fields are defined in Table 3-2 along with serialization requirements for modification of each field and the state of the field after an interruption.

**Figure 3-2. Processor Status Register (PSR)**



The PSR instructions and their serialization requirements are defined in Table 3-1. These instructions explicitly read or write portions of the PSR. Other instructions also read and write portions of the PSR as described in Table 3-2 and Table 5-2.

**Table 3-1. Processor Status Register Instructions**

Mnemonic	Description	Operation	Instr. Type	Serialization Required
<code>sum imm</code>	Set user mask from immediate	$PSR\{5:0\} \leftarrow PSR\{5:0\}   imm$	M	implicit
<code>rum imm</code>	Reset user mask from immediate	$PSR\{5:0\} \leftarrow PSR\{5:0\} \& \sim imm$	M	implicit
<code>mov psr.um = r<sub>2</sub></code>	Move to user mask	$PSR\{5:0\} \leftarrow GR[r_2]$	M	implicit
<code>mov r<sub>1</sub> = psr.um</code>	Move from user mask	$GR[r_1] \leftarrow PSR\{5:0\}$	M	none
<code>ssm imm</code>	Set system mask from immediate	$PSR\{23:0\} \leftarrow PSR\{23:0\}   imm$	M	data/inst <sup>a</sup>
<code>rsm imm</code>	Reset system mask from immediate	$PSR\{23:0\} \leftarrow PSR\{23:0\} \& \sim imm$	M	data/inst <sup>a</sup>
<code>mov psr.l = r<sub>2</sub></code>	Move to lower PSR	$PSR\{31:0\} \leftarrow GR[r_2]$	M	data/inst <sup>a</sup>
<code>mov r<sub>1</sub> = psr</code>	Move from PSR	$GR[r_1] \leftarrow PSR\{36:35,31:0\}^b$	M	none
<code>bsw.0, bsw.1</code>	Bank switch	$PSR\{44\} \leftarrow 0 \text{ or } 1$	B	implicit
<code>vmsw.0, vmsw.1</code>	Virtual machine switch	$PSR\{46\} \leftarrow 0 \text{ or } 1$	B	implicit
<code>rfi</code>	Return From Interruption	$PSR\{63:0\} \leftarrow IPSR$	B	implicit

- a. Based upon the resource being serialized, use data or instruction serialization.  
b. All other bits of the PSR read as zero.

The user mask, PSR{5:0}, can be set and cleared by the Set User Mask (`sum`), Reset User Mask (`rum`) and Move to User Mask (`mov psr.um=`) instructions at any privilege level. For user mask modifications by `sum`, `rum` and `mov`, the processor ensures all side effects are observed before subsequent instruction groups.

The system mask, PSR{23:0}, can be set and cleared by the Set System Mask (*ssm*) and Reset System Mask (*rsm*) instructions. Software must issue the appropriate serialization operation before dependent instructions. The system mask instructions are privileged.

The lower half of the PSR, PSR{31:0}, can be written with the Move to Lower PSR (*mov psr.l=*) instruction. Software must issue the appropriate serialization operation before dependent instructions. The Move to Lower PSR instruction is privileged.

The PSR can be read with the Move from PSR (*mov =psr*) instruction. Only PSR{36:35} and PSR{31:0} are written to the target register by Move from PSR. PSR{63:37} and PSR{34:32} can only be read after an interruption by reading the state in IPSR. The entire PSR is updated from IPSR by the Return from Interruption (*rfi*) instruction. An *rfi* also implicitly serializes the PSR. Both Move from PSR and Return from Interruption are privileged.

**Table 3-2. Processor Status Register Fields**

Field	Bits	Description	Interruption State	Serialization Required
User Mask = PSR{5:0}				
rv	0	reserved		
be	1	Big-Endian – When 1, data memory references are big-endian. When 0, data memory references are little endian. This bit is ignored for IA-32 data references, which are always performed little-endian. Instruction fetches are always performed little endian.	DCR.be	data <sup>a</sup>
up	2	User Performance monitor enable – When 1, performance monitors configured as user monitors are enabled to count events (including IA-32). When 0, user configured monitors are disabled. See <a href="#">“Performance Monitoring” on page 2:155</a> for details.	unchanged	data <sup>a</sup> inst <sup>b</sup>
ac	3	Alignment Check – When 1, all unaligned data memory references result in an Unaligned Data Reference fault. When 0, unaligned data memory references may or may not result in a Unaligned Data Reference fault. See <a href="#">“Memory Datum Alignment and Atomicity” on page 2:93</a> for details. Unaligned semaphore references also result in a Unaligned Data Reference fault, regardless of the state of PSR.ac. For IA-32 instructions, if PSR.ac is 1 an unaligned IA-32 data memory reference raises an IA_32_Exception(AlignmentCheck) fault. When 0, additional IA-32 control bits as defined in Section 10.6.7, “Memory Alignment” also generate alignment checks.	0	data <sup>a</sup>
mfl	4	Lower (f2..f31) floating-point registers written – This bit is set to one when an Intel Itanium instruction completes that uses register f2..f31 as a target register. This bit is sticky and only cleared by an explicit write of the user mask. When leaving the IA-32 instruction set, PSR.mfl is set to 1 if PSR.dfl is 0, otherwise PSR.mfl is unmodified.	unchanged	data <sup>a</sup>

**Table 3-2. Processor Status Register Fields (Continued)**

Field	Bits	Description	Interruption State	Serialization Required
mfh	5	Upper (f32 .. f127) floating-point registers written – This bit is set to one when an Intel Itanium instruction completes that uses register f32..f127 as a target register. This bit is sticky and only cleared by an explicit write of the user mask. PSR.mfh is unmodified by IA-32 instruction set execution.	unchanged	data <sup>a</sup>
System Mask = PSR{23:0}				
ic	13	Interrupt Collection – When 1 and an interruption occurs, the current state of the processor is loaded in IIP, IPSR, IIM and IFS; and additional registers defined in “ <a href="#">Interrupt Vector Descriptions</a> ” on page 2:165. When 0, IIP, IPSR, IIM and IFS are not modified on an interruption (see <a href="#">Table 8-1, “Writing of Interruption Resources by Vector”</a> on page 2:166 for details). When 0, speculative load exceptions result in deferred exception behavior, regardless of the state of the DCR and ITLB deferral bits. Processor operation is undefined if PSR.ic is 0 and a transition is made to execute IA-32 code.	0	inst/data <sup>c</sup>
i	14	Interrupt Bit – When 1 and executing Intel Itanium instructions, unmasked pending external interrupts will interrupt the processor by transferring control to the external interrupt handler. When 0, pending external interrupts do not interrupt the processor. The effect of clearing PSR.i via Reset System Mask ( $r_{sm}$ ) instructions is observed by the next instruction. Toggling PSR.i from one to zero via Move to PSR.I requires data serialization. When executing IA-32 instructions, external interrupts are enabled if PSR.i and (CFLG.if is 0 or EFLAG.if is 1). NMI interrupts are enabled if PSR.i is 1 regardless of EFLAG.if.	0	clear: implicit serialization set: data <sup>d</sup>
pk	15	Protection Key enable – When 1 and PSR.it is 1, instruction references (including IA-32) check for valid protection keys. When 1 and PSR.dt is 1, data references (including IA-32) check for valid protection keys. When 1 and PSR.rt is 1, protection key checks are enabled for register stack references. When 0, neither instruction, data, nor register stack references are checked for valid protection keys. When PSR.dt, PSR.rt or PSR.it are 0, PSR.pk is ignored for the corresponding reference.	unchanged	inst/data <sup>e</sup>
rv	12:6, 16	reserved		
dt	17	Data address Translation – When 1, virtual data addresses are translated and access rights checked. When 0, data accesses use physical addressing. PSR.dt must be 1 when entering IA-32 code, otherwise processor operation is undefined.	unchanged/0 <sup>j</sup>	inst/data <sup>c</sup>
dfl	18	Disabled Floating-point Low register set – When 1, a read or write access to f2 through f31 results in a Disabled Floating-Point Register fault. When 1, all IA-32 FP, Intel SSE and Intel MMX technology instructions raise a Disabled FP Register fault (regardless whether the instruction actually references f2-31).	0	data

**Table 3-2. Processor Status Register Fields (Continued)**

Field	Bits	Description	Interruption State	Serialization Required
dfh	19	Disabled Floating-point High register set – When 1, a read or write access to f32 through f127 results in a Disabled Floating-Point Register fault. When 1, a Disabled FP Register fault is raised on the first IA-32 target instruction following a <code>br.ia</code> or <code>rfi</code> , regardless whether f32-127 are referenced.	0	data
sp	20	Secure Performance monitors – Controls the ability of non-privileged code (including IA-32 code) to read non-privileged performance monitors. See <a href="#">Table 7-5 on page 2:158</a> for values returned by PMD read instructions. Also, when 0, PSR.up can be modified by user mask instructions; otherwise, PSR.up is unchanged by user mask instructions. When 1 or CFLAG.pce is 0, non-privileged IA-32 performance monitor reads (via <code>rdpmc</code> ) raise an <code>IA_32_Exception(GPFault)</code> .	0	data
pp	21	Privileged Performance monitor enable – When 1, monitors configured as privileged monitors are enabled to count events (including IA-32 events). When 0, privileged monitors are disabled. See <a href="#">“Performance Monitoring” on page 2:155</a> for details.	DCR.pp	inst/data <sup>e</sup>
di	22	Disable Instruction set transition – When 1, attempts to switch instruction sets via the IA-32 <code>jmppe</code> or <code>br.ia</code> instructions results in a Disabled Instruction Set Transition fault. This bit doesn't restrict instruction set transitions due to interruptions or <code>rfi</code> .	0	data
si	23	Secure Interval timer – When 1, the Interval Time Counter (ITC) register and the Resource Utilization Counter (RUC) are readable only by privileged code; non-privileged reads result in a Privileged Register fault. When 0, ITC and RUC are readable at any privilege level. System software can secure the ITC from non-privileged IA-32 access by setting either PSR.si or CFLAG.tsd to 1. When secured, an IA-32 <code>rdtsc</code> (read time stamp counter) instruction at any privilege level other than the most privileged raises an <code>IA_32_Exception(GPfault)</code>	0	data
PSR.I = PSR{31:0}				
db	24	Debug Breakpoint fault – When 1, data and instruction address breakpoints are enabled and can cause an Data/Instruction Debug fault. When 1, IA-32 instruction address breakpoints are enabled and can cause an <code>IA_32_Exception(Debug)</code> fault. When 1, IA-32 data address breakpoints are enabled and can cause an <code>IA_32_Exception(Debug) Trap</code> . When 0, address breakpoint faults and traps are disabled.	0	inst/data <sup>e</sup>
lp	25	Lower Privilege transfer trap – When 1, a Lower Privilege Transfer trap occurs whenever a taken branch lowers the current privilege level (numerically increases). This bit is ignored during IA-32 instruction set execution.	0	data

**Table 3-2. Processor Status Register Fields (Continued)**

Field	Bits	Description	Interruption State	Serialization Required
tb	26	Taken Branch trap – When 1, the successful completion of a taken branch results in a Taken Branch trap. <i>rfi</i> and interruptions can not raise a Taken Branch trap. When 1, successful completion of a taken IA-32 branch results in an IA_32_Exception(Debug) trap.	0	data
rt	27	Register stack Translation – When 1, register stack accesses are translated and access rights are checked. When 0, register stack accesses use physical addressing. PSR.dt is ignored for register stack accesses. The register stack engine must be in enforced lazy mode (RSC.mode = 00) when modifying this bit; otherwise, processor behavior is undefined. During IA-32 instruction execution this bit is ignored and the register stack is disabled.	unchanged	data
rv	31:28	reserved		
PSR{63:0}				
cpl <sup>f</sup>	33:32	Current Privilege Level –The current privilege level of the processor (including IA-32). Controls accessibility to system registers, instructions and virtual memory pages. A value of 0 is most privileged, a value of 3 is least privileged. Written by the <i>rfi</i> , <i>epc</i> , and <i>br.ret</i> instructions. PSR.cpl is unchanged by the <i>jmp</i> e and <i>br.ia</i> instructions. PSR.cpl cannot be updated by any IA-32 instructions.	0	rfi <sup>g</sup>
is	34	Instruction Set – When 0, Intel Itanium instructions are executing. When 1, IA-32 instructions are executing. Written by the <i>rfi</i> and <i>br.ia</i> instructions and the IA-32 <i>jmp</i> e instruction.	0	rfi <sup>g</sup> , br.ia <sup>h</sup>
mc	35	Machine Check abort mask – When 1, machine check aborts are masked. When 0, machine check aborts can be delivered (including IA-32 instruction set execution). Processor operation is undefined if PSR.mc is 1 and a transition is made to execute IA-32 code.	unchanged/1 <sup>i</sup>	rfi <sup>g</sup>
it	36	Instruction address Translation – When 1, virtual instruction addresses are translated and access rights checked. When 0, instruction accesses use physical addressing. PSR.it must be 1 when entering IA-32 code, otherwise processor operation is undefined.	unchanged/0 <sup>j</sup>	rfi <sup>g</sup>
id	37	Instruction Debug fault disable – When 1, Instruction Debug faults are disabled on the first restart instruction in the current bundle. <sup>k</sup> When PSR.id is 1 or EFLAG.rf is 1, IA-32 instruction debug faults are disabled for one IA-32 instruction. PSR.id and EFLAG.rf are set to 0 after the successful execution of each IA-32 instruction.	0	rfi <sup>g</sup>
da	38	Disable Data Access and Dirty-bit faults – When 1, Data Access and Dirty-Bit faults are disabled on the first restart instruction in the current bundle or for the first mandatory RSE reference following the <i>rfi</i> . <sup>k</sup> IA-32 Access/Dirty-bit faults are not affected by PSR.da. <sup>l</sup>	0	rfi <sup>g</sup>
dd	39	Data Debug fault disable – When 1, Data Debug faults are disabled on the first restart instruction in the current bundle or for the first mandatory RSE reference. <sup>k</sup> IA-32 Data Debug traps are not affected by PSR.dd. <sup>l</sup>	0	rfi <sup>g</sup>

**Table 3-2. Processor Status Register Fields (Continued)**

Field	Bits	Description	Interruption State	Serialization Required
ss	40	Single Step enable – When 1, a Single Step trap occurs following the successful execution of the first restart instruction in the current bundle. Instruction slots 0, 1, and 2 can be single stepped. When 1 or EFLAG.tf is 1, an IA_32_Exception(Debug) trap is taken after each IA-32 instruction.	0	rff <sup>g</sup>
ri	42:41	Restart Instruction – Set on an interruption, indicating the next instruction in the bundle to be executed. When the next instruction is the L+X instruction of an MLX, this field is set to the value 1. When restarting instructions with <i>rfi</i> , this field in IPSR specifies which instruction(s) in the bundle are restarted. The specified and subsequent instructions are restarted, all instructions prior to the restart point are ignored. 0 – restart execution at instruction slot 0 1 – restart execution at instruction slot 1 2 – restart execution at instruction slot 2 3 – reserved Except at an interruption and for the first restart instruction following an <i>rfi</i> , the value of this field is undefined. This field is set to 0 after any interruption from the IA-32 instruction set and is ignored when IA-32 instructions are restarted.	instruction pointer	rff <sup>g</sup>
ed	43	Exception Deferral – When 1, if the first restart instruction in the current bundle is a speculative load, the operation is forced to indicate a deferred exception by setting the load target register to NaT or NaTval. No memory references are performed, however any address post increments are performed. If the operation is a speculative advanced load, the ALAT entry corresponding to the load address and target register is purged. If the operation is an <i>lfetch</i> instruction, memory promotion is not performed, however any address post increments are performed. When 0, exception deferral is not forced on restarted speculative loads. If the first restart instruction is not a speculative load or <i>lfetch</i> instruction, this bit is ignored. <sup>kl</sup>	0	rff <sup>g</sup>
bn	44	register Bank – When 1, registers GR16 to GR31 for bank 1 are accessible. When 0, registers GR16 to GR31 for bank 0 are accessible. Written by <i>rfi</i> and <i>bsw</i> instructions.	0	implicit <sup>m</sup>
ia	45	Disable Instruction Access-bit faults – When 1, Instruction Access-Bit faults are disabled on the first restart instruction in the current bundle. <sup>k</sup> IA-32 Access-bit faults are not affected by PSR.ia. <sup>l</sup>	0	rff <sup>g</sup>
vm	46	Virtual Machine – When 1, an attempt to execute certain instructions results in a Virtualization fault. Implementation of this bit is optional. If the bit is not implemented, it is treated as a reserved bit. Written by the <i>rfi</i> and <i>vmsw</i> instructions.	0	rfi, vmsw: implicit <sup>n</sup>
rv	63:47	reserved		



- a. User mask bits are implicitly serialized if accessed via user mask instructions; `sum`, `rum`, and `move` to User Mask. If modified with system mask instructions; `rsm`, `ssm` and `move` to PSR.I, software must explicitly serialize to ensure side effects are observed before dependent instructions.
- b. User mask modification serialization is implicit only for monitoring data execution events. Software should issue instruction serialization operations before monitoring instruction events to achieve better accuracy.
- c. Requires instruction serialization to guarantee that VHPT walks initiated on behalf of an instruction reference observe the new value of this bit. Otherwise, data serialization is sufficient to guarantee that the new value is observed.
- d. The effect of masking external interrupts with `rsm` is observed by the next instruction. However, the processor does not ensure unmasking interrupts with `ssm` is immediately observed. Software can issue a data serialization operation to ensure the effects of setting PSR.i are observed before a given point in program execution.
- e. Requires instruction or data serialization, based on whether the dependent “use” is an instruction fetch access or data access.
- f. CPL can be modified due to interruptions, Return From Interruption (`rfi`), Enter Privilege Code (`epc`), and Branch Return (`br.ret`) instructions.
- g. Can only be modified by the Return From Interruption (`rfi`) instruction. `rfi` performs an explicit instruction and data serialization operation.
- h. Modification of the PSR.is bit by a `br.ia` instruction set is implicitly instruction serialized.
- i. PSR.mc is set to 1 after a machine check abort or INIT; otherwise, unmodified on interruptions.
- j. After an interruption this bit is normally unchanged, however after a PAL-based interruption this bit is set to 0.
- k. This bit is set to 0 after the successful execution of each instruction in a bundle except for `rfi` which may set it to 1.
- l. This bit is ignored when restarting IA-32 instructions and set to zero when `br.ia` or `rfi` successfully complete and before the first IA-32 instruction starts execution.
- m. After an interruption, `rfi`, or `bsw` the processor ensures register accesses are made to the new register bank. For interruptions, `rfi` and `bsw`, the processor ensures all register accesses and outstanding loads prior to the bank switch operate on the prior register bank.
- n. Can be modified by the Return From Interruption (`rfi`) and Virtual Machine Switch (`vmsw`) instructions. `rfi` performs an explicit instruction and data serialization operation. Modification of PSR.vm bit by the `vmsw` instruction is implicitly serialized.

### 3.3.3 Control Registers

Table 3-3 defines all registers in the control register name space along with serialization requirements to ensure side effects are observed by subsequent instructions. However, reads of a control register must be data serialized with prior writes to the same register. The serialization required column only refers to the side effects of the data value.

Writes to read-only registers (IVR, IRR0-3) result in an Illegal Operation fault, accesses to reserved registers result in a Privileged Operation fault. Accesses can only be performed by `mov` to/from instructions defined in Table 3-4 at privilege level 0; otherwise, a Privileged Operation fault is raised.

**Table 3-3. Control Registers**

	Register	Name	Description	Serialization Required
Global Control Registers	CR0	DCR	Default Control Register	inst/data
	CR1	ITM	Interval Timer Match register	data <sup>a</sup>
	CR2	IVA	Interrupt Vector Address	inst <sup>a</sup>
	CR3		reserved	
	CR4	ITO	Interval Timer Offset Register	data <sup>a</sup>
	CR5-7		reserved	
	CR8	PTA	Page Table Address	inst/data <sup>b</sup>
	CR9-15		reserved	

**Table 3-3. Control Registers (Continued)**

	Register	Name	Description	Serialization Required
Interruption Control Registers	CR16	IPSR	Interruption Processor Status Register	implied <sup>d</sup>
	CR17	ISR	Interruption Status Register	implied <sup>c</sup>
	CR18		reserved	
	CR19	IIP	Interruption Instruction Pointer	implied <sup>d</sup>
	CR20	IFA	Interruption Faulting Address	implied <sup>d</sup>
	CR21	ITIR	Interruption TLB Insertion Register	implied <sup>d</sup>
	CR22	IIPA	Interruption Instruction Previous Address	implied <sup>c</sup>
	CR23	IFS	Interruption Function State	implied <sup>d,e</sup>
	CR24	IIM	Interruption Immediate register	implied <sup>c</sup>
	CR25	IHA	Interruption Hash Address	implied <sup>c</sup>
	CR26	IIB0	Interruption Instruction Bundle 0	implied <sup>c</sup>
	CR27	IIB1	Interruption Instruction Bundle 1	implied <sup>c</sup>
	Reserved	CR28-63		reserved
Interrupt Control Registers	CR64	LID	Local Interrupt ID	data <sup>a</sup>
	CR65	IVR	External Interrupt Vector Register (read only)	data <sup>a</sup>
	CR66	TPR	Task Priority Register	data <sup>a</sup>
	CR67	EOI	End Of External Interrupt	data <sup>a</sup>
	CR68	IRR0	External Interrupt Request Register 0 (read only)	data <sup>a</sup>
	CR69	IRR1	External Interrupt Request Register 1 (read only)	data <sup>a</sup>
	CR70	IRR2	External Interrupt Request Register 2 (read only)	data <sup>a</sup>
	CR71	IRR3	External Interrupt Request Register 3 (read only)	data <sup>a</sup>
	CR72	ITV	Interval Timer Vector	data <sup>a</sup>
	CR73	PMV	Performance Monitoring Vector	data <sup>a</sup>
	CR74	CMCV	Corrected Machine Check Vector	data <sup>a</sup>
	CR75-79		reserved	reserved
	CR80	LRR0	Local Redirection Register 0	data <sup>a</sup>
	CR81	LRR1	Local Redirection Register 1	data <sup>a</sup>
	Reserved	CR82-127		reserved

- Serialization is needed to ensure external interrupt masking, new interval timer match values or new interruption table addresses are observed before a given point in program execution.
- Serialization is needed to ensure new values in PTA are visible to the hardware Virtual Hash Page Table (VHPT) walker before a dependent instruction fetch or data access.
- These registers are modified by the processor on an interruption or by an explicit move to these registers. There are no side effects when written.
- These registers are implied operands to the rfi and/or TLB insert instructions. The processor ensures writes in previous instruction groups are observed by rfi and/or TLB insert instructions in subsequent instruction groups. These registers are also modified by the processor on an interruption, subsequent reads return the results of the interruption. There are no other side effects.
- IFS written by a `cover` instruction followed by a move-from IFS is implicitly serialized.

**Table 3-4. Control Register Instructions**

Mnemonic	Description	Operation	Format
<code>mov cr<sub>3</sub> = r<sub>2</sub></code>	Move to control register	$CR[r_3] \leftarrow GR[r_2]$	M
<code>mov r<sub>1</sub> = cr<sub>3</sub></code>	Move from control register	$GR[r_1] \leftarrow CR[r_3]$	M

**Table 3-4. Control Register Instructions (Continued)**

Mnemonic	Description	Operation	Format
srlz.i, rfi	Serialize instruction references	Ensure side effects are observed by the instruction fetch stream	M
srlz.d	Serialize data references	Ensure side effects are observed by the execute and data streams	M

### 3.3.4 Global Control Registers

#### 3.3.4.1 Default Control Register (DCR – CR0)

The DCR specifies default parameters for PSR values on interruption, some additional global controls, and whether speculative load faults can be deferred. Figure 3-3 and Table 3-5 define and describe the DCR fields.

**Figure 3-3. Default Control Register (DCR – CR0)**



**Table 3-5. Default Control Register Fields**

Field	Bit	Description	Serialization Required
pp	0	Privileged Performance monitor default – On interruption, DCR.pp is loaded into PSR.pp.	data
be	1	Big-Endian default – When 1, Virtual Hash Page Table (VHPT) walker accesses are performed big-endian; otherwise, little-endian. On interruption, DCR.be is loaded into PSR.be.	inst
lc	2	IA-32 Lock Check enable – When 1, and an IA-32 atomic memory reference is defined as requiring a read-modify-write operation external to the processor under an external bus lock, an IA_32_Intercept(Lock) is raised. (IA-32 atomic memory references are defined to require an external bus lock for atomicity when the memory transaction is made to non-write-back memory or are unaligned across an implementation-specific non-supported alignment boundary.) When 0, and an IA-32 atomic memory reference is defined as requiring a read-modify-write operation external to the processor under external bus lock, the processor may either execute the transaction as a series of non-atomic transactions or perform the transaction with an external bus lock, depending on the processor implementation. Intel Itanium semaphore accesses ignore this bit. All unaligned Intel Itanium semaphore references generate an Unaligned Data Reference fault. All aligned Intel Itanium semaphore references made to memory that is neither write-back cacheable nor a NaTPage result in an Unsupported Data Reference fault.	data
dm	8	Defer TLB Miss faults only (VHPT data, Data TLB, and Alternate Data TLB faults) – When 1, and a TLB miss is deferred, lower priority Debug faults may still be delivered. A TLB miss fault, deferred or not, precludes concurrent Page not Present, Key Miss, Key Permission, Access Rights, or Access Bit faults. This bit is ignored by IA-32 instructions.	data
dp	9	Defer Page not Present faults only – When 1, and a Page not Present fault is deferred, lower priority Debug faults may still be delivered. A Page not Present fault, deferred or not, precludes concurrent Key Miss, Key Permission, Access Rights, or Access Bit faults. This bit is ignored by IA-32 instructions.	data

**Table 3-5. Default Control Register Fields (Continued)**

Field	Bit	Description	Serialization Required
dk	10	Defer Key Miss faults only – When 1, and a Key Miss fault is deferred, lower priority Access Bit, Access Rights or Debug faults may still be delivered. A Key Miss fault, deferred or not, precludes concurrent Key Permission faults. This bit is ignored by IA-32 instructions.	data
dx	11	Defer Key Permission faults only – When 1, and a Key Permission fault is deferred, lower priority Access Bit, Access Rights or Debug faults may still be delivered. This bit is ignored by IA-32 instructions.	data
dr	12	Defer Access Rights faults only – When 1, and an Access Rights fault is deferred, lower priority Access Bit or Debug faults may still be delivered. This bit is ignored by IA-32 instructions.	data
da	13	Defer Access Bit faults only – When 1, and an Access Bit fault is deferred, lower priority Debug faults may still be delivered. This bit is ignored by IA-32 instructions.	data
dd	14	Defer Debug faults – When 1, Data Debug faults on speculative loads are deferred. This bit is ignored by IA-32 instructions.	data
rv	7:3, 63:15	reserved	reserved

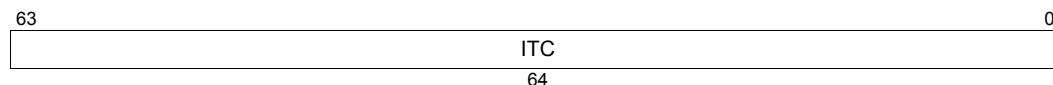
For the DCR exception deferral bits, when the bit is 1, and a speculative load results in the specified fault condition, and the speculative load’s code page exception deferral bit (ITLB.ed) is 1, the exception is deferred by setting the speculative load target register to NaT or NaTVal. Otherwise, the specified fault is taken on the speculative load. For a description of faults on speculative loads see [“Deferral of Speculative Load Faults” on page 2:105](#).

Since DCR.be also controls byte ordering of VHPT references that are the result of instruction misses, DCR.be requires instruction serialization. Other DCR bits require data serialization only.

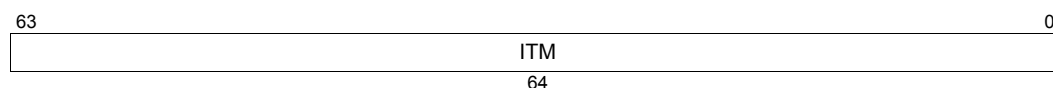
### 3.3.4.2 Interval Time Counter and Match Register (ITC – AR44 and ITM – CR1)

The Interval Time Counter (ITC) and Interval Timer Match (ITM) register support elapsed time notification, see [Figure 3-4](#) and [Figure 3-5](#).

**Figure 3-4. Interval Time Counter (ITC – AR44)**



**Figure 3-5. Interval Timer Match Register (ITM – CR1)**



The ITC is a free-running 64-bit counter that counts up at a fixed relationship to the input clock to the processor. The ITC may be clocked at a somewhat lower frequency than the instruction execution frequency. This clocking relationship is described in the PAL procedure PAL\_FREQ\_RATIOS on page 2:393. The ITC is guaranteed to be clocked at a constant rate, even if the instruction execution frequency may vary. The ITC counting rate is not affected by power management mechanisms.

A sequence of reads of the ITC is guaranteed to return ever-increasing values (except for the case of the counter wrapping back to 0) corresponding to the program order of the reads. Applications can directly sample the ITC for time-based calculations.

A 64-bit overflow condition can occur without notification. The ITC can be read at any privilege level if PSR.si is zero. The timer can be secured from non-privileged access by setting PSR.si to one. When secured, a read of the ITC by non-privileged code results in a Privileged Register fault. Writes to the ITC can only be performed at privilege level 0; otherwise, a Privileged Register fault is raised.

The IA-32 Time Stamp Counter (TSC) is similar to ITC. The ITC can be read by the IA-32 `rdtsc` (read time stamp counter) instruction. System software can secure the ITC from non-privileged IA-32 access by setting either PSR.si or CFLG.tsd to 1. When secured, an IA-32 read of the ITC at any privilege level other than the most privileged raises an `IA_32_Exception(GPfault)`.

When the value in the ITC is equal to the value in the ITM an Interval Timer Interrupt is raised. Once the interruption is taken by the processor and serviced by software, the ITC may not necessarily be equal to the ITM. The ITM is accessible only at privilege level 0; otherwise, a Privileged Operation fault is raised.

The interval counter can be written, for initialization purposes, by privileged code. The ITC is not architecturally guaranteed to be synchronized with any other processor's interval time counter in a multiprocessor system, nor is it synchronized with the wall clock. Software must calibrate interval timer ticks to wall clock time and periodically adjust for drift. In a multiprocessor system, a processor's ITC is not architecturally guaranteed to be clocked synchronously with the ITC's on other processors, and may not be clocked at the same nominal clock rate as ITC's on other processors. The platform firmware provides information on the clocking of processors in a multiprocessor system.

Modification of the ITC or ITM is not necessarily serialized with respect to instruction execution. Software can issue a data serialization operation to ensure the ITC or ITM updates and possible side effects are observed by a given point in program execution. Software must accept a level of sampling error when reading the interval timer due to various machine stall conditions, interruptions, bus contention effects, etc. Please see the processor-specific documentation for further information on the level of sampling error of the Itanium processor.

### **3.3.4.3 Resource Utilization Counter (RUC – AR45)**

The Resource Utilization Counter (RUC) is a 64-bit counter that counts up at a fixed relationship to the input clock to the processor, when the processor is active. Processors may be inactive due to hardware multi-threading. Virtual processors may be inactive when not scheduled to run by the VMM. (See [Section 11.7, "PAL Virtualization Support" on page 2:324](#) for details on virtual processors.)

The RUC is clocked such that, in a given time interval, the difference in the RUC values for all of the logical or virtual processors on a given physical processor add up to approximately the difference seen in the ITC on that physical processor for that same interval.

A sequence of reads of the RUC is guaranteed to return ever-increasing values (except for the case of the counter wrapping back to 0) corresponding to the program order of the reads. Applications can directly sample the RUC for active-running-time calculations.

A 64-bit overflow condition can occur without notification. The RUC can be read at any privilege level if PSR.si is zero. The timer can be secured from non-privileged access by setting PSR.si to one. When secured, a read of the RUC by non-privileged code results in a Privileged Register fault. Writes to the RUC can only be performed at privilege level 0; otherwise, a Privileged Register fault is raised.

Modification of the RUC is not necessarily serialized with respect to instruction execution. Software can issue a data serialization operation to ensure the RUC updates are observed by a given point in program execution. Software must accept a level of sampling error when reading the resource utilization counter due to various machine stall conditions, interruptions, bus contention effects, etc. Please see the processor-specific documentation for further information on the level of sampling error of the Itanium processor.

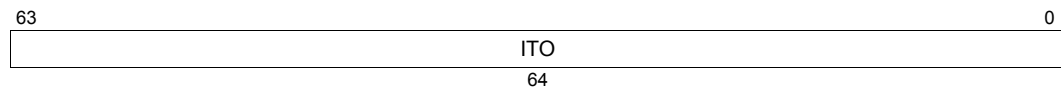
RUC should only be written by Virtual Machine Monitors; other Operating Systems should not write to RUC, but should only read it.

The RUC register is not supported on all processor implementations. Software can check CPUID register 4 to determine the availability of this feature. The RUC register is reserved when this feature is not supported.

#### 3.3.4.4 Interval Timer Offset (ITO – CR4)

The Interval Timer Offset (ITO) register allows virtual machine monitors to specify an offset to the Interval Timer Counter (ITC) for the virtual processor. The layout of the register is shown in [Figure 3-6](#). For details of the usage of this register in virtual environment, please refer to [Section 11.7.4.1.3, “Guest MOV-from-AR.ITC Optimization”](#) on page 2:337.

**Figure 3-6. Interval Timer Offset Register (ITO – CR4)**



The ITO register has no effects on instruction execution when PSR.vm is 0.

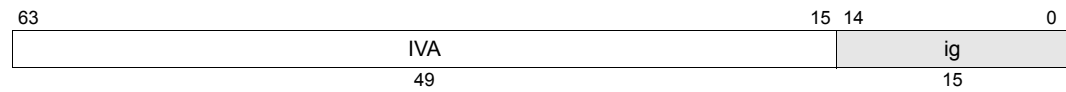
The ITO register does not affect the generation of interval timer interrupts, discussed in [Section 3.3.4.2, “Interval Time Counter and Match Register \(ITC – AR44 and ITM – CR1\)”](#).

The ITO register is not supported on all processor implementations. Software can call either PAL\_PROC\_GET\_FEATURES or PAL\_VP\_ENV\_INFO to determine the availability of this feature. The ITO register is reserved when this feature is not supported.

### 3.3.4.5 Interruption Vector Address (IVA – CR2)

The IVA specifies the location of the interruption vector table in the virtual address space, or the physical address space if PSR.it is 0, see [Figure 3-7](#). The size of the vector table is 32K bytes and is 32K byte aligned. The lower 15 bits of the IVA are ignored when written, reads return zeros. All upper 49 address bits of IVA must be implemented regardless of the size of the physical and virtual address space. If an unimplemented virtual or physical address (see [“Unimplemented Address Bits” on page 2:73](#)) is loaded into IVA, and an interruption occurs, processor behavior is unpredictable. See [“IVA-based Interruption Vectors” on page 2:113](#) for a description of an interruption table layout.

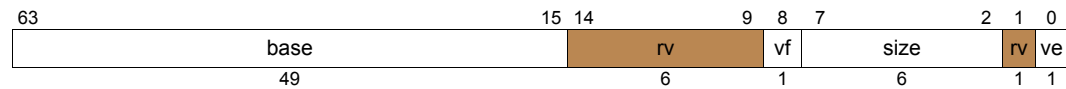
**Figure 3-7. Interruption Vector Address (IVA – CR2)**



### 3.3.4.6 Page Table Address (PTA – CR8)

The PTA anchors the Virtual Hash Page Table (VHPT) in the virtual address space. See [“Virtual Hash Page Table \(VHPT\)” on page 2:61](#) for a complete definition of the VHPT. Operating systems must ensure that the table is aligned on a natural boundary; otherwise, processor operation is undefined. See [Figure 3-8](#) and [Table 3-6](#) for the PTA field definitions.

**Figure 3-8. Page Table Address (PTA – CR8)**



**Table 3-6. Page Table Address Fields**

Field	Bits	Description
ve	0	VHPT Enable – When 1, the processor is enabled to walk the VHPT.
size	7:2	VHPT Size – VHPT table size in power of 2 increments, table size is $2^{\text{size}}$ bytes. Size generates a mask that is logically AND'ed with the result of the VHPT hash function. Minimum VHPT table size is 32K bytes; otherwise, a Reserved Register/Field fault is raised (see <a href="#">“Virtual Hash Page Table (VHPT)” on page 2:61</a> ). The maximum size is $2^{61}$ bytes for long format VHPTs, and $2^{52}$ bytes for short format VHPTs.
vf	8	VHPT Format – When 0, 8-byte short format entries are used, when 1, 32-byte long format entries are used.
base	63:15	VHPT Base virtual address – Defines the starting virtual address of the VHPT table. Base is logically OR'ed with the hash index produced by the VHPT hash function when referencing the VHPT. Base must be on $2^{\text{size}}$ boundary otherwise processor operation is undefined. All base address bits of PTA must be implemented regardless of the size of the physical and virtual address space. If an unimplemented virtual address (see <a href="#">“Unimplemented Address Bits” on page 2:73</a> ) is used by the processor as a page table base, all VHPT walks generate an Instruction/Data TLB miss (see <a href="#">“Translation Searching” on page 2:69</a> ).
rv	1, 14:9	reserved

### 3.3.5 Interruption Control Registers

Registers CR16 - CR27 record information at the time of an interruption (including from the IA-32 instruction set) and are used by handlers to process the interruption.

The interruption control registers can only be read or written while PSR.ic is 0; otherwise, an Illegal Operation fault is raised. These registers are only guaranteed to retain their values when PSR.ic is 0. When PSR.ic is 1, the processor does not preserve their contents.

The contents of the interruption control registers are defined only when the PSR.ic bit is cleared by an interruption. If the PSR.ic bit is explicitly cleared (e.g., by using `rsm`, or `mov` to PSR), then the contents of these registers are undefined. If the PSR.ic bit is explicitly set (e.g., by using `ssm`, or `mov` to PSR), then the contents of these registers are undefined until the PSR.ic bit has been serialized and an interruption occurs.

IIPA has special behavior in case of an `rfi` to a fault. Refer to “[Interruption Instruction Previous Address \(IIPA – CR22\)](#)” on page 2:40.

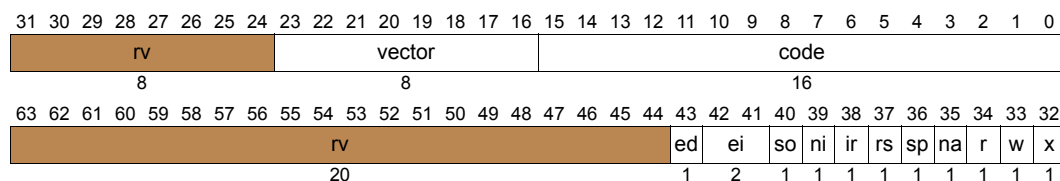
#### 3.3.5.1 Interruption Processor Status Register (IPSR – CR16)

On an interruption and if PSR.ic is 1, the IPSR receives the value of the PSR. The IPSR, IIP and IFS are used to restore processor state on a Return From Interruption (`rfi`). The IPSR has the same format as PSR, see “[Processor Status Register \(PSR\)](#)” on page 2:23 for details.

#### 3.3.5.2 Interruption Status Register (ISR – CR17)

The ISR receives information related to the nature of the interruption, and is written by the processor on all interruption events regardless of the state of PSR.ic, except for Data Nested TLB faults. The ISR contains information about the excepting instruction and its properties such as whether it was doing a read, write, execute, speculative, or non-access operation, see [Figure 3-9](#) and [Table 3-7](#). Multiple bits may be concurrently set in the ISR, for example, a faulting semaphore operation will set both ISR.r and ISR.w, and faults on speculative loads will set ISR.sp and ISR.r. Additional fault- or trap-specific information is available in ISR.code and ISR.vector. Refer to [Section 8.2, “ISR Settings”](#) for complete definition of the ISR field settings.

**Figure 3-9. Interruption Status Register (ISR – CR17)**





**Table 3-7. Interruption Status Register Fields**

Field	Bits	Description
code	15:0	Interruption Code – 16 bit code providing additional information specific to the current interruption. For IA-32 specific exceptions and software interrupts, contains the IA-32 interruption error code or zero.
vector	23:16	IA-32 exception/interception vector number. For IA-32 exceptions and software interrupts, contains the IA-32 vector number (e.g., GPFault has a vector number of 13). See <a href="#">Chapter 9, “IA-32 Interruption Vector Descriptions”</a> for details.
x	32	Execute exception – Interruption is associated with an instruction fetch (including IA-32).
w	33	Write exception – Interruption is associated with a write operation. Both ISR.r and ISR.w are set for IA-32 read-modify-write instructions.
r	34	Read exception – Interruption is associated with a read operation. Both ISR.r and ISR.w are set for IA-32 read-modify-write instructions.
na	35	Non-access exception – See <a href="#">Section 5.5.2, “Non-access Instructions and Interruptions”</a> on page 2:103. This bit is always 0 for interruptions taken in the IA-32 instruction set.
sp	36	Speculative load exception – Interruption is associated with a speculative load instruction. This bit is always 0 for interruptions taken in the IA-32 instruction set.
rs	37	Register Stack – Interruption is associated with a mandatory RSE fill or spill. This bit is always 0 for interruptions taken in the IA-32 instruction set.
ir	38	Incomplete Register frame – The current register frame is incomplete when the interruption occurred. This bit is always 0 for interruptions taken in the IA-32 instruction set.
ni	39	Nested Interruption – Indicates that PSR.ic was 0 or in-flight when the interruption occurred. This bit is always 0 for interruptions taken in the IA-32 instruction set.
so	40	IA-32 Supervisor Override – Indicates the fault occurred during an IA-32 instruction set supervisor override condition (the processor was performing a data memory accesses to the IDT, GDT, LDT or TSS segments) or an IA-32 data memory access at a privilege level of zero. This bit is always 0 for interruptions taken while executing Intel Itanium instructions.
ei	42:41	Excepting Instruction – 0 – exception due to instruction in slot 0 1 – exception due to instruction in slot 1 2 – exception due to instruction in slot 2 For faults and external interrupts, ISR.ei is equal to IPSR.ri. For traps, ISR.ei defines the slot of the excepting instruction. Traps on the L+X instruction of an MLX set ISR.ei to 2. This field is always 0 for interruptions taken in the IA-32 instruction set.
ed	43	Exception Deferral – this bit is set to the value of the TLB exception deferral bit (TLB.ed) for the instruction page containing the faulting instruction. If a translation does not exist or instruction translation is disabled, or if the interruption is caused by a mandatory RSE spill or fill, ISR.ed is set to 0. This bit is always 0 for interruptions taken in the IA-32 instruction set.
rv	31:24, 63:44	reserved

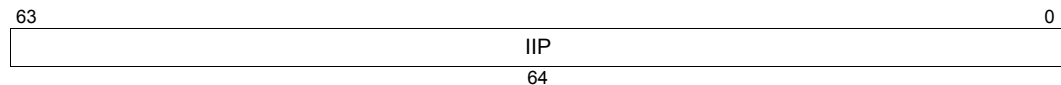
### 3.3.5.3 Interruption Instruction Bundle Pointer (IIP – CR19)

On an interruption and if PSR.ic is 1, the IIP receives the value of IP. IIP contains the virtual address (or physical if instruction translations are disabled) of the next instruction bundle or the IA-32 instruction to be executed upon return from the interruption. For IA-32 instruction addresses, IIP is zero extended to 64-bits and specifies a byte granular address. For traps and interrupts, IIP points to the next instruction to execute. For faults, IIP points to the faulting instruction. As shown in

Figure 3-10, all 64-bits of the IIP must be implemented regardless of the size of the physical and virtual address space supported by the processor model (see “Unimplemented Address Bits” on page 2:73). IIP also receives byte-aligned IA-32 instruction pointers. The IIP, IPSR and IFS are used to restore processor state on a Return From Interruption instruction (*rfi*). See “Interrupt Vector Descriptions” on page 2:165 for usages of the IIP.

An *rfi* to Itanium architecture-based code (IPSR.is is 0) ignores IIP{3:0}, an *rfi* to IA-32 code (IPSR.is is 1) ignores IIP{63:32}. Ignored bits are assumed to be zero.

**Figure 3-10. Interruption Instruction Bundle Pointer (IIP – CR19)**



Control transfers to unimplemented addresses (see “Unimplemented Address Bits” on page 2:73) result in an Unimplemented Instruction Address trap or fault. When the trap or fault is delivered, IIP is written as follows:

- If the trap is taken for an unimplemented virtual address, IIP is written in one of two ways, depending on the implementation: 1) IIP may be written with the implemented virtual address bits IP{63:61} and IP{IMPL\_VA\_MSB:0} only. Bits IIP{60:IMPL\_VA\_MSB+1} are set to IP{IMPL\_VA\_MSB}, i.e., sign-extended. 2) IIP may be written with the full, unimplemented virtual address from IP.
- If the trap is taken for an unimplemented physical address, IIP is written in one of two ways, depending on the implementation: 1) IIP may be written with the physical addressing memory attribute bit IP{63} and the implemented physical address bits IP{IMPL\_PA\_MSB:0} only. Bits IIP{62:IMPL\_PA\_MSB+1} are set to 0. 2) IIP may be written with the full, unimplemented physical address from IP.

When an *rfi* is executed with an unimplemented address in IIP (an unimplemented virtual address if IPSR.it is 1, or an unimplemented physical address if IPSR.it is 0), and an Unimplemented Instruction Address trap is taken, an implementation may optionally leave IIP unchanged (preserving the unimplemented address in IIP).

**Note:** Since IP{3:0} are always 0 when executing Itanium architecture-based code, IIP{3:0} will always be 0 when any interruption is taken from Itanium architecture-based code, with the exception of an Unimplemented Instruction Address trap on an *rfi*, where IIP may optionally be preserved as whatever value it held before executing the *rfi*.

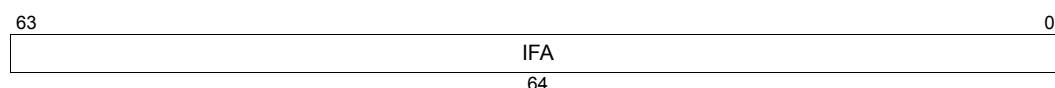
### 3.3.5.4 Interruption Faulting Address (IFA – CR20)

On an interruption and if PSR.ic is 1, the IFA receives the virtual address (or physical address if translations are disabled) that raised a fault. IFA reports the faulting address for both instruction and data memory accesses (including IA-32). For faulting data references (including IA-32), IFA points to the first byte of the faulting data memory operand. IFA reports a byte granular address. For faulting instruction references (including IA-32), IFA contains the 16-byte aligned bundle address (IFA{3:0} are zero) of the faulting instruction. For faulting IA-32 instructions, IIP points to the first byte of the IA-32 instruction, and is byte granular. In the event of an IA-32 instruction spanning a virtual page boundary, IA-32 instruction fetch faults are reported as either (1) for faults on the first page, IFA is set to the bundle address (IFA{3:0}=0) of the

faulting instruction and IIP points to the first byte of the faulting instruction, or (2) for faults on the second page, IFA contains the bundle address of the second virtual page and IIP points to the first byte of the faulting IA-32 instruction.

The IFA also specifies a translation’s virtual address when a translation entry is inserted into the instruction or data TLB. See [“Interrupt Vector Descriptions” on page 2:165](#) and [“Translation Insertion Format” on page 2:53](#) for usages of the IFA. As shown in [Figure 3-11](#), all 64-bits of the IFA must be implemented regardless of the size of the virtual and physical space supported by the processor model (see [“Unimplemented Address Bits” on page 2:73](#)). In some implementations, a mov to IFA instruction may raise an Unimplemented Data Address fault if an unimplemented virtual address is used.

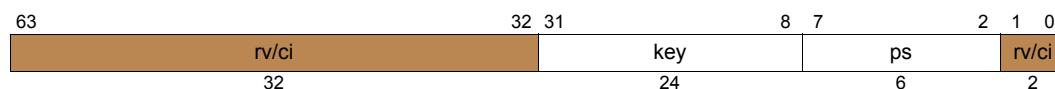
**Figure 3-11. Interruption Faulting Address (IFA – CR20)**



### 3.3.5.5 Interruption TLB Insertion Register (ITIR – CR21)

The ITIR receives default translation information from the referenced virtual region register on a virtual address translation fault. See [“Interrupt Vector Descriptions” on page 2:165](#) for the fault conditions that set the ITIR. The ITIR provides additional virtual address translation parameters on an insertion into the instruction or data TLB. See [“Translation Instructions” on page 2:60](#) for ITIR usage information. [Figure 3-12](#) and [Table 3-8](#) define the ITIR fields.

**Figure 3-12. Interruption TLB Insertion Register (ITIR)**



**Table 3-8. ITIR Fields**

Field	Bits	Description
rv/ci	63:32, 1:0	Reserved / Check on Insert – On a read these fields may return zeros or the value last written to them. If a non-zero value is written, a Reserved Register/Field fault may be raised on the mov to ITIR instruction. If not, a subsequent TLB insert will raise a Reserved Register/Field fault depending on other parameters to the insert. See <a href="#">“Translation Insertion Format” on page 2:53</a> . On an instruction or data translation fault, these fields are set to zero.
ps	7:2	Page Size – On a TLB insert, specifies the size of the virtual to physical address mapping. If an unsupported page size is written, a Reserved Register/Field fault may be raised on the mov to ITIR instruction. If not, a subsequent TLB insert will raise a Reserved Register/Field fault. See <a href="#">“Translation Insertion Format” on page 2:53</a> . On an instruction or data translation fault, this field is set to the accessed region’s page size (RR.ps).
key	31:8	Protection Key – On a TLB insert specifies a protection key that uniquely tags translations to a protection domain. If non-zero values are written to unimplemented protection key bits, a Reserved Register/Field fault may be raised on the mov to ITIR instruction. If not, a subsequent TLB insert will raise a Reserved Register/Field fault depending on other parameters to the insert. See <a href="#">“Translation Insertion Format” on page 2:53</a> . On an instruction or data translation fault, this field is set to the accessed Region Identifier (RR.rid).

### 3.3.5.6 Interruption Instruction Previous Address (IIPA – CR22)

For Itanium instructions, IIPA records the last successfully executed instruction bundle address. For IA-32 instructions, IIPA records the byte granular virtual instruction address zero extended to 64-bits of the faulting or trapping IA-32 instruction. In the case of a fault, IIPA does not report the address of the last successfully executed IA-32 instruction, but rather the address of the faulting IA-32 instruction. IIPA preserves bits 3:0 for byte aligned IA-32 instruction addresses.

The IIPA can be used by software to locate the address of the instruction bundle or IA-32 instruction that raised a trap or the instruction executed prior to a fault or interruption. In the case of a branch related trap, IIPA points to the instruction bundle which contained the branch instruction that raised the trap, while IIP points to the target of the branch.

When an instruction successfully executes without a fault, and the PSR.ic bit was 1 prior to instruction execution, it becomes the “last successfully executed instruction.” On interruptions, IIPA contains the address of the last successfully executed instruction bundle or IA-32 instruction, if PSR.ic was 1 prior to the interruption. Note that execution of an `rfi` instruction with PSR.ic equal to 0, but which sets PSR.ic to 1 does not update IIPA, since PSR.ic was zero prior to instruction execution.

When PSR.ic is one, accesses to IIPA cause an Illegal Operation fault. When PSR.ic is zero, IIPA is not updated by hardware and can be read and written by software. This permits low-level code to preserve IIPA across interruptions.

If the PSR.ic bit is explicitly cleared, e.g., by using `rsm`, then the contents of IIPA are undefined. Only when the PSR.ic bit is cleared by an interruption is the value of IIPA defined. It may point at the instruction which caused a trap, or at the instruction just prior to a faulting instruction, at an earlier instruction that became defined by some prior interruption, or by a move to IIPA instruction when PSR.ic was zero.

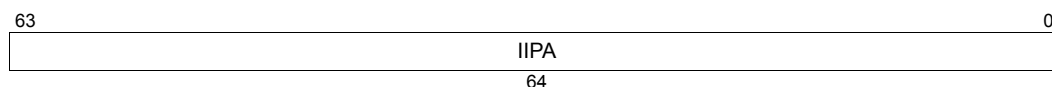
If the PSR.ic bit is explicitly set, e.g., by using `ssm`, then the contents of IIPA are undefined until the PSR.ic bit has been serialized and an interruption occurs.

During instruction set transitions the following boundary cases exist:

- On faults taken on the first IA-32 instruction after a `br.ia` or `rfi`, IIPA records the faulting IA-32 instruction address.
- On `br.ia` traps, IIPA records the address of the trapping instruction bundle.
- On faults taken on the first Itanium instruction after leaving the IA-32 instruction set, due to a `jmp` or interruption, IIPA contains the address of the `jmp` instruction or the interrupted IA-32 instruction.
- On `jmp` Data Debug, Single Step and Taken Branch traps, IIPA contains the address of the `jmp` instruction.

As shown in [Figure 3-13](#), all 64-bits of the IIPA must be implemented regardless of the size of the physical and virtual address space supported by the processor model (see “Unimplemented Address Bits” on page 2:73).

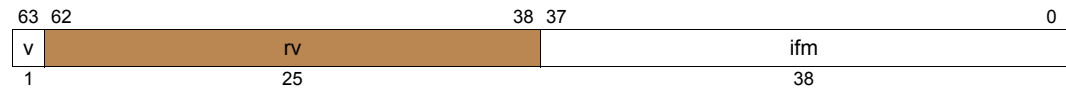
**Figure 3-13. Interruption Instruction Previous Address (IIPA – CR22)**



### 3.3.5.7 Interruption Function State (IFS – CR23)

The IFS register is used to reload the current register stack frame (CFM) on a Return From Interruption (*rfi*). If the IFS is accessed while PSR.ic is 1, an Illegal Operation fault is raised. The IFS can only be accessed at privilege level 0; otherwise, a Privileged Operation fault is raised. The IFS.v bit is cleared on interruption if PSR.ic is 1. All other fields are undefined after an interruption. If PSR.ic is 0, the *cover* instruction copies CFM to IFS.ifm and sets IFS.v to 1. See Figure 3-14 and Table 3-9 for the IFS field definitions.

**Figure 3-14. Interruption Function State (IFS – CR23)**



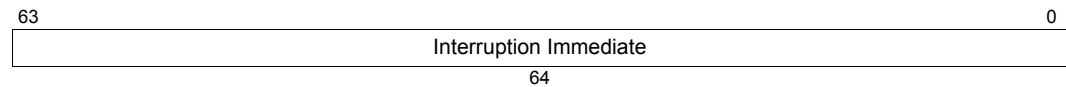
**Table 3-9. Interruption Function State Fields**

Field	Bits	Description
ifm	37:0	Interruption Frame Marker
v	63	Valid bit, cleared to 0 on interruption if PSR.ic is 1.
rv	62:38	reserved

### 3.3.5.8 Interruption Immediate (IIM – CR24)

If PSR.ic is 1, the IIM (Figure 3-15) records the zero-extended immediate field encoded in *chk.a*, *chk.s*, *fchkf* or *break* instruction faults. The *break.b* instruction always writes a zero value and ignores its immediate field. The IA\_32\_Intercept vector writes all 64-bits of IIM to indicate the cause of the intercept. See Table 8-1 on page 2:166 for the value of IIM in other situations. For the purpose of resource dependency, IIM is written as a result of the fault, not by the instruction itself.

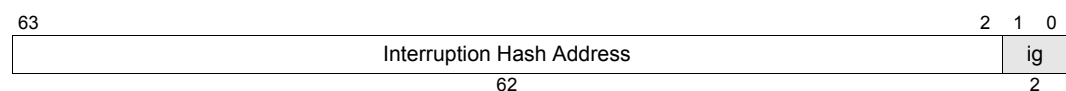
**Figure 3-15. Interruption Immediate (IIM – CR24)**



### 3.3.5.9 Interruption Hash Address (IHA – CR25)

The IHA (Figure 3-16) is loaded with the address of the Virtual Hash Page Table (VHPT) entry the processor referenced or would have referenced to resolve a translation fault. The IHA is written on interruptions by the processor when PSR.ic is 1. Refer to “VHPT Hashing” on page 2:65 for complete details. See Table 8-1 on page 2:166 for the value of IHA in other situations. All upper 62 address bits of IHA must be implemented regardless of the size of the virtual address space supported by the processor model (see “Unimplemented Address Bits” on page 2:73). The virtual address written to IHA by the processor is guaranteed to be an implemented virtual addresses on all processor models; however, if the address referenced by the VHPT is an unimplemented virtual address, the value of IHA is undefined.

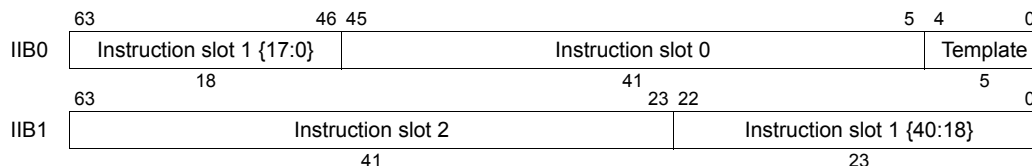
**Figure 3-16. Interruption Hash Address (IHA – CR25)**



### 3.3.5.10 Interruption Instruction Bundle Registers (IIB0-1 – CR26, 27)

On an interruption and if PSR.ic is 1, the IIB registers receive the 16-byte instruction bundle corresponding to the interruption. The bundle reported in the IIB registers is the bundle exactly as it was fetched for execution of the instruction which raised the interruption. [Figure 3-17](#) shows the format of the IIB0 and IIB1 registers. For details on instruction bundle format, see [Section 3.3, “Instruction Encoding Overview” on page 1:38](#).

**Figure 3-17. Interruption Instruction Bundle Registers (IIB0-1, – CR26, 27)**



If the interruption is a fault, the IIB registers record the instruction bundle pointed to by IIP. If the interruption is a trap, the IIB registers record the instruction bundle pointed to by IIPA.

The IIB registers only provide valid interruption bundle information on certain IVA-based faults and traps. Please refer to [Table 8-1, “Writing of Interruption Resources by Vector” on page 2:166](#) and corresponding interruption vector pages in [Section 8.3, “Interruption Vector Definition” on page 2:166](#) for information on which faults and traps these registers are valid. For faults and traps that indicate IIB is not valid, updates to the register may occur, but the information is undefined.

For IA-32 interruptions, instruction bundle information is not provided and the values in IIB registers are undefined.

The IIB registers are not supported on all processor implementations. Software can call PAL\_PROC\_GET\_FEATURES to determine the availability of this feature, see [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)” on page 2:446](#) for details. The IIB registers are reserved when this feature is not supported.

## 3.3.6 External Interrupt Control Registers

The external interrupt control registers (CR64-81) are defined in [“External Interrupt Control Registers” on page 2:121](#). They are used to prioritize and deliver external interrupts, send inter-processor interrupts to other processors and assign interrupt vectors for locally generated processor interrupts.

## 3.3.7 Banked General Registers

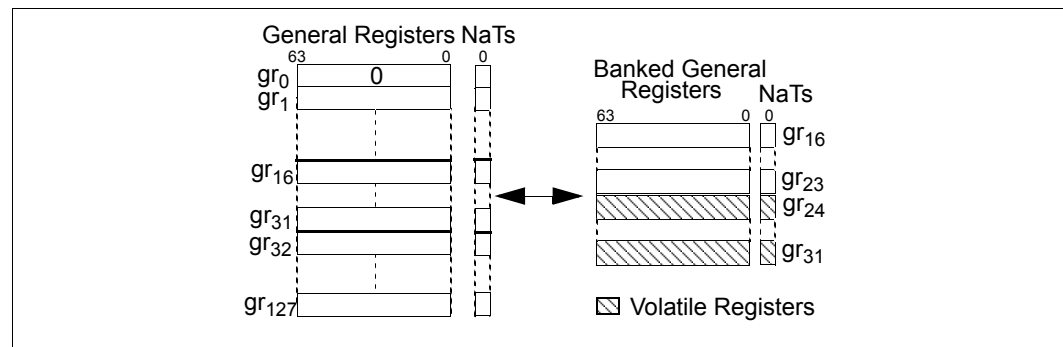
Banked general registers (see [Figure 3-18](#)) provide immediate register context for low-level interruption handlers (e.g., speculation and TLB miss handlers). Upon interruption, the processor switches 16 general purpose registers (GR16 to GR31) to register bank 0, register bank 1 contents are preserved.

When PSR.bn is 1, bank 1 for registers GR16 to GR31 is selected; when 0, bank 0 for registers GR16 to GR31 is selected. Banks are switched in the following cases:

- An interruption selects bank 0,
- `rfi` switches to the bank specified by `IPSR.bn`, or
- `bsw` switches to the specified bank.

On an interruption or bank switch, the processor ensures all prior register accesses (reads and writes) are performed to the prior register bank. Data values in banked registers are preserved across bank switches and both banks maintain NaT values when loaded from general registers. Registers from both banks cannot be addressed at the same time. However, non-banked general registers (GR0-15, and GR32-127) are accessible regardless of the state of `PSR.bn`.

**Figure 3-18. Banked General Registers**



Whether the ALAT register target tracking mechanism (see [“Data Speculation” on page 1:63](#)) distinguishes between the two register banks is implementation dependent; from the ALAT’s perspective, GR16 in bank 0 may be the same register as GR16 in bank 1 in some implementations.

Operating systems should ensure that IA-32 and Itanium architecture-based application code is executed within register bank 1. If IA-32 or Itanium architecture-based application code executes out of register bank 0, the application register state (including IA-32) will be lost on any interruption. During interruption processing the operating system uses register bank 0 as the initial working register context.

Usage of these additional registers is determined by software conventions. However, registers GR24 to GR31, of bank 0, are not preserved when `PSR.ic` is 1; operating system code can not rely on register values being preserved unless `PSR.ic` is 0. While `PSR.ic` is 1, processor-specific firmware may use these registers for machine check or firmware interruption handling at any point regardless of the state of `PSR.i`. If `PSR.ic` is 0, GR24 to GR31 can be used as scratch registers for low-level interruption handlers. Registers GR16 to GR23 are always preserved; operating system code can rely on the values being preserved.

## 3.4 Processor Virtualization

Processors in the Itanium Processor Family may optionally implement a mechanism to support processor virtualization. This includes an additional PSR.vm bit (see [Section 3.3.2, “Processor Status Register \(PSR\)”](#)), which, when 1, causes certain instructions to take a Virtualization fault (see [Section 5.6, “Interruption Priorities”](#) and [“Virtualization vector \(0x6100\)”](#) on page 2:209).

The set of instructions which are virtualized by PSR.vm are listed in [Table 3-10](#) below.

**Table 3-10. Virtualized Instructions**

Class	Virtualized Instructions
All privileged instructions	<code>itc.i, itc.d, itr.i, itr.d, ptc.l, ptc.g, ptc.ga, ptc.e, ptr, tak, tpa, mov rr, mov pkr, mov cr, mov ibr, mov dbr, mov pmc, mov to pmd, ssm, rsm, mov psr, rfi, bsw</code>
Some non-privileged instructions (virtualized at all privilege levels)	<code>thash, ttag, mov from cpuid, probe<sup>a</sup></code>
Some non-privileged instructions (virtualized at privilege level 0)	<code>cover, probe<sup>a</sup></code>
Reading AR[ITC] or AR[RUC] with PSR.si==1 (virtualized at all privilege levels)	<code>mov from ar.itc, mov from ar.ruc</code>
Instructions which write privileged registers	<code>mov to ar.itc, mov to ar.ruc</code>

a. Virtualization of the `probe` instruction is configurable, see [Section 11.7.4.2.8, “Probe Instruction Virtualization”](#) on page 2:344 for details.

Processors which support processor virtualization must provide an implementation-dependent mechanism for disabling the `vmsw` instruction. When enabled, the `vmsw` instruction functions as described on the `vmsw` instruction page. When disabled, the `vmsw` instruction always raises a Virtualization fault when executed at the most privileged level.

Processors which support processor virtualization may provide an implementation-dependent mechanism to disable virtual machine features, see [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)”](#) on page 2:446 for details.

Processor virtualization is largely invisible to system software, and therefore its effects on virtualized instructions are not discussed in this document, except on the instruction description pages themselves.

### §



This chapter defines operating system resources to translate 64-bit virtual addresses into physical addresses, 32-bit virtual addressing, virtual aliasing, physical addressing, memory ordering and properties of physical memory. Register state defined to support virtual memory management is defined in [Chapter 3](#), while [Chapter 5](#) provides complete information on virtual memory faults.

**Note:** Unless otherwise noted, references to “interruption” in this chapter refer to IVA-based interruptions. See [“Interruption Definitions” on page 2:95](#).

The following key features are supported by the virtual memory model.

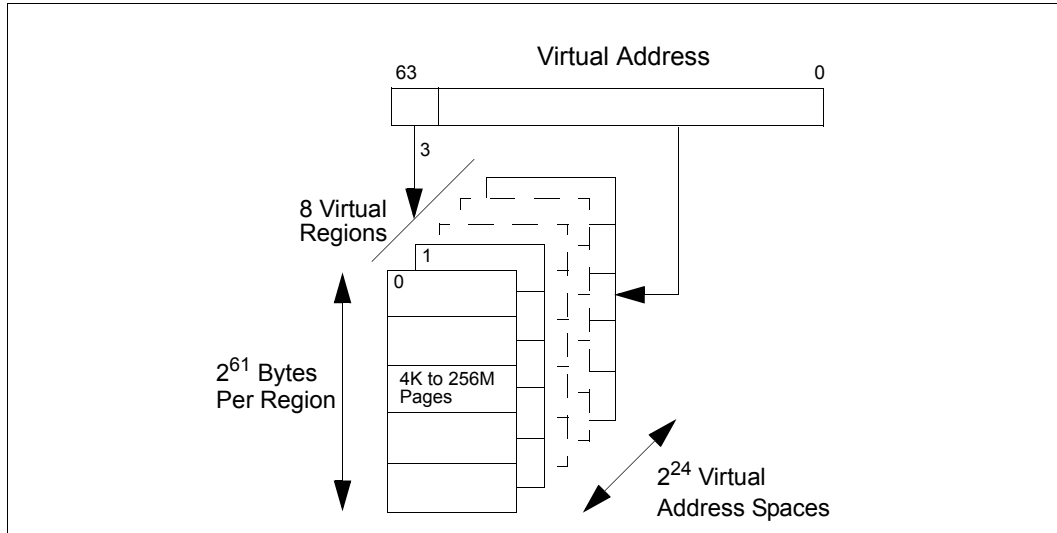
- Virtual Regions are defined to support contemporary operating system Multiple Address Space (MAS) models of placing each process within a unique address space. Region identifiers uniquely tag virtual address mappings to a given process.
- Protection Domain mechanisms support the Single Address Space (SAS) model, where processes co-exist within the same virtual address space.
- Translation Lookaside Buffer (TLB) structures are defined to support high-performance paged virtual memory systems. Software TLB fill and protection handlers are utilized to defer translation policies and protection algorithms to the operating system.
- A Virtual Hash Page Table (VHPT) is designed to augment the performance of the TLB. The VHPT is an extension of the processor’s TLB that resides in memory and can be automatically searched by the processor. A particular operating system page table format is not dictated. However, the VHPT is designed to mesh with two common translation structures: the virtual linear page table and hashed page table. Enabling of the VHPT and the size of the VHPT are completely under software control.
- Sparse 64-bit virtual addressing is supported by providing for large translation arrays (including multiple levels of hierarchy similar to a cache hierarchy), efficient translation miss handling support, multiple page sizes, pinned translations, and mechanisms to promote sharing of TLB and page table resources.

## 4.1 Virtual Addressing

As seen by Itanium architecture-based application programs, the virtual addressing model is fundamentally a 64-bit flat linear virtual address space. 64-bit general registers are used as pointers into this address space. IA-32 32-bit virtual linear addresses are zero extended into the 64-bit virtual address space.

As shown in [Figure 4-1](#), the 64-bit virtual address space is divided into eight  $2^{61}$  byte virtual regions. The region is selected by the upper 3-bits of the virtual address. Associated with each virtual region is a region register that specifies a 24-bit region identifier (unique address space number) for the region. Eight out of the possible  $2^{24}$  virtual address spaces are concurrently accessible via the 8 region registers. The region identifier can be considered the high order address bits of a large 85-bit global address space for a single address space model, or as a unique ID for a multiple address space model.

**Figure 4-1. Virtual Address Spaces**



By assigning sequential region identifiers, regions can be coalesced to produce larger 62-, 63- or 64-bit spaces. For example, an operating system could implement a 62-bit region for process private data, 62-bit region for I/O, and a 63-bit region for globally shared data. Default page sizes and translation policies can be assigned to each virtual region.

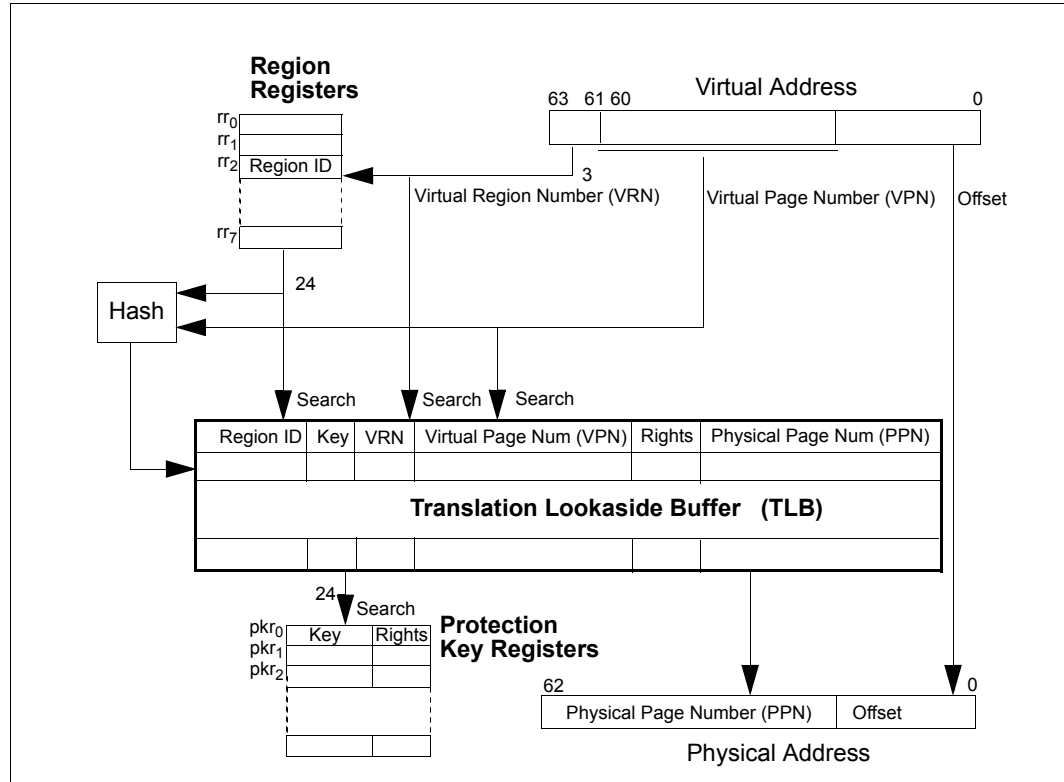
Figure 4-2 shows the process of mapping a virtual address into a physical address. Each virtual address is composed of three fields: the Virtual Region Number, the Virtual Page Number, and the page offset. The upper 3-bits select the Virtual Region Number (VRN). The least-significant bits form the page offset. The Virtual Page Number (VPN) consists of the remaining bits. The VRN bits are not included in the VPN. The page offset bits are passed through the translation process unmodified. Exact bit positions for the page offset and VPN bits vary depending on the page size used in the virtual mapping.

On a memory reference (any reference other than an insert or purge), the VRN bits select a Region Identifier (RID) from 1 of the 8 region registers, the TLB is then searched for a translation entry with a matching VPN and RID value. The VRN may optionally be used when searching for a matching translation on memory references (references other than inserts and purges – see Section 4.1.1.4, “Purge Behavior of TLB Inserts and Purges”). If a matching translation entry is found, the entry’s physical page number (PPN) is concatenated with the page offset bits to form the physical address. Matching translations are qualified by page-granular privilege level access right checks and optional protection domain checks by verifying the translation’s key is contained within a set of protection key registers and read, write, execute permissions are granted.

If the required translation is not resident in the TLB, the processor may optionally search the VHPT structure in memory for the required translation and install the entry into the TLB. If the required entry cannot be found in the TLB and/or VHPT, the processor raises a TLB Miss fault to request that the operating system supply the translation. After the operating system installs the translation in the TLB and/or VHPT, the faulting instruction can be restarted and execution resumed.

Virtual addressing for instruction references are enabled when PSR.it is 1, data references when PSR.dt is 1, and register stack accesses when PSR.rt is 1.

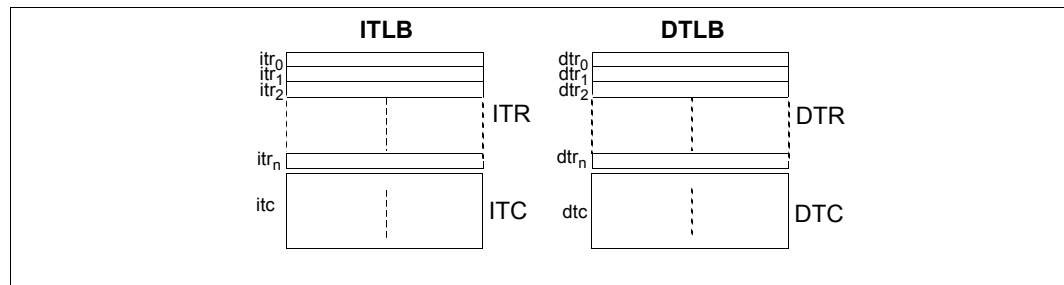
**Figure 4-2. Conceptual Virtual Address Translation for References**



### 4.1.1 Translation Lookaside Buffer (TLB)

The processor maintains two architectural TLBs as shown in Figure 4-3, the Instruction TLB (ITLB) and Data TLB (DTLB). Each TLB services translation requests for instruction and data memory references (including IA-32), respectively. The Data TLB also services translation requests for references by the RSE and the VHPT walker. The TLBs are further divided into two sub-sections; Translation Registers (TR) and Translation Cache (TC).

**Figure 4-3. TLB Organization**



In the remainder of this document, the term TLB refers to the combined instruction, data, translation register, and translation cache structures.

The TLB is a local processor resource; installation of a translation or local processor purges do not affect other processor's TLBs. Global TLB purges are provided to purge translations from all processors within a TLB coherence domain in a multiprocessor system.

#### 4.1.1.1 Translation Registers (TR)

The Translation Register (TR) section of the TLB is a fully-associative array defined to hold translations that software directly manages. Software can explicitly insert a translation into a TR by specifying a register slot number. Translations are removed from the TRs by specifying a virtual address, page size and a region identifier. Translation registers allow the operating system to "pin" critical virtual memory translations in the TLB. Examples include I/O spaces, kernel memory areas, frame buffers, page tables, sensitive interruption code, etc. Instruction fetches for interruption handlers are performed using virtual addresses; therefore, virtual address ranges containing software translation miss routines and critical interruption sequences should be pinned or else additional TLB faults may occur. Other virtual mappings may be pinned for performance reasons.

Entries are placed into a specific TR slot with the Insert Translation Register (`itr`) instruction. Once a translation is inserted, the processor will not replace the translation to make room for other translations. Local translations can only be removed by software issuing the Purge Translation Register (`ptr`) instruction.

TR inserts and purges may cause other TR and/or TC entries to be removed (refer to Section 4.1.1.4, "Purge Behavior of TLB Inserts and Purges" for details). Prior to inserting a TR entry, software must ensure that no overlapping translation exists in any TR (including the one being written); otherwise, a Machine Check abort may be raised, or the processor may exhibit other undefined behavior. Translation register entries may be removed by the processor due to hardware or software errors. In the presence of an error, the processor can remove TR entries; notification is raised via a Machine Check abort.

There are at least 8 instruction and 8 data TR slots implemented on all processor models. Please see the processor-specific documentation for further information on the number of translation registers implemented on the Itanium processor. Translation registers support all implemented page sizes and must be implemented in a single-level fully-associative array. Any register slot can be used to specify any virtual address mapping. Translation registers are not directly readable.

In some processor models, translation registers are physically implemented as a subsection of the translation cache array. Valid TR slots are ignored for purposes of processor replacement on an insertion into the TC. However, invalid TR slots (unused slots) may be used as TC entries by the processor. As a result, software inserts into previously invalid TR entries may invalidate a TC entry in that slot.

Implementations may also place a floating boundary between TR and TC entries within the same structure where any entry above the boundary is considered a TC and any entry below the boundary a TR. To maximize TC resources, software should allocate contiguous translation registers starting at slot 0 and continuing upwards.

#### 4.1.1.2 Translation Cache (TC)

The Translation Cache (TC) is an implementation-specific structure defined to hold the large working set of dynamic translations for memory references (including IA-32). Please see the processor-specific documentation for further information on Itanium processor TC implementation details. The processor directly controls the replacement policy of all TC entries.

Entries are installed by software into the translation cache with the Insert Data Translation Cache (`itc.d`) and Insert Instruction Translation Cache (`itc.i`) instructions. The Purge Translation Cache Local (`ptc.l`) instruction purges all ITC/DTC entries in the local processor that match the specified virtual address range and region identifier. Purges of all ITC/DTC entries matching a specified virtual address range and region identifier among all processors in a TLB coherence domain can be globally performed with the Purge Translation Cache Global (`ptc.g`, `ptc.ga`) instruction. The TLB coherence domain covers at least the processors on the same local bus on which the purge was broadcast. Propagation between multiple TLB coherence domains is platform dependent. Software must handle the case where a purge does not propagate to all processors in a multiprocessor system. Translation cache purges do not invalidate TR entries.

All the entries in a local processor's ITC and DTC can be purged of all entries with a sequence of Purge Translation Cache Entry (`ptc.e`) instructions. A `ptc.e` does not propagate to other processors.

In all processor models, the translation cache has at least 1 instruction and 1 data entry in addition to the specified 8 instruction and 8 data translation registers. Implementations are free to implement translation cache arrays of larger sizes. Implementations may also choose to implement additional hierarchies for increased performance. At least one translation cache level is required to support all implemented page sizes. Additional hierarchy levels may or may not be performance optimized for the preferred page size specified by the virtual region, may be set-associative or fully associative, and may support a limited set of page sizes. Please see the processor-specific documentation for further information on the Itanium processor implementation details of the translation cache.

The translation cache is managed by both software and hardware. In general, software cannot assume any entry installed will remain, nor assume the lifetime of any entry since replacement algorithms are implementation specific. The processor may discard or replace a translation at any point in time for any reason (subject to the forward progress rules below). TC purges may remove more entries than explicitly requested. In the presence of a processor hardware error, the processor may remove TC entries and optionally raise a Corrected Machine Check Interrupt.

In order to ensure forward progress for Itanium architecture-based code, the following rules must be observed by the processor and software.

- Software may insert multiple translation cache entries per TLB fault, provided that only the last installed translation is required for forward progress.
- The processor may occasionally invalidate the last TC entry inserted. The processor must eventually guarantee visibility of the last inserted TC entry to all references while `PSR.ic` is zero. The processor must eventually guarantee visibility of the last inserted TC entry until an `rfi` sets `PSR.ic` to 1 and at least one instruction is executed with `PSR.ic` equal to 1, and completes without a fault or interrupt. The last

inserted TC entry may be occasionally removed before this point, and software must be prepared to re-insert the TC entry on a subsequent fault. For example, eager or mandatory RSE activity, speculative VHPT walks, or other interruptions of the restart instruction may displace the software-inserted TC entry, but when software later re-inserts the same TC entry, the processor must eventually complete the restart instruction to ensure forward progress, even if that restart instruction takes other faults which must be handled before it can complete. If PSR.ic is set to 1 by instructions other than `rfi`, the processor does not guarantee forward progress.

- If software inserts an entry into the TLB with an overlapping entry (same or larger size) in the VHPT, and if the VHPT walker is enabled, forward progress is not guaranteed. See [“VHPT Searching” on page 2:62](#).
- Software may only make references to memory with physical addresses or with virtual addresses which are mapped with TRs, or to addresses mapped by the just-inserted translation, between the insertion of a TC entry, and the execution of the instruction with PSR.ic equal to 1 which is dependent on that entry for forward progress. Software may also make repeated attempts to execute the same instruction with PSR.ic equal to 1. If software makes any other memory references than these, the processor does not guarantee forward progress.
- Software must not defeat forward progress by consistently displacing a required TC entry through a global or local translation cache purge.

IA-32 code has more stringent forward progress rules that must be observed by the processor and software. IA-32 forward progress rules are defined in [Section 10.6.3, “IA-32 TLB Forward Progress Requirements” on page 2:261](#).

The translation cache can be used to cache TR entries if the TC maintains the instruction vs. data distinction that is required of the TRs. A data reference cannot be satisfied by a TC entry that is a cache of an instruction TR entry, nor can an instruction reference be satisfied by a TC entry that is a cache of a data TR entry. This approach can be useful in a multi-level TLB implementation.

#### **4.1.1.3 Unified Translation Lookaside Buffers**

Some processor models may merge the ITC and DTC into a unified translation cache. The minimum number of unified entries is 2 (1 for instruction, and 1 for data). Processors may service instruction fetch memory references with TC entries originally installed into the DTC and service data memory references with translations originally installed in the ITC. To ensure consistent operation across processor implementations, software is recommended to not install different translations into the ITC or DTC for the same virtual region and virtual address. ITC inserts may remove DTC entries. DTC inserts may remove ITC entries. TC purges remove ITC and DTC entries.

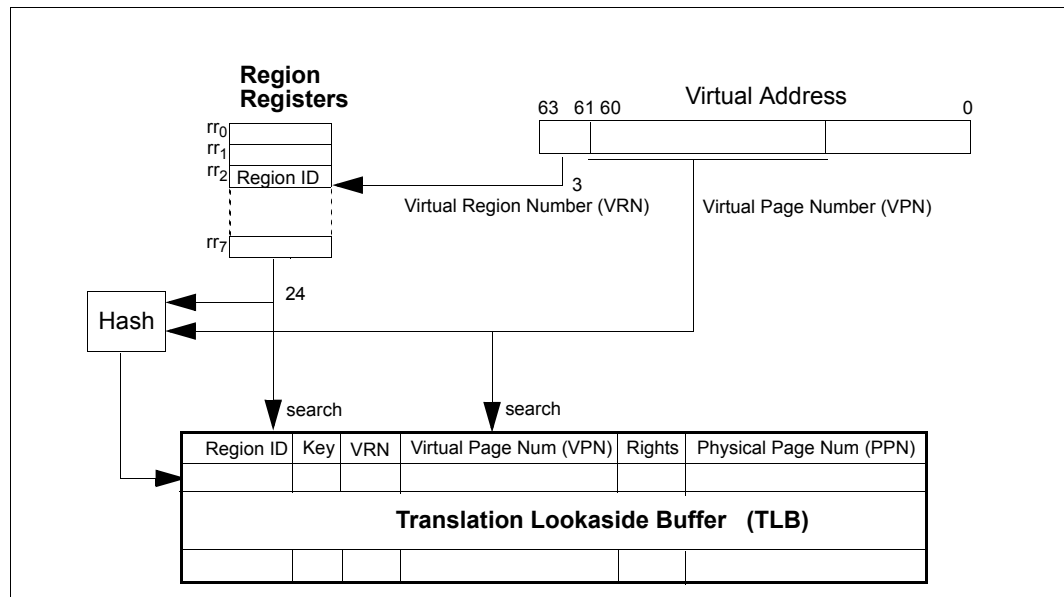
Instruction and data translation registers cannot be unified. DTR entries cannot be used by instruction references and ITR entries cannot be used by data references. ITR inserts and purges do not remove DTR entries. DTR inserts and purges do not remove ITR entries.

#### 4.1.1.4 Purge Behavior of TLB Inserts and Purges

Translations contained in the translation caches (TC) and translation registers (TR) are maintained in a consistent state by ensuring that TLB insertions remove existing overlapping entries before new TR or TC entries are installed. Similarly, TLB purges that partially or fully overlap with existing translations may remove all overlapping entries. In this context, “overlap” refers to two translations with the same region identifier (but not necessarily identical virtual region numbers), and with partially or fully overlapping virtual address ranges (determined by the virtual address and the page size). Examples are: two 4K-byte pages at the same virtual address, or an 8K-byte page at virtual address 0x2000 and a 4K-byte page at 0x3000.

As described in Section 4.1, “Virtual Addressing” on page 2:45, each TLB may contain a VRN field, and virtual address bits {63:61} may be used as part of the match for memory references (references other than inserts and purges). This binding of a translation to the VRN implies that a lookup of a given virtual address (region identifier/VPN pair) in either the translation cache or translation registers may result in a TLB miss if a memory reference is made through a different VRN (even if the region identifiers in the two region registers are identical). Some processor models may also omit the VRN field of the TLB, causing the TLB search on memory references to find an entry independent of VRN bits. However, all processor models are required, during translation cache purge and insert operations, to purge all possible translations matching the region identifier and virtual address regardless of the specified VRN.

**Figure 4-4. Conceptual Virtual Address Searching for Inserts and Purges**



A processor may overpurge translation cache entries; i.e., it may purge a larger virtual address range than required by the overlap. Since page sizes are powers of 2 in size and aligned on that same power of 2 boundary, purged entries can either be a superset of, identical to, or a subset of the specified purge range.

Table 4-1 define the purge behavior of different TLB insert and purge instructions. Table 4-2 describes the purge behavior for VHPT inserts.

**Note:** Please refer to Table 4-1 for footnotes in Table 4-2.

**Table 4-1. Purge Behavior of TLB Inserts and Purges**

Case	Insert?	Purge?	Machine Check?
<code>it[cr].[id] overlaps [ID]TC<sup>a</sup></code>	Must <sup>b</sup>	Must <sup>c</sup>	Must not <sup>d</sup>
<code>it[cr].[id] overlaps [DI]TC<sup>e</sup></code>	Must	May <sup>f</sup>	Must not
<code>it[cr].[id] overlaps [ID]TR</code>	May <sup>g</sup>	May	Must <sup>h</sup>
<code>it[cr].[id] overlaps [DI]TR</code>	Must	Must not <sup>i</sup>	Must not
<code>ptc.l overlaps [ID]TC</code>	N/A	Must	Must not
<code>ptc.l overlaps [ID]TR</code>		May	Must
<code>ptc.g (local) overlaps [ID]TC<sup>j</sup></code>		Must	Must not
<code>ptc.g (local) overlaps [ID]TR</code>		May	Must
<code>ptc.g (remote) overlaps [ID]TC</code>		Must	Must not
<code>ptc.g (remote) overlaps [ID]TR</code>		Must not	Must not
<code>ptc.e overlaps [ID]TC</code>		Must	Must not
<code>ptc.e overlaps [ID]TR</code>		Must not	Must not
<code>ptr.[id] overlaps [ID]TC</code>		Must	Must not
<code>ptr.[id] overlaps [DI]TC</code>		May	Must not
<code>ptr.[id] overlaps [ID]TR</code>		Must	Must not
<code>ptr.[id] overlaps [DI]TR</code>		Must not	Must not

- a. Bracketed notation is intended to specify TC and TR overlaps in the same stream, e.g. `itc.i` and `ITC`.
- b. Must Insert: requires that the translation specified by the operation is inserted into a TC or TR as appropriate. For `itc` and VHPT walker inserts, there is no guarantee to software that the entry will exist in the future, with the exception of the relevant forward-progress requirements specified in Section 4.1.1.2, "Translation Cache (TC)".
- c. Must Purge: requires that all partially or fully overlapped translations are removed prior to the insert or purge operation.
- d. Must not Machine Check: indicates that a processor does not cause a Machine Check abort as a result of the operation.
- e. Bracketed notation is intended to specify TC and TR overlaps in the opposite stream, e.g. `itc.i` and `DTC`.
- f. May Purge: indicates that a processor may remove partially or fully overlapped translations prior to the insert or purge operation. However, software must not rely on the purge.
- g. May Insert: indicates that the translation specified by the operation may be inserted into a TC. However, software must not rely on the insert.
- h. Must Machine Check: indicates that a processor will cause a Machine Check abort if an attempt is made to insert or purge a partially or fully overlapped translation. The Machine Check abort may not be delivered synchronously with the TLB insert or purge operation itself, but is guaranteed to be delivered, at the latest, on a subsequent instruction serialization operation.
- i. Must not Purge: the processor does not remove (or check for) partially or fully overlapped translations prior to the insert or purge operation. Software can rely on this behavior.
- j. `ptc.g` (and `ptc.ga`): two forms of global TLB purges are distinguished: local and remote. The local form indicates that the `ptc.g` or `ptc.ga` was initiated on the local processor. The remote form indicates that this is an incoming TLB shoot-down from a remote processor.



**Table 4-2. Purge behavior of VHPT Inserts**

Case	VRN bits used for TLB searching on VHPT insert						VRN bits not used for TLB searching on VHPT insert		
	VRN Match			No VRN Match			Insert?	Purge?	Machine Check?
	Insert?	Purge?	Machine Check?	Insert?	Purge?	Machine Check?			
[ID]VHPT overlaps [ID]TC <sup>a</sup>	Must <sup>b</sup>	Must <sup>c</sup>	Must not <sup>d</sup>	Must	May	Must not	Must	Must	Must not
[ID]VHPT overlaps [DI]TC <sup>e</sup>	Must	May <sup>f</sup>	Must not	Must	May	Must not	Must	May	Must not
[ID]VHPT overlaps [ID]TR	May <sup>g</sup>	May	Must <sup>h</sup>	May	Must not <sup>i</sup>	May	May	Must not	Must
[ID]VHPT overlaps [DI]TR	Must	Must not	Must not	Must	Must not	Must not	Must	Must not	Must not

The VHPT walker's inserts into the TC follow purge-before-insert rules similar to those for software inserts. VHPT walker inserts into the DTC behave similar to *itc.d*; VHPT walker inserts into the ITC behave similar to *itc.i*. If an instruction reference results in a VHPT walk that misses in the data TLB, the DTC insert for the translation for the VHPT acts similar to an *itc.d*.

As described in Section 4.1, “Virtual Addressing” on page 2:45, processors may optionally use VRN bits when searching for a matching translation for a memory reference (references other than inserts and purges). In processors which do use VRN bits for such searches, VHPT inserts optionally may also use VRN bits in searching for overlapping entries. Thus, if a VHPT insertion overlaps a translation in the TC, but the VRN of the address being inserted does not match the VRN of the existing TC translation, the purge of the existing TC entry is optional. If a VHPT insertion overlaps a translation in a TR, but the VRN of the address being inserted does not match the VRN of the TR translation, the VHPT insertion is allowed, and a machine check is optional. In processors which do not use VRN bits when searching for a matching translation for a memory reference, the behavior of VHPT inserts is identical to that of software inserts (see Table 4-1, “Purge Behavior of TLB Inserts and Purges” on page 2:52).

If a VHPT insert overlaps with an existing TR entry and the VRN of the insertion matches the VRN of the existing TR entry (for example, if the translation being inserted is for a large page which overlaps with a small page translation in the TR), the VHPT insertion can be done, but a machine check must be raised. Software must not create overlapping translations in the VHPT that are larger than a currently existing TR translation. The behavior of VHPT inserts is summarized in Table 4-2.

#### 4.1.1.5 Translation Insertion Format

Figure 4-5 shows the register interface to insert entries into the TLB. TLB insertions are performed by issuing the Insert Translation Cache (*itc.d*, *itc.i*) and Insert Translation Registers (*itr.d*, *itr.i*) instructions. The first 64-bit field containing the physical address, attributes and permissions is supplied by a general purpose register operand. Additional protection key and page size information is supplied by the Interruption TLB Insertion Register (ITIR). The Interruption Faulting Address register (IFA) specifies the virtual address for instruction and data TLB inserts. ITIR and IFA are defined in “Control Registers” on page 2:29. The upper 3 bits of IFA (VRN bits{63:61}) select a virtual region register that supplies the RID field for the TLB entry. The RID of the selected region is tagged to the translation as it is inserted into the TLB.

Reserved fields or encodings are checked as follows:

- The GR[r] value is checked when a TLB insert instruction is executed, and if reserved fields or reserved encodings are used, a Reserved Register/Field fault is raised on the TLB insert instruction. If GR[r]{0} is zero (not-present Translation Insertion Format), the rest of GR[r] is ignored.
- The RR[vrn] value is checked when a mov to RR instruction is executed, and if reserved fields or reserved encodings are used, a Reserved Register/Field fault is raised on the mov to RR instruction.
- The ITIR value is checked either when a mov to ITIR instruction is executed, or when a TLB insert instruction is executed, depending on the processor implementation. If reserved fields or reserved encodings are used, a Reserved Register/Field fault is raised on the mov to ITIR or TLB insert instruction. In implementations where ITIR is checked on a TLB insert instruction, ITIR{63:32} and ITIR{31:8} may be ignored if GR[r]{0} is zero (not-present Translation Insertion Format).
- The IFA value is checked either when a mov to IFA instruction is executed, or when a TLB insert instruction is executed, depending on the processor implementation. If an unimplemented virtual address is used, an Unimplemented Data Address fault is raised on the mov to IFA or TLB insert instruction.

Software must issue an instruction serialization operation to ensure installs into the ITLB are observed by dependent instruction fetches and a data serialization operation to ensure installs into the DTLB are observed by dependent memory data references.

**Figure 4-5. Translation Insertion Format**

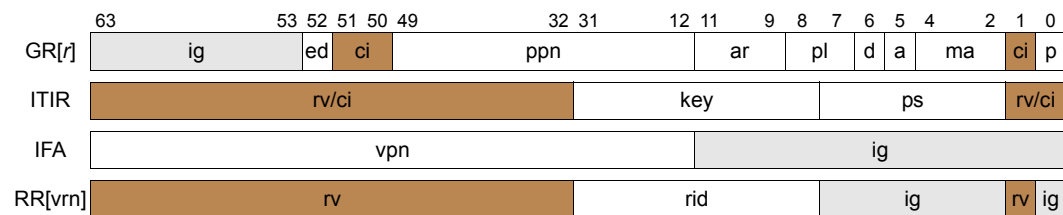


Table 4-3 describes all the translation interface fields.

**Table 4-3. Translation Interface Fields**

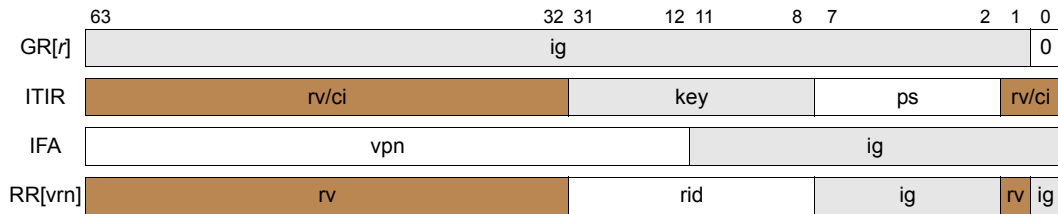
TLB Field	Source Field	Description
ci	GR[r]{1,51:50}	Checked on Insert – Checked on a TLB insert instruction. If reserved fields or encodings are used, a Reserved Register/Field fault is raised on the TLB insert instruction.
rv/ci	ITIR{1:0,63:32}	Reserved/Checked on Insert – Depending on implementation, may be reserved (checked on a mov to ITIR instruction) or checked on a TLB insert instruction. If reserved fields or encodings are used, a Reserved Register/Field fault is raised on the mov to ITIR or TLB insert instruction. In implementations where ITIR is checked on a TLB insert instruction, ITIR{63:32} may be ignored if GR[r]{0} is zero (not-present Translation Insertion Format).
rv	RR[vrn]{1,63:32}	Reserved – Checked on a mov to RR instruction. If reserved fields or encodings are used, a Reserved Register/Field fault is raised on the mov to RR instruction.

**Table 4-3. Translation Interface Fields (Continued)**

TLB Field	Source Field	Description
p	GR[r]{0}	Present bit – When 0, references using this translation cause an Instruction or Data Page Not Present fault. Most other fields are ignored by the processor, see <a href="#">Figure 4-6</a> for details. This bit is typically used to indicate that the mapped physical page is not resident in physical memory. The present bit is not a valid bit. For each TLB entry, the processor maintains an additional hidden valid bit indicating if the entry is enabled for matching.
ma	GR[r]{4:2}	Memory Attribute – describes the cacheability, coherency, write-policy and speculative attributes of the mapped physical page. See <a href="#">“Memory Attributes” on page 2:75</a> for details.
a	GR[r]{5}	Accessed Bit – When 0 and PSR.da is 0, data references to the page cause a Data Access Bit fault. When 0 and PSR.ia is 0, instruction references to the page cause an Instruction Access Bit fault. When 0, IA-32 references to the page cause an Instruction or Data Access Bit fault. This bit can trigger a fault on reference for tracing or debugging purposes. The processor does not update the Accessed bit on a reference.
d	GR[r]{6}	Dirty Bit – When 0 and PSR.da is 0, Intel Itanium store or semaphore references to the page cause a Data Dirty Bit fault. When 0, IA-32 store or semaphore references to the page cause a Data Dirty Bit fault. The processor does not update the Dirty bit on a store or semaphore reference.
pl	GR[r]{8:7}	Privilege Level – Specifies the privilege level or promotion level of the page. See <a href="#">“Page Access Rights” on page 2:56</a> for complete details.
ar	GR[r]{11:9}	Access Rights – page granular read, write and execute permissions and privilege controls. See <a href="#">“Page Access Rights” on page 2:56</a> for details.
ppn	GR[r]{49:12}	Physical Page Number – Most significant bits of the mapped physical address. Depending on the page size used in the mapping, some of the least significant PPN bits are ignored.
ig	GR[r]{63:53} IFA{11:0}, RR[vrn]{0,7:2}	available – Software can use these fields for operating system defined parameters. These bits are ignored when inserted into the TLB by the processor.
ed	GR[r]{52}	Exception Deferral – For a speculative load that results in an exception, the speculative load’s instruction page TLB.ed bit is one of the conditions which determines whether the exception must be deferred. See <a href="#">“Deferral of Speculative Load Faults” on page 2:105</a> for complete details. This bit is ignored in the data TLB for data memory references and for IA-32 memory references.
ps	ITIR{7:2}	Page Size – Page size of the mapping. For page sizes larger than 4K bytes the low-order bits of PPN and VPN are ignored. Page sizes are defined as $2^{Ps}$ bytes. See <a href="#">“Page Sizes” on page 2:57</a> for a list of supported page sizes.
key	ITIR{31:8}	Protection Key – Uniquely tags the translation to a protection domain. If a translation’s Key is not found in the Protection Key Registers (PKRs), access is denied and a Data or Instruction Key Miss fault is raised. See <a href="#">“Protection Keys” on page 2:59</a> for complete details. In implementations where ITIR is checked on a TLB insert instruction, ITIR{31:8} may be ignored if GR[r]{0} is zero (not-present Translation Insertion Format).
vpn	IFA{63:12}	Virtual Page Number – Depending on a translation’s page size, some of the least-significant VPN bits specified are ignored in the translation process. VPN{63:61} (VRN) selects the region register.
rid	RR[VRN].rid	Virtual Region Identifier – On TLB inserts the Region Identifier selected by VPN{63:61} (VRN) is used as additional match bits for subsequent accesses and purges (much like vpn bits).

The format in [Figure 4-6](#) is defined for not-present translations (P-bit is zero).

**Figure 4-6. Translation Insertion Format – Not Present**



### 4.1.1.6 Page Access Rights

Page granular access controls use 4 levels of privilege. Privilege level 0 is the most privileged and has access to all privileged instructions; privilege level 3 is least privileged. Access (including IA-32) to a page is determined by the TLB.ar and TLB.pl fields, and by the privilege level of the access, as defined in Table 4-4. RSE fills and spills obtain their privilege level from RSC.pl; all other accesses (including IA-32) obtain their privilege level from PSR.cpl. Within each cell, “-” means no access, “R” means read access, “W” means write access, “X” means execute access, and “Pn” means promote PSR.cpl to privilege level “n” when an Enter Privileged Code (epc) instruction is executed.

**Table 4-4. Page Access Rights**

TLB.ar	TLB.pl	Privilege Level <sup>a</sup>				Description
		3	2	1	0	
0	3	R	R	R	R	read only
	2	-	R	R	R	
	1	-	-	R	R	
	0	-	-	-	R	
1	3	RX	RX	RX	RX	read, execute
	2	-	RX	RX	RX	
	1	-	-	RX	RX	
	0	-	-	-	RX	
2	3	RW	RW	RW	RW	read, write
	2	-	RW	RW	RW	
	1	-	-	RW	RW	
	0	-	-	-	RW	
3	3	RWX	RWX	RWX	RWX	read, write, execute
	2	-	RWX	RWX	RWX	
	1	-	-	RWX	RWX	
	0	-	-	-	RWX	
4	3	R	RW	RW	RW	read only / read, write
	2	-	R	RW	RW	
	1	-	-	R	RW	
	0	-	-	-	RW	
5	3	RX	RX	RX	RWX	read, execute / read, write, exec
	2	-	RX	RX	RWX	
	1	-	-	RX	RWX	
	0	-	-	-	RWX	

**Table 4-4. Page Access Rights (Continued)**

TLB.ar	TLB.pl	Privilege Level <sup>a</sup>				Description
		3	2	1	0	
6	3	RWX	RW	RW	RW	read, write, execute / read, write
	2	–	RWX	RW	RW	
	1	–	–	RWX	RW	
	0	–	–	–	RW	
7	3	X	X	X	RX	exec, promote <sup>b</sup> / read, execute
	2	XP2	X	X	RX	
	1	XP1	XP1	X	RX	
	0	XP0	XP0	XP0	RX	

a. RSC.pl, for RSE fills and spills; PSR.cpl for all other accesses.

b. User execute only pages can be enforced by setting PL to 3.

Software can verify page level permissions by the `probe` (regular\_form `probe` or `probe.fault`) instruction, which checks accessibility to a given virtual page by verifying privilege levels, page level read and write permission, and protection key read and write permission.

Execute-only pages (TLB.ar 7) can be used to promote the privilege level on entry into the operating system. User level code would typically branch into a promotion page (controlled by the operating system) and execute the Enter Privileged Code (`epc`) instruction. When `epc` successfully promotes, the next instruction group is executed at the target privilege level specified by the promotion page. A procedure return branch type (`br.ret`) can demote the current privilege level.

#### 4.1.1.7 Page Sizes

A range of page sizes are supported to assist software in mapping system resources and improve TLB/VHPT utilization. Typically, operating systems will select a small range of fixed page sizes to implement virtual memory algorithms. Larger pages may be statically allocated. For example, large areas of the virtual address space may be reserved for operating system kernels, frame buffers, or memory-mapped I/O regions. Software may also elect to pin these translations, by placing them in the translation registers.

Table 4-5 lists insertable and purgeable page sizes that are supported by all processor models. Insertable page sizes can be specified in the translation cache, the translation registers, the region registers and the VHPT. Insertable page sizes can also be used as parameters to TLB purge instructions (`ptc.l`, `ptc.g`, `ptc.ga` or `ptr`). Page sizes that are purgeable only may only be used as parameters to TLB purge instructions.

Processors may also support additional insertable and purgeable page sizes. Please see the processor-specific documentation for further information on the page sizes supported by the Itanium processor.

**Table 4-5. Architected Page Sizes**

	Page Sizes										
	4k	8k	16k	64k	256k	1M	4M	16M	64M	256M	4G
Insertable	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes	-
Purgeable	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes	yes

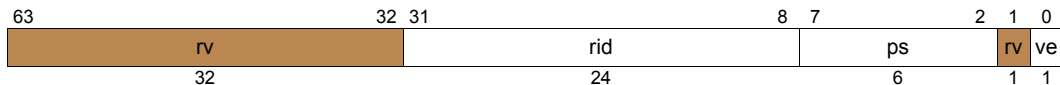
Page sizes are encoded in translation entries and region registers as a 6-bit encoded page size field. Each field specifies a mapping size of  $2^N$  bytes, thus a value of 12 represents a 4K-byte page. If unimplemented page sizes are specified to an `itc`, `itr` or `mov` to region register instruction, a Reserved Register/Field fault is raised. If unimplemented page sizes are specified for a TLB purge instruction an implementation may raise a Machine Check abort, may under-purge translations up to ignoring the request, or may over-purge translations up to removal of all entries from the translation cache. If unimplemented page sizes are specified by a `ptc.g` or `ptc.ga` broadcast from another processor, an implementation may under-purge translations up to ignoring the request, or may over-purge translations up to removal of all entries from the translation cache. However, it must not raise a Machine Check abort.

Virtual and physical pages are aligned on the natural boundary of the page. For example, 4K-byte pages are aligned on 4K-byte boundaries, and 4 M-byte pages on 4 M-byte boundaries.

### 4.1.2 Region Registers (RR)

Associated with each of the 8 virtual regions is a privileged Region Register (RR). Each register contains a Region Identifier (RID) along with several other region attributes, see Figure 4-7. The values placed in the region register by the operating system can be viewed as a collection of process address space identifiers.

**Figure 4-7. Region Register Format**



Regions support multiple address space operating systems by avoiding the need to flush the TLB on a context switch. Sharing between processes is promoted by mapping common global or shared region identifiers into the region register working set of multiple processes. All IA-32 memory references are through region register 0.

Table 4-6 describes the region register fields. Region Identifier (`rid`) bits 0 through 17 must be implemented on all processor models. Some processor models may implement additional bits. Additional implemented bits must be contiguous and start at bit 18. Unimplemented bits are reserved. Please see the processor-specific documentation for further information on the size of the Region Identifier implemented on the Itanium processor.

**Table 4-6. Region Register Fields**

Field	Bits	Description
<code>rv</code>	1,63:32	reserved
<code>ve</code>	0	VHPT Walker Enable – When 1, the VHPT walker is enabled for the region. When 0, disabled.

**Table 4-6. Region Register Fields (Continued)**

Field	Bits	Description
ps	7:2	Preferred page Size – Selects the virtual address bits used in hash functions for set-associative TLBs or the VHPT. Encoded as $2^{ps}$ bytes. The processor may make significant performance optimizations for the specified preferred page size for the region. <sup>a</sup>
rid	31:8	Region Identifier – During TLB inserts, the region identifier from the select region register is used to tag translations to a specific address space. During TLB/VHPT lookups, the region identifier is used to match translations and to distribute hash indexes among VHPT and TLB sets.

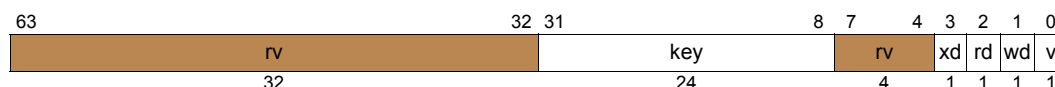
a. For more details on the usage of this field, See “VHPT Hashing” on page 2:65.

Software must issue an instruction serialization operation to ensure writes into the region registers are observed by dependent instruction fetches and issue a data serialization operation for dependent memory data references.

### 4.1.3 Protection Keys

Protection Keys provide a method to restrict permission by tagging each virtual page with a unique protection domain identifier. The Protection Key Registers (PKR) represent a register cache of all protection keys required by a process. The operating system is responsible for management and replacement policies of the protection key cache. Before a memory access (including IA-32) is permitted, the processor compares a translation’s key value against all keys contained in the PKRs. If a matching key is not found, the processor raises a Key Miss fault. If a matching Key is found, access to the page is qualified by additional read, write and execute protection checks specified by the matching protection key register. If these checks fail, a Key Permission fault is raised. Upon receipt of a Key Miss or Key Permission fault, software can implement the desired security policy for the protection domain. Figure 4-8 and Table 4-7 describe the protection key register format and protection key register fields.

**Figure 4-8. Protection Key Register Format**



**Table 4-7. Protection Register Fields**

Field	Bits	Description
v	0	Valid – When 1, the Protection Register entry is valid and is checked by the processor when performing protection checks. When 0, the entry is ignored.
wd	1	Write Disable – When 1, write permission is denied to translations in the protection domain.
rd	2	Read Disable – When 1, read permission is denied to translations in the protection domain.
xd	3	Execute Disable – When 1, execute permission is denied to translations in the protection domain.
key	31:8	Protection Key – uniquely tags translation to a given protection domain.
rv	7:4,63:32	reserved

Processor models have at least 16 protection key registers, and at least 18-bits of protection key. Some processor models may implement additional protection key registers and protection key bits. Unimplemented bits and registers are reserved. Key registers have at least as many implemented key bits as region registers have rid bits. Additional implemented bits must be contiguous and start at bit 18. Please see the processor-specific documentation for further information on the number of protection key registers and protection key bits implemented on the Itanium processor.

Software must issue an instruction serialization operation to ensure writes into the protection key registers are observed by dependent instruction fetches and a data serialization operation for dependent memory data references.

The processor ensures uniqueness of protection keys by checking new valid protection keys against all protection key registers during the move to PKR instruction. If a valid matching key is found in any PKR register, the processor invalidates the matching PKR register by setting PKR.v to zero, before performing the write of the new PKR register. The other fields in any matching PKR remain unchanged when it is invalidated.

Key Miss and Permission faults are only raised when memory translations are enabled (PSR.dt is 1 for data references, PSR.it is 1 for instruction references, PSR.rt is 1 for register stack references), and protection key checking is enabled (PSR.pk is one).

Data TLB protection keys can be acquired with the Translation Access Key ( $t_{ak}$ ) instruction. Instruction TLB key values are not directly readable. To acquire instruction key values software should make provisions to read memory structures.

#### 4.1.4 Translation Instructions

Table 4-8 lists translation instructions used to manage translations. Region registers, protection key registers and the TLBs are accessed indirectly; the register number is determined by the contents of a general register.

The processor does not ensure that modification of the translation resources is observed by subsequent instruction fetches or data memory references. Software must issue an instruction serialization operation before any dependent instruction fetch and a data serialization operation before any dependent data memory reference.

**Table 4-8. Translation Instructions**

Mnemonic	Description	Operation	Instr. Type	Serialization Requirement
mov rr[r <sub>3</sub> ] = r <sub>2</sub>	Move to region register	RR[GR[r <sub>3</sub> ]] = GR[r <sub>2</sub> ]	M	data/inst
mov r <sub>1</sub> = rr[r <sub>3</sub> ]	Move from region register	GR[r <sub>1</sub> ] = RR[GR[r <sub>3</sub> ]]	M	none
mov pkr[r <sub>3</sub> ] = r <sub>2</sub>	Move to protection key register	PKR[GR[r <sub>3</sub> ]] = GR[r <sub>2</sub> ]	M	data/inst
mov r <sub>1</sub> = pkr[r <sub>3</sub> ]	Move from protection key register	GR[r <sub>1</sub> ] = PKR[GR[r <sub>3</sub> ]]	M	none
itc.i r <sub>3</sub>	Insert instruction translation cache	ITC = GR[r <sub>3</sub> ], IFA, ITIR	M	inst



**Table 4-8. Translation Instructions (Continued)**

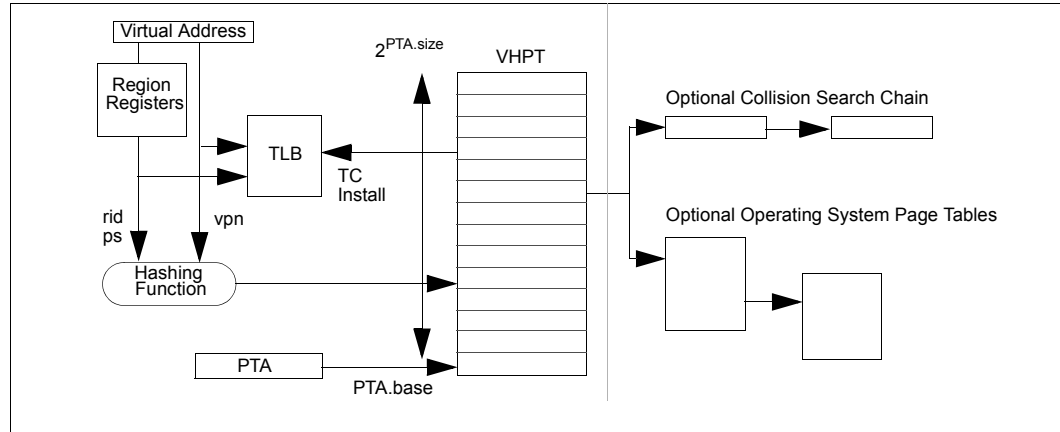
Mnemonic	Description	Operation	Instr. Type	Serialization Requirement
<code>itc.d r<sub>3</sub></code>	Insert data translation cache	DTC = GR[r <sub>3</sub> ], IFA, ITIR	M	data
<code>itr.i itr[r<sub>2</sub>] = r<sub>3</sub></code>	Insert instruction translation register	ITR[GR[r <sub>2</sub> ]] = GR[r <sub>3</sub> ], IFA, ITIR	M	inst
<code>itr.d dtr[r<sub>2</sub>] = r<sub>3</sub></code>	Insert data translation register	DTR[GR[r <sub>2</sub> ]] = GR[r <sub>3</sub> ], IFA, ITIR	M	data
<code>probe r<sub>1</sub> = r<sub>3</sub>, r<sub>2</sub></code>	Probe data TLB for translation		M	none
<code>probe.fault r<sub>3</sub>, imm<sub>2</sub></code>	Probe data TLB for translation		M	none
<code>ptc.l r<sub>3</sub>, r<sub>2</sub></code>	Purge a translation from local processor instruction and data translation cache		M	data/inst
<code>ptc.g r<sub>3</sub>, r<sub>2</sub></code>	Globally purge a translation from multiple processor's instruction and data translation caches		M	data/inst
<code>ptc.ga r<sub>3</sub>, r<sub>2</sub></code>	Globally purge a translation from multiple processor's instruction and data translation caches and remove matching entries from multiple processor's ALATs		M	data/inst
<code>ptc.e r<sub>3</sub></code>	Purge local instruction and data translation cache of all entries		M	data/inst
<code>ptr.i r<sub>3</sub>, r<sub>2</sub></code>	Purge instruction translation registers		M	inst
<code>ptr.d r<sub>3</sub>, r<sub>2</sub></code>	Purge data translation registers		M	data
<code>tak r<sub>1</sub> = r<sub>3</sub></code>	Obtain data TLB entry protection key		M	none
<code>thash r<sub>1</sub> = r<sub>3</sub></code>	Generate translation's VHPT hash address		M	none
<code>ttag r<sub>1</sub> = r<sub>3</sub></code>	Generate translation tag for VHPT		M	none
<code>tpa r<sub>1</sub> = r<sub>3</sub></code>	Translate a virtual address to a physical address		M	none

### 4.1.5 Virtual Hash Page Table (VHPT)

The VHPT is an extension of the TLB hierarchy designed to enhance virtual address translation performance. The processor's VHPT walker can optionally be configured to search the VHPT for a translation after a failed instruction or data TLB search. The VHPT walker provides significant performance enhancements by reducing the rate of flushing the processor's pipelines due to a TLB Miss fault, and by providing speculative translation fills concurrent to other processor operations.

The VHPT, resides in the virtual memory space and is configurable as either the primary page table of the operating system or as a single large translation cache in memory (see Figure 4-9). Since the VHPT resides in the virtual address space, an additional TLB miss can be raised when the VHPT is referenced. This property allows the VHPT to also be used as a linear page table.

**Figure 4-9. Virtual Hash Page Table (VHPT)**



The processor does not manage the VHPT or perform any writes into the table. Software is responsible for insertion of entries into the VHPT (including replacement algorithms), dirty/access bit updates, invalidation due to purges and coherency in a multiprocessor system. The processor does not ensure the TLBs are coherent with the VHPT memory image.

If software needs to control the entries inserted into the TLB more explicitly, or programs the VHPT with differing mappings for the same virtual address range, it may need to take additional action to ensure forward progress. See “VHPT Searching” on page 2:62.

#### 4.1.5.1 VHPT Configuration

The Page Table Address (PTA) register determines whether the processor is enabled to walk the VHPT, anchors the VHPT in the virtual address space, and controls VHPT size and configuration information. The VHPT can be configured as either a per-region virtual linear page table structure (8-byte short format) or as a single large hash page table (32-byte long format). No mixing of formats is allowed within the VHPT.

To implement a per-region linear page table structure an operating system would typically map the leaf page table nodes with small backing virtual translations. The size of the table is expanded to include all possible virtual mappings, effectively creating a large per-region flat page table within the virtual address space.

To implement a single large hash page table, the entire VHPT is typically mapped with a single large pinned virtual translation placed in the translation registers and the size of the table is reduced such that only a subset of all virtual mappings can be resident within the table. Operating systems can tune the size of the hash page table based on the size of physical memory and operating system performance requirements.

#### 4.1.5.2 VHPT Searching

When enabled, the processor’s VHPT walker searches the VHPT for a translation after a failed instruction or data TLB search. The VHPT walker checks only the specific VHPT entry addressed by the short- or the long-format hash function, as selected by PTA.vf. If additional TLB misses are encountered during the VHPT access, a VHPT Translation

fault is raised. If the region-based short-format VHPT entry contains no reserved bits or encodings, it is installed into the TLB, and the processor again attempts to translate the failed instruction or data reference. If the long-format VHPT entry's tag specifies the correct region identifier and virtual address, and the entry contains no reserved bits or encodings, it is installed into the TLB, and the processor again attempts to translate the failed instruction or data reference. Otherwise the processor raises a TLB Miss fault. The translation is installed into the TLB even if its VHPT entry is marked as not present (p=0). Software may optionally search additional VHPT collision chains (associativities) or search for translations within the operating system's primary page tables. Performance is optimized by placing frequently referenced translations within the VHPT structure directly searched by the processor.

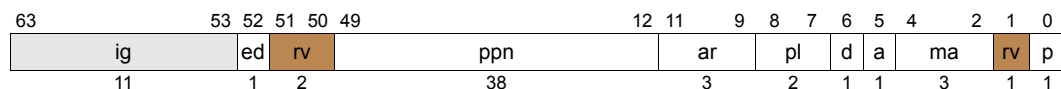
The VHPT walker is optional on a given processor model. Software can neither assume the presence of a VHPT walker, nor that the VHPT walker will find a translation in the VHPT. The VHPT walker can abort a search at any time for implementation-specific reasons, even if the required translation entry is in the VHPT. Operating systems must regard the VHPT walker strictly as a performance optimization and must be prepared to handle TLB misses if the walker fails.

VHPT walks may be done speculatively by the processor's VHPT walker. Additionally, VHPT walks triggered by non-speculatively-executed instructions are not required to be done in program order. Therefore, if the walker is enabled and if the VHPT contains multiple entries that map the same virtual address range, software must set up these entries such that any of them can be used in the translation of any part of this virtual address range. Additionally, if software inserts a translation into the TLB which is needed for forward progress, and this translation has a smaller page size than the translation which would have been inserted on a VHPT walk for the same address, then software may need to disable the VHPT walker in order to ensure forward progress, since this inserted translation may be displaced by a VHPT walk before it can be used.

### 4.1.5.3 Region-based VHPT Short Format

The region-based VHPT short format shown in [Figure 4-10](#) uses 8-byte VHPT entries to support a per-region linear page table configuration. To use the short-format VHPT, PTA.vf must be set to 0.

**Figure 4-10. VHPT Short Format**



See ["Translation Insertion Format" on page 2:53](#) for a description of all fields. The VHPT walker provides the following default values when entries are installed into the TLB.

- Virtual Page Number – implied by the position of the entry in the VHPT. The hashed short-format entry is considered to be the matching translation.
- Region Identifiers are not specified in the short format. To ensure uniqueness, software must provide unique VHPT mappings per region. Region identifiers obtained from the referenced region register are tagged with the translation when inserted into the TLB.
- Page Size – specified by the accessed region's preferred page size (RR[VA{63:61}].ps)

- Protection Key – specified by the accessed region identifier value (RR[VA{63:61}].rid). As a result, all implementations must ensure that the number of implemented key bits is greater than or equal to the number of implemented region identifier bits.

If a translation is marked as not present, ignored fields are usable by software as noted in Figure 4-11.

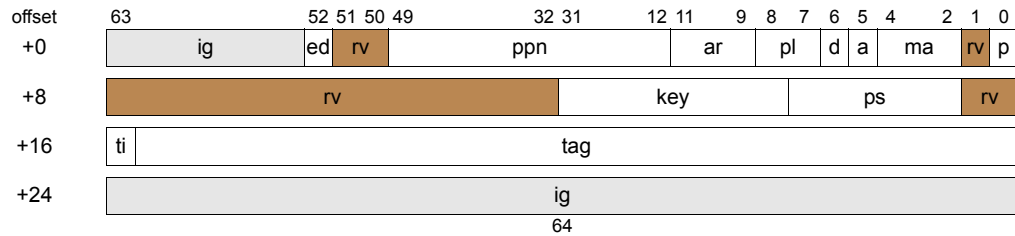
**Figure 4-11. VHPT Not-present Short Format**



#### 4.1.5.4 VHPT Long Format

The long-format VHPT uses 32-byte VHPT entries to support a single large virtual hash page table. To use the long-format VHPT, PTA.vf must be set to 1. The long format is a superset of the TLB insertion format, as noted in Figure 4-12, and specifies full translation information (including protection keys and page sizes). Additional fields are defined in Table 4-9. The long format is typically used to build the hash page table configuration.

**Figure 4-12. VHPT Long Format**

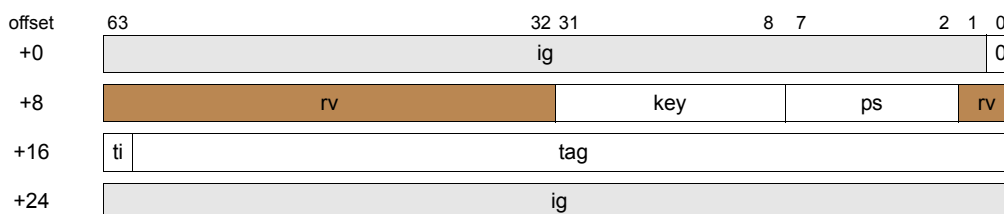


**Table 4-9. VHPT Long-format Fields**

Field	Offset	Description
tag	+16	Translation Tag – The tag, in conjunction with the VHPT hash index, is used to uniquely identify the translation. Tags are computed by hashing the virtual page number and the region identifier. See “VHPT Hashing” on page 2:65 for details on tag and hash index generation.
ti	+16	Tag Invalid Bit – If one, this bit of the tag indicates an invalid tag. On all processor implementations, the VHPT walker and the ttag instruction generate tags with the ti bit equal to 0. A VHPT entry with the ti bit equal to one will never be inserted into the processor’s TLBs. Software can use the ti bit to invalidate long-format VHPT entries in memory.
ig	+24	available – field for software use, ignored by the processor. Operating systems may store any value, such as a link address to extend collision chains on a hash collision.

If a translation is marked as not present, ignored fields are usable by software as noted in Figure 4-13. Also, in some implementations, +8{63:32} and +8{31:8} may be ignored as well.

**Figure 4-13. VHPT Not-present Long Format**



For multiprocessor systems, atomic updates of long-format VHPT entries may be ensured by software as follows:

- Before making multiple non-atomic updates to a VHPT entry in memory, software is required to set its ti bit to one.
- After making multiple non-atomic updates to a VHPT entry in memory, software may clear its ti bit to zero to re-enable tag matches.

The updates to the VHPT entry in memory must be constrained to be observable only after the store that sets the ti bit to one is observable. This can be accomplished with a `mfi` instruction, or by performing the updates to the VHPT entry with release stores. Similarly, the clearing of the ti bit must be constrained to be observable only after all of the updates to the VHPT entry are observable. This can be accomplished with a `mfi` instruction, or by performing the clear of the ti bit with a release store.

## 4.1.6 VHPT Hashing

The processor provides two methods for software to determine a VHPT entry’s address: the Translation Hash (`thash`) instruction, and the Interruption Hash Address (IHA) register defined on [page 2:41](#). The virtual address of the VHPT entry is placed in the IHA register when a VHPT Translation or TLB fault is delivered. In the long format, IHA can be used as a starting address to scan additional collision chains (associativities) defined by the operating system or to perform a search in software. The `thash` instruction is used to generate a VHPT entry’s address outside of interruption handlers and provides the same hash function that is used to calculate IHA.

`thash` produces a VHPT entry’s address for a given virtual address and region identifier, depending on the setting of the `PTA.vf` bit. When `PTA.vf=0`, `thash` returns the region-based short-format index as defined in [“Region-based VHPT Short-format Index” on page 2:65](#). When `PTA.vf=1`, `thash` returns the long-format hash as defined in [“Long-format VHPT Hash” on page 2:66](#). The `ttag` instruction is only useful for long-format hashing, and generates a 64-bit ti/tag identifier that the processor’s VHPT walker will check when it looks up a given virtual address and region identifier. Software should use the `ttag` instruction, and either the `thash` instruction or the IHA register when forming translation tags and hash addresses for the long-format VHPT. These resources encapsulate the implementation-specific long-format hashing functionality and improve performance.

### 4.1.6.1 Region-based VHPT Short-format Index

In the region-based short format, the linear page table for each region resides in the referenced region itself. As a result, the short-format VHPT consists of separate per-region page tables, which are anchored in each region by `PTA{60:15}`. For regions

in which the VHPT is enabled, the operating system is required to maintain a per-region linear page table. As defined in Figure 4-14, the VHPT walker uses the virtual address, the region’s preferred page size, and the PTA.size field to compute a linear index into the short-format VHPT.

**Figure 4-14. Region-based VHPT Short-format Index Function**

```
Mask = (1 << PTA.size) - 1;
VHPT_Offset = (VA{IMPL_VA_MSB:0} u>> RR[VA{63:61}].ps) << 3;
VHPT_Addr = (VA{63:61} << 61) |
  (((PTA{60:15} & ~Mask{60:15}) | (VHPT_Offset{60:15} &
    Mask{60:15})) << 15) |
  VHPT_Offset{14:0};
```

The size of the short-format VHPT (PTA.size) defines the size of the mapped virtual address space. The maximum architectural table size in the short format is  $2^{52}$  bytes per region. To map an entire region ( $2^{61}$  bytes) using 4Kbyte pages,  $2^{(61-12)} = 2^{49}$  pages must be mappable. A short-format VHPT entry is 8 bytes =  $2^3$  bytes large. As a result, the maximum table size is  $2^{(61-12+3)} = 2^{52}$  bytes per region. If the short format is used to map an address space smaller than  $2^{61}$ , a smaller short-format table (PTA.size < 52) can be used. Mapping of an address space of  $2^n$  with 4KByte pages requires a minimum PTA.size of (n-9).

In the short format, the `thash` instruction returns the region-based short-format index defined in Figure 4-14. The `ttag` instruction is not used with the short format. VHPT translation and TLB miss faults write the IHA register with the region-based short-format index defined in Figure 4-14.

#### 4.1.6.2 Long-format VHPT Hash

The long-format VHPT is a single large contiguous hash table that resides in the region defined by PTA.base. As defined in Figure 4-15, the VHPT walker uses the virtual address, the region identifier, the region’s preferred page size, and the PTA.size field to compute a hash index into the long-format VHPT. PTA{63:15} defines the base address and the region of the long-format VHPT. PTA.size reflects the size of the hash table, and is typically set to a number significantly smaller than  $2^{64}$ ; the exact number is based on operating system performance requirements.

**Figure 4-15. VHPT Long-format Hash Function**

```
Mask = (1 << PTA.size) - 1;
HPN = VA{IMPL_VA_MSB:0} u>> RR[VA{63:61}].ps;
Hash_Index = tlb_vhpt_hash_long(HPN, RR[VA{63:61}].rid);
// model-specific hash function
VHPT_Offset = Hash_Index << 5;
VHPT_Addr = (PTA{63:61} << 61) |
  (((PTA{60:15} & ~Mask{60:15}) | (VHPT_Offset{60:15}
    & Mask{60:15})) << 15) | VHPT_Offset{14:0};
```

The long-format hash function (`tlb_vhpt_hash_long`) and long-format tag generation function are implementation specific. However, on all processor models the hash and tag functions must exclude the virtual region number (virtual address bits VA{63:61}) from the hash and tag computations. This ensures that a unique 85-bit global virtual address hashes to the same VHPT hash address, regardless of which region the address is mapped to. All processor implementations guarantee that the most significant bit of

the tag (ti bit) is zero for all valid tags. The hash index and tag together must uniquely identify a translation. The processor must ensure that the indices into the hashed table, the region's preferred page size, and the tag specified in an indexed entry can be used in a reverse hash function to uniquely regenerate the region identifier and virtual address used to generate the index and tag. This must be possible for all supported page sizes, implemented virtual addresses and legal values of region identifiers. A hash function is reversible if using the hash result and all but one input produces the missing input as the result of the reverse hash function. The easiest hash function and reverse hash function is a simple XOR of bits. To ensure uniqueness, software must follow these rules:

1. Software must use only one preferred page size for each unique region identifier at any given time; otherwise, processor operation is undefined.
2. All tags for translations within a given region must be created with the preferred page size assigned to the region; otherwise, processor operation is undefined.
3. Software is not allowed to have pages in the VHPT that are smaller than the preferred page size for the region; otherwise, processor operation is undefined. Software can specify a page with a page size larger than the preferred page size in the VHPT, but tag values for the entries representing that page size must be generated using the preferred page size assigned to that region.
4. To reuse a region identifier with a different preferred page size, software must first ensure that the VHPT contains no insertable translations for that rid, purge all translations for that rid from all processors that may have used it, and then update the region register with the new preferred page size.

### 4.1.7 VHPT Environment

The processor's VHPT walker can optionally be configured to search the VHPT for a translation after a failed instruction or data TLB search. The VHPT walker is enabled for different types of references under the following conditions:

- Data and non-access references (including IA-32): PTA.ve=1, and RR[VA{63:61}].ve=1, and PSR.dt=1.
- Instruction fetches (including IA-32): PTA.ve=1, and RR[VA{63:61}].ve=1, and PSR.dt=1, and PSR.it=1, and PSR.ic=1.
- RSE references: PTA.ve=1, and RR[VA{63:61}].ve=1, and PSR.dt=1, and PSR.rt=1.

If the walker is not enabled, and an attempt is made to reference the VHPT, an Alternate Instruction/Data TLB Miss fault is raised. The remainder of this section assumes that the VHPT is enabled.

Region registers must support all implemented page sizes so software can use IHA, `thash` and `ttag` to manage the VHPT. `thash` and `ttag` are defined to operate on all page sizes supported by the translation cache, regardless of the VHPT walker's supported page sizes. The PTA register must be implemented on processor models that do not implement a VHPT walker. Software must ensure PTA is initialized and serialized before issuing `ttag`, `thash`, before enabling the VHPT walker or issuing a reference that may cause a VHPT walk. The minimum VHPT size is 32KBytes (PTA.size=15), and

operating systems must ensure that the VHPT is aligned on the natural boundary of the structure; otherwise, processor operation is undefined. For example, a 64K-byte table must be aligned on a 64K-byte boundary.

VHPT walker references to the VHPT are performed at privilege level 0, regardless of the state of `PSR.cpl`. VHPT byte ordering is determined by the state of `DCR.be`. When `DCR.be=1`, VHPT walker references are performed using big-endian memory formats; otherwise, VHPT walker references are little-endian. A long-format VHPT reference is matched against the data break-point registers as a 32-byte reference.

The VHPT is accessed by the processor only if the VHPT is virtually mapped into cacheable memory areas. The walker may access the VHPT speculatively, i.e., references may be performed that are not required by an in-order execution of the program. Any VHPT or TLB faults encountered during a VHPT walker's search are not reported until the faulting translation is required by an in-order execution of the program. If the VHPT is mapped into non-cacheable memory areas the VHPT is not referenced, and all TLB misses result in an Instruction/Data TLB Miss fault.

The VHPT walker will abort the search and deliver an Instruction/Data TLB Miss fault if an attempt is made to install translations that have reserved bits or encodings, or if the translation mapping the VHPT would have taken one of the following faults: Data Page Not Present, Data NaT Page Consumption, Data Key Miss, Data Key Permission, Data Access Bit, or Data Debug. The VHPT walker may abort a search and deliver an Instruction/Data TLB Miss fault at any time for implementation-specific reasons.

The processor's VHPT walker is required to read and insert VHPT entries from memory atomically (an 8-byte atomic read-and-insert for short format, and a 32-byte atomic read-and-insert for long format). Some implementation strategies for achieving this atomicity are as follows:

- If the walker performs its VHPT read with multiple cache accesses which are not done as an atomic unit, and if an update to part of the entry that is being installed is made in-between these multiple reads, the walker must abort the insert and deliver an Instruction/Data TLB Miss.
- If the walker performs its VHPT read and the insertion of the entry into the TLB as separate actions, and not as an atomic unit, and if an update to part of the entry that is being installed is made in-between the read and the insert, the walker must either abort the insert and deliver an Instruction/Data TLB Miss, or ignore the update and install the complete old entry.
- If the purge address range of a TLB purge operation (`ptc.l`, `ptc.e`, local or remote `ptc.g` or `ptc.ga`, `ptr.i`, or `ptr.d`) overlaps the virtual address the walker is attempting to insert, then the walker must either abort the insert and deliver an Instruction/Data TLB Miss, or delay the purge operation until after the walker either completes the insertion or aborts the walk.

The RSE can only raise a VHPT fault on a mandatory RSE spill/fill operation as defined for successful execution of an `alloc`, `loadrs`, `flushrs`, `br.ret` or `rfi` instruction. Eager RSE operations may generate speculative VHPT walks provided encountered faults are not reported.

Data TLB Miss faults encountered during a VHPT walk are permitted and, when `PSR.ic=1`, are converted into a VHPT Translation fault as defined in the next section.



### 4.1.8 Translation Searching

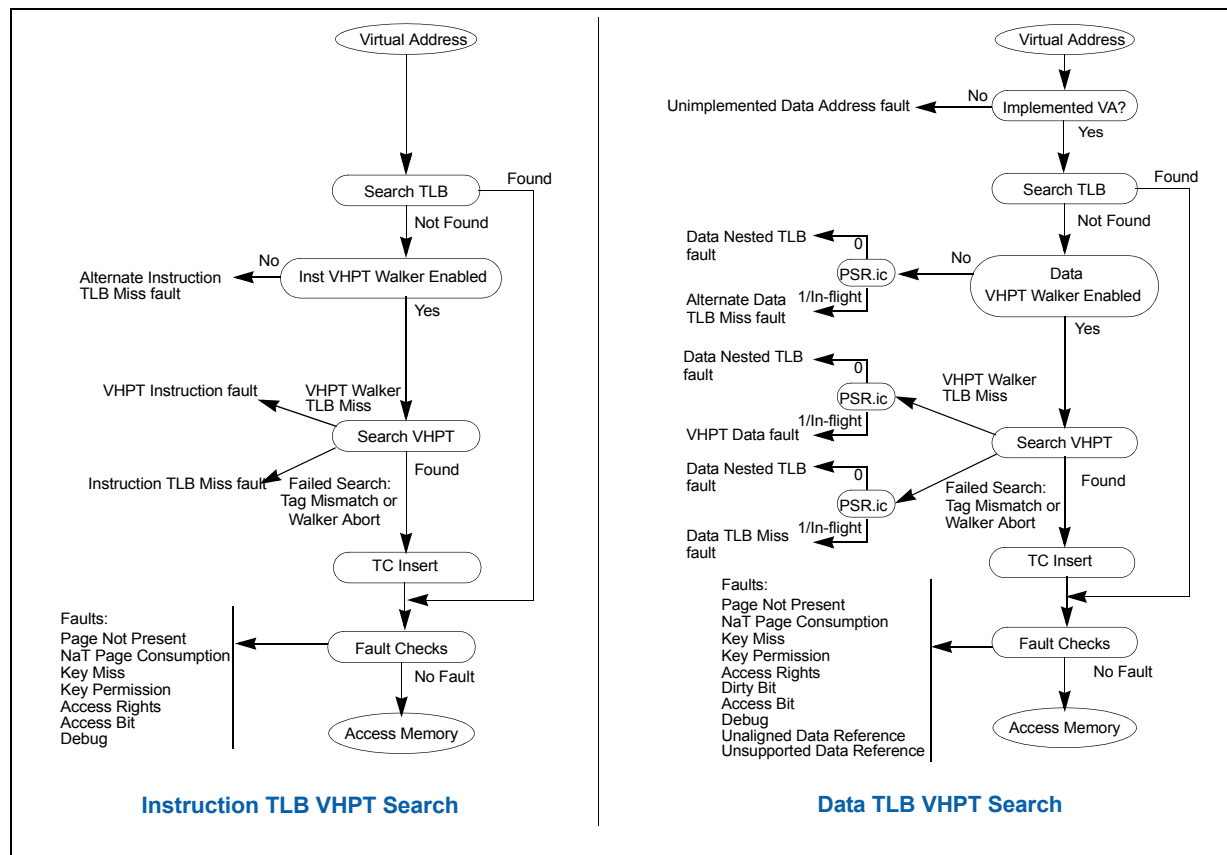
The general sequence of searching the TLB and VHPT is shown in [Figure 4-16](#). On a failed TLB search, if the VHPT walker is disabled for the referenced region an Alternate Instruction/Data TLB Miss fault is raised. If the VHPT walker is enabled for the referenced region, the VHPT is accessed to locate the missing translation. See [“VHPT Environment” on page 2:67](#). If additional TLB misses are encountered during the VHPT walker’s references, a VHPT Translation fault is raised. If the VHPT walker does not find the required translation in the VHPT or the search is aborted, an Instruction/Data TLB Miss fault is raised. Otherwise the entry is loaded into the ITC or DTC. Provided the above fault conditions are not detected, the processor may load the entry into the ITC or DTC even if an in-order execution of the program did not require the translation.

See [Table 4-1, “Purge Behavior of TLB Inserts and Purges,” on page 2:52](#) for the purge behavior of VHPT walker inserts.

After the translation entry is loaded, additional TLB faults are checked; these include in priority order: Page Not Present, NaT page Consumption, Key Miss, Key Permission, Access Rights, Access Bit, and Dirty Bit faults. [Table 4-10](#) describes the TLB and VHPT walker related faults.

On a failed TLB/VHPT search, the processor loads interruption registers and translation defaults as defined in [“Interruption Vector Descriptions” on page 2:165](#) defining the parameters of the translation fault. Provided the operating system accepts the defaults provided, only the physical address portion of a TLB entry need be provided on a TLB insert.

**Figure 4-16. TLB/VHPT Search**



**Table 4-10. TLB and VHPT Search Faults**

Fault	Description
VHPT Instruction/Data	Raised if there is an additional TLB miss when the VHPT walker attempts to access the VHPT. Typically used to construct leaf table mappings for linear page table configurations.
Alternate Instruction/Data TLB Miss	Raised when the VHPT walker is not enabled and an instruction or data reference causes a TLB miss. For example, the VHPT walker can be disabled within a given virtual region so region-specific translation algorithms can be utilized.

**Table 4-10. TLB and VHPT Search Faults (Continued)**

Fault	Description
Instruction/Data TLB Miss	<p>Raised when the VHPT walker is enabled, but the processor:</p> <ul style="list-style-type: none"> <li>• Cannot locate the required VHPT entry, or</li> <li>• The processor aborts the VHPT search for implementation-specific reasons, or</li> <li>• The VHPT walker is not implemented, or</li> <li>• The referenced region specifies a non-supported VHPT preferred page size, or</li> <li>• Reserved fields or unimplemented PPN bits are used in the translation, or</li> <li>• The hash address falls into unimplemented virtual address space, or</li> <li>• The hash address matches a data debug register.</li> </ul> <p>Instruction/Data TLB Miss handlers are essentially software walkers of the VHPT.</p>
Data Nested TLB	<p>Raised when a Data TLB Miss, Alternate Data TLB Miss, or VHPT Data Translation fault occurs and PSR.ic is 0 and not in-flight (e.g., fault within a TLB miss handler). Data Nested TLB faults enable software to avoid overheads for potential data TLB Miss faults.</p>
Instruction/Data Page Not Present	<p>The referenced translation's P-bit is 0.</p>
Instruction/Data NaT Page Consumption	<p>A non-speculative load, store, mandatory RSE load/store, execution on, or semaphore operation accesses a page marked with the physical memory attribute NaTPage. See <a href="#">“Not a Thing Attribute (NaTPage)” on page 2:86</a> for details.</p>
Instruction/Data Key Miss	<p>The referenced translation's permission key is not present in the set of valid protection key registers.</p>
Instruction/Data Key Permission	<p>The referenced translation is denied read, write, execute permissions by the matching protection key registers.</p>
Instruction/Data Access Rights	<p>Page granular read, write, execute and privilege level accesses are denied.</p>
Data Dirty Bit	<p>The referenced translation's Dirty bit is 0 on a store or semaphore operation.</p>
Instruction/Data Access Bit	<p>The referenced translation's Access bit is 0.</p>

### 4.1.9 32-bit Virtual Addressing

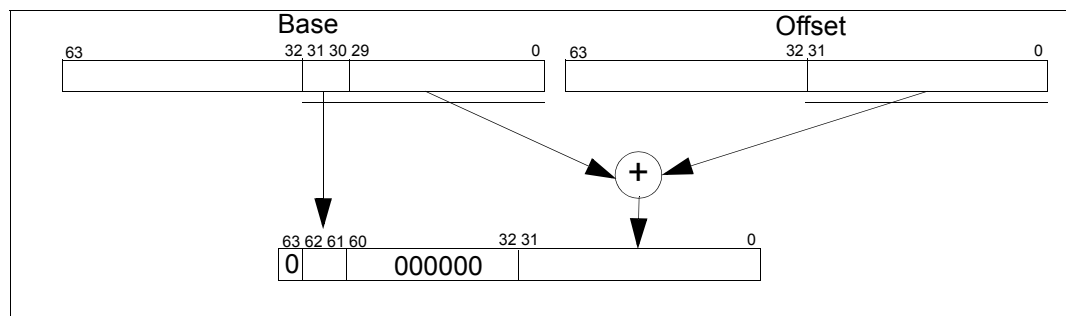
32-bit virtual data addressing is supported in the Itanium instruction set architecture by three models: zero-extension, sign-extension, and pointer “swizzling.” IA-32 memory references use the zero-extension model, all IA-32 32-bit virtual linear addresses are zero extended into the 64-bit virtual address space.

The zero-extension model performs address computations with the `add` and `shladd` instructions while software ensures that the upper 32-bits are always zeros. This model constrains 32-bit virtual addressing to virtual region zero. In this model, regions 1 to 7 are accessible only by 64-bit addressing.

In the sign-extension model, software ensures that the upper 32-bits of a virtual address are always equal to bit 31. Address computations use the `add`, `shladd`, and `sxt` instructions. This model splits the 32-bit address space into two halves that are spread into  $2^{31}$  bytes of virtual regions 0 and 7 within the 64-bit virtual address space. In this model, regions 2 to 6 are accessible only by 64-bit addressing.

The pointer “swizzling” model performs address computations with the `addp4`, and `shladdp4` instructions. These instructions generate a 32-bit address within the 64-bit virtual address space as shown in Figure 4-17. The 32-bit virtual address space is divided into 4 sections that are spread into  $2^{30}$  bytes of virtual regions 0 to 3 within the 64-bit virtual address space. In this model, regions 4 to 7 are accessible only by 64-bit addressing.

**Figure 4-17. 32-bit Address Generation using `addp4`**



In the pointer “swizzling” model, mappings within each region do not necessarily start at offset zero, since the upper 2-bits of a 32-bit address serve both as the virtual region number and an offset within each region. Virtual address bits {62:61} do not participate in the address addition, therefore some regions may be effectively larger than  $2^{30}$  bytes due to the addition of a 32-bit offset and lack of a carry into bits {62:61}. Note that the conversion is non-destructive: a converted 64-bit pointer can be used as a 32-bit pointer. Flat 31 or 32-bit address spaces can be constructed by assigning the same region identifier to contiguous region registers. Branches into another  $2^{30}$ -byte region are performed by first calculating the target address in the 32-bit virtual space and then converting to a 64-bit pointer by `addp4`. Otherwise, branch targets will extend above the  $2^{30}$  byte boundary within the originating region.

#### 4.1.10 Virtual Aliasing

Virtual aliasing (two or more virtual pages mapped to the same physical page) is functionally supported for memory references (including IA-32), however performance may be degraded on some processor models where the distance between virtual aliases is less than 1 MB. To avoid any possible performance degradation, software is advised to use aliases whose virtual addresses differ by an integer multiple of 1 MB. The processor ensures cache coherency and data dependencies in the presence of an alias. Stores using a virtual alias followed by a load with another alias to the same physical location see the effects of prior stores to the same physical memory location.

To support advanced loads in the presence of a virtual alias, the processor ensures that the Advanced Load Address Table (ALAT) is resolved using physical addresses and is coherent with physical memory. For details, please refer to “Detailed Functionality of the ALAT and Related Instructions” on page 1:65.

## 4.2 Physical Addressing

Objects in memory and I/O occupy a common 63-bit physical address space that is accessed using byte addresses. Accesses to physical memory and I/O may be performed via virtual addresses mapped to the 63-bit physical address space or by direct physical addressing. Current page table formats allow for mapping virtual addresses into 50 bits of physical address space (on processor implementations that support this many physical address bits). Future extensions to the page table formats will allow larger mappings, up to the full 63 bits of physical address space.

Physical addressing for instruction references (including IA-32) is enabled when PSR.it is 0, data references (including IA-32) when PSR.dt is 0, and register stack references when PSR.rt is 0.

While software views the physical addressing as being 63-bits, implementations may implement between 32 and 63 physical address bits. All processor models must implement a contiguous set of physical address bits starting at bit 32 and continuing upwards. Please see the processor-specific documentation for further information on the number of physical address bits implemented on the Itanium processor. Implementations must validate that memory references are performed to implemented physical address bits. Instruction references to unimplemented physical addresses result either in an Unimplemented Instruction Address trap on the last valid instruction, or in an Unimplemented Instruction Address fault on the instruction fetch of the unimplemented address. Data references to unimplemented physical addresses result in an Unimplemented Data Address fault. Memory references to unpopulated address ranges result in an asynchronous Machine Check abort, when the platform signals a transaction time-out. Exact machine check behavior is model specific.

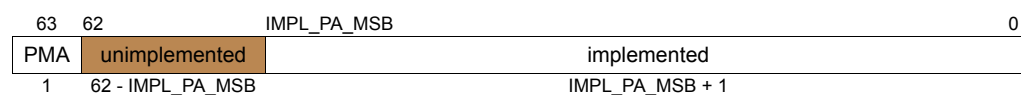
## 4.3 Unimplemented Address Bits

Based on the processor model, some physical and/or virtual address bits may not be implemented. Regardless of the number of implemented address bits, all general purpose, branch, control and application registers implement all 64 register bits on all processors. Similarly, regardless of the number of implemented address bits, data and instruction breakpoint registers must implement all 64 address bits and all 56 mask bits on all processors.

### 4.3.1 Unimplemented Physical Address Bits

As shown in [Figure 4-18](#), a 64-bit physical address consists of three fields: physical memory attribute (PMA), unimplemented and implemented bits.

**Figure 4-18. Physical Address Bit Fields**



All processor models implement at least 32 physical address bits, bits 0 to 31, plus the physical memory attribute bit. Additional implemented physical bits must be contiguous starting at bit 32. IMPL\_PA\_MSB is the implementation-specific position of the most

significant implemented physical address bit. In a processor that implements all physical address bits, IMPL\_PA\_MSB is 62. Please see the processor-specific documentation for further information on the number of physical address bits implemented on the Itanium processor.

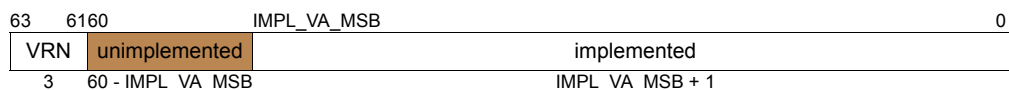
If unimplemented physical address bits are set by software, an Unimplemented Data Address fault is raised during the TLB insert instructions (*itc*, *itr*). Inserts performed by the VHPT walker, as noted in “VHPT Hashing” on page 2:65, abort the VHPT search if unimplemented or reserved fields are used. For translations marked as Not-Present (TLB.p is 0), the processor does not check the validity of PPN and some reserved bits as noted in Figure 4-6.

When a processor model does not implement all physical address bits, the missing bits are defined to be zero. Physical addresses in which bits  $PA\{62:\min(\text{IMPL\_PA\_MSB}+1,62)\}$  are not zero are considered “unimplemented” physical addresses on that processor model. Physical addresses are checked for correctness on use by ensuring that  $PA\{62:\min(\text{IMPL\_PA\_MSB}+1,62)\}$  bits are zero.

### 4.3.2 Unimplemented Virtual Address Bits

As shown in Figure 4-19, a 64-bit virtual address consists of three fields: virtual region number (VRN), unimplemented and implemented bits.

**Figure 4-19. Virtual Address Bit Fields**



All processor models provide three VRN bits in  $VA\{63:61\}$ . IMPL\_VA\_MSB is the implementation-specific bit position of the most significant implemented virtual address bit. In addition to the three VRN bits, all processor models implement at least 54 virtual address bits; i.e., the smallest IMPL\_VA\_MSB is 53. In a processor that implements all 64 virtual address bits IMPL\_VA\_MSB is 60. Please see the processor-specific documentation for further information on the number of virtual address bits implemented on the Itanium processor.

If the PSR.vm bit is implemented, and if PSR.vm is 1, then virtual addresses are treated as though one additional virtual address bit were unimplemented. If the PSR.vm bit is implemented, at least 55 virtual address bits must be implemented.

When a processor model does not implement all virtual address bits, the missing bits are defined to be a sign-extension of  $VA\{\text{IMPL\_VA\_MSB}\}$ . Virtual addresses in which bits  $VA\{60:\min(\text{IMPL\_VA\_MSB}+1,60)\}$  do not match  $VA\{\text{IMPL\_VA\_MSB}\}$  are considered “unimplemented” virtual addresses on that processor model. Virtual addresses are checked for correctness on use by ensuring that  $VA\{60:\min(\text{IMPL\_VA\_MSB}+1,60)\}$  bits are identical to  $VA\{\text{IMPL\_VA\_MSB}\}$ .

### 4.3.3 Instruction Behavior with Unimplemented Addresses

The use of an unimplemented address affects instruction execution as described in the bullet list below. If instruction address translation is enabled, an “unimplemented address” refers to an unimplemented virtual address. If instruction address translation is disabled, an “unimplemented address” refers to an unimplemented physical address.

- Non-speculative memory references (non-speculative loads, stores, and semaphores), the following non-access references: `fc`, `fc.i`, `tpa`, `lfetch.fault`, and `probe.fault`, and mandatory RSE operations to unimplemented addresses result in an Unimplemented Data Address fault.
- Virtual addresses used by instruction and data TLB purge/insert operations are checked, and if the base address (register `r3` of the purge, IFA for inserts) targets an unimplemented virtual address, a Unimplemented Data Address fault is raised. The page size of the insert or purge is ignored.
- Speculative loads from unimplemented addresses always return a NaT bit in the target register.
- A regular\_form `probe` instruction to an unimplemented address returns zero in the target register.
- A `tak` instruction to an unimplemented address returns one in the target register.
- A non-faulting `lfetch` to an unimplemented address is silently ignored.
- Eager RSE operations to unimplemented addresses do not fault.
- Execution of a taken branch, taken `chk`, or an `rfi` to an unimplemented address, or execution of a non-branching slot 2 instruction in a bundle at the upper edge of the implemented address space (where the next sequential bundle address would be an unimplemented address) results either in an Unimplemented Instruction Address trap on the branch, `chk`, `rfi` or non-branching slot 2 instruction, or in an Unimplemented Instruction Address fault on the fetch of the unimplemented address.
- When `ptc.g` or `ptc.ga` operations place a virtual address on the bus, the virtual address is sign-extended to a full 64-bit format. If an incoming `ptc.g` or `ptc.ga` presents a virtual address base that targets an unimplemented virtual address, the upper (unimplemented) virtual address bits are dropped, and the purge is performed with the truncated address.
- The behavior of executing `vmsw.1` in a bundle whose address will become unimplemented after `PSR.vm` is set to 1 is undefined.

## 4.4 Memory Attributes

When virtual addressing is enabled, memory attributes defining the speculative, cacheability and write-policies of the virtually mapped physical page are defined by the TLB. When physical addressing is enabled, memory attributes are supplied as described in “Physical Addressing Memory Attributes” on page 2:76.

### 4.4.1 Virtual Addressing Memory Attributes

For virtual memory references, the memory attribute field of each virtual translation describes physical memory properties as shown in [Table 4-11](#).

**Table 4-11. Virtual Addressing Memory Attribute Encodings**

Attribute	Mnemonic	ma	Cacheability	Write Policy	Speculation	Coherent <sup>a</sup> with Respect to
Write Back	WB	000	Cacheable	Write back	Non-sequential & speculative	WB, WBL
Write Coalescing	WC	110	Uncacheable	Coalescing		Not MP coherent <sup>b</sup>
Uncacheable	UC	100		Non-coalescing	Sequential & non-speculative	UC, UCE
Uncacheable Exported	UCE	101				
Reserved <sup>c</sup>		001				
Reserved		010 011				
NaTPage	NaTPage	111	Cacheable	N/A	Speculative	N/A

- a. The Coherency column in this table refers to multiprocessor coherence on normal, side-effect free memory. The data dependency rules defined in “Memory Access Ordering” on page 1:73 ensure uni-processor coherence for the memory attributes listed in each row.
- b. WC is not MP coherent w.r.t. any memory attribute, but is uni-processor coherent w.r.t. itself.
- c. This memory attribute is reserved for Software use.

The attribute UCE is identical to UC except when executing an `fetchadd` instruction. UCE enables the exporting of the `fetchadd` instruction outside the processor. Support for UCE is model-specific; see “Effects of Memory Attributes on Memory Reference Instructions” on page 2:86 for details.

Insert TLB instructions (`itc`, `itr`) that attempt to insert reserved memory attributes (Table 4-11) into the TLB raise Reserved Register/Field faults. External system operation is undefined if software inserts a memory attribute supported by the processor but not supported by the external system.

If software modifies the memory attributes for a page, it must follow the attribute transition requirements in Section 4.4.11, “Memory Attribute Transition” on page 2:88.

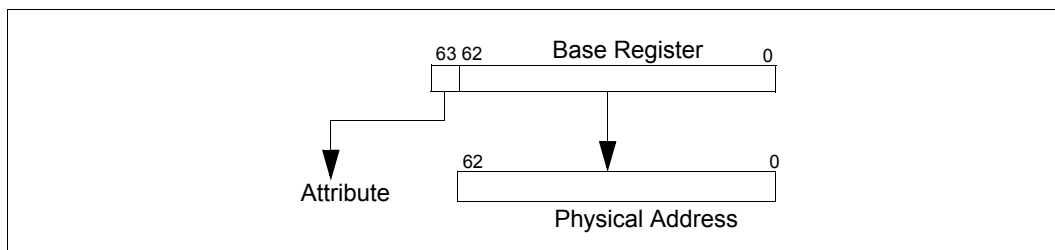
It is recommended that processor models report a Machine Check abort if the following memory attribute aliasing is detected:

- Cache hit on an uncacheable page, other than as the target of a local or remote flush cache (`fc`, `fc.i`) instruction (see “Effects of Memory Attributes on Memory Reference Instructions” on page 2:86).

## 4.4.2 Physical Addressing Memory Attributes

The selection of memory attributes for physical addressing is selected by bit 63 of the address contained in the address base register as shown in Figure 4-20 and Table 4-12.

**Figure 4-20. Physical Addressing Memory**





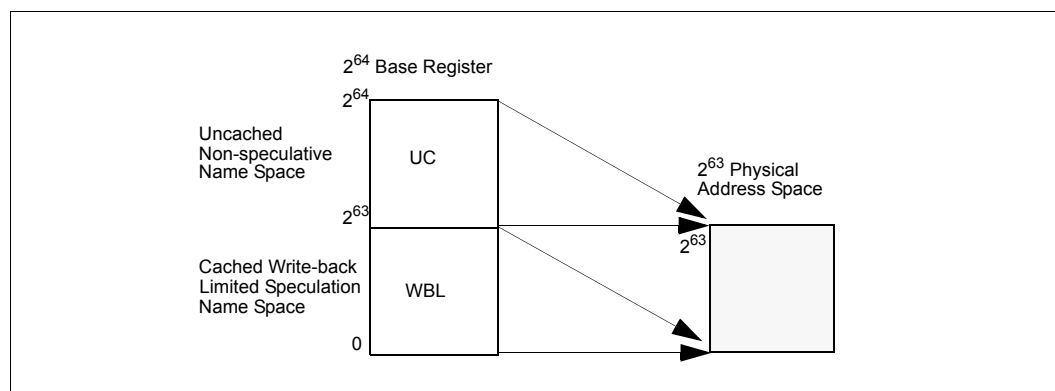
**Table 4-12. Physical Addressing Memory Attribute Encodings**

Bit{63}	Mnemonic	Cacheability	Write Policy	Speculation	Coherent <sup>a</sup> with respect to
0	WBL	Cacheable	Write Back	Non-sequential & limited speculation	WBL, WB
1	UC	Uncached	Non-coalescing	Sequential & non-speculative	UC, UCE

a. Coherency here refers to multiprocessor coherence on normal, side-effect free memory.

See “Speculation Attributes” on page 2:79 for a description of physical addressing limited speculation. Bit{63} is discarded when forming the physical address, effectively creating a write-back name space and an uncached name space as shown in Figure 4-21.

**Figure 4-21. Addressing Memory Attributes**



Software must use the correct name space when using physical addressing; otherwise, I/O devices with side-effects may be accessed speculatively. Physical addressing accesses are ordered only if ordered loads or ordered stores are used. Otherwise, physical addressing memory references are unordered.

### 4.4.3 Cacheability and Coherency Attribute

A page can be either **cacheable** or **uncacheable**. If a page is marked cacheable, the processor is permitted to allocate a local copy of the corresponding physical memory in all levels of the processor memory/cache hierarchy. Allocation may be modified by the cache control hints of memory reference instructions.

A page which is cached is coherent with memory; i.e., the processor and memory system ensure that there is a consistent view of memory from each processor. Processors support multiprocessor cache coherence based on physical addresses between all processors in the coherence domain (tightly coupled multiprocessors). Coherency is supported in the presence of virtual aliases, although software is recommended to use aliases which are an integer multiple of 1 MB apart to avoid any possible performance degradation.

Processors are not required to maintain coherency between processor local instruction and data caches for Itanium architecture-based code; i.e., locally initiated Itanium stores may not be observed by the local instruction cache. Processors are required to

maintain coherency between processor local instruction and data caches for IA-32 code. Instruction caches are also not required to be coherent with multiprocessor Itanium instruction set originated memory references. Instruction caches are required to be coherent with multiprocessor IA-32 instruction set originated memory references. The processor must ensure that transactions from other I/O agents (such as DMA) are physically coherent with the instruction and data cache.

For non-cacheable references the processor provides no coherency mechanisms; the memory system must ensure that a consistent view of memory is seen by each processor. See [“Coalescing Attribute” on page 2:78](#) for a description of coherency for the coalescing memory attribute.

#### 4.4.4 Cache Write Policy Attribute

Write-back cacheable pages need only modify the processor’s copy of the physical memory location; written data need only be passed to the memory system when the processor’s copy is displaced, or a Flush Cache (`fc`) instruction is issued to flush a virtual address. A cache line can only be written back to memory if a store, semaphore (successful or not), the `ld.bias`, a mandatory RSE store, or a `.excl` hinted lfetch instruction targeting that line has executed without a fault. These events enable write-backs. A synchronized `fc` instruction disables subsequent write-backs (after the line has been flushed).

As described in [“Invalidating ALAT Entries” on page 1:67](#), platform visible removal of cache lines from a processor’s caches (e.g., cache line write-backs or platform visible replacements) cause the corresponding ALAT entries to be invalidated.

#### 4.4.5 Coalescing Attribute

For uncacheable pages, the **coalescing** attribute informs the processor that multiple stores to this page may be collected in a coalescing buffer and issued later as a single larger merged transaction. The processor may accumulate stores for an indefinite period of time. Multiple pending loads may also be coalesced into a single larger transaction which is placed in a coalescing buffer. Coalescing is a performance hint for the processor; a processor may or may not implement coalescing.

A processor with multiple coalescing buffers must provide a flush policy that flushes buffers at roughly equal rate even if some buffers are only partially full. The processor may make coalesced buffer flushes visible in any order. Furthermore, individual bytes within a single coalesced buffer may be flushed and made visible in any order.

Stores (including IA-32), which are coalesced, are performed out of order; coalescing may occur in both the space and time domains. For example, a write to bytes 4 and 5 and a write to bytes 6 and 7 may be coalesced into a single write of bytes 4, 5, 6, and 7. In addition, a write of bytes 5 and 6 may be combined with a write of bytes 6 and 7 into a single write of bytes 5, 6, and 7.

Any release operation (regardless of whether it references a page with a coalescing memory attribute), or any fence type instruction, forces write-coalesced data to be flushed and made visible prior to the instruction itself becoming visible. (See [Table 4-15 on page 2:83](#) for a list of release and fence instructions.) Any IA-32 serializing instruction, or access to an uncached memory type, forces write-coalesced data to

become flushed and made visible prior to itself becoming visible. Even though IA-32 stores and loads are ordered, the write-coalesced data is not flushed unless the IA-32 stores or loads are to uncached memory types.

The Flush Cache ( $fc, fc.i$ ) instruction flushes all write-coalesced data whose address is within at least 32 bytes of the 32-byte aligned address specified by the Flush Cache ( $fc, fc.i$ ) instruction, forcing the data to become visible. The Flush Cache ( $fc, fc.i$ ) instruction may also flush additional write-coalesced data. The Flush Write buffers ( $fwb$ ) instruction is a “hint” to the processor to expedite flushing (visibility) of any pending stores held in the coalescing buffer(s), without regard to address.

No indication is given when the flushing of the stores is completed. An  $fwb$  instruction does not ensure ordering of coalesced stores, since later stores may be flushed before prior stores. To ensure prior coalesced stores are made visible before later stores, software must issue a release operation between stores.

The processor may at any time flush coalesced stores in any order before explicitly requested to do so by software.

Coalesced pages are not ensured to be coherent with other processors’ coalescing buffers or caches, or with the local processor’s caches. Loads to coalesced memory pages by a processor see the results of all prior stores by the same processor to the same coalesced memory page. Memory references made by the coalescing buffer (e.g., buffer flushes) have an unordered non-sequential memory ordering attribute. See [“Sequentiality Attribute and Ordering” on page 2:82](#).

Data that has been read or prefetched into a coalescing buffer prior to execution of an Itanium acquire or fence type instruction is invalidated by the acquire or fence instruction. (See [Table 4-15](#) for a list of acquire and fence instructions.)

#### 4.4.6 Speculation Attributes

For present pages (TLB.p=1) which are marked with a **speculative** or a NaTPage memory attribute, the processor may prefetch instructions (including IA-32), perform address generation and perform load accesses (including IA-32) without resolving prior control dependencies, including predicates, branches and interruptions. A page should only be marked speculative if accesses to that page have no side-effects. For example, many memory-mapped I/O devices have side-effects associated with reads and should be marked non-speculative. If a page is marked speculative, a processor can read any location in the page at any time independent of a programmer’s intentions or control flow changes. As a result, software is required, at all times, to maintain valid page table attributes for the ppn, ps and ma fields of all present translations whose memory attribute is speculative or NaTPage. (For example, software should not insert into the TLB, nor create in the VHPT, mappings whose memory attribute is WB, WC or NaTPage unless the entire corresponding physical address range is populated. Placing such mappings in the VHPT or inserting such mappings in the TLB could result in machine check aborts.) High-performance operation is only attainable on speculative pages. The speculative attribute is a hint; a processor may behave non-speculatively.

Prefetches are enabled if a speculative translation exists. Prefetches are asynchronous data and instruction memory accesses that appear logically to initiate and finish between some pair of instructions. This access may not be visible to subsequent flush cache ( $fc, fc.i$ ) and/or TLB purge instructions. This behavior is implementation-dependent.

The processor will not initiate memory references (16-byte instruction bundle fetches, IA-32 instruction fetches, RSE fills and spills, VHPT references, and data memory accesses) to non-speculative pages until all previous control dependencies (predicates, branches, and exceptions) are resolved; i.e., the memory reference is required by an in-order execution of the program. Additionally, for references to non-speculative pages, the processor:

- May not generate any memory access for a control or data speculative data reference.
- Will generate exactly one memory access for each aligned, non-speculative data reference. (Misaligned data references may cause multiple memory accesses, although these accesses are guaranteed to be non-overlapping – each byte will be accessed exactly once.)
- May generate multiple 16-byte memory accesses (to the same address) for each 16-byte instruction bundle fetch reference.

To ensure virtual and physical accesses to non-speculative pages are performed in program order and only once per program order occurrence, the rules in [Table 4-13](#) and [Table 4-14](#) are defined. Software should also ensure that RSE spill/fill transactions are not performed to non-speculative memory that may contain I/O devices; otherwise, system behavior is undefined.

**Table 4-13. Permitted Speculation**

Memory Attribute	Load (ld) <sup>a</sup>	Speculative Load (ld.s) <sup>b</sup>	Advanced Load (ld.a)	Speculative Advanced Load (ld.sa)	Hardware-generated Speculative References <sup>c</sup>
Speculative	Yes	Yes	Yes	Yes	Yes
Non-speculative	Yes	Always Fail	Always Fail	Always Fail	Prohibited
Limited Speculation	Yes	Always Fail	Yes	Always Fail	Limited <sup>d</sup>

a. Includes the faulting form of line prefetch ( $lfetch.fault$ ).

b. Includes the non-faulting form of line prefetch ( $lfetch$ ), which does not cause a cache fill if the memory attribute is non-speculative or limited speculation.

c. Hardware-generated speculative references include non-demand instruction prefetches (including IA-32), hardware-generated data prefetch references, and eager RSE memory references.

d. The processor may only issue hardware-generated speculative references to a 4K-byte physical page if it is a verified page.

**Table 4-14. Register Return Values on Non-faulting Advanced/Speculative Loads**

Memory Attribute	Speculative Load (ld.s)		Advanced Load (ld.a)		Speculative Advanced Load (ld.sa)	
	Success	Failure	Success	Failure	Success	Failure
Speculative	Value	NaT <sup>a</sup>	Value	N/a	Value	NaT <sup>a</sup>
Non-speculative	N/A	NaT <sup>b</sup>	N/A	Zero <sup>c</sup>	N/A	NaT <sup>b</sup>
Limited Speculation	N/A	NaT <sup>b</sup>	Value	N/a	N/a	NaT <sup>b</sup>

- a. Speculative or speculative advanced loads that cause deferred exceptions result in failed speculation. The processor aborts the reference. If the target of the load is a GR, the processor sets the register's NaT bit to one. If the target of the load is an FR, the processor sets the target FR to NaTVal. The processor performs all other side-effects (such as post-increment).
- b. Speculative or speculative advanced loads to limited or non-speculative memory pages result in failed speculation. The processor aborts the reference. If the target of the load is a GR, the processor sets the register's NaT bit to 1. If the target of the load is an FR, the processor sets the target FR to NaTVal. The processor performs all other side-effects (such as post-increment).
- c. Advanced loads to non-speculative memory pages always fail. The processor aborts the reference, sets the target register to zero, and performs all other side-effects (such as post-increment).

#### 4.4.6.1 Limited Speculation and the WBL Physical Addressing Attribute

Processors are allowed to reference limited speculation pages (WBL pages) speculatively, in order to increase performance, but this speculation is limited to prevent speculative references to 4Kbyte physical pages for which there is no actual memory (which would cause spurious machine checks).

Processors must not make hardware-generated speculative references to a given WBL 4Kbyte page until a **verified reference** has been made. Processors may optionally implement storage to hold the addresses of WBL 4Kbyte pages for which verified references have been made, and may make subsequent hardware-generated speculative references to these pages. Such pages are termed **verified pages**.

A verified reference is an instruction or data reference made to the page by an in-order execution of the program; that is, a reference which would have been made had the instructions from the program been fetched and executed one at a time. A hardware-generated speculative reference does not constitute a verified reference. Hardware-generated speculative references include:

- Instruction fetches when the processor has not yet determined whether prior branches were predicted correctly
- Instruction fetches when the processor has not yet determined whether prior instructions will raise faults or traps
- Data references by instructions when the processor has not yet determined whether prior branches were predicted correctly
- Data references by instructions when the processor has not yet determined whether prior instructions will raise faults or traps
- Hardware-generated instruction prefetch references
- Hardware-generated data prefetch references
- Eager RSE data references

For an instruction fetch to constitute a verified reference, it must only be determined that an in-order execution of the program requires that the IP point to this address, independent of whether the instruction at this address will subsequently take a fault or interrupt.

For a data reference to constitute a verified reference, the instruction must meet one of the following requirements:

- It executes without any fault or interrupt
- It takes an Unaligned Data Reference fault
- It takes a Data Debug fault

- It takes an External interrupt, but if it had not taken an External interrupt, it would have met one of the above qualifications (execute without fault, take an Unaligned Data Reference fault, or take a Data Debug fault)

Data-speculative loads are treated the same as normal loads, and if an in-order execution of the program requires the execution of a data speculative load, it constitutes a verified reference. Control-speculative loads to limited-speculation pages always defer and thus never constitute verified references.

It is not necessary for a processor to determine whether a reference will complete without generating a machine check for it to be a verified reference. If software actually references a physical address which will cause a machine check, hardware may generate multiple speculative references to the same page, potentially causing multiple machine checks.

Processors may access verified pages normally, as they would WB pages, including the use of caching, pipelining and hardware-generate speculative references to improve performance.

Calling the PAL\_PREFETCH\_VISIBILITY procedure forces the processor to clear the storage holding the addresses of verified pages.

#### 4.4.7 Sequentiality Attribute and Ordering

Memory ordering is defined in [Section 4.4.7, “Memory Access Ordering” on page 1:73](#). This section defines additional ordering rules for non-cacheable memory, cache synchronization (`sync.i`) and global TLB purge operations (`ptc.g`, `ptc.ga`).

As described in [Section 4.4.7, “Memory Access Ordering” on page 1:73](#), read-after-write, write-after-write, and write-after-read dependencies to the same memory location (memory dependency) are performed in program order by the processor. Otherwise, all other memory references may be performed in any order unless the reference is specifically marked as ordered. No ordering exists between instruction accesses and data accesses or between any two instruction accesses. IA-32 memory references follow a stronger processor consistency memory model. See [“IA-32 Memory Ordering” on page 2:265](#) for IA-32 memory ordering details. Explicit ordering takes the form of a set of Itanium instructions: ordered load and check load (`ld.acq`, `ld.c.clr.acq`), ordered store (`st.rel`), semaphores (`cmpxchg`, `xchg`, `fetchadd`), memory fence (`mf`), synchronization (`sync.i`) and global TLB purge (`ptc.g`, `ptc.ga`). The `sync.i` instruction is used to maintain an ordering relationship between instruction and data caches on local and remote processors. The global TLB purge instructions maintain multiprocessor TLB coherence.

For VHPT walks, visibility is defined by the memory read(s) which retrieves translation information, and the associated insertion of the translation into the TLB. VHPT walks are performed asynchronously with respect to program execution, and each walker VHPT read (which appears as though it were performed atomically) is made visible at some single point in the program order. Ordering constraints from [Table 4-15](#) do not prevent VHPT walks from becoming visible.

Table 4-15 defines a set of “Orderable Instructions” that follow one of four ordering semantics: **unordered**, **release**, **acquire** or **fence**. The table defines the ordering semantics and the instructions of each category. Only these Itanium instructions can be used to establish multiprocessor ordering relations.

In the following discussion, the terms **previous** and **subsequent** are used to refer to the program specified order. The term **visible** is used to refer to all architecturally visible effects of performing an instruction. For memory accesses and semaphores this involves at least reading or writing memory. For `mf.a`, visibility is defined by platform acceptance of previous memory accesses. Visibility of `sync.i` is defined by visibility of previous flush cache (`fc`, `fc.i`) operations. For ALAT lookups (`ld.c`, `chk.a`), visibility is determination of ALAT hit or miss. For global TLB purge operations, visibility is defined by removal of an address translation from the TLBs on all processors in the TLB coherence domain. Global TLB purge instructions (`ptc.g` and `ptc.ga`) follow release semantics on the local processor. They are also broadcast to all other processors in the TLB coherence domain. On each such remote processor, a point is chosen in its program-order execution and a local TLB purge operation is inserted at that point; this local TLB purge operation follows release semantics, except with respect to global purge instructions being executed by that remote processor. For local TLB purge operations, visibility is defined by removal of an address translation on the local processor. Local TLB purge instructions (`ptc.l`, `ptc.e`) ensure that all prior stores are made locally visible before the actual purge operation is performed.

**Table 4-15. Ordering Semantics and Instructions**

Ordering Semantics	Description	Orderable Intel® Itanium® Instructions
Unordered	Unordered instructions may become visible in any order.	<code>ld</code> , <code>ld.s</code> , <code>ld.a</code> , <code>ld.sa</code> , <code>ld.fill</code> , <code>ldf</code> , <code>ldf.s</code> , <code>ldf.sa</code> , <code>ldf.fill</code> , <code>ldfp</code> , <code>ldfp.s</code> , <code>ldfp.sa</code> , <code>st</code> , <code>st.spill</code> , <code>stf</code> , <code>stf.spill</code> , <code>mf.a</code> , <code>sync.i</code> , <code>ld.c</code> , <code>chk.a</code>
Release	Release instructions guarantee that all previous orderable instructions are made visible prior to being made visible themselves.	<code>cmp8xchg16.rel</code> , <code>cmpxchg.rel</code> , <code>fetchadd.rel</code> , <code>st.rel</code> , <code>ptc.g</code> , <code>ptc.ga</code>
Acquire	Acquire instructions guarantee that they are made visible prior to all subsequent orderable instructions.	<code>cmp8xchg16.acq</code> , <code>cmpxchg.acq</code> , <code>fetchadd.acq</code> , <code>xchg</code> , <code>ld.acq</code> , <code>ld.c.clr.acq</code>
Fence	Fence instructions combine the release and acquire semantics into a bi-directional fence; i.e., they guarantee that all previous orderable instructions are made visible prior to any subsequent orderable instruction being made visible.	<code>mf</code>

Itanium memory accesses to **sequential** pages occur in program order with respect to all other sequential pages in the same peripheral domain, but are not necessarily ordered with respect to non-sequential page accesses. A peripheral domain is a platform-specific collection of uncacheable addresses. An I/O device is normally contained in a peripheral domain and all sequential accesses from one processor to that device will be ordered with respect to each other. Sequentiality ensures that uncacheable, non-coalescing memory references from one processor to a peripheral domain reach that domain in program order. Sequentiality does not imply visibility.

Inter-Processor Interrupt Messages (8-byte stores to a Processor Interrupt Block address, through a UC memory attribute) are exceptions to the sequential semantics. IPI's are not ordered with respect to other IPI's directed at the same processor. Further, fence operations do not enforce ordering between two IPI's. See [Section 5.8.4.2, "Interrupt and IPI Ordering"](#) on page 2:130.

[Table 4-16](#) defines the ordering between unordered, release, acquire and fence type operations to sequential and non-sequential pages. [Table 4-16](#) defines the minimal ordering requirements; an implementation may enforce more restrictive ordering than required by the architecture. The actual mechanism for enforcing memory access ordering is implementation dependent.

**Table 4-16. Ordering Semantics**

First Operation		Second Operation						
		Fence	Non-sequential			Sequential <sup>a</sup>		
			Acquire	Release	Unordered	Acquire	Release	Unordered
	Fence	O	O	O	O	O	O	O
Non-sequential	Acquire	O	O	O	O	O	O	O
	Release	O	–	O	–	–	O	–
	Unordered	O	–	O	–	–	O	–
Sequential <sup>a</sup>	Acquire	O	O	O	O	OS	OS	OS
	Release	O	–	O	–	S	OS	S
	Unordered	O	–	O <sup>b</sup>	– <sup>c</sup>	S <sup>d</sup>	OS <sup>e</sup>	S

- a. Except for IPI.
- b. "O" indicates that the first and second operation become visible in program order.
- c. A dash indicates no ordering is implied.
- d. "S" indicates that the first and the second operation reach a peripheral domain in program order.
- e. "OS" implies that both "O" and "S" ordering relations apply.

[Table 4-16](#) establishes an order between operations on a particular processor. For operations to cacheable write-back memory the order established by these rules is observed by all observers in the coherence domain.

For example, when this sequence is executed on a processor:

```
st [a]
st.rel [b]
```

and a second processor executes this sequence:

```
ld.acq [b]
ld [a]
```

if the second processor observes the store to [b], it will also observe the store to [a].

Unless an ordering constraint from [Table 4-16](#) prevents a memory read<sup>1</sup> from becoming visible, the read may be satisfied with values found in a store buffer (or any logically equivalent structure). These values need not be globally visible even when the operation that created the value was a `st.rel`. This local bypassing behavior may make

---

1. This includes all types of loads (`ld` and `ld.acq`), and RSE memory reads. Note, however, that the read operation of semaphores cannot be satisfied with values found in a store buffer.



accesses of different sizes but with overlapping memory references appear to complete non-atomically. To ensure that a memory write is globally observed prior to a memory read, software must place an explicit fence operation between the two operations.

Aligned `st.rel` and semaphore operations<sup>1</sup> from multiple processors to cacheable write-back memory become visible to all observers in a single total order (i.e., in a particular interleaving; if it becomes visible to any observer, then it is visible to all observers), except that for `st.rel` each processor may observe (via `ld` or `ld.acq`) its own update prior to it being observed globally.

The Itanium architecture ensures this single total order only for aligned `st.rel` and semaphore operations to cacheable write-back memory. Other memory operations<sup>2</sup> from multiple processors are not required to become visible in any particular order, unless they are constrained w.r.t. each other by the ordering rules defined in [Table 4-16](#).

Ordering of loads is further constrained by data dependency. That is, if one load reads a value written by an earlier load by the same processor (either directly or transitively, through either registers or memory), then the two loads become visible in program order.

For example, when this sequence is executed on a processor:

```
st [a] = data
st.rel [b] = a
```

and a second processor executes this sequence:

```
ld x = [b]
ld y = [x]
```

if the second processor observes the store to `[b]`, it will also observe the store to `[a]`.

Also for example, when this sequence is executed on a processor:

```
st [a]
st.rel [b] = 'new'
```

and a second processor executes this sequence:

```
ld x = [b]
cmp.eq p1 = x, 'new'
(p1) ld y = [a]
```

if the second processor observes the store to `[b]`, it will also observe the store to `[a]`.

And for example, when this sequence is executed on a processor:

```
st [a]
st.rel [b] = 'new'
```

and a second processor executes this sequence:

- 
1. Both acquire and release semaphore forms
  2. e.g. unordered stores, loads, `ld.acq`, or memory operations to pages with attributes other than write-back cacheable.

```

        ld x = [b]
        cmp.eq p1 = x, 'new'
(p1)   br target
        ...
target:
        ld y = [a]

```

if the second processor observes the store to [b], it will also observe the store to [a].

The flush cache (*fc*, *fc.i*) instruction follows data dependency ordering. *fc* and *fc.i* are ordered only with respect to previous and subsequent load, store, or semaphore instructions to the same line, regardless of the specified memory attribute. Subsequent memory operations to the same line need not wait for prior *fc* or *fc.i* completion before being globally visible. *fc* and *fc.i* are not ordered with respect to memory operations to different lines. *mf* does not ensure visibility of *fc* and *fc.i* operations. Instead, the *sync.i* instruction synchronizes *fc* and *fc.i* instructions, and the *sync.i* is made visible using an *mf* instruction.

#### 4.4.8 Not a Thing Attribute (NaTPage)

A NaTPage attribute prevents non-speculative references to a page, and ensures that speculative references to the page always defer the Data NaT Page Consumption fault. However, as described in ["Speculation Attributes" on page 2:79](#), the processor may issue memory references to a NaTPage. As a result, all NaTPages must be backed by a valid physical page.

Speculative or speculative advanced loads to pages marked as a NaTPage cause the deferred exception indicator (NaT or NaTVal) to be written to the load target register, and the memory reference is aborted. However, all other effects of the load instruction such as post-increment are performed. Instruction fetches, loads, stores and semaphores (including IA-32), but except for Itanium speculative loads, pages marked as NaTPage raise a NaT Page Consumption fault.

A speculative reference to a page marked as NaTPage may still take lower priority faults, if not explicitly deferred in the DCR. See ["Deferral of Speculative Load Faults" on page 2:105](#).

#### 4.4.9 Effects of Memory Attributes on Memory Reference Instructions

Memory attributes affect the following Itanium instructions.

- *ldfe*, *stfe*: Hardware support for 10-byte memory accesses to a page that is neither a cacheable page with write-back write policy nor a NaTPage is optional. On processor implementations that do not support such accesses, an Unsupported Data Reference Fault is raised when an unsupported reference is attempted. For extended floating-point loads the fault is delivered only on the normal, advanced, and check load flavors (*ldfe*, *ldfe.a*, *ldfe.c.nc*, *ldfe.c.clr*). Control speculative flavors of the *ldfe* instruction that target pages that are not cacheable with write-back policy always defer the fault. Refer to ["Deferral of Speculative Load Faults" on page 2:105](#) for details.
- *cmpxchg* and *xchg*: These instructions are only supported to cacheable pages with write-back write policy. *cmpxchg* and *xchg* accesses to NaTPages causes a Data NaT

Page Consumption fault. `cmpxchg` and `xchg` accesses to pages with other memory attributes cause an Unsupported Data Reference fault.

- `fetchadd`: The `fetchadd` instruction can be executed successfully only if the access is to a cacheable page with write-back write policy or to a UCE page. `fetchadd` accesses to NaTPages cause a Data NaT Page Consumption fault. Accesses to pages with other memory attributes cause an Unsupported Data Reference fault. When accessing a cacheable page with write-back write policy, atomic fetch and add operation is ensured by the processor cache-coherence protocol. For highly contended semaphores, the cache line transactions required to guarantee atomicity can limit performance. In such cases, a centralized “fetch and add” semaphore mechanism may improve performance. If supported by the processor and the platform, the UCE attribute allows the processor to “export” the `fetchadd` operation to the platform as an atomic “fetch and add.” Effects of the exported `fetchadd` are platform dependent. If exporting of `fetchadd` instructions is not supported by the processor, a `fetchadd` instruction to a UCE page takes an Unsupported Data Reference fault.
- Flush Cache Instructions – `fc` instructions must always be “broadcast” to other processors, independent of the memory attribute in the local processor. It is legal to use an uncacheable memory attribute for any valid address when used as a flush cache (`fc`) instruction target. This behavior is required to enable transitions from one memory attribute to another and in case different memory attributes are associated with the address in another processor.
- Prefetch instructions – `lfetch` and any implicit prefetches to pages that are not cacheable are suppressed. No transaction is initiated. This allows programs to issue prefetch instructions even if the program is not sure the memory is cacheable.

#### 4.4.10 Effects of Memory Attributes on Advanced/Check Loads

The ALAT behavior of advanced and check loads is dependent on the memory attribute of the page referenced by the load. These behaviors are required; advanced and check load completers are not hints.

All speculative pages have identical behavior with respect to the ALAT. Advanced loads to speculative pages always allocate an ALAT entry for the register, size, and address tuple specified by the advanced load. Speculative advanced loads allocate an ALAT entry if the speculative load is successful (i.e., no deferred exception); if the speculative advanced load results in a deferred exception, any matching ALAT entry is removed and no new ALAT entry is allocated. Check loads with clear completers (`ld.c.clr`, `ld.c.clr.acq`, `ldf.c.clr`) remove a matching ALAT entry on ALAT hit and do not change the state of the ALAT on ALAT miss. Check loads with no-clear completers (`ld.c.nc`, `ldf.c.nc`) allocate an ALAT entry on ALAT miss. On ALAT hit, the ALAT is unchanged if an exact ALAT match is found (register, address, and size); a new ALAT entry with the register, address, and size specified by the no-clear check load may be allocated if a partial ALAT match is found (match on register).

Advanced loads (speculative or non-speculative variants) to non-speculative pages always remove any matching ALAT entry. Check loads to non-speculative pages that miss the ALAT never allocate an ALAT entry, even in the case of a no-clear check load. ALAT hits on check loads to non-speculative pages (which can occur if a previous advanced load referenced that page via a speculative memory attribute) result in

undefined behavior; when changing an existing page from speculative to non-speculative (or vice-versa), software should ensure that any ALAT entries corresponding to that page are invalidated.

Limited speculation pages behave like non-speculative pages with respect to speculative advanced loads, and behave like speculative pages with respect to all other advanced and/or check loads.

Table 4-17 describes the ALAT behavior of advanced and check loads for the different speculation memory attributes.

**Table 4-17. ALAT Behavior on Non-faulting Advanced/Check Loads**

Memory Attribute	Id.sa Response		Id.a Response	Id.c.clr, Id.c.clr.acq, Id.f.c.clr Response		Id.c.nc, Id.f.c.nc Response	
	No NaT	NaT		ALAT Hit	ALAT Miss	ALAT Hit	ALAT Miss
speculative	alloc	remove	alloc	remove	nop	unchanged <sup>a</sup>	alloc
non-speculative	N/A	remove	remove	undefined	nop	undefined	must not alloc
limited speculation	N/A	remove	alloc	remove	nop	unchanged <sup>a</sup>	alloc

a. May allocate a new ALAT entry if size and/or address are different than the corresponding Id.a or Id.sa whose ALAT entry was matched.

## 4.4.11 Memory Attribute Transition

If software modifies the memory attributes for a page, it must perform explicit actions to ensure that subsequent reads and writes using the new attribute will be coherent with prior reads and writes that were performed with the old attribute. Processors may have separate buffers for coalescing, uncacheable and cacheable references, and these buffers need not be coherent with each other.

### 4.4.11.1 Virtual Addressing Memory Attribute Transition

To change a virtually-addressed page from one attribute to another, software must perform the following sequence. (The address of the page whose attribute is being modified is referred to as "X").

**Note:** This sequence is ONLY required if the new mapping and the old mapping do not have the same memory attribute.

On the processor initiating the transition, perform the following steps 1-3:

1. `PTE[X].p = 0 // Mark page as not present`

This prevents any processors from reading the old mapping (with the old attribute) from the VHPT after this point.

2. `ptc.ga [X] ;; // Global shutdown and ALAT invalidate  
// for the entire page`

This removes the mapping from all processor TC's in the coherence domain, and it forces all processors to flush any pending WC or UC stores from write buffers.

3. 

```
mf ;; // Ensure visibility of ptc.ga to local data stream
srlz.i ;; // Ensure visibility of ptc.ga to local instruction stream
```

After step 3, no processor in the coherence domain will initiate new memory references or prefetches to the old translation. Note, however, that memory references or prefetches initiated to the old translation prior to step 2 may still be in progress after step 3. These outstanding memory references and prefetches may return instructions or data which may be placed in the processor cache hierarchy; this behavior is implementation-specific.

If the new memory attribute is an uncacheable attribute, and if the old attribute was cacheable (or if it is not known at this point in the code sequence what the old attribute was), then software must drain any current prefetches and ensure that any cached data from the page is removed from caches. To do this, software must perform steps 4-10. If the new memory attribute is cacheable, then software may skip steps 4-10, and go straight to step 11.

4. Call PAL\_PREFETCH\_VISIBILITY

Call PAL\_PREFETCH\_VISIBILITY with the input argument *trans\_type* equal to zero to indicate that the transition is for virtual memory attributes. The return argument from this procedure informs the caller if this procedure call is needed on remote processors or not. If this procedure call is not needed on remote processors, then software may skip the IPI in step 5 and go straight to step 6 below.

5. Using the IPI mechanism defined in [“Inter-processor Interrupt Messages” on page 2:128](#) to reach all processors in the coherence domain, perform step 4 above on all processors in the coherence domain, and wait for all PAL\_PREFETCH\_VISIBILITY calls to complete on all processors in the coherence domain before continuing.

After steps 4 and 5, no more new instruction or data prefetches will be made to page “X” by any processor in the coherence domain. However, processor caches in the coherence domain may still contain “stale” data or instructions from prior prefetch or memory references to page “X.”

6. Insert a temporary UC translation for page “X.”

7. 

```
fc [X] // flush all processor caches in the coherence domain
fc [X+32]
fc [X+64]
... // ... for all of page “X” (page size = ps)
fc [X+ps-32] ;;
```

```
// Ensure cache flushes are also seen by processors' instruction
fetch
sync.i ;;
```

After step 7, all flush cache instructions initiated in step 7 are visible to all processors in the coherence domain, i.e., no processor in the coherence domain will respond with a cache hit on a memory reference to an address belonging to page “X.”

8. Purge the temporary UC translation from the TLB

9. Call PAL\_MC\_DRAIN
10. Using the IPI mechanism defined in “Inter-processor Interrupt Messages” on page 2:128 to reach all processors in the coherence domain, perform step 9 above on all processors in the coherence domain, and wait for all PAL\_MC\_DRAIN calls to complete on all processors in the coherence domain before continuing.

This further guarantees that any cache lines containing addresses belonging to page [X] have been evicted from all caches in the coherence domain and forced onto the bus. Note that this operation does not ensure that the cache lines have been written back to memory.

11. Insert the new mapping with the new memory attribute

#### 4.4.11.2 Physical Addressing Attribute Transition – Disabling Prefetch/Speculation and Removing Cacheability

When a verified reference is made to a physical address with the WBL attribute, the 4K page containing that address becomes speculatively accessible. This allows the processor that made the verified reference to subsequently make speculative references to this page. (See the description of limited speculation in Section 4.4.6.1, “Limited Speculation and the WBL Physical Addressing Attribute” on page 2:81.)

If the same physical memory is then to be accessed with the UC attribute, software must first cause all such 4K pages to no longer be verified pages and flush any cached copies from the cache. Otherwise, an uncacheable reference may hit in cache, causing a Machine Check abort.

On the processor initiating the transition, perform the following steps:

1. Call PAL\_PREFETCH\_VISIBILITY

Call PAL\_PREFETCH\_VISIBILITY with the input argument *trans\_type* equal to one to indicate that the transition is for physical memory attributes. This PAL call terminates the processor’s rights to make speculative references to any limited speculation pages (i.e., it causes all WBL pages to no longer be verified pages – see the discussion on limited speculation in Section 4.4.6.1.)

The return argument from this procedure informs the caller if this procedure call is needed on remote processors or not. If this procedure call is not needed on remote processors, then software may skip the IPI in step 2 and go straight to step 3 below.

2. Using the IPI mechanism defined in “Inter-processor Interrupt Messages” on page 2:128 to reach all processors in the coherence domain, perform step 1 above on all processors in the coherence domain, and wait for all PAL\_PREFETCH\_VISIBILITY calls to complete on all processors in the coherence domain before continuing.

On the processor initiating the disabling process, continue the sequence:

3. 

```
fc [X] // flush all processor caches in the coherence domain
fc [X+32]
fc [X+64]
... // ... for all of page “X” (page size = ps)
fc [X+ps-32] ;;
```

```
// Ensure cache flushes are also seen by processors' instruction
fetch
sync.i ;;
```

After step 3, all flush cache instructions initiated in step 3 are visible to all processors in the coherence domain, i.e., no processor in the coherence domain will respond with a cache line hit on a memory reference to an address belonging to page "X."

4. Call PAL\_MC\_DRAIN.
5. Using the IPI mechanism defined in ["Inter-processor Interrupt Messages" on page 2:128](#) to reach all processors in the coherence domain, perform step 4 above on all processors in the coherence domain, and wait for all PAL\_MC\_DRAIN calls to complete on all processors in the coherence domain before continuing.

This further guarantees that any cache lines containing addresses belonging to page [X] have been evicted from all caches in the coherence domain and forced onto the bus. Note that this operation does not ensure that the cache lines have been written back to memory.

This sequence ensures that speculation and prefetch are disabled for all WBL pages, that all outstanding prefetches have completed, and that the caches have been flushed. It may also be necessary to take additional platform-dependent steps to ensure that all cache write-back transactions have completed to memory before re-configuring physical memory.

#### 4.4.11.3 Memory OLD Attribute Transition Sequence

In order to safely delete a memory range online (memory OLD), all speculative reference and prefetches to that range must be halted and all cache lines returned to the memory being deleted. If this is not done, an MCA could occur if data were to be delivered back to the memory controller after the memory had been removed. Software must perform the sequence shown below to ensure that no MCAs occur.

Before performing the memory OLD sequence shown below, all memory in the range being deleted belonging to firmware (PAL and SAL) must be evacuated, and control of the range given to the OS. If firmware cannot be evacuated from the range, then OLD cannot be done.

On the processor performing the memory OLD operation, perform the following:

1. Remove all mappings to all memory pages in this memory range from the page table. (PTE[X].p=0)
2. For each page which has a mapping in TLB, perform one of the following steps:
  - a. If there are any translations in TRs, perform `ptr.d` or `ptr.i`, depending on whether the translation is for code or data. If it is not known, do both. (This invalidates all TRs, and as a side effect, the mapping from all TCs on the processor.)
  - b. If there are no translations in TRs, perform a `ptc.ga`. (This removes mapping from all TC's and forces processors to flush any pending WC or UC stores from write buffers.)

3. Execute:

```
mf ;;  
srlz.i ;;
```

(This ensures visibility of `ptr.d`, `ptr.i`, or `ptc.ga` to both data and instruction stream, so that no new prefetches will be done to the old translations.)

4. Call `PAL_PREFETCH_VISIBILITY` with the input argument `trans_type` equal to one to indicate that the transition is for all memory attributes. This PAL call terminates the processor's rights to make speculative references to any limited speculation pages (i.e., it causes all WBL pages to no longer be verified pages – see the discussion on limited speculation in [Section 4.4.6.1, "Limited Speculation and the WBL Physical Addressing Attribute" on page 2:81.](#)). It also ensures all prefetches in flight have been completed. The return argument from this procedure informs the caller if this procedure call is needed on remote processors or not. If this procedure call is not needed on remote processors, and step 2.b was used above, then software may skip the IPI in step 5 and go straight to step 6 below.
5. If step 2.a was performed, or if the `PAL_PREFETCH_VISIBILITY` return argument indicated the call must be made on other processors in the coherency domain, then use the IPI mechanism defined in [Section 5.8.4.1, "Inter-processor Interrupt Messages" on page 2:128](#) to reach all processors in the coherency domain. If step 2a was performed, then steps 2 through 4 must be performed on all processors in the coherency domain. Otherwise, only step 4 must be performed. Wait for all `PAL_PREFETCH_VISIBILITY` calls to complete on all processors in the coherency domain before continuing. After step 5, no more new instruction or data prefetches will be made to page "X" by any processor in the coherency domain. However, processor caches in the coherency domain may still contain "stale" data or instructions from prior prefetch or memory references to page "X."
6. Perform one of the following steps:
  - a. Call `PAL_CACHE_FLUSH` with input parameters `cache_type=3` and `operation.inv=1`, or
  - b. On the processor where the OLD was initiated, perform the sequence:
    - i. If the sequence is to be executed with `PSR.dt=1`, then insert a temporary translation for the memory range with the "UC" memory attribute.
    - ii. Execute the following instruction sequence:

```
fc [X] // flush all processor caches in the coherence domain  
fc [X+32]  
fc [X+64]  
... // ... for the memory range being OLDed  
fc [X+ps-32] ;;  
// Ensure cache flushes are also seen  
// by processors' instruction fetch  
sync.i ;;
```
    - iii. If the sequence had been run with `PSR.dt=1`, then remove the temporary translation inserted in step 6.b.i.
7. Call `PAL_MC_DRAIN`.

**Note:** If the memory range being OLDed is much larger than the caches being flushed, option 6.a. may be significantly faster.



8. If `PAL_CACHE_FLUSH` is used to flush caches, it must also be called on all processors in the coherency domain. In any case, `PAL_MC_DRAIN` must be called on all processors. Using the IPI mechanism defined in [Section 5.8.4.1, “Inter-processor Interrupt Messages” on page 2:128](#) to reach all processors in the coherence domain, perform step 6.a, if necessary, and step 7 above in that order on all processors in the coherence domain, and wait for all `PAL_MC_DRAIN` calls to complete on all processors in the coherence domain before continuing. This further guarantees that any cache lines containing addresses belonging to page [X] have been evicted from all caches in the coherence domain and forced onto the platform fabric. Note that this operation does not ensure that the cache lines have been written back to memory.
9. Perform whatever platform dependent actions are necessary to flush any platform caches of any copies of the memory being OLDed and to force all cache lines back to the memory being OLDed. (Note: Refer to platform specific documentation.)

This sequence ensures that speculation and prefetching is disabled for the memory range, regardless of WB or WBL attribute, that all in-flight prefetches are completed, and that all caches lines are returned to memory.

## 4.5 Memory Datum Alignment and Atomicity

All Itanium instruction fetches, aligned load, store and semaphore operations (including IA-32) are atomic, except for floating-point extended memory references (`ldfe`, `stfe`, and IA-32 10-byte memory references) to non-write-back cacheable memory. In some processor models, aligned 10-byte Itanium floating-point extended memory references to non-write-back cacheable memory may raise an Unsupported Data Reference fault. See [“Effects of Memory Attributes on Memory Reference Instructions” on page 2:86](#) for details. Loads are allowed to be satisfied with values obtained from a store buffer (or any logically equivalent structure) where architectural ordering permits, and values loaded may appear to be non-atomic. For details, refer to [“Sequentiality Attribute and Ordering” on page 2:82](#).

Load pair instructions are performed atomically under the following conditions: a 16-byte aligned load integer/double pair is performed as an atomic 16-byte memory reference. An 8-byte aligned load single pair is performed as an atomic 8-byte memory reference.

An aligned `ld16` or `st16` instruction is performed as an atomic 16-byte memory reference. For these instructions, the address specified must be 16-byte aligned. Unaligned `ld16` and `st16` instructions result in an Unaligned Data Reference fault regardless of the state of `PSR.ac`.

Aligned Itanium data memory references never raise an Unaligned Data Reference fault. Minimally, each Itanium instruction and its corresponding template are fetched together atomically. Itanium unordered loads can use the store buffer for data values. See [“Sequentiality Attribute and Ordering” on page 2:82](#) for details.

When `PSR.ac` is 1, any Itanium data memory reference that is not aligned on a boundary the size of the operand results in an Unaligned Data Reference fault; e.g., 1, 2, 4, 8, 10, and 16-byte datums should be aligned on 1, 2, 4, 8, 16, and 16-byte

boundaries respectively to avoid generation of an Unaligned Data Reference fault. When PSR.ac is 1, any IA-32 data memory reference that is not aligned on a boundary the size of the operand results in an IA\_32\_Exception(AlignmentCheck) fault.

**Note:** 10-byte and floating-point load double pair datum alignment is 16-bytes. The alignment of long format 32-byte VHPT references is always 32-bytes.

Unaligned Itanium semaphore references (`cmpxchg`, `xchg`, `fetchadd`) result in an Unaligned Data Reference fault regardless of the state of PSR.ac. For the `cmp8xchg16` instruction, the address specified must be 8-byte aligned.

When PSR.ac is 0, Itanium data memory references that are not aligned may or may not result in an Unaligned Data Reference fault based on the implementation. The level of unaligned memory support is implementation specific. However, all implementations will raise an Unaligned Data Reference fault if the datum referenced by an Itanium instruction spans a 4K aligned boundary, and many implementations will raise an Unaligned Data Reference fault if the datum spans a cache line. Implementations may also raise an Unaligned Data Reference fault for any other unaligned Itanium memory reference. Software is strongly encouraged to align data values to avoid possible performance degradation for both IA-32 and Itanium architecture-based code. When PSR.ac is 0 and IA-32 alignment checks are also disabled, no fault is raised regardless of alignment for IA-32 data memory references.

Unaligned advanced loads are supported, though a particular implementation may choose not to allocate an ALAT entry for an unaligned advanced load. Additionally, the ALAT may “pessimistically” allocate an entry for an unaligned load by allocating a larger entry than the natural size of the datum being loaded, as long as the larger entry completely covers the unaligned address range (e.g. a `ld4.a` to address 0x3 may allocate an 8-byte entry starting at address 0x0). Stores (unaligned or otherwise) may also pessimistically invalidate unaligned ALAT entries.

## §

**Interruptions** are events that occur during instruction processing, causing the flow control to be passed to an interruption handling routine. In the process, certain processor state is saved automatically by the processor. Upon completion of interruption processing, a return from interruption (`rfi`) is executed which restores the saved processor state. Execution then proceeds with the interrupted instruction.

From the viewpoint of response to interruptions, the processor behaves as if it were not pipelined. That is, it behaves as if a single Itanium instruction (along with its template) is fetched and then executed; or as if a single IA-32 instruction is fetched and then executed. Any interruption conditions raised by the execution of an instruction are handled at execution time, in sequential instruction order. If there are no interruptions, the next Itanium instruction and its template, or the next IA-32 instruction, are fetched.

This chapter describes both the IA-32 and Itanium interruption mechanisms as well as the interactions between them. The descriptions of the Itanium interruption vectors and IA-32 exceptions, interruptions, and intercepts are in [Chapter 8](#).

## 5.1 Interruption Definitions

Depending on how an interruption is serviced, interruptions are divided into: IVA-based interruptions and PAL-based interruptions.

- **IVA-based interruptions** are serviced by the operating system. IVA-based interruptions are vectored to the Interruption Vector Table (IVT) pointed to by CR2, the IVA control register (see [“IVA-based Interruption Vectors” on page 2:113](#)).
- **PAL-based interruptions** are serviced by PAL firmware, system firmware, and possibly the operating system. PAL-based interruptions are vectored through a set of hardware entry points directly into PAL firmware (see [Chapter 11, “Processor Abstraction Layer”](#)).

Interruptions are divided into four types: Aborts, Interrupts, Faults, and Traps.

- **Aborts**  
A processor has detected a Machine Check (internal malfunction), or a processor reset. Aborts can be either synchronous or asynchronous with respect to the instruction stream. The abort may cause the processor to suspend the instruction stream at an unpredictable location with partially updated register or memory state. Aborts are PAL-based interruptions.
  - **Machine Checks (MCA)**  
A processor has detected a hardware error which requires immediate action. Based on the type and severity of the error the processor may be able to recover from the error and continue execution. The PALE\_CHECK entry point is entered to attempt to correct the error.
  - **Processor Reset (RESET)**  
A processor has been powered-on or a reset request has been sent to it. The

PALE\_RESET entry point is entered to perform processor and system self-test and initialization.

- **Interrupts**

An external or independent entity (e.g., an I/O device, a timer event, or another processor) requires attention. Interrupts are asynchronous with respect to the instruction stream. All previous instructions (including IA-32) appear to have completed. The current and subsequent instructions have no effect on machine state. Interrupts are divided into Initialization interrupts, Platform Management interrupts, and External interrupts. Initialization and Platform Management interrupts are PAL-based interruptions; external interrupts are IVA-based interruptions.

- **Initialization Interrupts (INIT)**

A processor has received an initialization request. The PALE\_INIT entry point is entered and the processor is placed in a known state.

- **Platform Management Interrupts (PMI)**

A platform management request to perform functions such as platform error handling, memory scrubbing, or power management has been received by a processor. The PALE\_PMI entry point is entered to service the request. Program execution may be resumed at the point of interruption. PMIs are distinguished by unique vector numbers. Vectors 0 through 3 are available for platform firmware use and are present on every processor model. Vectors 4 through 15 are reserved for processor firmware use. See [Section 11.5, "Platform Management Interrupt \(PMI\)"](#) on page 2:310 for details.

- **External Interrupts (INT)**

A processor has received a request to perform a service on behalf of the operating system. Typically these requests come from I/O devices, although the requests could come from any processor in the system including itself. The External Interrupt vector is entered to handle the request. External Interrupts are distinguished by unique vector numbers in the range 0, 2, and 16 through 255. These vector numbers are used to prioritize external interrupts. Two special cases of External Interrupts are Non-Maskable Interrupts and External Controller Interrupts.

- **Non-Maskable Interrupts (NMI)**

Non-Maskable Interrupts are used to request critical operating system services. NMIs are assigned external interrupt vector number 2.

- **External Controller Interrupts (ExtINT)**

External Controller Interrupts are used to service Intel 8259A-compatible external interrupt controllers. ExtINTs are assigned locally within the processor to external interrupt vector number 0.

- **Faults**

The current Itanium or IA-32 instruction which requests an action which cannot or should not be carried out, or system intervention is required before the instruction is executed. Faults are synchronous with respect to the instruction stream. The processor completes state changes that have occurred in instructions prior to the faulting instruction. The faulting and subsequent instructions have no effect on machine state. Faults are IVA-based interruptions.

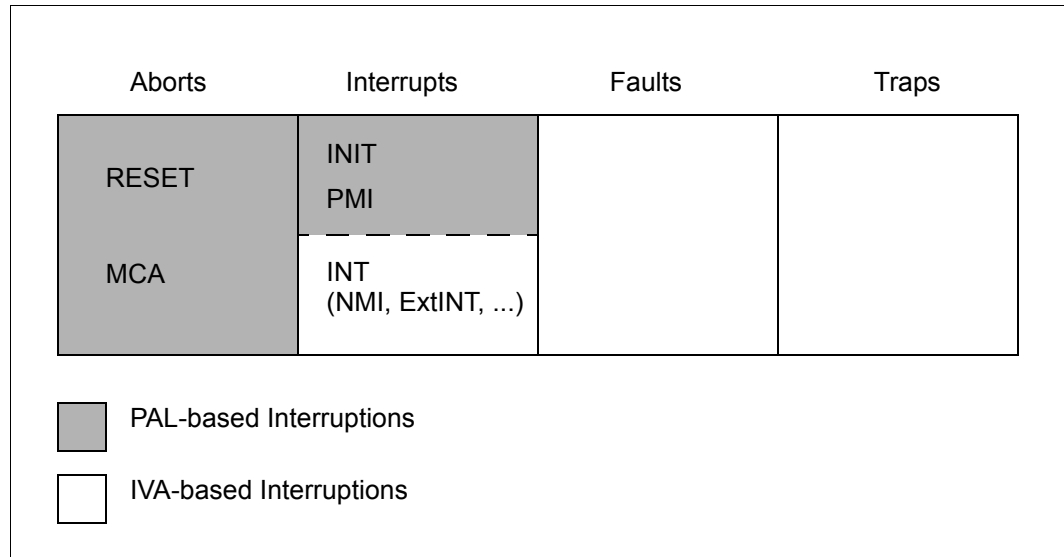
- **Traps**

The IA-32 or Itanium instruction just executed requires system intervention. Traps are synchronous with respect to the instruction stream. The trapping instruction

and all previous instructions are completed. Subsequent instructions have no effect on machine state. Traps are IVA-based interruptions.

Figure 5-1 summarizes the above classification.

**Figure 5-1. Interruption Classification**



Unless otherwise indicated, the term “interruptions” in the rest of this chapter refers to IVA-based interruptions. PAL-based interruptions are described in detail in [Chapter 11](#).

## 5.2 Interruption Programming Model

When an interruption event occurs, hardware saves the minimum processor state required to enable software to resolve the event and continue. The state saved by hardware is held in a set of interruption resources, and together with the interruption vector gives software enough information to either resolve the cause of the interruption, or surface the event to a higher level of the operating system. Software has complete control over the structure of the information communicated, and the conventions between the low-level handlers and the high-level code. Such a scheme allows software rather than hardware to dictate how to best optimize performance for each of the interruptions in its environment. The same basic mechanisms are used in all interruptions to support efficient low-level fault handlers for events such as a TLB fault, speculation fault, or a key miss fault.

On an interruption, the state of the processor is saved to allow a software handler to resolve the interruption with minimal bookkeeping or overhead. The banked general registers (see [“Efficient Interruption Handling” on page 2:102](#)) provide an immediate set of scratch registers to begin work. For low-level handlers (e.g., TLB miss) software need not open up register space by spilling registers to either memory or control registers.

Upon an interruption, asynchronous events such as external interrupt delivery are disabled automatically by hardware to allow software to either handle the interruption immediately or to safely unload the interruption resources and save them to memory. Software will either deal with the cause of the interruption and `rfi` back to the point of the interruption, or it will establish a new environment and spill processor state to memory to prepare for a call to higher-level code. Once enough state has been saved (such as the IIP, IPSR, and the interruption resources needed to resolve the fault) the low-level code can re-enable interruptions by restoring the PSR.ic bit and then the PSR.i bit. (See ["Re-enabling External Interrupt Delivery" on page 2:120.](#)) Since there is only one set of interruption resources, software must save any interruption resource state the operating system may require prior to unmasking interrupts or performing an operation that may raise a synchronous interruption (such as a memory reference that may cause a TLB miss).

The PSR.ic (interruption state collection) bit supports an efficient nested interruption model. Under normal circumstances the PSR.ic bit is enabled. When an interruption event occurs, the various interruption resources are overwritten with information pertaining to the current event. Prior to saving the current set of interruption resources, it is often advantageous in a miss handler to perform a virtual reference to an area which may not have a translation. To prevent the current set of resources from being overwritten on a nested fault, the PSR.ic bit is cleared on any interruption. This will suppress the writing of critical interruption resources if another interruption occurs while the PSR.ic bit is cleared. If a data TLB miss occurs while the PSR.ic bit is zero, then hardware will vector to the Data Nested TLB fault handler.

For a complete description of interruption resources (IFA, IIP, IPSR, ISR, IIM, IIPA, ITIR, IHA, IFS, IIB0-1) see ["Control Registers" on page 2:29.](#)

## 5.3 Interruption Handling during Instruction Execution

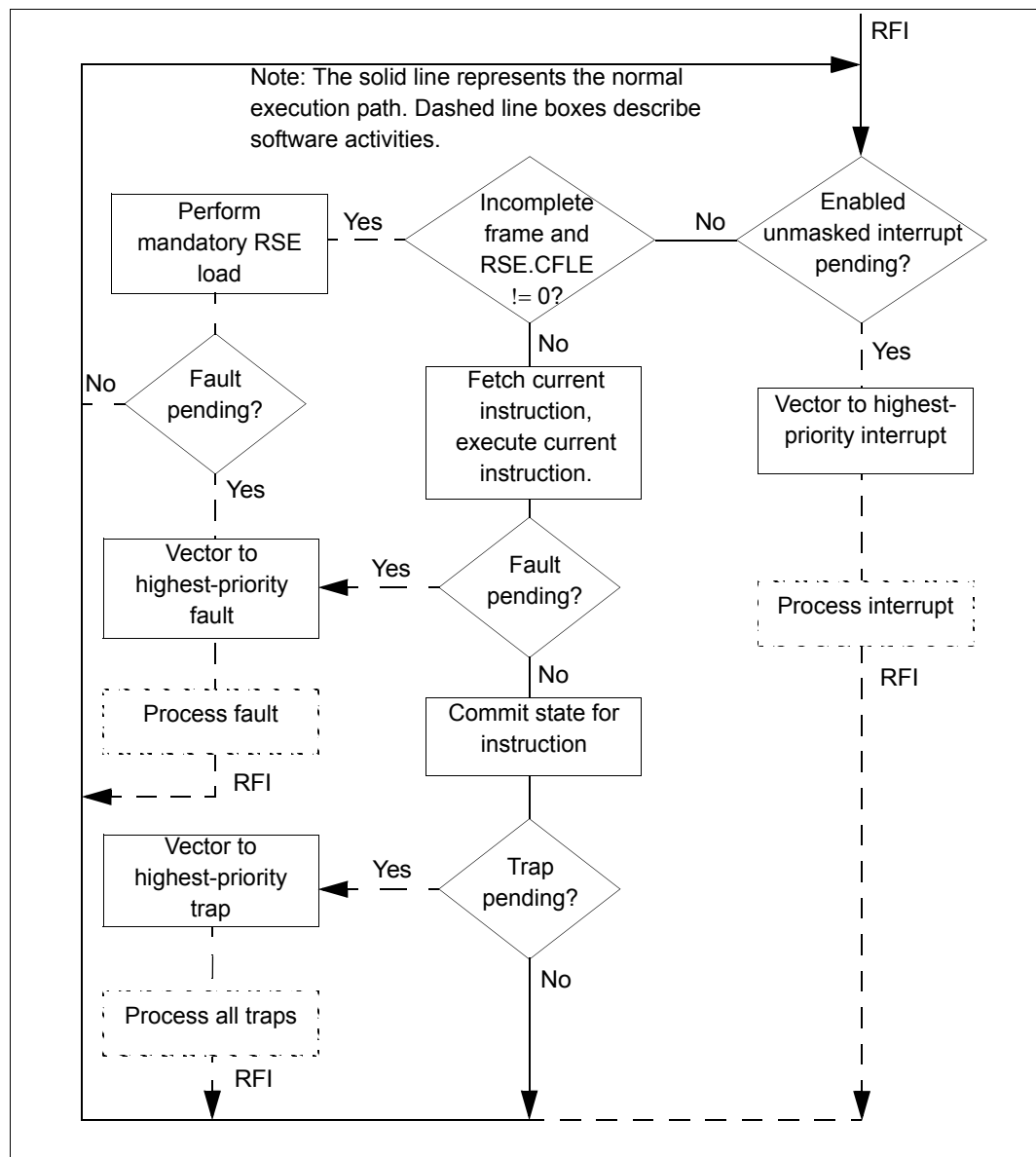
Execution of Itanium instructions involves calculating the address of the current bundle from the region registers and the IP and then fetching, decoding, and executing instructions in that bundle. Execution of IA-32 instructions involves calculating the 64-bit linear address of the current instruction from the EIP, code segment descriptors, and region registers and then fetching, decoding, and executing the IA-32 instruction. (See Section 3.4).

The execution process involves performing the events listed below. The values of the PSR bits are the values that exist before the instruction is executed (except for the case of instructions that are immediately preceded by a mandatory RSE load which clears the PSR.da and PSR.dd bits). Changes to the PSR bits only affect subsequent instructions, and are only guaranteed to be visible by the insertion of the appropriate serializing operation. See ["Serialization" on page 2:17.](#) Execution flow is shown in [Figure 5-2.](#)

1. Resets are always enabled, and may occur anytime during instruction execution.
2. If the PSR.mc bit is 0 then machine check aborts may occur.
3. The processor checks for enabled pending INITs and PMIs, and for enabled unmasked pending external interrupts.

4. For Itanium architecture-based code, the processor checks for a valid register stack frame.
  - If incomplete and RSE Current Frame Load Enable (RSE.CFLE) is set, then perform a mandatory RSE load and start again at step one. The mandatory load operation may fault. A non-faulting mandatory RSE load will clear PSR.da and PSR.dd.
  - If valid, then clear RSE.CFLE.
5. If the processor implements the Unimplemented Instruction Address (UIA) fault, instead of a UIA trap, it will check the instruction address and take the UIA fault if the instruction pointer (IP) falls outside of the implemented range.

**Figure 5-2. Interruption Processing**



6. For IA-32 code, IA-32 instruction addresses are checked for possible instruction

breakpoint faults. The IA-32 effective instruction address (EIP) is converted into a 64-bit virtual linear address IP and IA-32 defined code segmentation and code fetch faults are checked and may result in a fault.

7. When PSR.is is 0, the bundle is fetched using the IP. When PSR.is is 1, an IA-32 instruction is fetched using IP.
  - If the PSR.it bit is 1, virtual address translation of the instruction address is performed. Address translation may result in a fault.
  - If the PSR.pk bit is 1, access key checking is enabled and may result in a fault.
  - For Itanium instructions the IBR registers are checked for possible instruction breakpoint faults.
  - The fetched instruction is decoded and executed.
  - For IA-32 code, the fetched IA-32 instruction is checked to see if the opcode is an illegal opcode, results in an instruction intercept or the opcode bytes are longer than 15 bytes resulting in an fault.
  - If a fault occurs during execution, the processor completes all effects of the instructions prior to the faulting instruction, and does not commit the effect of the faulting instruction and all subsequent instructions. It then takes the interruption for the fault. IIP is loaded with the IP of the bundle or IA-32 instruction which contains the instruction that caused the fault.
  - The PSR.dd, PSR.id, PSR.ia, PSR.da, and PSR.ed bits are set to 0 after an Itanium instruction is successfully executed without raising a fault. The PSR.da and PSR.dd bits are also set to 0 after the execution of each mandatory RSE memory reference that does not raise a fault. PSR.da, PSR.ia, PSR.dd, and PSR.ed bits are cleared before the first IA-32 instruction starts execution after a `br.ia` or `rfi` instruction. EFLAG.rf and PSR.id bits are set to 0 after an IA-32 instruction is successfully executed.
  - If an `rfi` instruction is in the current bundle, then on the execution of `rfi`, the value from the IIP is copied into the IP, the value from IPSR is copied into the PSR, and the RSE.CFLE is set. On an `rfi` if IFS.v is set, then IFS.pfm is copied into CFM and the register stack BOF is decremented by CFM.sof. The following Itanium or IA-32 instruction is executed based on the new IP and PSR values.
8. Traps are handled after execution is complete.
  - If the processor reports unimplemented instruction addresses with an Unimplemented Instruction Address trap (rather than with an Unimplemented Instruction Address fault) and the instruction just completed set the instruction pointer (IP) to an unimplemented address, an Unimplemented Instruction Address trap is taken.
  - If the instruction just completed was an Itanium floating-point instruction which raised a trap, a Floating-point trap is taken.
  - For IA-32 instructions, if Data Breakpoint traps are enabled and one or more data breakpoint registers matched during execution of the instruction, a Data Breakpoint trap is taken.
  - If the PSR.lp bit is 1, and an Itanium branch lowers the privilege level, then a Lower-Privilege Transfer trap is taken.
  - If the PSR.tb bit is 1 and a branch (including IA-32) occurred during execution, then a Taken Branch trap occurs.
  - If no other trap was taken and the PSR.ss bit is 1, then a Single Step trap occurs.



- If more than one trap is triggered (such as Unimplemented Instruction Address trap, Lower-Privilege Transfer trap, and Single Step trap) the highest priority trap is taken. The ISR.code contains a bit vector with one bit set for each trap triggered.

A sequential execution model is presented in the preceding description. Implementations are free to use a variety of performance techniques such as pipelined, speculative, or out-of-order execution provided that, to the programmer, the illusion that instructions are executed sequentially is preserved.

## 5.4 PAL-based Interruption Handling

PAL-based interruption handling requires the processor to transfer control to the PAL firmware. The PAL firmware will execute handling code and set up the architected exit state before transferring control to the SAL firmware. See [Chapter 11, “Processor Abstraction Layer”](#) for more details on the architected exit state between the PAL and SAL firmware layers for PAL-based interruption handling.

It is strongly recommended that software ensure that, if machine check aborts are masked (PSR.mc), external interrupts are also masked (PSR.i). This will avoid cases where a corrected machine check interrupt (a lower priority interrupt) is handled before a machine check abort, which would cause an escalation in machine check abort severity when machine check aborts are unmasked.

## 5.5 IVA-based Interruption Handling

IVA-based interruption handling is implemented as a fast context switch. On IVA-based interruptions, instruction and data translation is left unchanged, the endian mode is set to the system default, and delivery of most PSR-controlled interruptions is disabled (including delivery of asynchronous events such as external interrupts). The processor is responsible for saving only a minimal amount of state in the interruption resource registers prior to vectoring to the Itanium architecture-based software handler.

When an interruption occurs, the processor takes the following actions:

1. If PSR.ic is 0:
  - IPSR, IIP, IIPA, IIB0-1, and IFS.v are unchanged.
  - Interruption-specific resources IFA, IIM, and IHA are unchanged.

If PSR.ic is 1:

- PSR is saved in IPSR. If PSR is in-flight, IPSR will get the most recent in-flight value of PSR (i.e., PSR is serialized by the processor before it is written into IPSR). For Itanium traps, the value written to IPSR.ri is the next instruction slot that would have been executed if there had been no trap. For all other interruptions, the value written to IPSR.ri is the instruction slot on which the interruption occurred (1 for interruptions on the L+X instruction of an MLX). For interruptions in the IA-32 instruction set, IPSR.ri is set to 0.
- IP is written into IIP. For faults and external interrupts, the saved IP is the IP at which the interruption occurred. For traps, the saved IP is the value after the execution of the IA-32 or Itanium instruction which caused the trap. For

branch-related traps, IIP is written with the target of the branch; for all other traps, IIP is written with the address of the bundle or IA-32 instruction containing the next sequential instruction.

- IIPA receives the IP of the last successfully executed Itanium instruction. For IA-32 instructions, IIPA receives the IP of the faulting or trapping IA-32 instruction.
- The interruption resources IFA, IIB0-1, IIM, IHA, and ITIR are written with information specific to the particular fault, trap, or interruption taken. These registers serve as parameters to each of the interruption vectors. The IFS valid bit (IFS.v) is cleared. All other bits in the IFS are undefined.

If PSR.ic is in-flight:

- Interruption state may or may not be collected in IIP, IPSR, IIPA, ITIR, IFA, IIM, IIB0-1 and IHA.
  - The value of IFS (including IFS.v) is undefined.
2. ISR bits are overwritten on all interruptions except for a Data Nested TLB fault. The instruction slot which caused the interruption is saved in ISR.ei (2 for traps, 1 for other interruptions, on the L+X instruction of an MLX). For IA-32 code, ISR.ei is set to 0. If PSR.ic is 0 or in-flight when the interruption occurs, ISR.ni is set to 1. Otherwise, ISR.ni is set to 0. ISR.ni is always 0 for interruptions taken in IA-32 code.
  3. The defined bits in the PSR are set to zero except as follows:
    - PSR.up, PSR.mfl, PSR.mfh, PSR.pk, PSR.dt, PSR.rt, PSR.mc, and PSR.it are unchanged for all interruptions.
    - PSR.be is set to the value of the default endian bit (DCR.be). If DCR.be is in-flight at the time of interruption, PSR.be may receive either the old value of DCR.be or the in-flight value.
    - PSR.pp is set to the value of the default privileged performance monitor bit (DCR.pp). If DCR.pp is in-flight at the time of interruption, PSR.pp may receive either the old value of DCR.pp or the in-flight value.

Since PSR.cpl is set to zero, the processor will execute at the most privileged level.

4. RSE.CFLE is set to zero.
5. IP gets the appropriate IVA vector for the interruption. If IVA is in-flight at the time of interruption, IP receives either the vector specified by the old IVA value or the vector specified by the in-flight value.
6. The processor performs an instruction serialization and execution of Itanium instructions begins at the IP obtained in step 5 above. The instruction serialization event ensures that all previous control register changes and side effects due to such changes are visible to the first instruction of the interruption handler.

### 5.5.1 Efficient Interruption Handling

A set of 16 banked registers are provided by the processor to assist in the efficient processing of low-level Itanium interruptions and instruction emulation. These registers allow a low-level routine to have immediate access to a small set of static registers without having to save and restore their contents to memory at the start and end of each handler. The extra bank of registers exists in the same name space as the normal

registers, overlapping GR16 to GR31. Which set of physical registers are accessed through GR16 to GR31 is determined by the PSR.bn bit. On an interruption this bit is forced to zero allowing access to the alternate set of 16 registers which can be used as scratch space or to hold predetermined values. Software can return to the original set of 16 GRs by setting the PSR.bn bit to one with `bsw` instruction. The `rfi` instruction may also restore the PSR.bn bit to the value at the time of the interruption which is held in the IPSR. Eight additional registers (KR0-KR7) can be used to hold latency critical information for a handler. These application registers (KR0-KR7) can be read but not written by non-privileged code.

When the processor handles an interruption event the current stack frame remains unchanged and the IFS valid bit is cleared. The remaining contents of IFS are undefined. While the interruption handler is running, the register stack engine (RSE) may spill/fill registers to/from the backing store if eager RSE stores/loads are enabled. The RSE will not load or store registers in the current frame (except as required on a `br.ret` or `rfi` in order to load the contents of the frame before continuing execution). For most low-level interruptions the current frame will not be modified. High-performance interruption handlers will not need to perform any register stack manipulation. For example, a TLB miss handler does not need access to any registers in the interrupted frame. An `rfi` instruction after an interruption and before a `cover` operation will also leave the frame marker unchanged (desired behavior for a low-level interruption handler). When an interruption handler falls off the fast path it is required to issue a `cover` instruction so that the interrupted frame can become part of backing store. See [“Switch from Interrupted Context” on page 2:148.](#)

It may be desirable to emulate a faulting instruction in the interruption handler and `rfi` back to the next sequential instruction rather than resuming at the faulting instruction. Some Itanium instructions can be emulated without having to read the bundle from memory, through knowledge of the vector, software convention, and information from the ISR (e.g., emulation of `tpa`). However, most Itanium instructions will require reading the bundle from memory and decoding the operation (e.g., an unaligned load). To correctly emulate an unaligned load, the bundle is read from memory using the value in the IIP which contains the bundle address. The instruction within the bundle that caused the interruption is determined by the ISR.ei field. Once the operation is decoded and emulation completes, the effect of the faulting instruction must be nullified when control is returned to the point of the fault.

An Itanium instruction is skipped by adjusting PSR.ri and possibly IIP prior to performing the `rfi` to the interrupted bundle. This is done by incrementing IPSR.ri by the number of slots this instruction occupies (usually 1). If the resulting IPSR.ri is 3, then reset IPSR.ri to 0 and advance IIP by 1 bundle (16 bytes). Emulating X-unit instructions requires setting IPSR.ri to 0 and setting IIP to the next bundle (X-unit instructions take up two instruction slots). IPSR, IIP, and IFS.pfm (if valid) will be restored on an `rfi` to the PSR, IP, and CFM registers.

## 5.5.2 Non-access Instructions and Interruptions

The non-access Itanium instructions are: `fc`, `fc.i`, `lfetch`, `probe`, `probe.fault`, `tpa`, and `tak`. These instructions reference the TLB but do not directly read or write memory. They are distinguished from normal load/store instructions since an operating system may wish to handle an interruption raised by a non-access instruction differently.

These non-access Itanium instructions can cause interruptions: *fc*, *fc.i*, *lfetch.fault*, *probe*, *probe.fault*, *tpa*, and *tak*. (*tak* can cause interruptions only for non-TLB reasons.) *ISR.code* will be set to indicate which non-access instruction caused the interruption. See [Table 5-1](#) for *ISR* field settings for non-access instructions.

**Table 5-1. ISR Settings for Non-access Instructions**

Instruction	ISR Fields			
	code{3:0}	na	r	w
<i>tpa</i>	0	1	0	0
<i>fc</i> , <i>fc.i</i>	1	1	1	0
<i>probe</i>	2	1	0 or 1 <sup>a</sup>	0 or 1 <sup>a</sup>
<i>tak</i>	3	1	0	0
<i>lfetch</i> , <i>lfetch.fault</i>	4	1	1	0
<i>probe.fault</i>	5	1	0 or 1 <sup>a</sup>	0 or 1 <sup>a</sup>

a. Sets *r* or *w* or both to 1 depending on the *probe* form.

### 5.5.3 Single Stepping

The processor can single step through a series of instructions by enabling the single step *PSR.ss* bit. This is accomplished by setting the *IPSR.ss* bit and performing an *rfi* back to the instruction to be single stepped over. When single stepping, the processor will execute one IA-32 instruction or one Itanium instruction pointed to by the *IPSR.ri* field.

After single stepping Itanium instruction slot 2 (*IPSR.ri* = 2) or when the template is MLX and single stepping instruction slot 1 (*IPSR.ri* = 1), the IIP will point to the next bundle, and *IPSR.ri* will point to slot 0.

### 5.5.4 Single Instruction Fault Suppression

Four bits, *PSR.id*, *PSR.da*, *PSR.ia*, and *PSR.dd* are defined to suppress faults for one Itanium instruction or one mandatory RSE memory operation. The *PSR.id* bit is used to suppress the instruction debug fault for one IA-32 or Itanium instruction. This bit will be cleared in the *PSR* after the first successfully executed instruction. The *PSR.ia* bit is used to suppress the Instruction Access Bit fault for one Itanium instruction. This bit will be cleared in the *PSR* after the first successfully executed instruction. The *PSR.da* and *PSR.dd* bits are used to suppress Dirty-Bit, Data Access-Bit and Data Debug faults for one Itanium instruction, or for one mandatory RSE memory reference. The *PSR.da* and *PSR.dd* bits will be cleared in the *PSR* after the first instruction is executed without raising a fault, or after the first mandatory RSE memory reference that does not raise a fault completes. *PSR.da*, *PSR.ia* and *PSR.dd* are cleared before the first IA-32 instruction starts execution after a *br.ia* or *rfi* instruction. Software may set the *PSR.id*, *PSR.da*, *PSR.ia* and *PSR.dd* bits in the *IPSR* prior to an *rfi*. The *rfi* will restore the *PSR* from the *IPSR*. By using these disable bits, software may step over a debug or dirty/access event and continue execution.

## 5.5.5 Deferral of Speculative Load Faults

Speculative and speculative advanced loads can defer fault handling by suppressing the speculative memory reference, and by setting the deferred exception indicator (NaT bit or NaTVal) of the load target register. Other effects of the instruction (such as post increment) are performed. Additionally, software can suppress the memory reference of speculative and speculative advanced loads independent of any exception.

Deferral is the process of generating a deferred exception indicator and not performing the exception processing at the time of its detection (and potentially never at all). Once a deferred exception indicator is generated, it will propagate through all uses until the speculation is checked by using either a `chk.s` instruction, a `chk.a` instruction (for speculative advanced loads), or a non-speculative use. This causes the appropriate action to be invoked to deal with the exception.

Three different programming models are supported: **no-recovery**, **recovery** and **always-defer**. In the no-recovery model, only fatal exceptional conditions are deferred – these are conditions which cannot be resolved without either involving the program’s exception-handling code or terminating the program. In the recovery model, performance may be increased by deferring additional exceptional conditions. The recovery model is used only if the program provides additional “recovery” code to re-execute failed speculative computations. When a speculative load is executed with PSR.ic equal to 1, and ITLB.ed equal to 0, the no-recovery model is in effect. When PSR.ic is 1 and ITLB.ed is 1, the recovery model is in effect. The **always-defer** model is supported for use in system code which has PSR.ic equal to 0. In this model, all exceptional conditions which can be deferred are deferred. This permits speculation in environments where faulting would be unrecoverable.

In addition to the deferral of exceptional conditions, speculative loads may be deferred automatically by hardware based on implementation-dependent criteria, such as the detection of a cache miss. Such deferral is referred to as **spontaneous deferral**, and is done in order to increase performance. Spontaneous deferral is allowed only in the recovery model.

**Table 5-2. Programming Models**

PSR.ic	PSR.it	ITLB.ed	Model	DCR-based Deferral	Spontaneous Deferral
0	x	x	Always defer	No	No
1	0	x	No recovery	No	No
1	1	0	No recovery	No	No
1	1	1	Recovery	Yes	Yes

Speculative load exceptions are categorized into three groups:

- Ones which always raise a fault
- Ones which always defer
- Ones which always raise a fault in the no-recovery model, but can defer based on the speculative deferral control bits in the DCR control register, in the recovery model.

Aborts, external interrupts, RSE or instruction-fetch-related faults that happen to occur on a speculative load are always raised (since they are not related to the speculative load instruction). Illegal Operation faults and Disabled Floating-point Register faults that occur on a speculative load are always raised.

Processing of exception conditions for speculative and speculative advanced loads is done in three stages: qualification, deferral and prioritization.

During the execution of a load instruction, multiple exception conditions may be detected simultaneously. For non-speculative loads these exception conditions are prioritized and only the highest priority one raises a fault. For speculative loads, however, some exception conditions may be deferred. As a result, it is possible for lower priority exceptions, which are not also deferred, to raise a fault. For some exception conditions, though, other lower priority conditions are meaningless, and are said to be qualified, or precluded. Exception qualification is described in [Table 5-3](#).

**Table 5-3. Exception Qualification**

Exception Condition	Precluded by Concurrent Exception Condition	
Register NaT Consumption (NaT'ed address)	none	
Unimplemented Data Address	Register NaT Consumption	
Alternate Data TLB	Register NaT Consumption	Unimplemented Data Address
VHPT data	Register NaT Consumption	Unimplemented Data Address
Data TLB	Register NaT Consumption	Unimplemented Data Address
Data Page Not Present	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB
Data NaT Page Consumption	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB Data Page Not Present
Data Key Miss	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB Data Page Not Present
Data Key Permission	Register NaT Consumption Unimplemented Data Address VHPT data Data TLB	Alternate Data TLB Data Page Not Present Data Key Miss
Data Access Rights	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB Data Page Not Present
Data Access Bit	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB Data Page Not Present
Data Debug	Register NaT Consumption	Unimplemented Data Address
Unaligned Data Reference	Register NaT Consumption	Unimplemented Data Address
Unsupported Data Reference	Register NaT Consumption Unimplemented Data Address VHPT data	Data TLB Alternate Data TLB Data Page Not Present

After exception conditions are detected and qualified, the remaining exception conditions are checked for deferral. Deferral occurs after fault qualification and determines which memory access exceptions raised by speculative loads are automatically deferred by hardware.

Deferral is controlled by PSR.ed, PSR.it, PSR.ic, the speculative deferral control bits in the DCR, the exception deferral bit of the code page's instruction TLB entry (ITLB.ed), and the memory attribute of the referenced data page. The speculative load and speculative advanced load exception deferral conditions are as follows:

- When PSR.ic is 0 and regardless of the state of DCR, and ITLB.ed bits (see [Table 5-2](#)), all exception conditions related to the data reference are deferred.
- Regardless of the state of DCR, PSR.it, PSR.ic, and ITLB.ed bits, Unimplemented Data Address exception conditions and Data NaT Page Consumption exception conditions (caused by references to NaTPages) are always deferred.
- When PSR.it and ITLB.ed are both 1, and the appropriate DCR bit is 1 for the exception, the speculative load exception is deferred.
- When PSR.it and ITLB.ed are both 1, Unaligned Data Reference exception conditions are deferred.

The conditions for deferral are given in [Table 5-4](#). See also "Default Control Register (DCR - CR0)" on page 2:31.

**Table 5-4. Qualified Exception Deferral**

Qualified Exception	Deferred If
Register NaT Consumption (NaT'ed address)	always
Unimplemented Data Address	always
Alternate Data TLB	!PSR.ic    (PSR.it && ITLB.ed && DCR.dm)
VHPT data	!PSR.ic    (PSR.it && ITLB.ed && DCR.dm)
Data TLB	!PSR.ic    (PSR.it && ITLB.ed && DCR.dm)
Data Page Not Present	!PSR.ic    (PSR.it && ITLB.ed && DCR.dp)
Data NaT Page Consumption	always
Data Key Miss	!PSR.ic    (PSR.it && ITLB.ed && DCR.dk)
Data Key Permission	!PSR.ic    (PSR.it && ITLB.ed && DCR.dx)
Data Access Rights	!PSR.ic    (PSR.it && ITLB.ed && DCR.dr)
Data Access Bit	!PSR.ic    (PSR.it && ITLB.ed && DCR.da)
Data Debug	!PSR.ic    (PSR.it && ITLB.ed && DCR.dd)
Unaligned Data Reference	!PSR.ic    (PSR.it && ITLB.ed)
Unsupported Data Reference	always

The conditions for spontaneous deferral are given in [Table 5-5](#). See the PAL\_PROC\_GET\_FEATURES – Get Processor Dependent Features (17) procedure for details on enabling/disabling spontaneous deferral.

**Table 5-5. Spontaneous Deferral**

Implementation-dependent condition may optionally be deferred if
(PSR.ic && PSR.it && ITLB.ed && spontaneous_deferral_enabled())

After checking for deferral, execution of a speculative load instruction proceeds as follows:

- When PSR.ed is 1, then a deferred exception indicator (NaT bit or NaTVal) is written to the load target register, regardless of whether it has an exception or not and regardless of the state of DCR, PSR.it, PSR.ic and the ITLB.ed bits.
- If PSR.ed is 0 and there is at least one exception condition which is neither precluded nor deferred, then a fault is taken corresponding to the highest-priority

exception condition which is neither precluded nor deferred. Prioritization of non-deferred speculative load faults follows the same interruption priorities as non-speculative instruction faults (Table 5-6 on page 2:109). However, deferred speculative load faults do not take part in the prioritization. As a result, depending on DCR settings, a lower priority fault may be taken, even if a higher priority exception condition exists, but is deferred.

- If PSR.ed is 0 and there are exception conditions, but all are either precluded or deferred, then a deferred exception indicator (NaT bit or NaTVal) is written to the load target register.
- If PSR.ed is 0, and there are no exception conditions, and if the memory attribute of the referenced page is uncacheable or limited speculation, then a deferred exception indicator (NaT bit or NaTVal) is written to the load target register. See “Speculation Attributes” on page 2:79..
- If PSR.ed is 0, and there are no exception conditions, and if spontaneous deferral is enabled and permitted by the programming model, then a deferred exception indicator (NaT bit or NaTVal) may optionally be written to the load target register.
- Otherwise, the load executes normally.

If automatic hardware deferral is not enabled, software may still choose to defer exception processing (for speculative loads) at the time of the fault. If the code page has its ITLB.ed bit equal to 1, then the operating system may choose to defer a non-fatal exception. It is expected that the operating system will always defer fatal exceptions. To assist software in the deferral of non-fatal or fatal exceptions, the system architecture provides three additional resources: ISR.sp, ISR.ed, and PSR.ed.

ISR.sp indicates whether the exception was the result of a speculative or speculative advanced load. The ISR.ed bit captures the code page ITLB.ed bit, and allows deferral of a non-fatal exception due to a speculative load. If both the ISR.sp and ISR.ed bit are 1 on an interruption, then the operating system may defer a non-fatal exception by using the PSR.ed bit to perform the action of hardware deferral for one executed instruction. Software may use the same PSR.ed mechanism to defer fatal speculative load exceptions.

## 5.6 Interruption Priorities

Table 5-6 contains a complete list of the architecture defined interruptions (including IA-32), grouped according to type (aborts, interrupts, faults and traps), instruction set, and listed in priority order. Interruptions are delivered in priority order. If more than one instruction detects an interruption within a bundle, the interruption occurring in the lowest numbered instruction slot is raised. Lower priority faults and traps are discarded. Lower priority interrupts are held pending.

The shaded interruptions are disabled if the instruction generating the interruption is predicated off. All other interruptions are either “bundle related” (so the predicate bits do not affect them) or are caused by instructions that cannot be predicated off. Incomplete Register frame (IR) faults 7 through 18 are identical in behavior to faults 45, 51 through 62 (exclusive of 60) except they are of a higher priority. IR faults 7 through 18 can only be caused by mandatory RSE load operations that result from `br.ret`, or `rfi` instructions, but not from `loadrs` instructions (for details see Section 6.6, “RSE Interruptions” on page 2:144).



**Table 5-6. Interruption Priorities**

Type	Instr. Set	Interruption Name	Vector Name	IA-32 Class <sup>a</sup>			
Aborts	IA-32, Intel Itanium	1 Machine Reset (RESET)	PALE_RESET vector	N/A			
		2 Machine Check (MCA)	PALE_CHECK vector				
Interrupts		3 Initialization Interrupt (INIT)	PALE_INIT vector	N/A			
		4 Platform Management Interrupt (PMI)	PALE_PMI vector				
		5 External Interrupt (INT)	External Interrupt vector	N/A			
		6 Virtual External Interrupt (VINT)	Virtual External Interrupt vector				
Faults	Intel Itanium	7 IR Unimplemented Data Address fault	General Exception vector	N/A			
		8 IR Data Nested TLB fault	Data Nested TLB vector				
		9 IR Alternate Data TLB fault	Alternate Data TLB vector				
		10 IR VHPT Data fault	VHPT Translation vector				
		11 IR Data TLB fault	Data TLB vector				
		12 IR Data Page Not Present fault	Page Not Present vector				
		13 IR Data NaT Page Consumption fault	NaT Consumption vector				
		14 IR Data Key Miss fault	Data Key Miss vector				
		15 IR Data Key Permission fault	Key Permission vector				
		16 IR Data Access Rights fault	Data Access Rights vector				
		17 IR Data Access Bit fault	Data Access-Bit vector				
		18 IR Data Debug fault	Debug vector				
		19 Unimplemented Instruction Address fault <sup>b</sup>	Lower-Privilege Transfer Trap vector				
		Faults	IA-32		20 IA-32 Instruction Breakpoint fault	IA-32 Exception vector (Debug)	A
					21 IA-32 Code Fetch fault <sup>c</sup>	IA-32 Exception vector (GPFault)	
			IA-32, Intel Itanium		22 Alternate Instruction TLB fault	Alternate Instruction TLB vector	
					23 VHPT Instruction fault	VHPT Translation vector	
					24 Instruction TLB fault	Instruction TLB vector	
					25 Instruction Page Not Present fault	Page Not Present vector	
26 Instruction NaT Page Consumption fault	NaT Consumption vector						
27 Instruction Key Miss fault	Instruction Key Miss vector						
28 Instruction Key Permission fault	Key Permission vector						
29 Instruction Access Rights fault	Instruction Access Rights vector						
30 Instruction Access Bit fault	Instruction Access-Bit vector						
	Intel Itanium	31 Instruction Debug fault	Debug vector				
	IA-32	32 IA-32 Instruction Length > 15 bytes	IA-32 Exception vector (GPFault)	B			
		33 IA-32 Invalid Opcode fault	IA-32 Intercept vector (Instruction)				
		34 IA-32 Instruction Intercept fault	IA-32 Intercept vector (Instruction)				
	Intel Itanium	35 Illegal Operation fault <sup>d</sup>	General Exception vector				
		36 Illegal Dependency fault	General Exception vector				
		37 Break Instruction fault	Break Instruction vector				
		38 Privileged Operation fault	General Exception vector				

**Table 5-6. Interruption Priorities (Continued)**

Type	Instr. Set	Interrupt Name	Vector Name	IA-32 Class <sup>a</sup>
	IA-32, Intel Itanium	39 Disabled Floating-point Register fault	Disabled FP-Register vector	B
		40 Disabled Instruction Set Transition fault	General Exception vector	
	IA-32	41 IA-32 Device Not Available fault	IA-32 Exception vector (DNA)	C
		42 IA-32 FP Error fault <sup>e</sup>	IA-32 Exception vector (FPError)	
	IA-32, Intel Itanium	43 Register NaT Consumption fault	NaT Consumption vector	
	Intel Itanium	44 Reserved Register/Field fault	General Exception vector	
		45 Unimplemented Data Address fault	General Exception vector	
		46 Privileged Register fault	General Exception vector	
		47 Speculative Operation fault	Speculation vector	
	IA-32	48 Virtualization fault	Virtualization vector	
		49 IA-32 Stack Exception	IA-32 Exception vector (StackFault)	
		50 IA-32 General Protection Fault	IA-32 Exception vector (GPFault)	
Faults	IA-32, Intel Itanium	51 Data Nested TLB fault	Data Nested TLB vector	C
		52 Alternate Data TLB fault <sup>f</sup>	Alternate Data TLB vector	
		53 VHPT Data fault <sup>f</sup>	VHPT Translation vector	
		54 Data TLB fault <sup>f</sup>	Data TLB vector	
		55 Data Page Not Present fault <sup>f</sup>	Page Not Present vector	
		56 Data NaT Page Consumption fault <sup>f</sup>	NaT Consumption vector	
		57 Data Key Miss fault <sup>f</sup>	Data Key Miss vector	
		58 Data Key Permission fault <sup>f</sup>	Key Permission vector	
		59 Data Access Rights fault <sup>f</sup>	Data Access Rights vector	
		60 Data Dirty Bit fault	Dirty-Bit vector	
		61 Data Access Bit fault <sup>f</sup>	Data Access-Bit vector	
	Intel Itanium	62 Data Debug fault <sup>f</sup>	Debug vector	
		63 Unaligned Data Reference fault <sup>f</sup>	Unaligned Reference vector	
	IA-32	64 IA-32 Alignment Check fault	IA-32 Exception vector (AlignmentCheck)	C
		65 IA-32 Locked Data Reference fault	IA-32 Intercept vector (Lock)	
		66 IA-32 Segment Not Present fault	IA-32 Exception vector (NotPresent)	
		67 IA-32 Divide by Zero fault	IA-32 Exception vector (Divide)	
		68 IA-32 Bound fault	IA-32 Exception vector (Bound)	
		69 IA-32 SSE Numeric Error fault	IA-32 Exception vector (StreamSIMD)	
	Intel Itanium	70 Unsupported Data Reference fault	Unsupported Data Reference vector	
		71 Floating-point fault	Floating-point Fault vector	
Traps	Intel Itanium	72 Unimplemented Instruction Address trap <sup>b,9</sup>	Lower-Privilege Transfer Trap vector	
		73 Floating-point trap	Floating-point Trap vector	
	Intel Itanium	74 Lower-Privilege Transfer trap	Lower-Privilege Transfer Trap vector	
		75 Taken Branch trap	Taken Branch Trap vector	
		76 Single Step trap	Single Step Trap vector	

**Table 5-6. Interruption Priorities (Continued)**

Type	Instr. Set	Interruption Name	Vector Name	IA-32 Class <sup>a</sup>
	IA-32	77 IA-32 System Flag Intercept trap	IA-32 Intercept vector (SystemFlag)	D
		78 IA-32 Gate Intercept trap	IA-32 Intercept vector (Gate)	
		79 IA-32 INTO trap	IA-32 Exception vector (Overflow)	
		80 IA-32 Breakpoint (INT 3) trap	IA-32 Exception vector (Break)	
		81 IA-32 Software Interrupt (INT) trap	IA-32 Interrupt vector (Vector#)	
		82 IA-32 Data Breakpoint trap	IA-32 Exception vector (Debug)	
		83 IA-32 Taken Branch trap	IA-32 Exception vector (Debug)	
		84 IA-32 Single Step trap	IA-32 Exception vector (Debug)	

- a. IA-32 Interruption Class, see [Section 5.6.1, “IA-32 Interruption Priorities and Classes” on page 2:111](#) for details
- b. Processor implementations may report unimplemented instruction addresses either with an Unimplemented Instruction Address trap on the taken branch, taken `chk`, or an `rfi` to an unimplemented address, or on a non-branching slot 2 instruction in a bundle at the upper edge of the implemented address space (where the next sequential bundle address would be an unimplemented address), or with an Unimplemented Instruction Address fault on the fetch of the unimplemented address.
- c. IA-32 Code Fetch faults include Code Segment Limit Violation and other Code Fetch checks defined in [Section 6.2.2.3.3, “IA-32 Environment Runtime Integrity Checks” on page 1:122](#).
- d. Illegal Operation faults can be taken for certain predicated off reserved opcodes. For details, refer to [Section 4.1, “Format Summary” on page 3:294](#).
- e. IA-32 FP Error fault conditions detected on an IA-32 FP instruction are reported as a fault on the next IA-32 FP instruction that performs an FWAIT operation.
- f. If not deferred.
- g. Unimplemented Instruction Address traps on emulated check instructions have a lower priority than Taken Branch trap and Single Step trap. See [“Speculation vector \(0x5700\)” on page 2:198](#).

## 5.6.1 IA-32 Interruption Priorities and Classes

Table 5-6 establishes a well defined priority between faults, traps and interrupts (including IA-32). However, IA-32 instruction set generated interruptions are divided into interruption classes. While priority among these IA-32 interruption classes is well defined by the table (except as noted below), interruption priority within each IA-32 interruption class is implementation dependent and may vary from processor to processor as defined below:

**Class A** – Faults from fetching an instruction. Priority of IA-32 Instruction Breakpoint, IA-32 Code Fetch (GPFault(0)), and Instruction TLB faults (Alternate Instruction TLB fault to Instruction Access Bit fault) may vary based on instruction alignment and page boundaries in a model-specific way. Faults are prioritized as defined in the table if the instruction does not span a virtual page. If an IA-32 instruction spans a virtual page, IA-32 Code Fetch faults (IA\_32\_Exception(GPFault)) due to code segment (CS) Limit violations can be raised above or below Instruction TLB faults as defined below:

- If the starting effective address of the IA-32 instruction exceeds the code segment limit, then the IA-32 Code Fetch fault has higher priority than any Instruction TLB faults. If the starting effective address of the IA-32 instruction is within the code segment limit, then Instruction TLB faults have higher priority for the starting effective address.
- If the IA-32 instruction spans a virtual page and the code segment limit is equal to the page boundary, the IA-32 Code Fetch fault has higher priority than any Instruction TLB faults on the second page. Otherwise if the code segment limit is

greater than the page boundary, any Instruction TLB faults on the second page have higher priority than the IA-32 Code Fetch fault.

**Class B** – Faults from decoding an instruction. Priority of IA-32 Instruction Length, IA-32 Invalid Opcode, and IA-32 Instruction Intercept, Disabled Floating Point Register, Disabled Instruction Set Transition, and Device Not Available faults are model specific. If the IA-32 instruction spans a virtual page, IA-32 Instruction Length >15 byte Faults (IA\_32\_Exception(GPFault)) can have higher priority than Instruction TLB faults as defined below:

- If the IA-32 prefix bytes on the first page are  $\geq 15$  bytes, an IA-32 Instruction >15 byte fault (GPFault) is taken first regardless of any Instruction TLB faults on the second page.
- If the IA-32 prefix bytes on the first page are  $< 15$  bytes, Instruction TLB faults on the second page may or may not have priority over any possible IA-32 Instruction Length fault.

**Class C** – Faults resulting from executing an instruction. Priority of faults is model specific and can vary across processor implementations. Most faults are related to data memory references, other fault priorities can vary due to model-specific differences across processor implementations. The memory fault priorities (IA-32 Stack Exception through Data Access Bit fault) defined in the table only apply to a single IA-32 data memory reference that does not cross a virtual page. If an IA-32 instruction requires multiple data memory references or a single data memory reference crosses a virtual page:

- If any given IA-32 instruction requires multiple data memory references, all possible faults are raised on the first data memory reference before any faults are checked on subsequent data memory references. This implies lower priority faults on an earlier memory reference will be raised before higher priority faults on a later data memory reference within a single IA-32 instruction. The order of data memory references initiated by an IA-32 instruction is implementation dependent and may vary from processor to processor. Software can not assume all higher priority data memory faults are raised before all lower priority data memory faults within a single IA-32 instruction.
- If a single IA-32 data memory reference crosses a virtual page, the processor checks for faults in a model-specific order: Any faults present on one page are checked and reported before any faults are checked and reported on the other page. This implies that a single data reference that crosses a virtual page can raise lower priority data memory faults on one page before higher priority data memory faults are raised on the other page. For example, Data Key Miss faults (lower priority) on the first page could be raised before a Data TLB Miss Fault (higher priority) on the second page. Software can not assume all higher priority data memory faults are raised before all lower priority data memory faults within a single IA-32 instruction.

**Class D** – Traps on the current IA-32 instruction. Trap conditions are reported concurrently on the same exception vector or via a trap code specifying all concurrent traps.

## 5.7 IVA-based Interruption Vectors

Table 5-7 contains the processor’s interruption vector table (IVT). The base of the IVT is held in the IVA control register. The size of the IVT is 32KB. The first 20 vectors are designed to provide more code space by allowing 64 bundles per vector (16 bytes per bundle) for performance-critical interruption handlers. The second 48 vectors provide 16 bundles per vector. Several vectors have more than one interruption associated with them. Information provided in the ISR allows the handler to distinguish which fault or trap caused the event.

Some vectors require additional software decoding to determine the cause of the interruption. Additional information for this decoding is provided in the ISR.code field. See Chapter 8, “Interruption Vector Descriptions” for a complete specification of the information supplied in the ISR for each of the vectors.

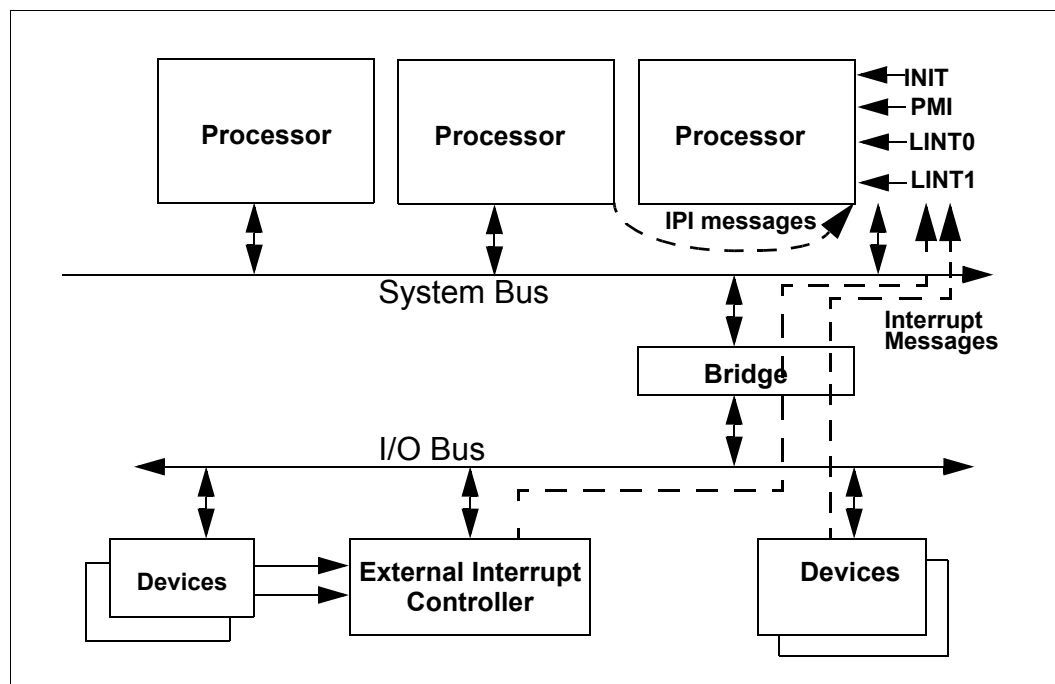
**Note:** PAL-based interruptions (RESET, MCA, INIT, and PMI) do not reference the IVT.

**Table 5-7. Interruption Vector Table (IVT)**

Offset	Vector Name	Interruption(s)	Page
0x0000	VHPT Translation vector	10, 23, 53	2:173
0x0400	Instruction TLB vector	24	2:175
0x0800	Data TLB vector	11, 54	2:176
0x0c00	Alternate Instruction TLB vector	22	2:177
0x1000	Alternate Data TLB vector	9, 52	2:178
0x1400	Data Nested TLB vector	8, 51	2:179
0x1800	Instruction Key Miss vector	27	2:180
0x1c00	Data Key Miss vector	14, 57	2:181
0x2000	Dirty-Bit vector	60	2:182
0x2400	Instruction Access-Bit vector	30	2:183
0x2800	Data Access-Bit vector	17, 61	2:184
0x2c00	Break Instruction vector	37	2:185
0x3000	External Interrupt vector	5	2:186
0x3400	Virtual External Interrupt vector	6	2:187
0x3800	Reserved		
0x3c00	Reserved		
0x4000	Reserved		
0x4400	Reserved		
0x4800	Reserved		
0x4c00	Reserved		
0x5000	Page Not Present vector	12, 25, 55	2:188
0x5100	Key Permission vector	15, 28, 58	2:189
0x5200	Instruction Access Rights vector	29	2:190
0x5300	Data Access Rights vector	16, 59	2:191
0x5400	General Exception vector	7, 35, 36, 38, 40, 44, 45, 46	2:192
0x5500	Disabled FP-Register vector	39	2:195
0x5600	NaT Consumption vector	13, 26, 43, 56	2:196
0x5700	Speculation vector	47	2:198
0x5800	Reserved for software use <sup>a</sup>		



**Figure 5-3. Interrupt Architecture Overview**



As defined in “[Interruption Definitions](#)” on page 2:95 there are three kinds of interrupts: initialization interrupts (INITs), platform management interrupts (PMIs), and external interrupts (INTs).

The processors and external interrupt controllers communicate over the processor’s system bus with an implementation-specific interrupt messaging protocol. Interrupts are generated by a number of different interrupt sources in the system:

- **External (I/O) devices** – Interrupt messages from any external source can be directed to any one processor by an external interrupt controller or by I/O devices capable of directly sending interrupt messages. An interrupt message informs the processor that an interrupt request is being made, and, in the case of PMIs and external interrupts, specifies a unique vector number for the interrupt. Interrupt messages are only issued on the “assertion edge” of an interrupt; “deassertion” of an interrupt does not result in an interrupt message.
- **Locally connected devices** – These interrupts originate on the processor’s interrupt pins (LINT, INIT, PMI)<sup>1</sup>, and are always directed to the local processor. The LINT pins can be connected directly to an Intel 8259A-compatible external interrupt controller. The LINT pins are programmable to be either **edge-sensitive** or **level-sensitive**, and for the kind of interrupt that gets generated. If programmed to generate external interrupts, the vector number is a programmed constant per LINT pin. Only the LINT pins connected to the processor can directly generate level-sensitive interrupts (See “[Edge- and Level-sensitive Interrupts](#)” on page 2:131). LINT pins cannot be programmed to generate level-sensitive PMIs or INITs. The INIT and PMI pins generate their corresponding interrupts. For PMI pins a PMI vector 0 interrupt is generated.

1. Processors are not required to support externally connected interrupt pins. Software can query the presence of the INIT, PMI, and LINT pins via the PAL\_PROC\_GET\_FEATURES procedure call.

- **Internal processor interrupts** – such as interval timer, performance monitoring, and corrected machine checks. These are always directed to the local processor. A unique vector number can be programmed for each source.
- **Other processors** – A processor can interrupt any individual processor, including itself, by sending an Inter-Processor Interrupt (IPI) message to a specific target processor. See [“Inter-processor Interrupt Messages” on page 2:128](#).

The destination of an interrupt message is any one processor in the system, and is specified by a unique processor identifier. A different destination can be specified for each interrupt. There is no mechanism to “broadcast” a single interrupt to all processors in the system.

The following terms are used in the interrupt definition:

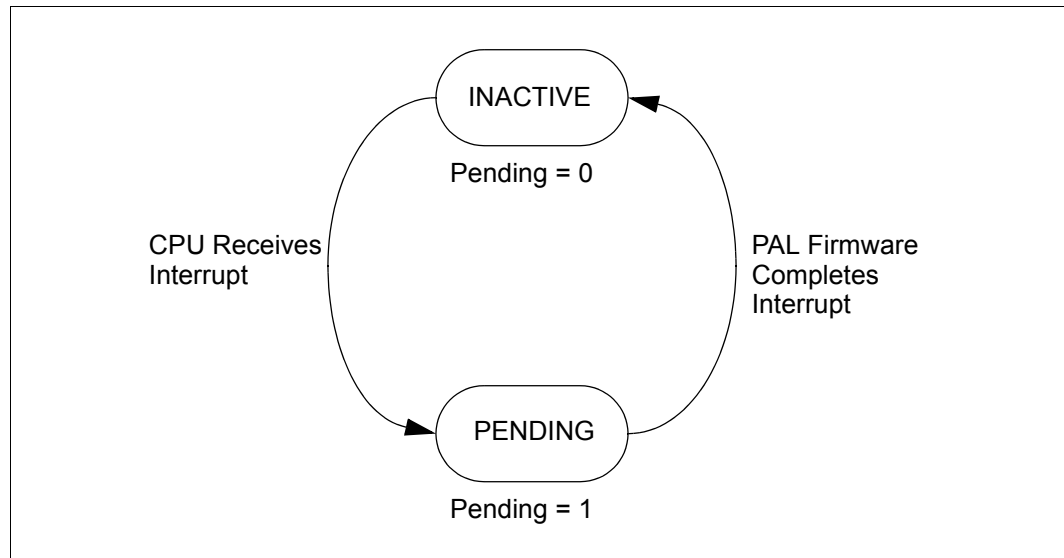
- The processor is said to **receive** an interrupt, if one of the processor’s interrupt pins is asserted, the processor detected an interrupt message bus transaction containing the processor’s unique identifier, or the processor detected an internal interrupt event.
- After receiving an interrupt, the processor internally holds the interrupt **pending**. The interrupt is said to be **pending** when it is received and held by the processor.
- For edge-sensitive interrupts, an external interrupt is held pending until the interrupt is acquired by software at which point it is said to be in-service. INITs and PMIs are held pending until the corresponding PAL vector is entered and PAL firmware clears the pending indication at which point they are said to be completed. For level-sensitive interrupts programmed through the LINT pins, the interrupt is held pending as long as the pin is asserted. Deassertion of a level-sensitive interrupt removes the pending indication (see [“Edge- and Level-sensitive Interrupts” on page 2:131](#)).
- The processor maintains an individual interrupt pending indication for INITs. Since external interrupts and PMIs are also signified by a unique interrupt **vector** number, the processor maintains individual pending indications per vector. An occurrence of an interrupt on a vector that is already marked as pending cannot be distinguished from previous interrupts on the same vector because the interrupts are pending in the same internal pending bit, and are therefore treated as “the same” interrupt occurrence.
- When interrupt delivery is enabled and the highest priority pending interrupt is unmasked (as defined below), the processor **accepts** the pending interrupt, interrupts the control flow of the processor and transfers control to the software interrupt handler.
- An external interrupt is said to be **in-service** when software **acquires** the interrupt vector from the processor by reading the IVR register (see [“External Interrupt Vector Register \(IVR – CR65\)” on page 2:123](#)). The processor then removes the pending indication for the interrupt vector. The processor maintains one in-service indicator for each unique vector number. Note that there are no in-service indicators for INITs and PMIs.
- Once an external interrupt is in-service it remains so until software indicates service for that external interrupt is **complete**. By writing to the EOI register (see [“End of External Interrupt Register \(EOI – CR67\)” on page 2:124](#)) software indicates that service for the highest-priority in-service external interrupt is complete. The processor then removes the in-service indication for the highest-priority external interrupt vector. INITs and PMIs are completed when PAL firmware clears the corresponding pending indication.



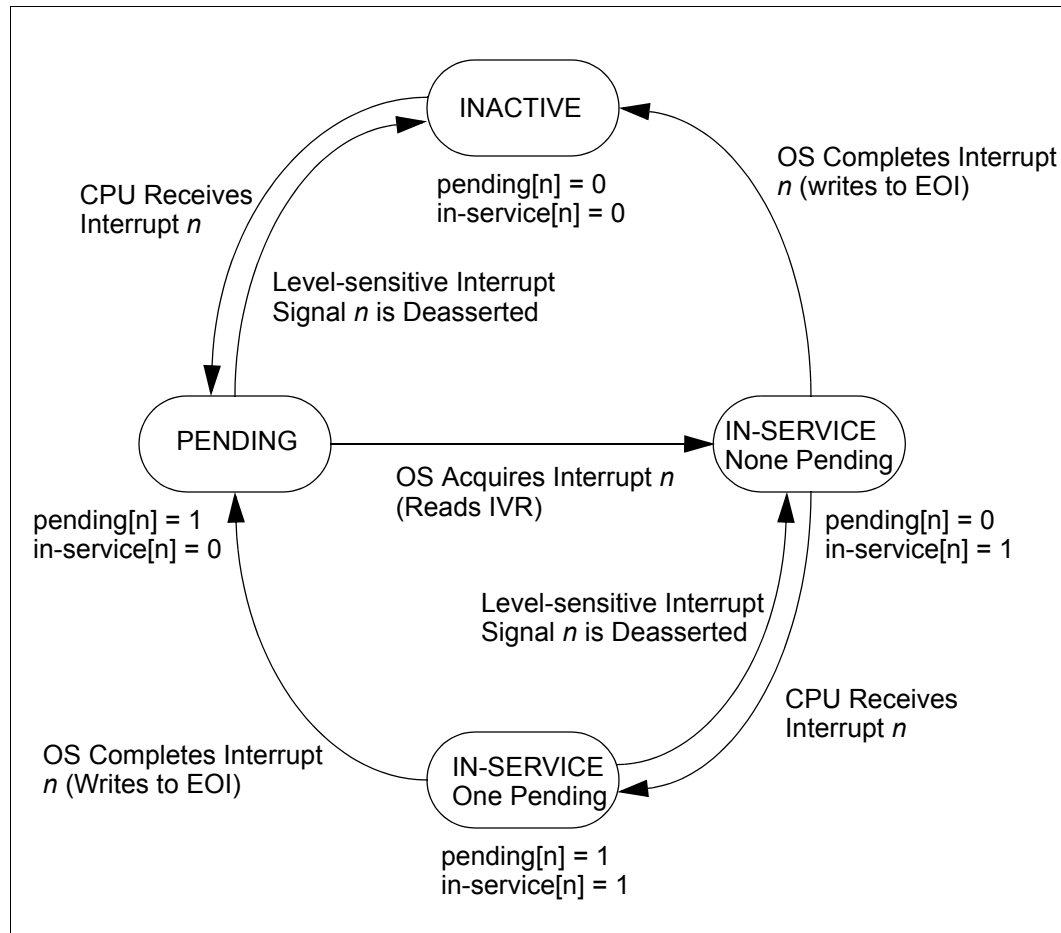
- The **priority** of interrupts is defined in [Table 5-8](#). Entry *A* is higher priority than interrupt *B*, if entry *A* appears at a higher location in the table than entry *B*. Interrupt priority is used to select interrupts that require urgent service over less urgent interrupt requests.
- Interrupt **delivery is enabled** when software programs the processor to accept any unmasked interrupt. INITs delivery is enabled when PSR.mc is 0. PMIs delivery is enabled when PSR.ic is 1. For Itanium architecture-based code execution, external interrupts delivery is enabled when PSR.i is 1.
- **Masking** applies only to external interrupts. Unmasked interrupts are those external interrupts of higher priority than the highest priority external interrupt vector currently in-service (if any) and whose priority level is higher than the current priority masking level specified by the TPR register (see “[Task Priority Register \(TPR – CR66\)](#)” on [page 2:123](#)). Masking conditions are defined in [Table 5-8](#). PSR.i does not affect masking of external interrupts.

[Figure 5-4](#) shows how this terminology is applied to the handling of a PAL-based interrupt. Similarly, [Figure 5-5](#) shows the handing of a vectored external interrupt *n*. Both figures show the different states and transitions interrupts go through.

**Figure 5-4. PAL-based Interrupt States**



**Figure 5-5. External Interrupt States**



### 5.8.1 Interrupt Vectors and Priorities

As indicated in [Table 5-6 on page 2:109](#), INITs have higher priority than PMIs, which in turn have higher priority than external interrupts. PMIs and external interrupts are further prioritized by vector number.

PMIs have a separate vector space from external interrupts. PMI vectors 0-3 can be used by platform firmware. PMI vectors 4 through 15 are reserved for use by processor firmware. Assertion of the processor's PMI pin, when present, results in PMI vector number 0. PMI vector priorities are described in [Section 11.5, "Platform Management Interrupt \(PMI\)" on page 2:310](#).

Each external interrupt (INT) in the system is distinguished from other external interrupts by a unique vector number. There are 256 distinct vector numbers in the range 0 - 255. Vector numbers 1 and 3 through 14 are reserved for future use. Vector number 0 (ExtINT) is used to service Intel 8259A-compatible external interrupt controllers. Vector number 2 is used for the Non-Maskable Interrupt (NMI). The remaining 240 external interrupt vector numbers (16 through 255) are available for general operating system use. [Table 5-8](#) summarizes the interrupt priority model.

**Table 5-8. Interrupt Priorities, Enabling, and Masking**

Priority	Priority Class	Interrupt	Vector Number	Interrupt Delivery Enabled	Interrupt Unmasked Condition
Highest	N/A	INIT	N/A	if PSR.mc is 0	Always
		PMI	0..3	if PSR.ic is 1	Always
		INT	2 (NMI)	if PSR.i is 1 <sup>a</sup>	Interrupt is higher priority than all in-service external interrupts
	0 (ExtINT)		TPR.mmi is 0, and interrupt is higher priority than all in-service external interrupts		
	15	240..255	TPR.mmi is 0, and interrupt is higher priority than all in-service external interrupts, and Vector Number{7:4} > TPR.mic		
	14	224..239			
	13	208..223			
	12	192..207			
	11	176..191			
	10	160..175			
	9	144..159			
	8	128..143			
	7	112..127			
	6	96..111			
	5	80..95			
	4	64..79			
3	48..63				
2	32..47				
Lowest	1	16..31			

a. For Itanium architecture-based code execution external interrupt delivery is enabled if PSR.i is 1. For IA-32 code execution external interrupt delivery is enabled if (PSR.i AND (!CFLAG.if OR EFLAG.if)) is true.

NMI (vector 2) has higher interrupt priority than ExtINT (vector 0), which has higher priority than external interrupt vectors 16 through 255.

External interrupts vectors 16 through 255 are divided into 15 interrupt priority classes. Sixteen different interrupt vectors share a single interrupt priority class, with class 1 being the lowest priority and class 15 being the highest. For these external interrupts, higher number external interrupts have priority over lower number external interrupts, including those within the same priority class.

Vector number 15 is used to indicate that the highest priority pending interrupt in the processor is at a priority level that is currently masked or there are no pending external interrupts. This encoding is referred to as a “spurious” interrupt.

## 5.8.2 Interrupt Enabling and Masking

Upon receiving an interrupt, the processor holds the interrupt pending internally until interrupt delivery is enabled and, in the case of external interrupts, the interrupt is unmasked. When all of the interrupt enabling and unmasking conditions are satisfied (see Table 5-8), the processor accepts the pending interrupt, interrupts the control flow of the processor, and transfers control to the External Interrupt handler for external interrupts, or to PAL firmware for INITs and PMIs.

**Note:** The TPR controls the masking of external interrupts. TPR is described in “Task Priority Register (TPR – CR66)” on page 2:123.

The processor provides nested interrupt priority support for external interrupt vectors 0, 2, and 16 through 255 by:

- Automatically masking external interrupts of equal or lower priority than the highest priority external interrupt currently in-service. This raises the in-service external interrupt masking level when each external interrupt begins service by an IVR read.
- Associating EOI writes with the highest priority in-service external interrupt, and removing the in-service indication for this external interrupt. This lowers the in-service masking level to that of the next highest priority currently in-service external interrupt (if any).

This mechanism allows software external interrupt handlers to be interrupted by higher priority external interrupts.

For example, assume software acquires an external interrupt vector 45 by reading IVR. During the service of this interrupt other external interrupts can still be received and are pended. If software sets PSR.i to a 1, pending external interrupts of equal or lower priority than 45 are masked. However, a higher priority pending external interrupt can be accepted by the processor (provided it is not masked by TPR.mmi or TPR.mic). Assuming external interrupt vector 80 is received by the processor, the processor will accept the interrupt by interrupting the control flow of the processor. During the service of this interrupt, external interrupts of equal or lower priority than vector 80 are masked. When EOI is issued by software, the processor will remove the in-service indication for external interrupt vector 80. External interrupt masking will then revert back to the next highest priority in-service external interrupt, vector 45. External interrupt vectors of equal or lower priority than vector 45 would remain masked until EOI is issued by software. The in-service indication for vector 45 is then removed by the write to EOI.

### 5.8.2.1 Re-enabling External Interrupt Delivery

When emerging from code in which external interrupt delivery is disabled and interruption state collection is turned off, the following minimal code sequence describes the architectural method with which to re-enable interruption collection and enable external interrupts:

```
    ssm PSR.ic          // enable interruption collection
    ;;
    srlz.d              // guarantee that interruption collection is enabled
    ssm PSR.i           // enable external interrupts
```

The processor does not ensure that enabling external interrupts is immediately observed after the `ssm PSR.i` instruction. Software must perform a data serialization operation after `ssm PSR.i` to ensure that external interrupt delivery is enabled prior to a given point in program execution.

### 5.8.2.2 External Interrupt Sampling

Assuming that external interrupt delivery is currently disabled (PSR.i is 0), the following minimal code sequence describes the architectural method with which to briefly open the external interrupt window for external interrupt sampling (typically PSR.ic is 1 to enable interruption collection):

```

    ssm PSR.i
    ;;
    srlz.d          // external interrupts may be sampled anywhere here
    ;;
    rsm PSR.i

```

The stop following the `srlz.d` instruction in the above code sequence is required to force the Reset System Mask (`rsm`) instruction into a subsequent instruction group. The stop guarantees that the `srlz.d` will open the external interrupt window for at least one cycle before the `rsm` instruction closes it again.

**Note:** In the above code sequence, the effect of disabling interrupts due to the `rsm` instruction is observed on the next instruction following the `rsm`.

### 5.8.2.3 Disabling of External Interrupt Delivery and `rsm`

When the current privilege level is zero, an `rsm` instruction whose mask includes `PSR.i` may cause external interrupt delivery to be disabled for an implementation-dependent number of instructions, even if the qualifying predicate for the `rsm` instruction is false. Architecturally, the extents of this delivery disable “window” are defined as follows:

1. External interrupt delivery may be disabled for any instructions in the same instruction group as the `rsm`, including those that precede the `rsm` in sequential program order, regardless of the value of the qualifying predicate of the `rsm` instruction.
2. If the qualifying predicate of the `rsm` is true, then external interrupt delivery is disabled immediately following the `rsm` instruction.
3. If the qualifying predicate of the `rsm` is false, then external interrupt delivery may be disabled until the next data serialization operation that follows the `rsm` instruction.

The delivery disable window is guaranteed to be no larger than defined by the above criteria, but it may be smaller, depending on the implementation.

When the current privilege level is non-zero, an `rsm` instruction whose mask includes `PSR.i` may briefly disable external interrupt delivery, regardless of the value of the qualifying predicate of the `rsm` instruction. However, the implementation guarantees that non-privileged code cannot lock out external interrupts indefinitely (e.g., via an arbitrarily long sequence of `rsm PSR.i` instructions with zero-valued qualifying predicates).

## 5.8.3 External Interrupt Control Registers

Software interacts with external interrupts by reading and writing the external interrupt control registers (CR64-81). These registers are summarized in [Table 5-9](#), and are used to prioritize and deliver external interrupts, and to assign external interrupt vectors for processor-internal interrupt sources such as interval timer, performance monitoring, and corrected machine check.

The external interrupt control registers can only be accessed at privilege level 0, otherwise a Privileged Operation fault is raised.

**Table 5-9. External Interrupt Control Registers**

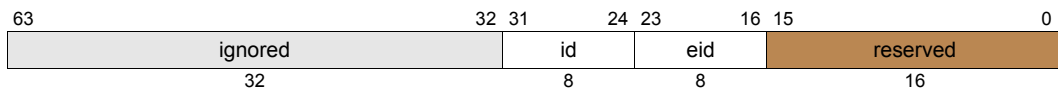
Register	Name	Description
CR64	LID	Local ID
CR65	IVR	External Interrupt Vector Register (read only)
CR66	TPR	Task Priority Register
CR67	EOI	End Of External Interrupt
CR68	IRR0	External Interrupt Request Register 0 (read only)
CR69	IRR1	External Interrupt Request Register 1 (read only)
CR70	IRR2	External Interrupt Request Register 2 (read only)
CR71	IRR3	External Interrupt Request Register 3 (read only)
CR72	ITV	Interval Timer Vector
CR73	PMV	Performance Monitoring Vector
CR74	CMCV	Corrected Machine Check Vector
CR80	LRR0	Local Redirection Register 0
CR81	LRR1	Local Redirection Register 1

**5.8.3.1 Local ID (LID – CR64)**

The LID register contains the processor’s local interrupt identifier. Two fields (*id* and *eid*) serve as the processor’s physical name for all interrupt messages (external interrupts, INITs, and PMIs). LID is loaded by firmware during platform initialization based on the processor’s physical location within the system. Processors receiving an interrupt message on the system interconnect may or may not compare their *id/eid* fields with the target address for the interrupt message, depending on the type of system interconnect. If this comparison is performed, then a match would indicate that the interrupt received was intended for this processor. In case of no comparison, processors use other system topology mechanisms to determine the correct target of the interrupt message.

The LID register fields are either read-only or read-write. Details of the programmability of these fields is communicated by PAL at PALE\_RESET handoff (see Section 11.2.2, “PALE\_RESET Exit State” on page 2:289 for details). Read-only LID bits always return a value of 0. Writes to read-only bits are ignored. To ensure that future arriving interrupts see the updated LID value by a given point in program execution, software must perform a data serialization operation after a LID write and prior to that point. The Local ID fields are defined in Figure 5-6 and Table 5-10.

**Figure 5-6. Local ID (LID – CR64)**



**Table 5-10. Local ID Fields**

Field	Bits	Description
id/eid	31:16	The low order bits of id correspond to a unique, geographically significant address of the processor on the local system bus. The eid field and the higher order bits of the id field correspond to a unique address of the local system bus within the entire system. These fields are initialized by platform firmware to an implementation-dependent value and should not be modified by software. The two fields corresponds to physical address bits{19:4} of the inter-processor interrupt message.

### 5.8.3.2 External Interrupt Vector Register (IVR – CR65)

A read of IVR returns the highest priority, pending, unmasked external interrupt vector, independent of the value of PSR.i. The external interrupt vector is an 8-bit encoded number. If there are no pending external interrupts or all external interrupts are currently masked, IVR returns the “spurious” interrupt indication (vector 15). IVR fields are shown in [Figure 5-7](#). See “Interrupt Unmasked Condition” column in [Table 5-8 on page 2:119](#) for masking conditions.

IVR reads also have two atomic side effects:

- The interrupt pending bit in IRR is cleared for the reported external interrupt vector. Subsequent IVR reads will not report the interrupt as pending unless a new interrupt was pended for the specified interrupt vector.
- The processor marks the interrupt vector as being in-service and masks all pending external interrupts with equal or lower priority until software writes the end-of-interrupt (EOI) register for the in-service interrupt.

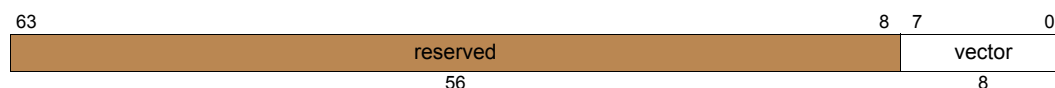
To ensure IVR side effects are observed by a given point in program execution (e.g., before the next IVR read, EOI write, or PSR.i write to enable external interrupt delivery), software must perform a data serialization operation after an IVR read and prior to that point. To ensure that the reported external interrupt vector is correctly masked before the next IVR read, software must perform a data serialization operation after a TPR or EOI write and prior to that IVR read.

Software must be prepared to service any possible external interrupt if it reads IVR, since IVR reads are destructive and removes the highest priority pending external interrupt (if any).

IVR is a read-only register; writes to IVR result in a Illegal Operation fault.

IVR reads do not issue an external INTA cycle. If the interrupt vector must be acquired from an Intel 8259A-compatible external interrupt controller, software should perform a load from the INTA byte. See “[Interrupt Acknowledge \(INTA\) Cycle](#)” on [page 2:130](#) for details.

**Figure 5-7. External Interrupt Vector Register (IVR – CR65)**



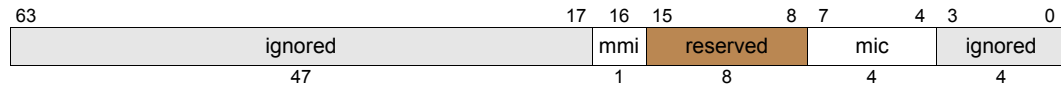
### 5.8.3.3 Task Priority Register (TPR – CR66)

The processor’s Task Priority Register (TPR) provides the ability to create additional masking of external interrupts based on a “priority class.” The 240 external interrupt vectors (16 - 255) are divided into 15 priority classes of 16 numerically contiguous interrupt vectors each. The value written in TPR.mic masks all external interrupts of equal or lower priority classes.

To ensure that new priority levels are established by a given point in program execution, software must perform a data serialization operation after a TPR write and prior to that point. For example, if PSR.i is subsequently set to 1, thus enabling interrupts, and the new priority levels need to be in place before this enabling, a data serialization must be performed prior to the setting of PSR.i. Similarly, if PSR.pp or

PSR.up is set to 1, potentially enabling performance monitor interrupts, and the new priority levels need to be in place before this enabling, a data serialization must be performed. (Note that there's no dependence between writing TPR and then changing the PSR for any other bits in the PSR than these.) A data serialization operation must be performed after TPR is written and before IVR is read to ensure that the reported IVR vector is correctly masked. The TPR fields are described in [Figure 5-8](#) and [Table 5-11](#).

**Figure 5-8. Task Priority Register (TPR – CR66)**



**Table 5-11. Task Priority Register Fields**

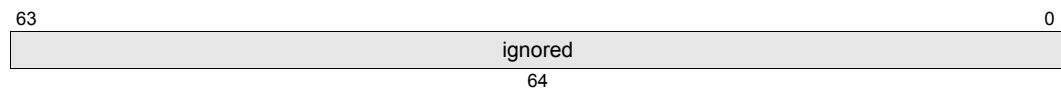
Field	Bits	Description
mic	7:4	Mask Interrupt Class: all external interrupt vectors of equal or lower priority classes than the TPR.mic field are masked. For example, if mic field is 4, interrupt priority classes 1, 2, 3, and 4 are masked. A TPR.mic value of 0 has no masking effect; a value of 15 will mask all external interrupt vectors in the range 16 - 255. TPR.mic has no effect on external interrupt vectors 0 and 2, INITs and PMIs. See <a href="#">“Processor Interrupt Block” on page 2:127..</a>
mmi	16	Mask Maskable Interrupts: When 1, masks all external interrupts other than NMI (vector 2). When 0, external interrupt vectors 16 - 255, are masked by the TPR.mic field.

#### 5.8.3.4 End of External Interrupt Register (EOI – CR67)

A write to the EOI (end-of-external interrupt) register, shown in [Figure 5-9](#), indicates that software has finished servicing the highest priority in-service external interrupt. The processor removes its internal in-service indication for the highest priority currently in-service external interrupt vector. Pending external interrupts are then masked by the next highest priority in-service external interrupt (if any).

Writes to EOI affect the local processor only, and do not propagate to other processors or external interrupt controllers.

**Figure 5-9. End of External Interrupt Register (EOI – CR67)**



EOI is a read-write register. Reads return 0. Data associated with the EOI writes is ignored.

To ensure that the previous in-service interrupt indication has been cleared by a given point in program execution, software must perform a data serialization operation after an EOI write and prior to that point. To ensure that the reported IVR vector is correctly masked before the next IVR read, software must perform a data serialization operation after an EOI write and prior to that IVR read.

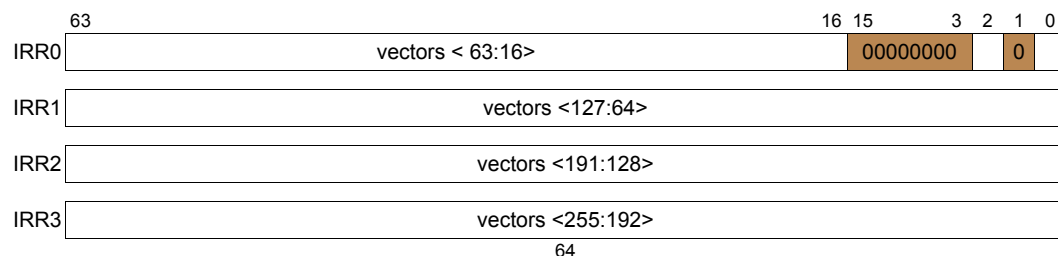


### 5.8.3.5 External Interrupt Request Registers (IRR0-3 – CR68,69,70,71)

Four 64-bit read-only External Interrupt Request Registers (IRR0-3, see [Figure 5-10](#)) provide the capability for software to determine the set of pending asynchronous external interrupts. IRR0 contains vectors <63:0> where vector 0 is in bit position 0, IRR1 contains vectors <127:64>, IRR2 contains vectors <191:128>, and IRR3 contains vectors <255:192>. A bit in the IRR, corresponding to the pending interrupt vector number, is set when the processor receives an external interrupt. The IRR bit is cleared when software reads the IVR and the vector number corresponding to the IRR bit value is returned in the IVR. The IRR bit is also cleared when a level-sensitive external interrupt signal is deasserted, effectively removing the pending interrupt.

Since IRR0-3 are read-only registers, writes to these registers result in Illegal Operation faults.

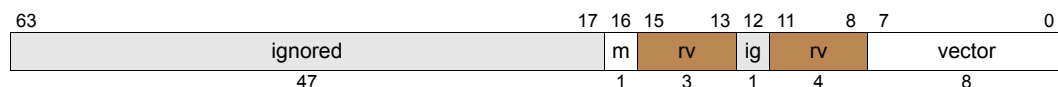
**Figure 5-10. External Interrupt Request Register (IRR0-3 – CR68, 69, 70, 71)**



### 5.8.3.6 Interval Timer Vector (ITV – CR72)

ITV specifies the external interrupt vector number for Interval Timer Interrupts. To ensure that subsequent interval timer interrupts reflect the new state of the ITV by a given point in program execution, software must perform a data serialization operation after an ITV write and prior to that point. See [Figure 5-11](#) and [Table 5-12](#) for the definitions of the ITV fields.

**Figure 5-11. Interval Timer Vector (ITV – CR72)**



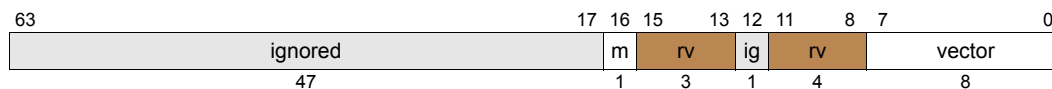
**Table 5-12. Interval Timer Vector Fields**

Field	Bits	Description
vector	7:0	External interrupt vector number to use when generating an Interval Timer interrupt. Vector values can be 0, 2 or 16-255. All other vectors are ignored and reserved for future use.
m	16	Mask: When 1, occurrences of Interval Timer interrupts are discarded and not pending. When 0, occurrences of Interval Timer interrupts are pending.

### 5.8.3.7 Performance Monitoring Vector (PMV – CR73)

PMV specifies the external interrupt vector number for Performance Monitoring overflow interrupts. To ensure that subsequent performance monitor interrupts reflect the new state of PMV by a given point in program execution, software must perform a data serialization operation after a PMV write and prior to that point. See Figure 5-12 and Table 5-13 for the definitions of the PMV fields.

**Figure 5-12. Performance Monitor Vector (PMV – CR73)**



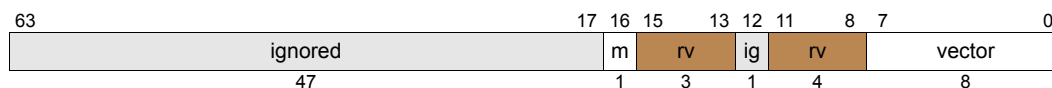
**Table 5-13. Performance Monitor Vector Fields**

Field	Bits	Description
vector	7:0	Vector number to use when generating a Performance Monitor interrupt. Vector values can be 0, 2, or 16-255. All other vectors are ignored and reserved for future use.
m	16	Mask: When 1, occurrences of Performance Monitor interrupts are discarded and not pended. When 0, occurrences of Performance Monitor interrupts are pended.

### 5.8.3.8 Corrected Machine Check Vector (CMCV – CR74)

CMCV specifies the external interrupt vector number for Corrected Machine Checks. To ensure that subsequent corrected machine check interrupts reflect the new state of CMCV by a given point in program execution, software must perform a data serialization operation after a CMCV write and prior to that point. See Figure 5-13 and Table 5-14 for the CMCV field definitions.

**Figure 5-13. Corrected Machine Check Vector (CMCV – CR74)**



**Table 5-14. Corrected Machine Check Vector Fields**

Field	Bits	Description
vector	7:0	Vector number to use when generating a Corrected Machine Check. Vector values can be 0, 2, or 16 - 255. All other vectors are ignored and reserved for future use.
m	16	Mask: When 1, occurrences of Corrected Machine Check interrupts are discarded and not pended. When 0, occurrences of Corrected Machine Check interrupts are pended.

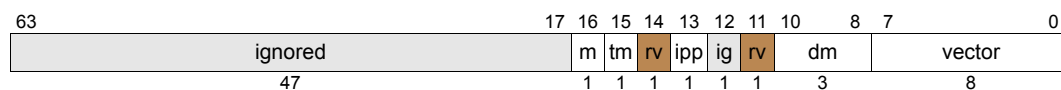
### 5.8.3.9 Local Redirection Registers (LRR0-1 – CR80,81)

Local Redirection Registers (LRR0-1) steer external signal-based interrupts that are directly connected to the local processor to a specific external interrupt vector. Processors may optionally support two direct external interrupt pins. When supported these external interrupt signals (pins) are referred to as Local Interrupt 0 (LINT0) and Local Interrupt 1 (LINT1). Software can query the presence of these pins via the PAL\_PROC\_GET\_FEATURES procedure call.

To ensure that subsequent interrupts from LINT0 and LINT1 reflect the new state of LRR prior to a given point in program execution, software must perform a data serialization operation after an LRR write and prior to that point. In the case when

LINT0 and LINT1 pins are absent, writes to LRR would have no effect, and reads from LRR would return 0. Software can query the presence of the LINT pins via the PAL\_PROC\_GET\_FEATURES procedure call. The LRR fields are defined in Figure 5-14 and Table 5-15.

**Figure 5-14. Local Redirection Register (LRR – CR80,81)**



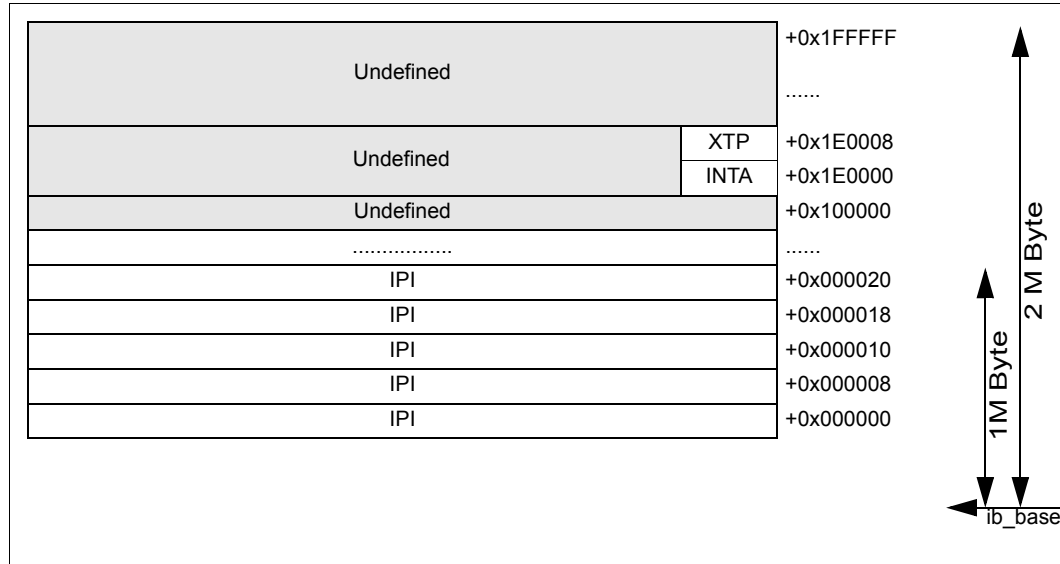
**Table 5-15. Local Redirection Register Fields**

Field	Bits	Description	
vector	7:0	External interrupt vector number to use when generating an interrupt for this entry. For INT delivery mode, allowed vector values are 0, 2, or 16-255. All other vectors are ignored and reserved for future use. For all other delivery modes this field is ignored.	
dm	10:8	000	INT – pend an external interrupt for the vector number specified by the vector field in LRR. Allowed vector values are 0, 2, or 16-255. All other vector numbers are ignored and reserved for future use.
		001	reserved
		010	PMI – pend a Platform Management Interrupt Vector number 0 for system firmware. The vector field is ignored.
		011	reserved
		100	NMI – pend a Non-Maskable Interrupt. This interrupt is pended at external interrupt vector number 2. The vector field is ignored.
		101	INIT – pend an Initialization Interrupt for system firmware. The vector field is ignored.
		110	reserved
		111	ExtINT – pend an Intel 8259A-compatible interrupt. This interrupt is delivered at external interrupt vector number 0. For details on servicing ExtINT external interrupts see “Interrupt Acknowledge (INTA) Cycle” on page 2:130. The vector field is ignored.
ipp	13	Interrupt Pin Polarity – specifies the polarity of the interrupt signal. When 0, the signal is active high. When 1, the signal is active low.	
tm	15	Trigger Mode – When 0, specifies edge sensitive interrupts. If the m field is 0, assertion of the corresponding LINT pin pends an interrupt for the specified vector corresponding to the dm field. The pending interrupt indication is cleared by software servicing the interrupt. When 1, specifies level sensitive interrupts. If the m field is 0, assertion of the corresponding LINT pin pends an external interrupt for the specified vector. Deassertion of the corresponding LINT pin clears the pending interrupt indication. The processor has undefined behavior if the dm and tm fields specify level sensitive PMIs or INITs.	
m	16	Mask – When 1, edge or level occurrences of the local interrupt pins are discarded and not pended. When 0, edge or level occurrences of local interrupt pins are pended.	

## 5.8.4 Processor Interrupt Block

Inter-Processor Interrupt (IPI) messages, Interrupt Acknowledge (INTA) cycles, and External Task Priority (XTP) cycles on the processor system bus are initiated by software by accessing a special physical memory range known as the “Processor Interrupt Block.” Figure 5-15 defines its memory layout. The entire 2 MByte Processor Interrupt Block is relocatable by a PAL firmware call and must be aligned on a 2 MByte boundary; by default, the block is located at physical address 0x0000 0000 FEE0 0000.

**Figure 5-15. Processor Interrupt Block Memory Layout**



The Inter-Processor Interrupt region occupies the lower half of the Processor Interrupt Block; by default its physical address range is 0x0000 0000 FEE0 0000 through 0x0000 0000 FEEF FFFF. A processor generates Inter-Processor Interrupts by performing an aligned 8-byte store to this memory region.

The Processor Interrupt Block does not support all forms of memory operations. Unsupported memory accesses result in undefined processor operation.

- When targeted at the inter-processor interrupt delivery region (lower half of the Processor Interrupt Block), the following memory operations are undefined: instruction fetch, RSE accesses, or memory read references (only writes are permitted), references other than aligned 8-byte accesses, and references through any memory attribute other than UC.
- When targeted at the upper half of the Processor Interrupt Block, the following memory operations are undefined: instruction fetches, references other than 1-byte accesses to the XTP byte and 1-byte read access to the INTA byte, and references through any memory attribute other than UC.

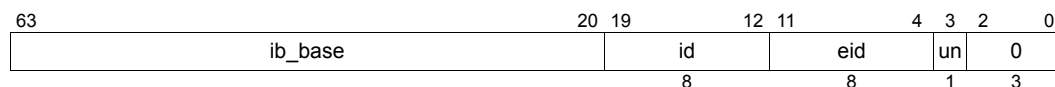
Any memory operation targeted at the lower half of the Processor Interrupt Block which does not correspond to any actual processor is undefined.

#### 5.8.4.1 Inter-processor Interrupt Messages

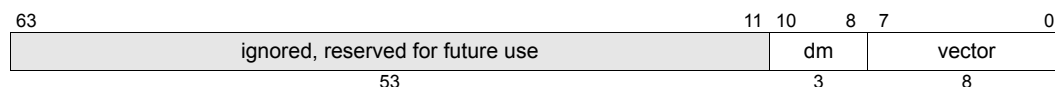
A processor can interrupt any individual processor, including itself, by issuing an inter-processor interrupt message (IPI). A processor generates an IPI by storing an 8-byte interrupt command to an 8-byte aligned address in the interrupt delivery region of the Processor Interrupt Block defined in “Processor Interrupt Block” on page 2:127. (If the address is not 8-byte aligned, the processor must either generate an Unaligned Data Reference Fault, see Section “Memory Datum Alignment and Atomicity” on page 2:93, or have undefined behavior). The address being stored to designates the target processor to receive the interrupt. The store address and data format of the

inter-processor interrupt message are defined in [Figure 5-16](#) and [Figure 5-17](#). The data fields are defined in [Table 5-17](#). The address processor identifier fields specify the target processor and are defined in [Table 5-16](#).

**Figure 5-16. Address Format for Inter-processor Interrupt Messages**



**Figure 5-17. Data Format for Inter-processor Interrupt Messages**



**Table 5-16. Address Fields for Inter-processor Interrupt Messages**

Field	Bits	Description
un	3	Unused. This field must be set to 0. Behavior of the inter-processor interrupt (IPI) message is undefined if this field is set to 1.
id/eid	19:4	Specify the target processor. See <a href="#">Table 5-10 on page 2:122</a> for a definition of these fields.
ib_base	63:20	Physical Base address of Processor Interrupt Block. This is a PAL relocatable physical address. The default is 0x0000 0000 FEE. See <a href="#">“Processor Interrupt Block” on page 2:127</a> . Based on the processor model some of the high order physical address bits may be reserved.

**Table 5-17. Data Fields for Inter-processor Interrupt Messages**

Field	Bits	Description	
vector	7:0	Vector number for the interrupt. For INT delivery, allowed vector values are 0, 2, or 16-255. All other vectors are ignored and reserved for future use. For PMI delivery, allowed PMI vector values are 0-3. All other PMI vector values are reserved for use by processor firmware.	
dm	10:8	000	INT – pend an external interrupt for the specified vector to the processor listed in the destination. Allowed vector values are 0, 2, or 16-255. All other vector numbers are ignored and reserved for future use.
		001	Reserved
		010	PMI – pend a PMI interrupt for the specified vector to the processor listed in the destination. Allowed PMI vector values are 0-3. All other PMI vector values are reserved for use by processor firmware. See <a href="#">Section 11.5, “Platform Management Interrupt (PMI)” on page 2:310</a> for details.
		011	Reserved
		100	NMI – pend an external interrupt as an NMI (vector 2) to the processor listed in the destination. The vector field is ignored.
		101	INIT – pend an Initialization Interrupt for platform firmware on the processor listed in the destination. The vector field is ignored.
		110	Reserved
	111	ExtINT – pend an Intel 8259A-compatible interrupt. This interrupt is delivered at external interrupt vector number 0. For details on servicing ExtINT external interrupts see <a href="#">“Interrupt Acknowledge (INTA) Cycle” on page 2:130</a> . The vector number field is ignored.	
ignored	63:11	Ignored, reserved for future use	

### 5.8.4.2 Interrupt and IPI Ordering

Interrupt messages from external device(s), or external interrupts routed to the processor's LINT pins, when present, may arrive at one or more processors and become pending in any order. No ordering is enforced by the processor or the platform.

As observed by a receiving processor, IPIs emitted from the same issuing processor may be pended in any order, even when the receiving processor and the issuing processor are the same.

As observed by a receiving processor, IPIs are pended after all prior loads and stores emitted by the same issuing processor are visible if and only if the IPI is issued with a `st.rel` (or preceded by an `mf`), even when the receiving processor and the issuing processor are the same. For all other cases, no ordering is implied between IPI transactions and prior cacheable or uncached memory references.

As observed by a receiving processor, no ordering is implied between IPIs and subsequent loads/stores from the same issuing processor, even when the receiving processor and the issuing processor are the same. Subsequent loads or stores may become visible before an IPI is seen as pending. Data or instruction serialization operations, memory fences (`mf` or `mf.a`), or `st.rel` do not ensure an IPI is pending at the target processor (including self) by a given point in program execution on the local processor.

### 5.8.4.3 Interrupt Acknowledge (INTA) Cycle

Intel 8259A-compatible external interrupt controllers can not issue interrupt messages and therefore do not specify an external interrupt vector number when the interrupt request is generated. When accepting an external interrupt, software must inspect the vector number supplied by the IVR register. If the vector matches the vector number assigned to the external controller (can be `ExtINT`, or any other vector number based on software convention), software must acquire the actual external interrupt vector number from the external interrupt controller by issuing a 1-byte load from the INTA Byte.

The INTA Byte is located within the upper half of the Processor Interrupt Block, at offset `0x1E0000` from the base. A single byte load from the INTA address causes the processor to emit the INTA cycle on the processor system bus. An Intel 8259A-compatible external interrupt controller must respond with the actual interrupt vector number as the data to be loaded. If two INTA cycles are required by the external interrupt controller, the platform must provide this functionality. Any memory operation to the INTA address other than a single byte load is undefined.

Software must issue an EOI to the local processor, to clear the interrupt in-service indication for the vector associated with the external interrupt controller.

### 5.8.4.4 External Task Priority (XTP) Cycle

Some model-specific system configurations support an External Task Priority Register (XTPR) per processor in external bus logic. A processor's XTPR can be modified by storing one byte of data to the processor's XTP Byte address. This generates a special bus transaction required to change the processor's XTPR within the system. Please refer to system-specific documentation for XTPR bit format and field definitions. The

processor does not interpret any data stored to the XTP Byte address and all data bits are passed to the external system unmodified. Any memory operation to the XTP address other than a single byte store is undefined.

XTPR is written by operating system code to notify the system that the processor's current task priority has been changed. Based on this task priority information, system implementations can steer interrupt messages from the I/O subsystems to the processors that have registered the lowest task priority levels. The XTPR register is a system performance hint and need not be updated by operating system code nor be implemented in all system configurations. If the system does not implement the XTPR, it must still accept a processor's XTP cycle and discard it. Operating system code can issue XTPR updates regardless of external system support.

### 5.8.5 Edge- and Level-sensitive Interrupts

The processor's LINT pins, when present, directly support edge and level sensitive interrupts, however all other interrupt sources are edge sensitive. A single external interrupt messages is issued only on the assertion of an interrupt by external interrupt controllers or devices, deassertion of an external interrupt sends no interrupt message to the processor. Since the processor removes the pending interrupt when the interrupt is serviced, the processor guarantees exactly-one interrupt acceptance for each external interrupt message. By definition external interrupt messages are edge sensitive.

Level sensitive external interrupts can be supported using edge sensitive interrupt messages as follows:

- Software services the external interrupt generated by an edge interrupt message.
- Software removes the external interrupt request from the requesting device, the device should then deassert its interrupt request line.
- To avoid spurious external interrupts, it is highly recommended that software issue a dummy read from the device to ensure that the interrupt request has been actually been removed before the interrupt is resampled in the next step.
- Software issues a command to the external interrupt controller to resample the interrupt (typically an external interrupt controller end-of-interrupt command). The external interrupt controller must issue another interrupt message back to the processor if service is still required by the processor for a given vector number. For example, if there are other devices still requiring service that are attached to the same level sensitive interrupt request line.

#### §





The register stack engine (RSE) moves registers between the register stack and the backing store in memory without explicit program intervention. The RSE operates concurrently with the processor and can take advantage of unused memory bandwidth to dynamically issue register spill and fill operations. In this manner, the latency of register spill/fill operations can be overlapped with useful program work. The basic principles of the register stack are discussed in [Section 4.1, “Register Stack” on page 1:47](#). This chapter presents the internal state, the programming model and the interruption behavior of the register stack engine.

## 6.1 RSE and Backing Store Overview

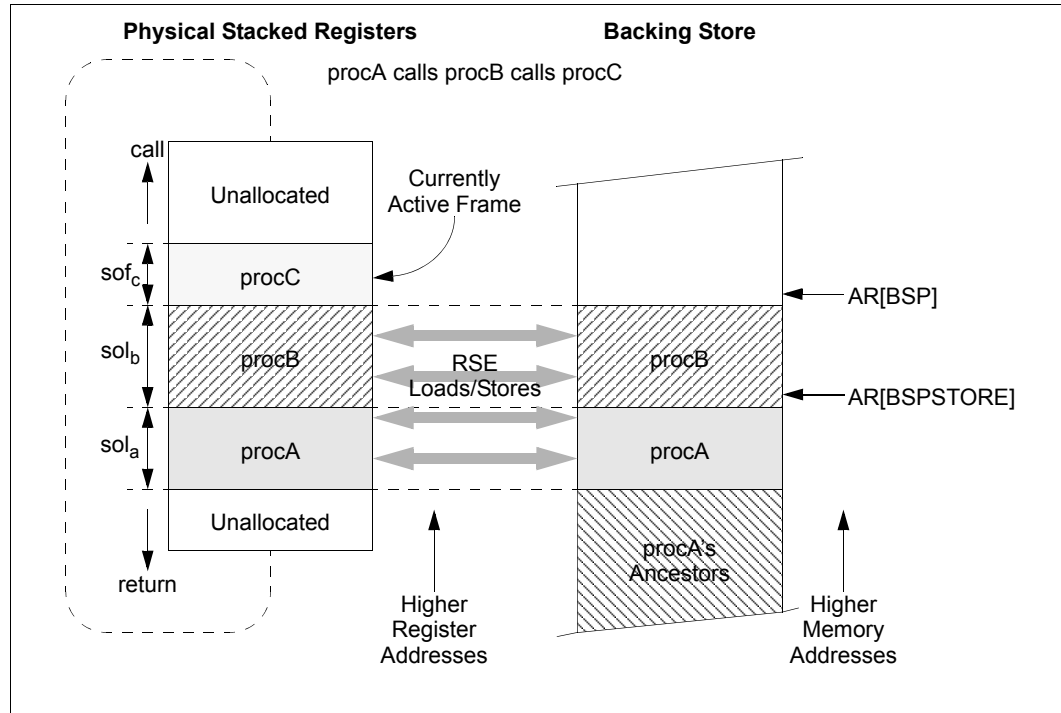
The register stack frames are mapped onto a set of physical registers which operate as a circular buffer containing the most recently created frames. The RSE spills and fills these physical registers to/from a backing store in memory. The RSE moves registers between the physical register stack and the backing store without explicit program intervention. As indicated in [Figure 6-1](#), the RSE operates on the physical stacked registers outside of the currently active frame (as defined by CFM). These registers contain the frames of the parent procedures of the current procedure.

As shown in [Figure 6-1](#), the backing store is organized as a stack in memory that grows from lower to higher addresses. The Backing Store Pointer (BSP) application register contains the address of the first (lowest) memory location reserved for the current frame (i.e., the location at which GR32 of the current frame will be spilled). RSE spill/fill activity occurs at addresses below what is contained in the BSP since the RSE spills/fills the frames of the current procedure’s parents. The BSPSTORE application register contains the address at which the next RSE spill will occur. The address register which corresponds to the next RSE fill operation, the BSP load pointer, is not architecturally visible. The addresses contained in BSP and BSPSTORE are always aligned to an 8-byte boundary. The backing store contains the local area of each frame. The output area is not spilled to the backing store (unless it later becomes part of a callee’s local area). Within each stack frame, lower-addressed registers are stored at lower memory addresses. RSE spills of NaTed stacked general registers are subject to the same memory update constraints as software spills (st8.spill) of NaTed static general registers (see [“Register Spill and Fill” on page 1:62](#)).

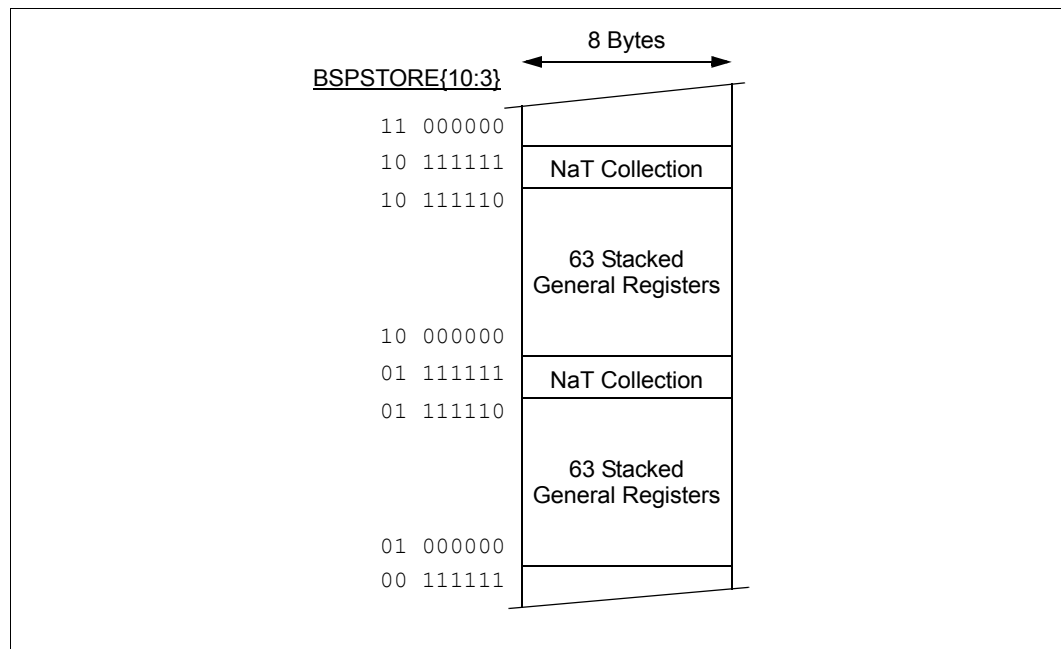
The RSE also spills/fills the NaT bits corresponding to the stacked registers. The NaT bits corresponding to the static subset must be spilled/filled as necessary by software. The NaT bits are the 65th bit of each general register. The NaT bits for the stacked subset are spilled/filled in groups of 63 corresponding to 63 consecutive physical stacked registers. When the RSE spills a register to the backing store the corresponding NaT bit is copied to the RSE NaT collection (RNAT) application register. Whenever bits 8:3 of BSPSTORE are all ones, the RSE stores RNAT to the backing store. As shown in [Figure 6-2](#), this results in a backing store memory image in which every 63 register values are followed by a collection of NaT bits. Bit 0 of the NaT collection corresponds to the first (lowest addressed) of the 63 register values; bit 62 corresponds to the 63rd register value. Bit 63 of the NaT collection is always written as zero. When the RSE fills

a stacked register from the backing store it also fills the register's NaT bit. Whenever bits 8:3 of the RSE backing store load pointer are all ones, the RSE reloads a NaT collection from the backing store. Bit 63 of the NaT collection is ignored when read from the backing store.

**Figure 6-1. Relationship Between Physical Registers and Backing Store**



**Figure 6-2. Backing Store Memory Format**



The RSE operates concurrently and asynchronously with respect to instruction execution by taking advantage of unused memory bandwidth to dynamically perform register spill and fill operations. The algorithm employed by the RSE to determine whether and when to spill/fill is implementation dependent. Software can not depend on the spill/fill algorithm. To ensure that the processor and RSE activities do not interfere with each other, software should not access stacked registers outside of the current stack frame. The architecture guarantees register stack integrity by faulting on writes to out-of-frame registers. Reads from out-of-frame registers may interact with RSE operations and return undefined data values. However, out-of-frame reads are required to propagate NaT bits.

The operation of the RSE is controlled by the Register Stack Configuration (RSC) application register. Activity between the processor and the RSE is synchronized only when `alloc`, `flushrs`, `loadrs`, `br.ret`, or `rfi` instructions actually require registers to be spilled or filled, or when software explicitly requests RSE synchronization by executing a `mov to/from RSC`, `BSPSTORE` or `RNAT` application register instruction.

## 6.2 RSE Internal State

Table 6-1 describes architectural state that is maintained by the register stack engine. The RSE internal state elements described here are not directly exposed to the programmer as architecturally visible registers. As a consequence, RSE internal state does not need to be preserved across context switches or interruptions. Instead, it is modified as the side-effect of register stack-related instructions. To describe the effects of these instructions a complete definition of the RSE internal state is essential. To distinguish them from architecturally visible resources, all RSE internal state elements are prefixed with "RSE." Other RSE related resources are architecturally visible and are exposed to software as application registers: `RSC`, `BSP`, `BSPSTORE`, and `RNAT`.

**Table 6-1. RSE Internal State**

Name	Description	Corresponds To
RSE.N_STACKED_PHYS	Number of Stacked Physical registers: Implementation dependent size of the stacked physical register file.	
RSE.BOF	Bottom-of-frame register number: Physical register number of GR32.	AR[BSP]
RSE.StoreReg	RSE Store Register number: Physical register number of next register to be stored by RSE.	AR[BSPSTORE]
RSE.LoadReg	RSE Load Register number: Physical register number one greater than the next register to load (modulo the number of stacked physical registers).	RSE.BspLoad
RSE.BspLoad	Backing Store Pointer for memory loads: 64-bit Backing Store Address 8 bytes greater than the next address to be loaded by the RSE.	RSE.BspLoad
RSE.RNATBitIndex	RSE NaT Collection Bit Index: 6-bit wide RNAT Collection Bit Index (defines which RNAT collection bit gets updated)	AR[BSPSTORE][8:3]
RSE.CFLE	RSE Current FrameLoad Enable: Control bit that permits the RSE to load registers in the current frame after a <code>br.ret</code> or <code>rfi</code> .	

**Table 6-1. RSE Internal State (Continued)**

Name	Description	Corresponds To
RSE.ndirty	Number of dirty registers on the register stack	
RSE.ndirty_words	Number of dirty words on the register stack plus corresponding number of NaT collection registers	AR[BSP] - AR[BSPSTORE]

## 6.3 Register Stack Partitions

The processor's physical register file provides at least 96 stacked registers. The actual number of stacked registers (RSE.N\_STACKED\_PHYS) is implementation dependent and must be an even multiple of 16. Figure 6-3 illustrates the circular nature of the physical register file, and shows the correspondence of the registers to the backing store. Figure 6-3 also shows the four partitions of the stacked register file:

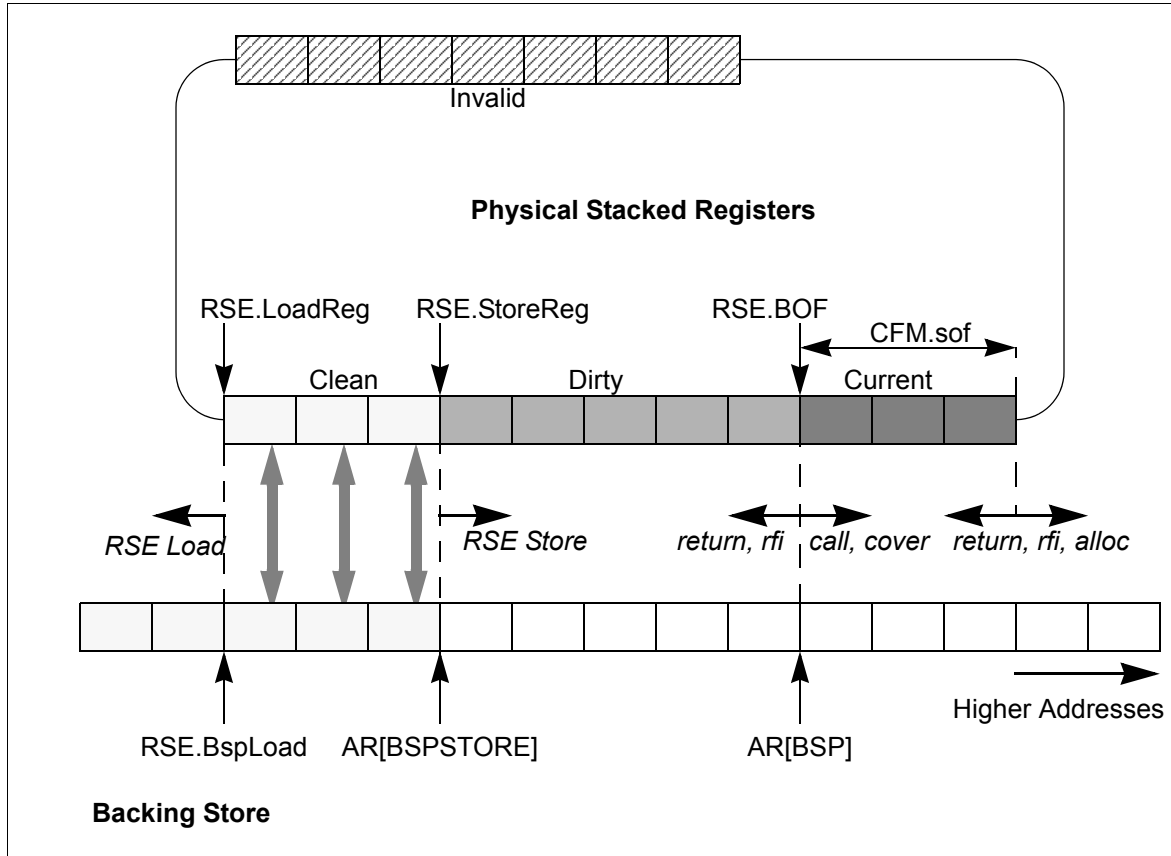
**Clean** partition (lightly-shaded): registers that contain values from parent procedure frames. The registers in this partition have been successfully spilled to the backing store by the RSE and their contents have not been modified since they were written to the backing store.

**Dirty** partition (medium-shaded): registers that contain values from parent procedure frames. The registers in this partition have not yet been spilled to the backing store by the RSE. The number of registers contained in the dirty partition (distance between RSE.StoreReg and RSE.BOF) is referred to as RSE.ndirty.

**Current** frame (shaded dark): stacked registers allocated for computation. The position of the current frame in the physical stacked register file is defined by the Bottom-of-frame register (RSE.BOF). The number of registers in the current frame is defined by the size of frame field in the current frame marker (CFM.sof).

**Invalid** partition (diagonally striped): registers outside the current frame that do not contain values from parent procedure frames. They are immediately available for allocation into the current frame or for RSE load operations.

**Figure 6-3. Four Partitions of the Register Stack**



The boundaries between the four register stack partitions are defined by the current frame marker (CFM) and three physical register numbers: a load, store and bottom-of-frame register number. As described in Table 6-1 each of these physical register numbers has a corresponding 64-bit backing store memory address pointer. (For example, AR[BSP] always contains the address where GR[32] of the current frame will be stored.)

Figure 6-3 also shows the effects of various instructions on the partition boundaries. RSE loads use invalid registers. RSE stores use dirty registers. Eager RSE loads and stores grow the clean partition. A `br.call`, `brl.call`, or `cover` instruction can increase the bottom-of-frame pointer (RSE.BOF) which moves registers from the current frame to the dirty partition. An `alloc` may shrink or grow the current frame by updating CFM.sof. A `br.ret` or `rfi` instruction may shrink or grow the current frame by updating both the bottom-of-frame pointer (RSE.BOF) and CFM.sof.

## 6.4 RSE Operation

The register stack backing store is organized as a stack in memory that grows from lower addresses towards higher addresses. The top of the backing store stack is defined by the Backing Store Pointer (BSP) application register, which points to the first memory location reserved for the current frame. The RSE load and store activities take

place at lower addresses, defined relative to BSP by the sizes of the clean and dirty partitions. Although the stack is conceptually infinite in both directions, the effective base of the stack is expected to be the first memory location of the first page allocated to the backing store.

To allow the highest possible degree of concurrent execution, the processor and the RSE operate independently of each other during normal program execution. The RSE distinguishes between **mandatory** and **eager** load/store operations. Mandatory load/store operations occur as the result of `alloc`, `flushrs`, `loadrs`, `br.ret` or `rfi` instructions. Eager operations occur when the RSE is speculatively working ahead of program execution, and it is not known whether this register spill/fill is actually required by the program.

When the RSE works in the background, it issues eager RSE spill and fill operations to extend the size of the clean partition in both directions—by decreasing the RSE load pointer and loading values from the backing store into invalid registers (eager RSE load), and by saving dirty registers to the backing store and increasing the RSE store pointer (eager RSE store). Allocation of a sufficiently large frame (using `alloc`) or execution of a `flushrs` instruction may cause the RSE to suspend program execution and issue mandatory RSE stores until the required number of registers have been spilled to the backing store. Similarly a `br.ret` or `rfi` back to a sufficiently large frame or execution of a `loadrs` instruction may cause the RSE to suspend program execution and issue mandatory RSE loads until the required number of registers have been restored from the backing store. The RSE only operates in the foreground and suspends program execution whenever forward progress of the program actually requires registers to be spilled or filled.

Table 6-2 describes the RSE operation instructions and state modifications.

**Table 6-2. RSE Operation Instructions and State Modification**

Affected State	Instruction			
	<code>alloc</code> $r_1 = ar.pfs, i, l, o, r^a$	<code>br.call<sup>a</sup></code> , <code>brl.call<sup>a</sup></code>	<code>br.ret<sup>a</sup></code>	<code>rfi</code> when <code>CR[IFS].v = 1</code>
AR[BSP]{63:3}	unchanged	AR[BSP]{63:3} + CFM.sol + (AR[BSP]{8:3} + CFM.sol)/63	AR[BSP]{63:3} - AR[PFS].pfm.sol - (62-AR[BSP]{8:3} + AR[PFS].pfm.sol)/63	AR[BSP]{63:3} - CR[IFS].ifm.sof - (62-AR[BSP]{8:3} + CR[IFS].ifm.sof)/63
AR[PFS]	unchanged	AR[PFS].pfm = CFM AR[PFS].pec = AR[EC] AR[PFS].ppl = PSR.cpl	unchanged	unchanged
GR[r <sub>i</sub> ]	AR[PFS]	N/A	N/A	N/A
CFM	CFM.sof = $i+1+o$ CFM.sol = $i+1$ CFM.sor = $r \gg 3$	CFM.sof -= CFM.sol CFM.sol = 0 CFM.sor = 0 CFM.rrb.gr = 0 CFM.rrb.fr = 0 CFM.rrb.pr = 0	AR[PFS].pfm or <sup>b</sup> CFM.sof = 0 CFM.sol = 0 CFM.sor = 0 CFM.rrb.gr = 0 CFM.rrb.fr = 0 CFM.rrb.pr = 0	CR[IFS].ifm

a. These instructions have undefined behavior with an incomplete frame. See “RSE Behavior with an Incomplete Register Frame” on page 2:146.

b. Normal `br.ret` instructions restore CFM with AR[PFS].pfm. However, if a bad PFS value is read by the `br.ret` instruction, all CFM fields are set to zero. See “Bad PFS used by Branch Return” on page 2:143.

## 6.5 RSE Control

The RSE can be controlled at all privilege levels by means of three instructions (`cover`, `flushrs`, and `loadrs`) and by accessing four application registers (`mov to/from RSC`, `BSP`, `BSPSTORE` and `RNAT`). This section first presents each of the RSE application registers, and then discusses the three RSE control instructions.

### 6.5.1 Register Stack Configuration Register

The layout of the Register Stack Configuration application register (RSC) is defined in [Section 3.1.8.2, “Register Stack Configuration Register \(RSC – AR 16\)” on page 1:29](#). This section describes the semantics of the mode, the privilege level and the byte order fields of the RSC. The `loadrs` field is described as part of the `loadrs` instruction in [Section 6.5.4, “RSE Control Instructions” on page 2:142](#).

**RSE Mode:** Two mode bits in the RSC register determine when the RSE generates register spill or fill operations. When both mode bits are zero (enforced lazy mode) the RSE issues only mandatory loads and stores (when an `alloc`, `br.ret`, `flushrs` or `rfi` instruction requires registers to be spilled or filled). Bit 0 of the `RSC.mode` field enables eager RSE stores and bit 1 enables eager RSE loads. [Table 6-3](#) defines all four possible RSE modes. Please see the processor-specific documentation for further information on the RSE modes implemented by the Itanium processor.

**Table 6-3. RSE Modes (RSC.mode)**

Mode	RSE Loads	RSE Stores	RSC.mode
Enforced Lazy	Mandatory only	Mandatory only	00
Store Intensive	Mandatory only	Eager and Mandatory	01
Load Intensive	Eager and Mandatory	Mandatory only	10
Eager	Eager and Mandatory	Eager and Mandatory	11

The algorithm that decides whether and when to speculatively perform eager register spill or fill operations is implementation dependent. Software may not make any assumptions about the RSE load/store behavior when the `RSC.mode` is non-zero. Furthermore, access to the `BSPSTORE` and `RNAT` application registers and the execution of the `loadrs` instructions require `RSC.mode` to be zero (enforced lazy mode). If `loadrs`, move to/from `BSPSTORE` or move to/from `RNAT` are executed when `RSC.mode` is non-zero an Illegal operation fault is raised. Eager spill/fill of the `RNAT` register to/from the backing store is only permitted if the RSE is in store/load intensive or eager mode. In enforced lazy mode, the RSE may spill/fill the `RNAT` register only if a subsequent mandatory register spill/fill is required.

**RSE Privilege Level:** When address translation is enabled (`PSR.rt` is one), the RSE operates at a privilege level defined by two privilege level bits in the Register Stack Configuration register (`RSC.pl`). All privilege level checks for RSE virtual accesses are performed using the privilege level in `RSC.pl`. When the RSC is written, the privilege level bits are clipped to the current privilege level of the process, i.e., the numerical maximum of the current privilege level and the privilege level in the source register is written to `RSC.pl`.

Protection is also checked based on the current entries in the data TLB. The RSE always remains coherent with respect to the data TLB. If a translation that is being used by the RSE is changed or purged, the RSE will immediately begin using the new translation or suffer a TLB miss. Only mandatory loads and stores can cause RSE memory related faults. Details on RSE fault delivery are described in [“RSE Interruptions”](#) Although eager RSE loads and stores do not cause interruptions they can, under certain conditions, cause a VHPT walk and TLB insert. Details on when RSE loads and stores can cause a VHPT walk are described in [“VHPT Environment” on page 2:67](#).

The RSE expects its backing store to be mapped to cacheable speculative memory. If RSE spill/fill transactions are performed to non-speculative memory that may contain I/O devices, system behavior is unpredictable.

**RSE Byte Order:** Because the RSE runs asynchronously with the processor, it may be running on behalf of a context with a different byte order from the current one. Consequently, the RSE defines its own byte ordering bit: RSC.be. When RSC.be is zero, registers are stored in little-endian byte order (least significant bytes to lower addresses). When RSC.be is one, registers are stored in big-endian byte order (most significant bytes to lower addresses). RSC.be also determines the byte order of NaT collections spilled/filled by the RSE. RSC.be may be written by code at any privilege level. Changes to RSC.be should only be made by software when RSC.mode is zero. Failure to do so results in undefined backing store contents.

## 6.5.2 Register Stack NaT Collection Register

As described in [Section 6.1, “RSE and Backing Store Overview” on page 2:133](#), the RSE is responsible for saving and restoring NaT bits associated with the stacked registers to and from the backing store. The RSE writes its NaT collection register (the RNAT application register) to the backing store whenever  $BSPSTORE\{8:3\} = 0x3F$  (1 NaT collection for every 63 registers). The RNAT acts as a temporary holding area for up to 63 unsaved NaT bits. The RSE NaT collection bit index (RSE.RNATBitIndex) determines which bit of the RNAT register receives the NaT bit of a spilled register as the result of an RSE store. The six-bit wide RSE.RNATBitIndex is always equal to  $BSPSTORE\{8:3\}$ . As a result,  $RNAT\{x\}$  corresponds to the register saved at

$$\text{concatenate}(BSPSTORE\{63:9\}, x\{5:0\}, 0\{2:0\}) .$$

The RSE never saves partial NaT collections to the backing store, so software must save and restore the RNAT application register when switching the backing store pointer. RSE.RNATBitIndex determines which RNAT bits are valid. Bits  $RNAT\{RSE.RNATBitIndex:0\}$  contain defined values, and bits  $RNAT\{62:RSE.RNATBitIndex+1\}$  contain undefined values. Bit 63 of the RNAT application register always reads as zero. Writes to bit 63 of the RNAT application register are ignored. The execution of RSE control instructions `mov` to `BSPSTORE` and `loadrs` as well as an RSE spill of the RNAT register cause the contents of the RNAT register to become undefined. The RNAT application register can only be accessed when RSC.mode is zero. If RSC.mode is non-zero, accessing the RNAT application register results in an Illegal Operation fault.



### 6.5.3 Backing Store Pointer Application Registers

The RSE defines two Backing Store Pointer application registers: BSPSTORE and BSP. Since the RSE backing store pointers are always 8-byte aligned, bits {2:0} of the backing store pointers always read as zero. When writing the BSPSTORE application register, bits {2:0} in the presented address are ignored.

The RSE Backing Store Pointer for memory stores (BSPSTORE) is a 64-bit application register that provides the main interface to the three RSE backing store memory pointers: BSP, BSPSTORE and RSE.BspLoad. The BSPSTORE application register can only be accessed when RSC.mode is zero. If RSC.mode is non-zero, accessing BSPSTORE results in an Illegal Operation fault.

Reading BSPSTORE (`mov` from BSPSTORE application register) returns the address of the next RSE store.

Writing BSPSTORE (`mov` to BSPSTORE application register) has side-effects on all three RSE pointers and the NaT collection process. The operation is defined as follows: the BSPSTORE and RSE.BspLoad pointers are both set to the address presented, which forces the size of the clean partition to zero. Writes to the BSPSTORE application register do not change the size of the dirty partition: the BSP pointer is set to the address presented plus the size of the dirty partition plus the size of any intervening NaT collections. The dirty partition is preserved to allow software to change the backing store pointer without having to flush the register stack. Writing BSPSTORE causes the contents of the RNAT register to become undefined. Therefore software must preserve the contents of RNAT prior to writing BSPSTORE. After writing to BSPSTORE, the NaT collection bit index (RSE.RNATBitIndex) is set to bits {8:3} of the presented address. If an unimplemented address in BSPSTORE is used by a mandatory RSE spill or fill, an Unimplemented Data Address fault is raised.

The RSE Backing Store Pointer (BSP) is a 64-bit read-only application register. Writing BSP (`mov` to BSP application register) results in an Illegal Operation fault. Reads from BSP (`mov` from BSP application register) return the address of the top of the register stack in memory. This location is the backing store address to which the current GR32 would be written. Reading BSP does not have any side-effect on any of the internal RSE pointers or the NaT collection process. Therefore, BSP can be read regardless of the RSE mode, i.e., even when RSC.mode is non-zero. Since BSP is determined by BSPSTORE and the size of the dirty partition, it is possible for BSPSTORE to contain an implemented address and for BSP to contain an unimplemented address. BSP reads always return a full 64-bit (possibly unimplemented) address; only a subsequent data memory reference with an unimplemented address will cause an Unimplemented Data Address fault.

[Table 6-4](#) summarizes the effects of the three instructions that access the backing store pointer application registers.

**Table 6-4. Backing Store Pointer Application Registers**

Affected State	Instruction		
	Read BSP <code>mov r<sub>1</sub>=AR[BSP]</code>	Read BSPSTORE <code>mov r<sub>1</sub>=AR[BSPSTORE]</code>	Write BSPSTORE <sup>a</sup> <code>mov AR[BSPSTORE]=r<sub>2</sub></code>
GR[ <i>r</i> <sub>1</sub> ]	AR[BSP]	AR[BSPSTORE]	N/A
AR[BSP]{63:3}	Unchanged	Unchanged	(GR[ <i>r</i> <sub>2</sub> ]{63:3} + RSE.ndirty) + ((GR[ <i>r</i> <sub>2</sub> ]{8:3} + RSE.ndirty)/63)
AR[BSPSTORE]{63:3}	Unchanged	Unchanged	GR[ <i>r</i> <sub>2</sub> ]{63:3}
RSE.BspLoad {63:3}	Unchanged	Unchanged	GR[ <i>r</i> <sub>2</sub> ]{63:3}
AR[RNAT]	Unchanged	Unchanged	UNDEFINED
RSE.RNATBitIndex	Unchanged	Unchanged	GR[ <i>r</i> <sub>2</sub> ]{8:3}

a. Writing to AR[BSPSTORE] has undefined behavior with an incomplete frame. See “RSE Behavior with an Incomplete Register Frame” on page 2:146.

## 6.5.4 RSE Control Instructions

This section describes the RSE control instructions: `cover`, `flushrs` and `loadrs`. The effects of the three RSE control instructions on the RSE state are summarized in Table 6-5.

The `cover` instruction adds all registers in the current frame to the dirty partition, and allocates a zero-size current frame. As a result AR[BSP] is updated. `cover` clears the register rename base fields in the current frame marker CFM. If PSR.ic is zero, the original value of CFM is copied into CR[IFS].ifm and CR[IFS].v is set to one. The `cover` instruction must be the last instruction in an instruction group; otherwise, operation is undefined.

The `flushrs` instruction spills all dirty registers to the backing store. When it completes, RSE.ndirty is defined to be zero, and BSPSTORE equals BSP. Since `flushrs` may cause RSE stores, the RNAT application register is updated. A `flushrs` instruction must be the first instruction in an instruction group otherwise the results are undefined.

The `loadrs` instruction ensures that a specified portion of the backing store below the current BSP is present in the physical stacked registers. The size of the backing store section is specified in the `loadrs` field of the RSC application register (AR[RSC].loadrs). After `loadrs` completes, all registers and NaT collections between the current BSP and the tear-point (BSP-(RSC.loadrs{13:3} << 3)), and no more than that, are guaranteed to be present and marked as dirty in the stacked physical registers. When `loadrs` completes BSPSTORE and RSE.BspLoad are defined to be equal to the backing store tear-point address. All other physical stacked registers are marked invalid.

- If the tear-point specifies an address below RSE.BspLoad, the RSE issues mandatory loads to restore registers and NaT collections. All registers between the current BSP and the tear-point are marked dirty.
- If the RSE has already loaded registers beyond the tear-point when the `loadrs` instruction executes, the RSE marks clean registers below the tear-point as invalid and marks clean registers above the tear-point as dirty.
- If the tear-point specifies an address greater than BSPSTORE, the RSE marks clean and dirty registers below the tear-point as invalid (in this case dirty registers are lost).

**Table 6-5. RSE Control Instructions**

Affected State	Instruction		
	cover	flushrs <sup>a</sup>	loadrs <sup>a</sup>
AR[BSP]{63:3}	AR[BSP]{63:3} + CFM.sof + (AR[BSP]{8:3} + CFM.sof)/63	Unchanged	Unchanged
AR[BSPSTORE]{63:3}	Unchanged	AR[BSP]{63:3}	AR[BSP]{63:3} - AR[RSC].loadrs{13:3}
RSE.BspLoad{63:3}	Unchanged	Model specific <sup>b</sup>	AR[BSP]{63:3} - AR[RSC].loadrs{13:3}
AR[RNAT]	Unchanged	Updated	UNDEFINED
RSE.RNATBitIndex	Unchanged	AR[BSPSTORE]{8:3}	AR[BSPSTORE]{8:3}
CR[IFS]	if (PSR.ic == 0) { CR[IFS].ifm = CFM CR[IFS].v = 1}	Unchanged	Unchanged
CFM	CFM.sof = 0 CFM.sol = 0 CFM.sor = 0 CFM.rrb.gr = 0 CFM.rrb.fr = 0 CFM.rrb.pr = 0	Unchanged	Unchanged

a. These instructions have undefined behavior with an incomplete frame. See “RSE Behavior with an Incomplete Register Frame” on page 2:146.

b. In general, eager RSE implementations will preserve RSE.BspLoad during a `flushrs`. Lazy RSE implementations may set RSE.BspLoad to AR[BSPSTORE] after `flushrs` completes or faults.

By specifying a zero RSC.loadrs value `loadrs` can be used to invalidate all stacked registers outside the current frame. `loadrs` causes the contents of the RNAT register to become undefined. The NaT collection index is set to bits {8:3} of the new BSPSTORE. A `loadrs` instruction must be the first instruction in an instruction group otherwise the results are undefined. The following conditions cause `loadrs` to raise an Illegal Operation fault:

- If RSC.mode is non-zero.
- If both CFM.sof and RSC.loadrs are non-zero.
- If RSC.loadrs specifies more words to be loaded than will fit in the stacked physical register file (RSE.N\_STACKED\_PHYS).

### 6.5.5 Bad PFS used by Branch Return

On a `br.ret`, if the PFS application register defines an output area which is larger than the number of implemented stacked registers minus the size of dirty partition ((AR[PFS].sof - AR[PFS].sol) > (RSE.N\_STACKED\_PHYS - RSE.ndirty)), the return will not restore CFM with AR[PFS].pfm (normal behavior); instead, the return sets all fields in the CFM (of the procedure being returned to) to zero.

Typical procedure call and return sequences that preserve PFS values and that do not use `cover` or `loadrs` instructions will not encounter this situation.

The RSE will detect the above condition on a `br.ret`, and update its state as follows:

- The register rename base (RSE.BOF), AR[BSP], and AR[BSPSTORE] are updated as required by the return.

- The CFM (after the return) is forced to zero; i.e., all CFM fields (including CFM.sof and CFM.sol) are set to zero.
- The registers from the returned-from frame and the preserved registers from the returned-to frame are added to the invalid partition of the register stack.
- The dirty partition of the register stack is shrunk by AR[PFS].pfm.sol.
- The clean partition of the register stack remains unchanged. RSE.BspLoad and RSE.LoadReg remain unchanged.
- No other indication is given to software.

Since the size of the current frame is set to zero, the contents of some (possibly all) stacked GRs may be overwritten by subsequent eager RSE operations or by subsequent instructions allocating a new stack frame and then targeting a stacked GR. Therefore, explicit register stack management sequences that manipulate PFS, use the `cover` instruction, or use the `loadrs` instruction must avoid this situation by executing one of the two following code sequences prior to a `br.ret`:

- Use a `flushrs` instruction prior to the `br.ret`. This preserves all dirty registers to memory, and sets RSE.ndirty to zero, which avoids the condition.
- Use a `loadrs` instruction with an AR[RSC].loadrs value in the following range:  

$$\text{AR[RSC].loadrs} \leq 8 * (\text{ndirty\_max} + ((62 - \text{AR[BSP]\{8:3\}} + \text{ndirty\_max}) / 63)),$$
 where  $\text{ndirty\_max} = (\text{RSE.N\_STACKED\_PHYS} - (\text{AR[PFS].sof} - \text{AR[PFS].sol}))$

This adjusts the size of the dirty partition appropriately to avoid the condition. A `loadrs` with `RSC.loadrs=0` works on all processor models, regardless of the number of implemented stacked physical registers. Note that `loadrs` may cause registers in the dirty partition to be lost.

## 6.6 RSE Interruptions

Although the RSE runs asynchronously to processor execution, RSE related interruptions are delivered synchronously with the instruction stream. These RSE interruptions are a direct consequence of register stack-related instructions such as: `alloc`, `br.ret`, `rfi`, `flushrs`, `loadrs`, or `mov` to/from `BSP`, `BSPSTORE`, `RSC`, `PFS`, `IFS`, or `RNAT`. Register spills and fills that are executed by the RSE in the background (eager RSE loads or stores) do not raise interruptions. If a faulting/trapping register spill or fill operation is required for software to make forward progress (mandatory RSE load or store) then the RSE will raise an interruption.

Mandatory RSE stores occur in the context of `alloc` and `flushrs` instructions only. Any faults raised by these instructions are delivered on the issuing instruction. Faults raised by mandatory RSE loads caused by a `loadrs` are delivered on the issuing instruction. Mandatory RSE loads which fault while restoring the frame for a `br.ret` or `rfi` deliver the fault on the target instruction, and the `ISR.ir` (incomplete register frame) bit is set. When a mandatory RSE load faults, `AR[BSPSTORE]` points to a backing store location above the faulting address reported in `CR[IFA]`. This allows handlers that service RSE load faults to use the backing store switch routine described in ["Switch from Interrupted Context" on page 2:148](#).

The `br.ret` and the `rfi` instructions set the RSE Current Frame Load Enable bit (RSE.CFLE) to one if the register stack frame being returned to is not entirely contained in the stacked register file. This enables the RSE to restore registers for the current

frame of the target instruction. When RSE.CFLE is set, instruction execution is stalled until the RSE has completely restored the current frame or an interruption occurs. This is the only time that the RSE issues any memory traffic for the current frame. Interruption delivery clears RSE.CFLE which allows an interruption handler to execute in the presence of an incomplete frame (e.g., to handle the fault raised by the mandatory RSE load). The RSE.CFLE bit is RSE internal state and is not architecturally visible.

Table 6-6 summarizes RSE raised interruptions.

**Table 6-6. RSE Interruption Summary**

Instruction	Interruption	Description
<code>alloc</code>	Illegal Operation fault	Malformed <code>alloc</code> immediate.
<code>alloc</code>	Reserved Register/Field fault	<code>alloc</code> instruction which attempted to change the size of the rotating region when one or more of the RRB values in CFM were non-zero.
<code>alloc</code> , <code>flushrs</code> , <code>loadrs</code>	Unimplemented Data Address fault Data Nested TLB fault Alternate Data TLB fault VHPT Data fault Data TLB fault Data Page Not Present fault Data NaT Page Consumption fault Data Key Miss fault Data Key Permission fault Data Access Rights fault Data Dirty Bit fault Data Access Bit fault Data Debug fault	AR[BSPSTORE] contains an unimplemented address.  AR[BSPSTORE] pointed to a NaTVal data page.
<code>br.call</code> , <code>brl.call</code>	No RSE related interruptions	
<code>br.ret</code>	No RSE load related faults	RSE load related faults are delivered on target instruction.
<code>rfi</code>	No RSE related interruptions	RSE load related faults are delivered on target instruction.
Target of <code>br.ret</code> or <code>rfi</code>	IR Unimplemented Data Address fault IR Data Nested TLB fault  IR Alternate Data TLB fault IR VHPT Data TLB fault IR Data TLB fault IR Data Page Not Present fault IR Data NaT Page Consumption fault IR Data Key Miss fault IR Data Key Permission fault IR Data Access Rights fault IR Data Access Bit fault IR Data Debug fault	Mandatory RSE load targeted an unimplemented address.  <code>br.ret</code> with <code>PSR.ic = 0</code> or <code>rfi</code> executed when <code>IPSR.ic = 0</code> .  RSE.BspLoad pointed at a NaTPage.

## 6.7 RSE Behavior on Interruptions

When the processor raises an interruption, the current register stack frame remains unchanged. If PSR.ic is one, the valid bit in the Interruption Function State register (IFS.v) is cleared. When the IFS.v bit is clear, the contents of the interruption frame marker field (IFS.ifm) are undefined.

While an interruption handler is running and the RSE is in store/load intensive or eager mode, the RSE continues spilling/filling registers to/from the backing store on behalf of the interrupted context as long as the registers are not part of the current frame as defined by CFM.

A sequence of mandatory RSE loads or stores (from `alloc`, `br.ret`, `flushrs`, `loadrs` and `rfi`) can be interrupted by an external interrupt.

When PSR.ic is 0, faults taken on mandatory RSE operations may not be recoverable.

## 6.8 RSE Behavior with an Incomplete Register Frame

The current register frame is considered **incomplete** when one of the mandatory RSE loads after a `br.ret` or a `rfi` faults, leaving `BSPSTORE` pointing to a location above `BSP` (i.e., `RSE.ndirty_words` is negative). The frame becomes complete when `RSE.ndirty_words` becomes non-negative, either by executing a cover instruction, or by handling the fault and completing the original sequence of mandatory RSE loads.

When the current frame is incomplete the following instructions have undefined behavior: `alloc`, `br.call`, `brl.call`, `br.ret`, `flushrs`, `loadrs`, and `move` to `BSPSTORE`. Software must guarantee that the current frame is complete before executing these instructions.

## 6.9 RSE and ALAT Interaction

The ALAT (see “Data Speculation” on page 1:63) uses physical register addresses to track advanced loads. `RSE.BOF` may only change as the result of a `br.call` (by `CFM.sol`), `cover` (by `CFM.sof`), `br.ret` (by `AR[PFM].sol`) or `rfi` (by `CR[IFS].ifm.sof` when `CR[IFS].v = 1`). This ensures, for ALAT invalidation purposes, that hardware does not update virtual to physical register address mapping, unless explicitly instructed to do so by software.

When software performs backing store switches that could cause program values to be placed in different physical registers, then the ALAT must be explicitly invalidated with the `invala` instruction. Typically this happens as part of a process or thread context switch, `longjmp` or `call stack unwind`, when software re-writes `AR[BSPSTORE]`, but cannot guarantee that `RSE.BOF` was preserved.

A stacked register is said to be **deallocated** when an `alloc`, `br.ret`, or `rfi` instruction changes the top of the current frame such that the register is no longer part of the current frame. Once a stacked register is deallocated, its value, its corresponding NaT bit, and its ALAT state are undefined. If that register is subsequently made part of the

current frame again (either via another `alloc` instruction, or via a `br.ret` or `rfi` to a previous frame that contained that register), the value stored in the register, the NaT bit for the register, and the corresponding ALAT entry for the register remain undefined.

RSE stores do not invalidate ALAT entries. Therefore, software cannot use the ALAT to trace RSE stores to the backing store.

**Note:** While an implementation is allowed to remove entries from the ALAT at any time, performance considerations strongly encourage not invalidating ALAT entries due to RSE stores.

## 6.10 Backing Store Coherence and Memory Ordering

RSE loads and stores are coherent with respect to the processor's data cache at all times. The backing store below `BSPSTORE` is defined to be consistent with the register stack (the memory image contains consecutive register values and NaT collections). Addresses below `BSPSTORE` are not modified by the RSE until `br.ret`, `rfi` or a move to `BSPSTORE` causes `BSP` to drop below the original `BSPSTORE` value. The RSE never writes to a memory address greater than or equal to `BSP`.

In order for software to modify a value in the backing store and guarantee that it be loaded by the RSE, software must first place the RSE into enforced lazy mode (`RSC.mode=0`), and read `BSP` and `BSPSTORE` to determine the location of the RSE store pointer. If the location to be modified lies between `BSPSTORE` and `BSP`, software must issue a `flushrs`, update the backing store location in memory, and issue a `loadrs` instruction with the `RSC.loadrs` set to zero (this invalidates the current contents of the physical stacked registers, except the current frame, which forces the RSE to reload registers from the backing store). If the location to be modified lies below `BSPSTORE`, unnecessary memory traffic can be avoided as follows: software must read the `RNAT` application register, update the backing store location in memory, rewrite `BSPSTORE` with the original value, and then rewrite `RNAT`.

RSE loads and stores are weakly ordered. The `flushrs` and `loadrs` instructions do not include an implicit memory fence. Turning on and off the RSE does not affect memory ordering. To ensure ordering of RSE loads and stores on a multiprocessor system, software is required to issue explicit memory fence (`mf`) instructions.

## 6.11 RSE Backing Store Switches

The implementation of system calls, operating system context switches, user-level thread packages, debugging software, and certain types of exception handling (e.g., `setjmp/longjmp`, structured exception handling and call stack unwinding) require explicit user-level control of the RSE and/or knowledge of the backing store format in memory. Therefore, the RSE and the backing store can be controlled at all privilege levels.

Three RSE backing store switches are described here:

1. Switching from an interrupted context (as part of exception handler or interrupt bubble-up code)
2. Returning to a previously interrupted context

3. Non-preemptive, synchronous backing store switch (covers system calls, user-level thread and operating system context switches)

Failure to follow these sequences may result in undefined RSE and processor behavior.

### 6.11.1 Switch from Interrupted Context

To switch from the backing store of an interrupted context to a new backing store:

1. Read and save the RSC and PFS application registers.
2. Issue a `cover` instruction for the interrupted frame.
3. Read and save the IFS control register.
4. Place RSE in enforced lazy mode by clearing both RSC.mode bits.
5. Read and save the BSPSTORE and RNAT application registers.
6. Write BSPSTORE with the new backing store address.
7. Read and save the new BSP to calculate the number of dirty registers.
8. Select the desired RSE setting (mode, privilege level and byte order).

### 6.11.2 Return to Interrupted Context

To return to the backing store of an interrupted context:

1. Allocate a zero-sized frame.
2. Subtract the BSPSTORE value written in step 6 of Section 6.11.1, "Switch from Interrupted Context" from the BSP value read in step 7 of [Section 6.11.1, "Switch from Interrupted Context" on page 2:148](#), and deposit the difference into RSC.loadrs along with a zero into RSC.mode (to place the RSE into enforced lazy mode).
3. Issue a `loadrs` instruction to insure that any registers from the interrupted context which were saved on the new stack have been loaded into the stacked registers.
4. Restore BSPSTORE from the interrupted context (saved in step 5 of Section 6.11.1, "Switch from Interrupted Context").
5. Restore RNAT from the interrupted context (saved in step 5 of Section 6.11.1, "Switch from Interrupted Context").
6. Restore PFS and IFS from the interrupted context (saved in steps 1 and 3 of Section 6.11.1, "Switch from Interrupted Context").
7. Restore RSC from the interrupted context (saved in step 1 of Section 6.11.1, "Switch from Interrupted Context"). This restores the setting of the RSE mode bits as well as privilege level and byte order.
8. Issue an `rfi` instruction (IFS.ifm will become CFM).

### 6.11.3 Synchronous Backing Store Switch

A non-preemptive, synchronous backing store switch at any privilege level can be accomplished as follows:



1. Read and save the RSC, BSP and PFS application registers.
2. Issue a `flushrs` instruction to flush the dirty registers to the backing store.
3. Place RSE in enforced lazy mode by clearing both RSC.mode bits.
4. Read and save the RNAT application register.
5. Invalidate the ALAT using the `invala` instruction when switching from code that does not restore RSE.BOF to its original setting. A different RSE.BOF will cause program values in the new context to be placed in different physical registers. See ["RSE and ALAT Interaction" on page 2:146](#) for details.
6. Write the new context's BSPSTORE (was BSP after `flushrs` when switching out).
7. Write the new context's PFS and RNAT.
8. Write the new context's RSC which will set the RSE mode, privilege level and byte order.

## 6.12 RSE Initialization

At processor reset the RSE is defined to be in enforced lazy mode, i.e., the RSC.mode bits are both zero. The RSE privilege level (RSC.pl) is defined to be zero. RSE.BOF points to physical register 32. The values of AR[PFS].pfm and CR[IFS].ifm are undefined. The current frame marker (CFM) is set as follows: sof=96, sol=0, sor=0, rrb.gr=0, rrb.fr=0, and rrb.pr=0. This gives the processor access to 96 stacked registers.

The RSE performs no spill/fill operations until either an `alloc`, `br.ret`, `rfi`, `flushrs` or `loadrs` require a mandatory RSE operation, or software explicitly enables eager RSE operations. Software must provide the RSE with a valid backing store address in the BSPSTORE application register prior to causing any RSE spill/fill operations. Failure to initialize BSPSTORE results in undefined behavior.

### §



Processors based on the Itanium architecture provide comprehensive debugging and performance monitoring facilities for both IA-32 and Itanium instructions. This chapter describes the debug registers, performance monitoring registers and their programming models. The debugging facilities include several data and instruction breakpoint registers, single step trap, breakpoint instruction fault, taken branch trap, lower privilege transfer trap, instruction and data debug faults. The performance monitoring facilities include two sets of registers to configure and collect various performance-related statistics.

## 7.1 Debugging

Several Data Breakpoint Registers (DBR) and Instruction Breakpoint Registers (IBR) are defined to hold address breakpoint values for data and instruction references. In addition the following debugging facilities are supported:

- **Single Step trap** – When PSR.ss is 1, successful execution of each Itanium instruction results in a Single Step trap. When PSR.ss is 1 or EFLAG.tf is 1, successful execution of each IA-32 instruction results in an IA\_32\_Exception(Debug) single step trap. After the trap, IIP and IPSR.ri point to the next instruction to be executed. IIPA and ISR.ei point to the trapped instruction. See “[Single Stepping](#)” for complete single stepping behavior.
- **Break Instruction fault** – execution of a `break` instruction results in a Break Instruction fault. IIM is loaded with the immediate operand from the instruction. IIM values are defined by software convention. `break` can be used for profiling, debugging and entry into the operating system (although Enter Privileged Code (`epc`) is recommended since it has lower overhead). Execution of the IA-32 INT 3 (`break`) instruction results in a IA\_32\_Exception(Break) trap.
- **Taken Branch trap** – When PSR.tb is 1, a Taken Branch trap occurs on every taken Itanium branch instruction. When PSR.tb is 1, a IA\_32\_Exception(Debug) taken branch trap occurs on every taken IA-32 branch instruction (CALL, Jcc, JMP, RET, LOOP). This trap is useful for debugging and profiling. After the trap, IIP and IPSR.ri point to the branch target instruction and IIPA and ISR.ei point to the trapping branch instruction.
- **Lower Privilege Transfer trap** – When PSR.lp bit is 1, and an Itanium branch demotes the privilege level (numerically higher), a Lower Privilege Transfer trap occurs. This trap allows for auditing of privilege demotions, for example to remove permissions which were granted to higher privilege code. After the trap, IIP and IPSR.ri point to the branch target and IIPA and ISR.ei point to the trapping branch instruction. IA-32 instructions can not raise this trap.
- **Instruction Debug faults** – When PSR.db is 1, any Itanium instruction memory reference that matches the parameters specified by the IBR registers results in an Instruction Debug fault. Instruction Debug faults are reported even if Itanium instructions are nullified due to a false predicate. If PSR.id is 1, Itanium Instruction Debug faults are disabled for one instruction. The successful execution of an Itanium instruction clears PSR.id. When PSR.db is 1, any IA-32 instruction memory

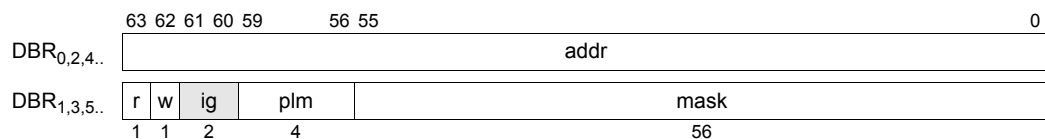
reference that matches the parameters specified by the IBR registers results in an IA\_32\_Exception(Debug) fault. If PSR.id is 1 or EFLAG.rf is 1, IA-32 Instruction Debug faults are disabled for one instruction. The successful execution of an IA-32 instruction clears the PSR.id and EFLAG.rf bits.

- **Data Debug faults** – When PSR.db is 1, any Itanium data memory reference that matches the parameters specified by the DBR registers results in a Data Debug fault. Data Debug faults are only reported if the qualifying predicate is true. Data Debug faults can be deferred on speculative loads by setting DCR.dd to 1. If PSR.dd is 1, Data Debug faults are disabled for one instruction or one mandatory RSE memory reference. When PSR.db is 1, any IA-32 data memory reference that matches the parameters specified by the DBR registers results in a IA\_32\_Exception(Debug) trap. IA-32 data debug events are traps, not faults as defined for the Itanium instruction set. The reported trap code returns the match status of the first 4 DBR registers that matched during the execution of the IA-32 instruction. See “IA-32 Trap Code” on page 2:213 for trap code details. Zero, one or more DBR registers may be reported as matching.

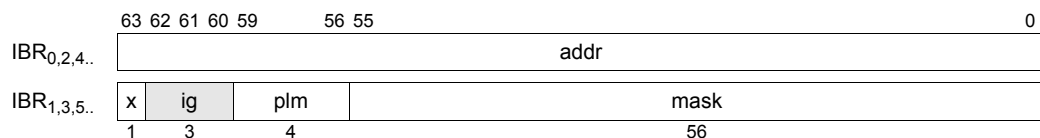
### 7.1.1 Data and Instruction Breakpoint Registers

Instruction or data memory addresses that match the Instruction or Data Breakpoint Registers (IBR/DBR) shown in Figure 7-1 and Figure 7-2 and Table 7-1 result in an Instruction or Data Debug fault. IA-32 Instruction or data memory addresses that match the Instruction or Data Breakpoint Registers (IBR/DBR) result in an IA\_32\_Exception(Debug) fault or trap. Even numbered registers contain breakpoint addresses, odd registers contain breakpoint mask conditions. At least 4 data and 4 instruction register pairs are implemented on all processor models. Implemented registers are contiguous starting with register 0.

**Figure 7-1. Data Breakpoint Registers (DBR)**



**Figure 7-2. Instruction Breakpoint Registers (IBR)**



When executing Itanium instructions, instruction and data memory addresses presented for matching are always in the implemented address space. Programming an unimplemented physical address into an IBR/DBR guarantees that physical addresses presented to the IBR/DBR will never match. Similarly, programming an unimplemented virtual address into an IBR/DBR guarantees that virtual addresses presented to the IBR/DBR will never match.

**Table 7-1. Debug Breakpoint Register Fields (DBR/IBR)**

Field	Bits	Description
addr	63:0	Match Address – 64-bit virtual or physical breakpoint address. Addresses are interpreted as either virtual or physical based on PSR.dt, PSR.it or PSR.rt. Data breakpoint addresses trap on load, store, semaphore, and mandatory RSE memory references. For Intel Itanium instruction set references, IBR.addr{3:0} is ignored in the address match. For IA-32 instruction references, IBR.addr{31:0} are used in the match and IBR.addr{63:32} must be zero to match. All 64 bits are implemented on all processors regardless of the number of implemented address bits.
mask	55:0	Address Mask – determines which address bits in the corresponding address register are compared in determining a breakpoint match. Address bits whose corresponding mask bits are 1, must match for the breakpoint to be signaled, otherwise the address bit is ignored. Address bits{63:56} for which there are no corresponding mask bits are always compared. For IA-32 instruction references, IBR.mask{55:32} are ignored. All 56 bits are implemented on all processors regardless of the number of implemented address bits.
plm	59:56	Privilege Level Mask – enables data breakpoint matching at the specified privilege level. Each bit corresponds to one of the four privilege levels, with bit 56 corresponding to privilege level 0, bit 57 with privilege level 1, etc. A value of 1 indicates that the debug match is enabled at that privilege level.
w	62	Write match enable – When DBR.w is 1, any non-nullified mandatory RSE store, IA-32 or Intel Itanium store, semaphore, probe.w.fault or probe.rw.fault to an address matching the corresponding address register causes a breakpoint.
r	63	Read match enable – When DBR.r is 1, any non-nullified IA-32 or Intel Itanium load, mandatory RSE load, semaphore, lfetch.fault, probe.r.fault or probe.rw.fault to an address matching the corresponding address register causes a breakpoint. When DBR.r is 1, a VHPT access that matches the DBR (except those for a <code>task</code> instruction) will cause an Instruction/Data TLB Miss fault. If DBR.r and DBR.w are both 0, that data breakpoint register is disabled.
x	63	Execute match enable – When IBR.x is 1, execution of an IA-32 instruction or Intel Itanium instruction in a bundle at an address matching the corresponding address register causes a breakpoint. If IBR.x is 0, that instruction breakpoint register is disabled. Instruction breakpoints are reported even if the qualifying predicate is false.
ig	62:60	Ignored

Four privileged instructions, defined in [Table 7-2](#), allow access to the debug registers. Register access is indirect, where the debug register number is determined by the contents of a general register. DBR/IBR registers can only be accessed at privilege level 0, otherwise a Privileged Operation fault is raised.

**Table 7-2. Debug Instructions**

Mnemonic	Description	Operation	Instr Type	Serialization Required
mov dbr[r <sub>3</sub> ] = r <sub>2</sub>	Move to data breakpoint register	DBR[GR[r <sub>3</sub> ]] ← GR[r <sub>2</sub> ]	M	data
mov r <sub>1</sub> = dbr[r <sub>3</sub> ]	Move from data breakpoint register	GR[r <sub>1</sub> ] ← DBR[GR[r <sub>3</sub> ]]	M	none
mov ibr[r <sub>3</sub> ] = r <sub>2</sub>	Move to instruction breakpoint register	IBR[GR[r <sub>3</sub> ]] ← GR[r <sub>2</sub> ]	M	inst
mov r <sub>1</sub> = ibr[r <sub>3</sub> ]	Move from instruction breakpoint register	GR[r <sub>1</sub> ] ← IBR[GR[r <sub>3</sub> ]]	M	none
break imm	Breakpoint Instruction fault	if (PSR.ic) IIM ← imm fault(Breakpoint_Instruction)	B/I/M	none

Changes to debug registers and PSR are not necessarily observed by following instructions. Software should issue a data serialization operation to ensure modifications to DBR, PSR.db, PSR.tb and PSR.lp are observed before a dependent instruction is executed. For register changes to IBR and PSR.db that affect fetching of subsequent instructions, software must issue an instruction serialization operation.

On some implementations, a hardware debugger may use two or more of these registers pairs for its own use. When a hardware debugger is attached, as few as 2 DBR pairs and as few as 2 IBR pairs may be available for software use. Software should be prepared to run with fewer than the implemented number of IBRs and/or DBRs if the software is expected to be debuggable with a hardware debugger. When a hardware debugger is not attached, at least 4 IBR pairs and 4 DBR pairs are available for software use.

Any debug registers used by an attached hardware debugger are allocated from the highest register numbers first (e.g. if only 2 DBR pairs are available to software, the available registers are DBR[0-3]).

**Note:** When a hardware debugger is attached and is using two or more debug registers pairs, the processor does not forcibly partition the registers between software and hardware debugger use; that is, the processor does not prevent software from reading or modifying any of the debug registers being used by the hardware debugger. However, if software modifies any of the registers being used by the hardware debugger, processor and/or hardware debugger operation may become undefined, or the processor and/or hardware debugger may crash.

## 7.1.2 Debug Address Breakpoint Match Conditions

For virtual memory accesses, breakpoint address registers contain the virtual addresses of the debug breakpoint. For physical accesses, the addresses in these registers are treated as a physical address. Software should be aware that debug registers configured to fault on virtual references, may also fault on a physical reference if translations are disabled. Likewise a debug register configured for physical references can fault on virtual references that match the debug breakpoint registers.

The range of addresses detected by the DBR and IBR registers for memory references by Itanium instructions is defined as:

- Instruction and single or multi-byte aligned data memory references that access any memory byte specified by the IBR/DBR address and mask fields results in an Instruction/Data Debug fault regardless of datum size. Implementations must only report a Debug fault if the specified aligned byte(s) are referenced.
- Floating-point load double/integer pair, floating-point spill/fill and 10-byte operands are treated as 16-byte datums for breakpoint matching, if the accesses are aligned. Floating-point load single pair operands are treated as 8-byte datums for breakpoint matching, if the accesses are aligned.
- If data memory references are unaligned, multi-byte memory references that access any memory byte specified by DBR address and mask fields result in a breakpoint Data Debug fault regardless of datum size. Processor implementations may also report additional breakpoint Data Debug faults for addresses not specifically specified by the DBR registers. Debugging software should perform a byte by byte breakpoint analysis of each address accessed by multi-byte unaligned datums to detect true breakpoint conditions.

- The `cmp8xchg16` operands are treated as 16-byte datums for both read and write breakpoint matching, even though this instruction only reads 8 bytes.

Address breakpoint Data Debug faults are not reported for the Flush Cache (`fc`, `fc.i`), regular\_form `probe`, non-faulting `lfetch`, insert TLB (`itc`, `itr`), purge TLB (`ptc`, `ptr`), or translation access (`thash`, `ttag`, `tak`, `tpa`) instructions. Accesses by the RSE to a debug region are checked, but the Data Debug fault is not reported until a subsequent `alloc`, `br.ret`, `rfi`, `loadrs`, or `flushrs` which requires that the faulting load or store actually occur.

The range of addresses detected by the DBR and IBR registers for IA-32 memory references is defined as:

- Instruction memory references where the first byte of the IA-32 instruction match the IBR address and mask fields results in an `IA_32_Exception(Debug)` fault. Subsequent bytes of a multiple byte IA-32 instruction are not compared against the IBR registers for breakpoints. The upper 32-bits of the IBR `addr` field must be zero to detect IA-32 instruction memory references.
- IA-32 single or multi-byte data memory references that access any memory byte specified by the DBR address and mask fields results in an `IA_32_Exception(Debug)` trap regardless of datum size and alignment. The processor ensures that all data breakpoint traps are precisely reported. Data breakpoint traps are reported if and only if any byte in the IA-32 data memory reference matches the DBR address and mask fields. No spurious data breakpoint events are generated for IA-32 data memory operands that are unaligned, nor are breakpoints reported if no bytes of the operand lie within the address range specified by the DBR address and mask fields.

## 7.2 Performance Monitoring

Performance monitors allow processor events to be monitored by programmable counters or give an external notification (such as a pin or transaction) on the occurrence of an event. Monitors are useful for tuning application, operating system and system performance. Two sets of performance monitor registers are defined. Performance Monitor Configuration (PMC) registers are used to control the monitors. Performance Monitor Data (PMD) Registers either provide data values from the monitors, or hold data values used by the PMU. The performance monitors can record performance values from either the IA-32 or Itanium instruction set.

As shown in [Figure 7-3](#), all processor implementations provide at least four performance counters (PMC/PMD[4]..PMC/PMD[7] pairs), and four performance counter overflow status registers (PMC[0]..PMC[3]). Performance monitors are also controlled by bits in the processor status register (PSR), the default control register (DCR) and the performance monitor vector register (PMV). Processor implementations may provide additional implementation-dependent PMC and PMD registers to increase the number of “generic” performance counters (PMC/PMD pairs). The remainder of the PMC and PMD register set is implementation dependent.

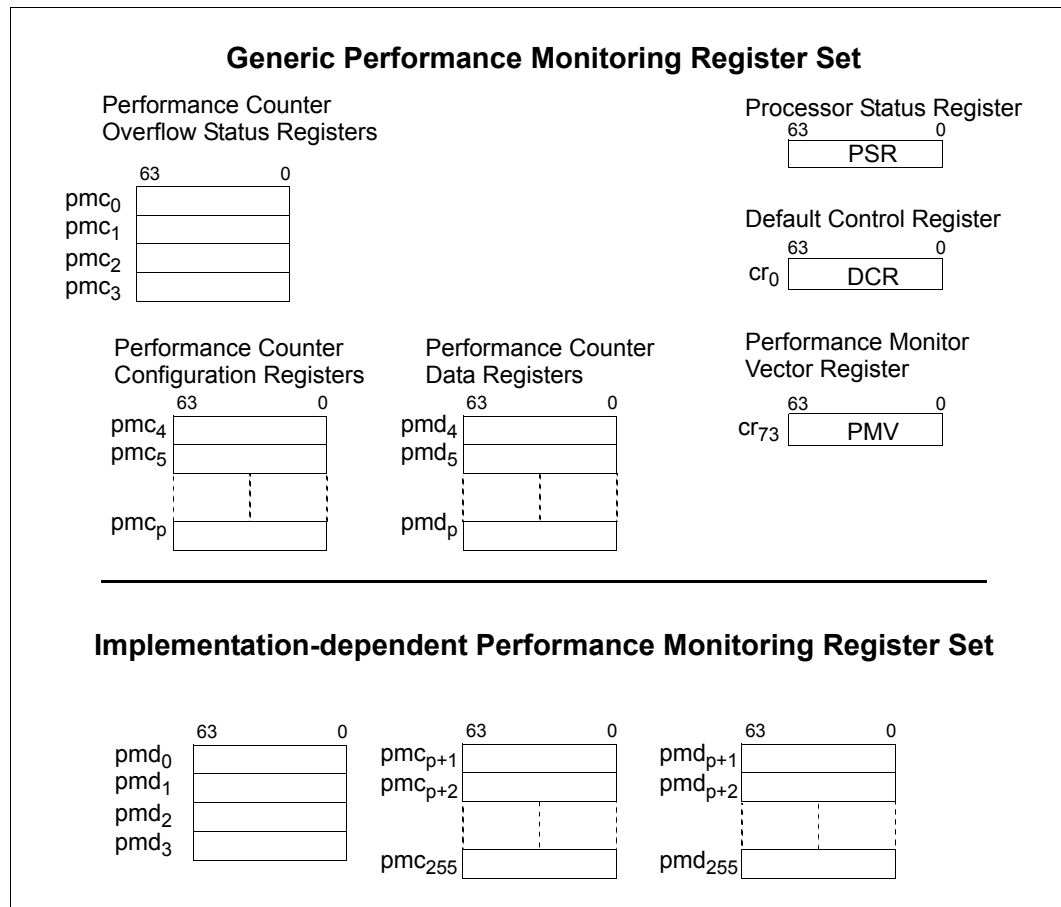
Event collection for implementation-dependent performance monitors is not specified by the architecture. Enabling and disabling functions are implementation dependent. For details, consult processor-specific documentation.

Processor implementations may not populate the entire PMC/PMD register space. Reading of an unimplemented PMC or PMD register returns zero. Writes to unimplemented PMC or PMD registers are ignored; i.e., the written value is discarded.

Writes to PMD and PMC and reads from PMC are privileged operations. At non-zero privilege levels, these operations result in a Privileged Operation fault, regardless of the register address.

Reading of PMD registers by non-zero privilege level code is controlled by PSR.sp. When PSR.sp is one, PMD register reads by non-zero privilege level code return zero.

**Figure 7-3. Performance Monitor Register Set**



### 7.2.1 Generic Performance Counter Registers

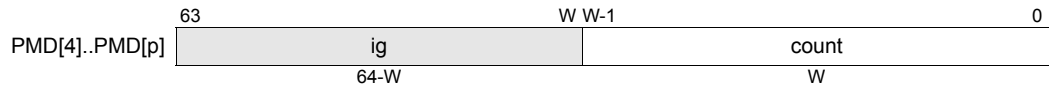
Generic performance counter registers are PMC/PMD pairs that contiguously populate the PMC/PMD name space starting at index 4. At least 4 performance counter register pairs (PMC/PMD[4]..PMC/PMD[7]) are implemented in all processor models. Each counter can be configured to monitor events for any combination of privilege levels and one of several event metrics. The number of performance counters is implementation specific. The figures and tables use the symbol "p" to represent the index of the last implemented generic PMC/PMD pair. The bit-width W of the counters is also implementation specific.



A counter overflow interrupt occurs when the counter wraps; i.e., a carry out from bit W-1 is detected. Counter overflow interrupts are edge-triggered; that is, the event of a counter incrementing and causing carry out from bit W-1 thus setting the overflow bit and the freeze bit, generates one PMU interrupt. Provided that software does not clear the freeze bit, while either or both of PSR.up and pp are 1, without also clearing the overflow bit (before or concurrent with the write to the freeze bit), no further interrupts are generated based on the fact that the carry out had been earlier detected.

Figure 7-4 and Figure 7-5 show the fields in PMD and PMC respectively, while Table 7-3 and Table 7-4 describe the fields in PMD and PMC respectively.

**Figure 7-4. Generic Performance Counter Data Registers (PMD[4]..PMD[p])**

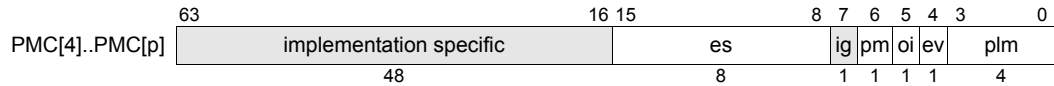


**Table 7-3. Generic Performance Counter Data Register Fields**

Field	Bits	Description
ig	63:W	Writes are ignored. Reads return 0.
count	W-1:0	Event Count. The counter is defined to overflow when the count field wraps (carry out from bit W-1).

Some implementations do not treat the upper, unimplemented bits of PMDs as ignored bits on reads, but rather return a copy of bit W-1 in these bit positions so that count values appear as if they were sign extended. Subsequent implementations will return 0 for these bits on reads.

**Figure 7-5. Generic Performance Counter Configuration Register (PMC[4]..PMC[p])**



**Table 7-4. Generic Performance Counter Configuration Register Fields (PMC[4]..PMC[p])**

Field	Bits	Description
plm	3:0	Privilege Level Mask – controls performance monitor operation for a specific privilege level. Each bit corresponds to one of the 4 privilege levels, with bit 0 corresponding to privilege level 0, bit 1 with privilege level 1, etc. A bit value of 1 indicates that the monitor is enabled at that privilege level. Writing zeros to all plm bits effectively disables the monitor. In this state, the corresponding PMD register(s) do not preserve values, and the processor may choose to power down the monitor.
ev	4	External visibility – When 1, an external notification (such as a pin or transaction) may be provided, dependent upon implementation, whenever the monitor overflows. Overflow occurs when a carry out from bit W-1 is detected.
oi	5	Overflow interrupt – If 1, when the monitor overflows, a Performance Monitor Interrupt is raised and the performance monitor freeze bit (PMC[0].fr) is set. If 0, no interrupt is raised and the performance monitor freeze bit (PMC[0].fr) remains unchanged. Overflow occurs when a carry out from bit W-1 is detected. See <a href="#">“Performance Monitor Overflow Status Registers (PMC[0]..PMC[3])”</a> for details on configuring interrupt vectors.

**Table 7-4. Generic Performance Counter Configuration Register Fields (PMC[4]..PMC[p]) (Continued)**

Field	Bits	Description
pm	6	Privileged monitor – When 0, the performance monitor is configured as a user monitor, and enabled by PSR.up. When PMC.pm is 1, the performance monitor is configured as a privileged monitor, enabled by PSR.pp, and the corresponding PMD can only be read by privileged software.
ig	7	ignored
es	15:8	Event select – selects the performance event to be monitored. Actual event encodings are implementation dependent. Some processor models may not implement all event select (es) bits. At least one bit of es must be implemented on all processors. Unimplemented es bits are ignored.
implem. specific	63:16	Implementation-specific bits – Reads from implemented bits return implementation-dependent values. For portability, software should write what was read; i.e., software may not use these bits as storage. Hardware will ignore writes to unimplemented bits.

Event collection is controlled by the Performance Monitor Configuration (PMC) registers and the processor status register (PSR). Four PSR fields (PSR.up, PSR.pp, PSR.cpl and PSR.sp) and the performance monitor freeze bit (PMC[0].fr) affect the behavior of all generic performance monitor registers. Finer, per monitor, control of generic performance monitors is provided by two PMC register fields (PMC[i].plm, PMC[i].pm). Event collection for a generic monitor is enabled under the following constraints:

- Generic Monitor Enable[i] = (not PMC[0].fr) and PMC[i].plm[PSR.cpl] and ((not (PMC[i].pm) and PSR.up) or (PMC[i].pm and PSR.pp))

Generic performance monitor data registers (PMD[i]) can be configured to be user readable (useful for user level sampling and tracing user level processes) by setting the PMC[i].pm bit to 0. All user-configured monitors can be started and stopped synchronously by the user mask instructions (`rum` and `sum`) by altering PSR.up. User-configured monitors can be secured by setting PSR.sp to 1. A user-configured secured monitor continues to collect performance values; however, reads of PMD, by non-privileged code, return zeros until the monitor is unsecured.

Monitors configured as privileged (PMC[i].pm is 1) are accessible only at privilege level 0; otherwise, reads return zeros. All privileged monitors can be started and stopped synchronously by the system mask instructions (`rsm` and `ssm`) by altering PSR.pp. [Table 7-5](#) summarizes the effects of PSR.sp, PMC[i].pm, and PSR.cpl on reading PMD registers.

Updates to generic PMC registers and PSR bits (up, pp, is, sp, cpl) require implicit or explicit data serialization prior to accessing an affected PMD register. The data serialization ensures that all prior PMD reads and writes as well as all prior PMC writes have completed.

**Table 7-5. Reading Performance Monitor Data Registers**

PSR.sp	PMC[i].pm	PSR.cpl	PMD Reads Return
0	0	0	PMD value
0	1	0	PMD value
1	0	0	PMD value
1	1	0	PMD value
0	0	>0	PMD value

**Table 7-5. Reading Performance Monitor Data Registers (Continued)**

PSR.sp	PMC[i].pm	PSR.cpl	PMD Reads Return
0	1	>0	0
1	0	>0	0
1	1	>0	0

Generic PMD counter registers may be read by software without stopping the counters. Under normal counting conditions (PMC[0].fr is zero and has been serialized), the processor guarantees that a sequence of reads of a given PMD will return non-decreasing values corresponding to the program order of the reads. Under frozen count conditions (PMC[0].fr is one and has been serialized), the counters are unchanging and ordering is irrelevant. When the freeze bit is in-flight, whether counters count events and reads return non-decreasing values is implementation dependent. Instruction serialization is required to ensure that the behavior specified by PMC[0].fr is observed.

Software must accept a level of sampling error when reading the counters due to various machine stall conditions, interruptions, and bus contention effects, etc. The level of sampling error is implementation specific. More accurate measurements can be obtained by disabling the counters and performing an instruction serialize operation for instruction events or data serialize operation for data events before reading the monitors. Other (non-counter) implementation-dependent PMD registers can only be read reliably when event monitoring is frozen (PMC[0].fr is one).

For accurate PMD reads of disabled counters, data serialization (implicit or explicit) is required between any PMD read and a subsequent `ssm` or `sum` (that could toggle PSR.up or PSR.pp from 0 to 1), or a subsequent `epc`, demoting `br.ret` or branch to IA-32 (`br.ia`) (that could affect PSR.cpl or PSR.is). Note that implicit post-serialization semantics of `sum` do not meet this requirement.

Table 7-6 defines the instructions used to access the PMC and PMD registers.

**Table 7-6. Performance Monitor Instructions**

Mnemonic	Description	Operation	Instr Type	Serialization Required
<code>mov pmd[r<sub>3</sub>] = r<sub>2</sub></code>	Move to performance monitor data register	$PMD[GR[r_3]] \leftarrow GR[r_2]$	M	data/inst
<code>mov r<sub>1</sub> = pmd[r<sub>3</sub>]</code>	Move from performance monitor data register	$GR[r_1] \leftarrow PMD[GR[r_3]]$	M	none <sup>a</sup>
<code>mov pmc[r<sub>3</sub>] = r<sub>2</sub></code>	Move to performance monitor configure register	$PMC[GR[r_3]] \leftarrow GR[r_2]$	M	data/inst
<code>mov r<sub>1</sub> = pmc[r<sub>3</sub>]</code>	Move from performance monitor configure register	$GR[r_1] \leftarrow PMC[GR[r_3]]$	M	none

a. When the freeze bit is in-flight, whether counters count events and reads return non-decreasing values is implementation dependent. Instruction serialization is required to ensure that the behavior specified by PMC[0].fr is observed.

## 7.2.2 Performance Monitor Overflow Status Registers (PMC[0]..PMC[3])

Performance monitor interrupts may be caused by an overflow from a generic performance monitor or an implementation-dependent event from a model-specific monitor. The four performance monitor overflow registers (PMC[0]...PMC[3]) shown in Figure 7-6 indicate which monitor caused the interruption.

Each of the 252 overflow bits in the performance monitoring overflow status registers (PMC[0]...PMC[3]) corresponds to a generic performance counter pair or to an implementation-dependent monitor. For generic performance counter pairs, overflow status bit  $PMC[i/64]\{i\%64\}$  corresponds to generic counter pair PMC[i]/PMD[i], where  $4 < i \leq p$ , and p is the index of the last implemented generic PMC/PMD pair.

There are currently two criteria for generating a performance monitor interrupt:

1. A generic performance counter pair (PMC[n]/PMD[n]) overflows and its overflow interrupt bit (PMC[n].oi) is 1.
2. An implementation-dependent monitor wants to report an event with an interruption.

If any of these criteria are met, the processor will:

- Set the corresponding overflow status bit in PMC[0]..PMC[3] to 1, and
- Raise a Performance Monitor interrupt, and
- Set the freeze bit (PMC[0].fr) which suspends event monitoring.

PMU interrupts are generated by events, such as the overflowing of a generic counter pair which is configured to interrupt on overflow. Each such event generates one interrupt. Provided that software does not clear the freeze bit, while either or both of PSR.up and .pp are 1, before clearing the overflow bits, writes to PMCs and PMDs by software do not generate interrupts, nor cause a monitor which had generated an interrupt to generate a second interrupt. (For overflow bits in PMC 0, even if either or both of PSR.up and .pp are 1, software may clear the overflow bits and the freeze bit with a single write to PMC 0 without causing any additional interrupts to be generated.)

Software may restore PMU state which has the freeze bit equal to 1 and one or more overflow bits equal to 1 without generating any interrupts provided that it ensures either that:

- both PSR.up and .pp are zero during the restore, or
- the freeze bit is a 1 (and serialized) before any overflow bits are set to 1

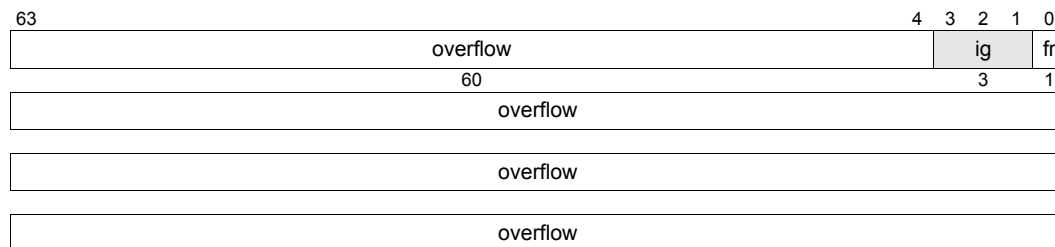
When the PMU is disabled by writing a 0 into PSR.up and .pp and serializing this write, the PMU cannot generate any interrupts and no SW writes to any PMU state can cause any interrupts.

When a generic performance counter pair (PMC[n]/PMD[n]) overflows and its overflow interrupt bit (PMC[n].oi) is 0, the corresponding overflow status register bit is set to 1. However, in this case of counter overflow without interrupt, the freeze bit in the PMC[0] is left unchanged, and event monitoring continues.

If control register bit PMV.m is one, a performance monitoring interrupt is disabled from being pended. When PMV.m is zero, the interruption is received and held pending. (Further masking by the PSR.i, TPR and in-service masking can keep the interrupt from being raised.) Figure 7-6 shows the Performance Monitor Overflow Status registers.

Implementation dependent PMD registers (0-3) cannot report events in the overflow registers; those 4 bit positions are used for other purposes.

**Figure 7-6. Performance Monitor Overflow Status Registers (PMC[0]..PMC[3])**



Under frozen count conditions when PMC[0].fr is one (either by a performance counter overflow, or an explicit software write and serialization), the processor suspends all event monitoring, i.e. counters do not increment and overflow bits as well as model-specific monitoring are frozen. Normal counting conditions are restored by software writing a zero to the freeze bit and serializing to resume event monitoring. When the freeze bit is in-flight, whether counters count events and reads return non-decreasing values is implementation dependent. Instruction serialization is required to ensure that the behavior specified by PMC[0].fr is observed.

**Table 7-7. Performance Monitor Overflow Register Fields (PMC[0]...PMC[3])**

Register	Field	Bits	Description
PMC[0]	fr	0	Performance Monitor “freeze” bit. This bit is volatile state, i.e., it is set by the processor whenever: <ul style="list-style-type: none"> <li>a generic performance monitor overflow occurs and its overflow interrupt bit (PMC[n].oi) is set to one.</li> <li>a model-specific performance monitor signals an interrupt.</li> </ul> The freeze bit can also be set by software to enable or disable all event monitoring. If the freeze bit is one, event monitoring is disabled. If the freeze bit is zero, event monitoring is enabled. If the freeze bit is in-flight, event monitoring behavior is implementation dependent.
PMC[0]	ig	3:1	Ignored
PMC[0]..PMC[3]	overflow	implemented monitors	Bit vector indicating which performance monitor overflowed. Overflow status bits are sticky, they are set to 1 by the processor if the corresponding PMD overflows; otherwise they are left unchanged. Multiple overflow status bits may be set, independent of whether counter overflow causes an interrupt or not.
		unimplemented monitors	Ignored

Multiple overflow bits may be set to 1, if counters overflow concurrently. The overflow bits and the freeze bit are sticky; i.e., the processor sets them to 1 but never resets them to 0. It is software's responsibility to reset the overflow and freeze bits.

The overflow status bits are populated only for implemented counters. Overflow bits of unimplemented counters read as zero and writes are ignored.

### 7.2.3 Performance Monitor Events

The set of monitored events is implementation-specific. All processor models are required to provide at least two events:

1. The number of retired instructions. These are defined as all instructions which execute without a fault, including nops and those which were predicated off. Generic counters configured for this event count only when the processor is in the NORMAL or LOW-POWER state (see [Figure 11-8 on page 2:314](#)).
2. The number of processor clock cycles. Generic counters configured for this event count only when the processor is in the NORMAL or LOW-POWER state (see [Figure 11-8 on page 2:314](#)).

Events may be monitorable only by a subset of the available counters. PAL calls provide an implementation-independent interface that provides information on the number of implemented counters, their bit-width, the number and location of other (non-counter) monitors, etc.

### 7.2.4 Implementation-independent Performance Monitor Code Sequences

This section describes implementation-independent code sequences for servicing overflow interrupts and context switches of the performance monitors. For forward compatibility, the code sequences outlined in [Section 7.2.4.1](#) and [Section 7.2.4.2](#) use PAL-provided implementation-specific information to collect/preserve data values for all implemented counters.

#### 7.2.4.1 Performance Monitor Interrupt Service Routine

When a generic performance counter pair (PMC[n]/PMD[n]) overflows and its overflow interrupt bit (PMC[n].oi) is 1, or an implementation-dependent monitor wants to report an event with an interruption, then the processor:

- Sets the corresponding overflow status bit in PMC[0]..PMC[3] to one,
- Raises a Performance Monitor Interrupt, and
- Sets the freeze bit in PMC[0] which suspends event monitoring.

Event monitoring remains frozen until software clears the freeze bit. When the freeze bit is in-flight, whether counters count events and reads return non-decreasing values is implementation dependent. Instruction serialization is required to ensure that the behavior specified by PMC[0].fr is observed. Performance monitor interrupts may be caused by an overflow of any of the counters. The processor indicates which performance monitor overflowed in the performance monitor overflow status registers (PMC[0]..PMC[3]). If multiple counters overflow concurrently, multiple overflow bits will be set to one. For forward compatibility, event collection interrupt handlers must

follow the implementation-independent overflow interrupt service routine outlined in [Figure 7-7](#). Use of alternate context-switch sequences may be incompatible with future implementations.

If the outgoing context has an interrupt pending but has not yet invoked the performance monitor interrupt service routine, the interrupt may be delivered to the incoming context even if it is a non-monitored process. The interrupt service routine can recognize this kind of bogus interrupt by noticing that either: the freeze bit is zero or the context is not being monitored.

**Figure 7-7. Performance Monitor Interrupt Service Routine (Implementation Independent)**

```
//Assumes PSR.up and PSR.pp are switched to zero together
if ((PMC[0].fr==1) && (PSR.up == 1) || (PSR.pp == 1)){
    // freeze bit is set. Search for interrupt.
    for (i=0; i< 4; i++) {
        if (PMC[i] != 0) {
            startbit = (i==0) ? 4 : 0;
            for (j=startbit; j < 64 ; j++) {
                if (PMC[i]{j}) {
                    counter_id = 64*i + j;
                    if (counter_id > PAL_GENERIC_PMC_PMD_PAIRS) {
                        Implementation_Specific_Update(counter_id);
                    }
                    else { // Generic PMC/PMD counter
                        if (PMC[counter_id].oi)
                            ovflcount[counter_id] += 1;
                    }
                }
            } // scan overflow bits
        }
    }
}
// Either ignore bogus interrupt or clear PMC[3]..PMC[1]
for (i=3; i>=1; i--) { PMC[i] = 0; }
rfi
```

#### 7.2.4.2 Performance Monitor Context Switch

The context switch routine described in [Figure 7-8](#) defines the implementation-independent context switching of Itanium performance monitors. Using bit masks provided by PAL (`PALPMCmask`, `PALPMDmask`) the routine can generically save/restore the contents of all implementation-specific performance monitoring registers. If the outgoing context is monitored, then all PMC and PMD registers whose mask bit is set are preserved by software. But if the outgoing context is monitored and the context switch routine determines that the outgoing context has a pending performance monitor interrupt (by reading the freeze bit with the knowledge that it was not generated by software) then software also preserves the outgoing context's overflow status registers (`PMC[0]..PMC[3]`) before all PMC and PMD registers whose mask bit is set. Here, it is explicitly assumed that software tracks monitored processes and can determine whether a process is monitored prior to reading the freeze bit. The context switch handler then restores the performance monitor freeze bit which resets event collection for the new context. Sometime into the incoming (possibly unmonitored) context, the performance overflow interrupt service routine will run, but by looking at the status of the freeze bit software can determine whether this interrupt can be ignored (for details refer to [Section 7.2.4.1](#)).

When switching back to the original context (that originally caused the counter overflow), the previously saved freeze bit can be inspected. If it was set (meaning there was a pending performance monitor interrupt), then the context switch routine posts an interrupt message to the incoming context's processor at the performance monitor vector specified by the PMV register (see [Section 10.5.8, "Inter-processor Interrupts Layout and Example" on page 2:612](#)). This will result in a new performance monitor overflow interrupt in the correct context. Essentially, the interrupt message is "replaying" the overflow interrupt that was missed because of the context switch.

**Figure 7-8. Performance Monitor Overflow Context Switch Routine**

```

// in context or thread switch

if (outgoing process is monitored) {
  1. Turn-off counting and ignore interrupts for context switch
    of counters.
    1a)  if not already done, raise interrupt priority above
          perf. mon overflow vector
    1b)  read and preserve PSR.up, PSR.pp, PSR.sp
    1c)  clear PSR.up, clear PSR.pp
    1d)  srlz.d
  2. Preserve PMC/PMD contents
    2a)  For each PMC whose PALPMCmask bit is set, preserve PMC.
    2b)  For each PMD whose PALPMDmask bit is set, preserve PMD.
}

.... continue context switch .....

// Now in incoming process/thread
if (incoming process is monitored) {
  // Event counting is disabled because PSR.up and pp are both
  // zero (step 1c above).

  3. Restore PMC/PMD contents (inverse of step 4 above)
    3a)  For each PMC whose PALPMCmask bit is set, reload PMC.
    3b)  For each PMD whose PALPMDmask bit is set, reload PMD.

  4. Restore Interrupt State (inverse of step 2 and 1a above)
    4a)  if (PMC[0].fr) {
          send myself a performance monitor interrupt
          (store to interrupt address)
        }
    4b)  Restore PSR.up and PSR.pp
    4c)  srlz.d
    4d)  lower interrupt priority below perf. mon overflow
          vector
}

```

§



Chapter 5 describes the interruption mechanism and programming model for the Itanium architecture. This chapter describes the IVA-based interruption handlers. “[Interruption Vector Descriptions](#)” describes all the Itanium IVA-based interruption vectors and “[IA-32 Interruption Vector Definitions](#)” describes all of the IA-32 interrupt vectors. PAL-based interruptions are described in [Chapter 11, “Processor Abstraction Layer.”](#) Note that unless otherwise noted, references to “interruption” in this chapter refer to IVA-based interruptions. See “[Interruption Definitions](#)” on page 2:95.

## 8.1 Interruption Vector Descriptions

The section lists all the Itanium interruption vectors. It describes the interruption vectors and the parameters that are defined when the vector is entered.

If an interruption is independent of the executing instruction set (including IA-32), such as an external interrupt or TLB fault, common Itanium interruption vectors are used. For exceptions and intercept conditions that are specific to the IA-32 instruction set three IA-32 specific vectors are used; IA\_32\_Exception, IA\_32\_Interrupt, and IA\_32\_Intercept.

[Table 8-1](#) defines which interruption resources are written, are left unmodified, or are undefined for each interruption vector. The individual vector descriptions below list interruption-specific resources for each vector.

See “[IVA-based Interruption Handling](#)” on page 2:101 for details on how the processor handles an interruption. See “[Interruption Control Registers](#)” on page 2:36 for the definition of bit fields within the interruption resources.

## 8.2 ISR Settings

For each of the interruption vectors, a figure depicts the ISR setting. These figures show the value that hardware writes into the ISR for the corresponding interruption.

[Table 8-2](#) provides an overview of ISR settings for all of the interruption vectors.

For some of the vectors, certain bits will always be 0 (or 1) simply because no instruction that would set that bit differently can ever end up on that vector. For example, ISR.sp is always 0 in the Break Instruction vector because ISR.sp is only set by speculative loads, and speculative loads can never take a Break Instruction fault.

After interruption from the IA-32 instruction set, the following ISR bits will always be zero: ISR.ni, ISR.na, ISR.sp, ISR.rs, ISR.ir, ISR.ei, and ISR.ed.

ISR.code settings for non-access instructions are described in “[Non-access Instructions and Interruptions](#)” on page 2:103.

[Table 8-3](#) on page 2:170 provides an overview of ISR.code field on all Itanium traps.

## 8.3 Interruption Vector Definition

**Table 8-1. Writing of Interruption Resources by Vector**

Interruption Resource	IIP, IPSR, IIPA, IFS.v		IFA		ITIR		IHA		IIM		ISR		IIB0, IIB1	
	0	1	0	1	0	1	0	1	0	1	0	1	0	1
<b>PSR.ic at time of interruption</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>
<b>Alternate Data TLB vector</b>														
Alternate Data TLB fault	N/A <sup>a</sup>	W <sup>b</sup>	N/A	W	N/A	W	N/A	x <sup>c</sup>	N/A	x	N/A	W	N/A	W
IR Alternate Data TLB fault	N/A	W	N/A	W	N/A	W	N/A	x	N/A	x	N/A	W	N/A	x
<b>Alternate Instruction TLB vector</b>														
Alternate Instruction TLB fault	- <sup>d</sup>	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Break Instruction vector</b>														
Break Instruction fault	-	W	x	x	x	x	x	x	-	W	W	W	-	W
<b>Data Access Rights vector</b>														
Data Access Rights fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
IR Data Access Rights fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Data Access-Bit vector</b>														
Data Access Bit fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
IR Data Key Miss fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Data Key Miss vector</b>														
Data Key Miss fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
IR Data Key Miss fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Data Nested TLB vector</b>														
Data Nested TLB fault	-	N/A	-	N/A	-	N/A	-	N/A	x	N/A	-	N/A	-	N/A
IR Data Nested TLB fault	-	N/A	-	N/A	-	N/A	-	N/A	x	N/A	-	N/A	-	N/A
<b>Data TLB vector</b>														
Data TLB fault	N/A	W	N/A	W	N/A	W	N/A	W	N/A	x	N/A	W	N/A	W
IR Data TLB fault	N/A	W	N/A	W	N/A	W	N/A	W	N/A	x	N/A	W	N/A	x
<b>Debug vector</b>														
Data Debug fault	-	W	-	W	x	x	x	x	x	x	W	W	-	W
Instruction Debug fault	-	W	-	W	x	x	x	x	x	x	W	W	-	x
IR Data Debug fault	-	W	-	W	x	x	x	x	x	x	W	W	-	x
<b>Dirty-Bit vector</b>														
Data Dirty Bit fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
<b>Disabled FP-Register vector</b>														
Disabled Floating-Point Register fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>External Interrupt vector</b>														
External Interrupt	-	W	x	x	x	x	x	x	x	x	W	W	-	x
<b>Floating-point Fault vector</b>														
Floating-Point Exception fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>Floating-point Trap vector</b>														
Floating-Point Exception trap	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>General Exception vector</b>														
Disabled ISA Transition fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
Illegal Dependency fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
Illegal Operation fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
IR Unimplemented Data Address fault	-	W	x	x	x	x	x	x	x	x	W	W	-	x
Privileged Operation fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
Privileged Register fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W

**Table 8-1. Writing of Interruption Resources by Vector (Continued)**

Interruption Resource	IIP, IPSR, IIPA, IFS.v		IFA		ITIR		IHA		IIM		ISR		IIB0, IIB1	
	0	1	0	1	0	1	0	1	0	1	0	1	0	1
<b>PSR.ic at time of interruption</b>	0	1	0	1	0	1	0	1	0	1	0	1	0	1
Reserved Register/Field fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
Unimplemented Data Address fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>IA-32 Exception vector</b>	N/A	W	N/A	x	N/A	x	N/A	x	N/A	x	N/A	W	N/A	x
<b>IA-32 Intercept vector</b>	N/A	W	N/A	x	N/A	x	N/A	x	N/A	W	N/A	W	N/A	x
<b>IA-32 Interrupt vector</b>	N/A	W	N/A	x	N/A	x	N/A	x	N/A	x	N/A	W	N/A	x
<b>Instruction Access Rights vector</b>														
Instruction Access Rights fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Instruction Access-Bit vector</b>														
Instruction Access Bit fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Instruction Key Miss vector</b>														
Instruction Key Miss fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Instruction TLB vector</b>														
Instruction TLB fault	-	W	-	W	-	W	-	W	x	x	W	W	-	x
<b>Key Permission vector</b>														
Data Key Permission fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
Instruction Key Permission fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
IR Data Key Permission fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Lower-Privilege Transfer Trap vector</b>														
Unimplemented Instruction Address fault	-	W	x	W	x	x	x	x	x	x	W	W	-	x
Lower-Privilege Transfer trap	-	W	x	x	x	x	x	x	x	x	W	W	-	W
Unimplemented Instruction Address trap	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>NaT Consumption vector</b>														
Data NaT Page Consumption fault	-	W	-	W	x	x	x	x	x	x	W	W	-	W
Instruction NaT Page Consumption fault	-	W	-	W	x	x	x	x	x	x	W	W	-	x
IR Data NaT Page Consumption fault	-	W	-	W	x	x	x	x	x	x	W	W	-	x
Register NaT Consumption fault	-	W	-	x	x	x	x	x	x	x	W	W	-	W
<b>Page Not Present vector</b>														
Data Page Not Present fault	-	W	-	W	-	W	x	x	x	x	W	W	-	W
Instruction Page Not Present fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
IR Data Page Not Present fault	-	W	-	W	-	W	x	x	x	x	W	W	-	x
<b>Single Step Trap vector</b>														
Single Step trap	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>Speculation vector</b>														
Speculative Operation fault	-	W	x	x	x	x	x	x	-	W	W	W	-	W
<b>Taken Branch Trap vector</b>														
Taken Branch trap	-	W	x	x	x	x	x	x	x	x	W	W	-	W
<b>Unaligned Reference vector</b>														

**Table 8-1. Writing of Interruption Resources by Vector (Continued)**

Interruption Resource	IIP, IPSR, IIPA, IFS.v		IFA		ITIR		IHA		IIM		ISR		IIB0, IIB1	
	0	1	0	1	0	1	0	1	0	1	0	1	0	1
PSR.ic at time of interruption	-	W	-	W	x	x	x	x	x	x	W	W	-	W
<b>Unsupported Data Reference vector</b>														
Unsupported Data Reference fault	-	W	-	W	x	x	x	x	x	x	W	W	-	W
<b>VHPT Translation vector</b>														
IR VHPT Data fault	N/A	W	N/A	W	N/A	W	N/A	W	N/A	x	N/A	W	N/A	x
VHPT Data fault	N/A	W	N/A	W	N/A	W	N/A	W	N/A	x	N/A	W	N/A	W
VHPT Instruction fault	N/A	W	N/A	W	N/A	W	N/A	W	N/A	x	N/A	W	N/A	x
<b>Virtual External Interrupt vector</b>														
Virtual External Interrupt	-	W	x	x	x	x	x	x	x	x	W	W	-	x
<b>Virtualization vector</b>														
Virtualization fault	-	W	x	x	x	x	x	x	x	x	W	W	-	W

- a. "N/A" indicates that this cannot happen.
- b. "W" indicates that the resource is written with a new value.
- c. "x" indicates that the resource may or may not be written; whether it is written and with what value is implementation specific.
- d. "-" indicates that the resource is not written.

**Table 8-2. ISR Values on Interruption**

Vector / Interruption	ed	ei <sup>a</sup>	so	ni <sup>b</sup>	ir <sup>c</sup>	rs <sup>d</sup>	sp <sup>e</sup>	na <sup>f</sup>	r	w	x
<b>Alternate Data TLB vector</b>											
Alternate Data TLB fault	ed <sup>k</sup>	ri	so	ni <sup>l</sup>	0	rs	sp	na	r	w	0
IR Alternate Data TLB fault	0	ri	0	ni <sup>l</sup>	1	1	0	0	1	0	0
<b>Alternate Instruction TLB vector</b>											
Alternate Instruction TLB fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Break Instruction vector</b>											
Break Instruction fault	0	ri	0	ni	0	0	0	0	0	0	0
<b>Data Access Rights vector</b>											
Data Access Rights fault	ed <sup>k</sup>	ri	so	ni	0	rs	sp	na	r	w	0
IR Data Access Rights fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Data Access-Bit vector</b>											
Data Access Bit fault	ed <sup>k</sup>	ri	so	ni	0	rs	sp	na	r	w	0
IR Data Access Bit fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Data Key Miss vector</b>											
Data Key Miss fault	ed <sup>k</sup>	ri	so	ni	0	rs	sp	na	r	w	0
IR Data Key Miss fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Data Nested TLB vector<sup>g</sup></b>											
Data Nested TLB fault	-	-	-	-	-	-	-	-	-	-	-
IR Data Nested TLB fault	-	-	-	-	-	-	-	-	-	-	-
<b>Data TLB vector</b>											
Data TLB fault	ed <sup>k</sup>	ri	so	ni <sup>l</sup>	0	rs	sp	na	r	w	0
IR Data TLB fault	0	ri	0	ni <sup>l</sup>	1	1	0	0	1	0	0
<b>Debug vector</b>											
Data Debug fault	ed <sup>k</sup>	ri	0	ni	0	rs	sp	na	r	w	0

**Table 8-2. ISR Values on Interruption (Continued)**

Vector / Interruption	ed	ei <sup>a</sup>	so	ni <sup>b</sup>	ir <sup>c</sup>	rs <sup>d</sup>	sp <sup>e</sup>	na <sup>f</sup>	r	w	x
Instruction Debug fault	0	ri	0	ni	0	0	0	0	0	0	1
IR Data Debug fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Dirty-Bit vector</b>											
Data Dirty Bit fault	ed <sup>k</sup>	ri	so	ni	0	rs	0	na <sup>h</sup>	r	1	0
<b>Disabled FP-Register vector</b>											
Disabled Floating-Point Register fault	0	ri	0	ni	0	0	sp	0	r	w	0
<b>External Interrupt vector</b>											
External Interrupt	0	ri	0	ni	ir <sup>i</sup>	0	0	0	0	0	0
<b>Floating-point Fault vector</b>											
Floating-Point Exception fault	0	ri	0	ni	0	0	0	0	0	0	0
<b>Floating-point Trap vector</b>											
Floating-Point Exception trap	0	ei	0	ni	0	0	0	0	0	0	0
<b>General Exception vector</b>											
Disabled ISA Transition fault	0	ri	0	ni	0	0	0	0	0	0	0
Illegal Dependency fault	0	ri	0	ni	0	0	0	0	0	0	0
Illegal Operation fault	0	ri	0	ni	0	0	0	0	0	0	0
IR Unimplemented Data Address fault	0	ri	0	ni	1	1	0	0	1	0	0
Privileged Operation fault	0	ri	0	ni	0	0	0	na	0	0	0
Privileged Register fault	0	ri	0	ni	0	0	0	0	0	0	0
Reserved Register/Field fault	0	ri	0	ni	0	0	0	0	0	0	0
Unimplemented Data Address fault	0	ri	0	ni	0	rs	0	na <sup>j</sup>	r	w	0
<b>IA-32 Exception vector</b>	0	0	0	0	0	0	0	0	0	0	x
<b>IA-32 Intercept vector</b>	0	0	0	0	0	0	0	0	r	w	0
<b>IA-32 Interrupt vector</b>	0	0	0	0	0	0	0	0	0	0	0
<b>Instruction Access Rights vector</b>											
Instruction Access Rights fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Instruction Access-Bit vector</b>											
Instruction Access Bit fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Instruction Key Miss vector</b>											
Instruction Key Miss fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Instruction TLB vector</b>											
Instruction TLB fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Key Permission vector</b>											
Data Key Permission fault	ed <sup>k</sup>	ri	so	ni	0	rs	sp	na	r	w	0
Instruction Key Permission fault	0	ri	0	ni	0	0	0	0	0	0	1
IR Data Key Permission fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Lower-Privilege Transfer Trap vector</b>											
Unimplemented Instruction Address fault	0	ri	0	ni	ir	0	0	0	0	0	1
Lower-Privilege Transfer trap	0	ei	0	ni	ir	0	0	0	0	0	0
Unimplemented Instruction Address trap	0	ei	0	ni	ir	0	0	0	0	0	0
<b>NaT Consumption vector</b>											
Data NaT Page Consumption fault	0	ri	so	ni	0	rs	0	na	r	w	0
Instruction NaT Page Consumption fault	0	ri	0	ni	0	0	0	0	0	0	1
IR Data NaT Page Consumption fault	0	ri	0	ni	1	1	0	0	1	0	0
Register NaT Consumption fault	0	ri	0	ni	0	0	0	na	r	w	0

**Table 8-2. ISR Values on Interruption (Continued)**

Vector / Interruption	ed	ei <sup>a</sup>	so	ni <sup>b</sup>	ir <sup>c</sup>	rs <sup>d</sup>	sp <sup>e</sup>	na <sup>f</sup>	r	w	x
<b>Page Not Present vector</b>											
Data Page Not Present fault	ed <sup>k</sup>	ri	so	ni	0	rs	sp	na	r	w	0
Instruction Page Not Present fault	0	ri	0	ni	0	0	0	0	0	0	1
IR Data Page Not Present fault	0	ri	0	ni	1	1	0	0	1	0	0
<b>Single Step Trap vector</b>											
Single Step trap	0	ei	0	ni	ir	0	0	0	0	0	0
<b>Speculation vector</b>											
Speculative Operation fault	0	ri	0	ni	0	0	0	0	0	0	0
<b>Taken Branch Trap vector</b>											
Taken Branch trap	0	ei	0	ni	ir	0	0	0	0	0	0
<b>Unaligned Reference vector</b>											
Unaligned Data Reference fault	ed	ri	0	ni	0	0	sp	0	r	w	0
<b>Unsupported Data Reference vector</b>											
Unsupported Data Reference fault	ed	ri	0	ni	0	0	0	0	r	w	0
<b>VHPT Translation vector</b>											
IR VHPT Data fault	0	ri	0	ni <sup>l</sup>	1	1	0	0	1	0	0
VHPT Data fault	ed <sup>k</sup>	ri	so	ni <sup>l</sup>	0	rs	sp	na	r	w	0
VHPT Instruction fault	0	ri	0	ni	0	0	0	0	0	0	1
<b>Virtual External Interrupt vector</b>											
Virtual External Interrupt	0	ri	0	ni	ir <sup>m</sup>	0	0	0	0	0	0
<b>Virtualization vector</b>											
Virtualization fault	0	ri	0	ni	0	0	0	0	0	0	0

- a. ISR.ei is equal to IPSR.ri for all faults and external interrupts (1 for faults and interrupts on the L+X instruction of an MLX). For traps, ISR.ei points at the excepting instruction (2 for traps on the L+X instruction of an MLX).
- b. If ISR.ni is 1, the interruption occurred either when PSR.ic was 0 or was in-flight.
- c. ISR.ir captures the value of RSE.CFLE at the time of an interruption.
- d. ISR.rs is 1 for interruptions caused by mandatory RSE fills/spills and 0 for all others.
- e. ISR.sp is 1 for interruptions caused by speculative loads and zero for all others.
- f. ISR.na is 1 for interruptions caused by non-access instructions and zero for all others.
- g. ISR is not written.
- h. A faulting `probe.w.fault` or `probe.rw.fault` can cause a Dirty Bit fault on a non-access instruction.
- i. ISR.ir is 1 if an external interrupt was taken when mandatory RSE fills caused by a `br.ret` or `rfi` were re-loading the current register stack frame.
- j. A faulting `lfetch.fault` or `probe.fault` to an unimplemented address will set ISR.na to 1.
- k. ISR.ed is 0 if the interruption was caused by a mandatory RSE fill or spill.
- l. If PSR.ic was 0 when the interruption was taken, these faults do not occur, but a Data Nested TLB fault is taken.
- m. ISR.ir is 1 if an external interrupt was taken when mandatory RSE fills caused by a `br.ret` or `rfi` were re-loading the current register stack frame.

Table 8-3 provides the definition for the ISR.code field on all Itanium traps. Hardware will always deliver the highest priority enabled trap. Software must look at the ISR.code bit vector to determine if any lower priority trap occurred at the same time as the trap being processed.

**Table 8-3. ISR.code Fields on Intel® Itanium® Traps**

Field	Bit	Description
fp	0	Floating-Point Exception trap
lp	1	Lower-Privilege Transfer trap

**Table 8-3. ISR.code Fields on Intel® Itanium® Traps (Continued)**

Field	Bit	Description
tb	2	Taken Branch trap
ss	3	Single Step trap
ui	4	Unimplemented Instruction Address trap
fp trap code	7	IEEE O (overflow) exception (Parallel FP-LO)
fp trap code	8	IEEE U (underflow) exception (Parallel FP-LO)
fp trap code	9	IEEE I (inexact) exception (Parallel FP-LO)
fp trap code	10	FPA, Added one to significand when rounding (Parallel FP-LO)
fp trap code	11	IEEE O (overflow) exception (Normal or Parallel FP-HI)
fp trap code	12	IEEE U (underflow) exception (Normal or Parallel FP-HI)
fp trap code	13	IEEE I (inexact) exception (Normal or Parallel FP-HI)
fp trap code	14	FPA, Added one to significand when rounding (Normal or Parallel FP-HI).

**Table 8-4. Interruption Vectors Sorted Alphabetically**

Vector Name	Offset	Page
Alternate Data TLB vector	0x1000	2:178
Alternate Instruction TLB vector	0x0c00	2:177
Break Instruction vector	0x2c00	2:185
Data Access Rights vector	0x5300	2:191
Data Access-Bit vector	0x2800	2:184
Data Key Miss vector	0x1c00	2:181
Data Nested TLB vector	0x1400	2:179
Data TLB vector	0x0800	2:176
Debug vector	0x5900	2:200
Dirty-Bit vector	0x2000	2:182
Disabled FP-Register vector	0x5500	2:195
External Interrupt vector	0x3000	2:186
Floating-Point Fault vector	0x5c00	2:203
Floating-Point Trap vector	0x5d00	2:204
General Exception vector	0x5400	2:192
IA-32 Exception vector	0x6900	2:210
IA-32 Intercept vector	0x6a00	2:211
IA-32 Interrupt vector	0x6b00	2:212
Instruction Access Rights vector	0x5200	2:190
Instruction Access-Bit vector	0x2400	2:183
Instruction Key Miss vector	0x1800	2:180
Instruction TLB vector	0x0400	2:175
Key Permission vector	0x5100	2:189
Lower-Privilege Transfer Trap vector	0x5e00	2:205
NaT Consumption vector	0x5600	2:196
Page Not Present vector	0x5000	2:188
Single Step Trap vector	0x6000	2:208
Speculation vector	0x5700	2:198
Taken Branch Trap vector	0x5f00	2:207
Unaligned Reference vector	0x5a00	2:201

**Table 8-4.    **Interruption Vectors Sorted Alphabetically (Continued)****

<b>Vector Name</b>	<b>Offset</b>	<b>Page</b>
Unsupported Data Reference vector	0x5b00	2:202
VHPT Translation vector	0x0000	2:173
Virtual External Interrupt vector	0x3400	2:187
Virtualization vector	0x6100	2:209



**Name**            **VHPT Translation vector (0x0000)**

**Cause**            The hardware VHPT walker encountered a TLB miss while attempting to reference the virtually addressed hashed page table for a memory reference (including IA-32).

Interruptions on this vector:

- IR VHPT Data fault
- VHPT Instruction fault
- VHPT Data fault

**Parameters**    IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IHA – The virtual address in the hashed page table which the hardware VHPT walker was attempting to reference.

ITIR – The ITIR contains default translation information for the virtual address contained in the IHA. The access key field within this register is set to the region id value from the region register selected by the virtual address in the IHA. The ITIR.ps field is set to the RR.ps field from the selected region register. All other fields are set to 0.

IIB0, IIB1 – If implemented, for VHPT Data faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR VHPT Data and VHPT Instruction faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

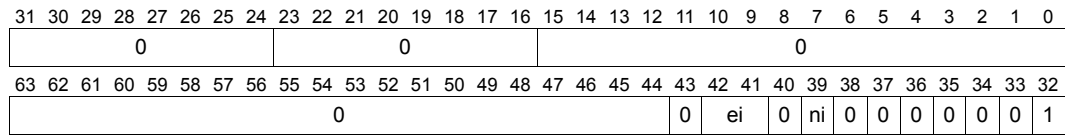
If the fault is due to a VHPT data fault for both original instruction and data references:

- IFA – The faulting address that the hardware VHPT walker was attempting to resolve.
- ISR – The ISR bits are set to reflect the original access on whose behalf the VHPT walker was operating. If the original operation was a non-access instruction then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0. For mandatory RSE fill or spill references, ISR.ed is always 0. The ISR.ni bit is 0 if PSR.ic was 1 when the interruption was taken, and is 1 if PSR.ic was in-flight. For IA-32 memory references the ISR.code, ni, ed, ei, ir, rs, sp, and na bits are always 0. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

If the fault is due to a VHPT instruction fault:

- IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits or, if the hardware VHPT walker was attempting to resolve a TLB miss, the virtual address of the translation.
- ISR – The ISR bits are set based on the original instruction fetch that the VHPT walker was attempting to resolve. The defined ISR bits are specified below. The ISR.ni bit is 0 if PSR.ic was 1 when the interruption was taken, and is 1 if PSR.ic was in-flight. For IA-32 memory references the ei and ni bits are always 0.



**Notes**

This fault can only occur when PSR.ic is 1 or in-flight, and the VHPT walker is enabled for the referenced region. Refer to ["VHPT Environment" on page 2:67](#) for details on VHPT enabling.

The original IFA address will be needed by the operating system page fault handler in the case where the page containing the VHPT entry has not yet been allocated. When the translation for the VHPT is available the handler must first move the address contained in the IHA to the IFA prior to the TLB insert.

Name **Instruction TLB vector (0x0400)**

Cause The instruction TLB entry needed by an instruction fetch (including IA-32) is absent, and the hardware VHPT walker could not find the translation in the VHPT, or the hardware VHPT walker is enabled but not implemented on this processor.

Interruptions on this vector:

Instruction TLB fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IHA – The virtual address of the hashed page table entry which corresponds to the reference that raised this fault.

ITIR – The ITIR contains default translation information for the original instruction address. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below. The ISR.ni bit is 0 if PSR.ic was 1 when the interruption was taken, and is 1 if PSR.ic was in-flight. The ISR.ei and ni bits are always 0 for IA-32 memory references.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																			
0								0								0																																		
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																			
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

Notes This fault can only occur when PSR.ic is 1 or in-flight, the VHPT hardware walker is enabled for the referenced region, the PSR.it bit is 1, and the fetched instruction bundle is to be executed. Refer to [“VHPT Environment”](#) on [page 2:67](#) for details on VHPT enabling.

The hardware VHPT walker may have failed due to an unimplemented page size, tag mismatch, illegal entry, or it may have terminated before reading the data. Software must be able to handle the case where the VHPT walker fails.

Name **Data TLB vector (0x0800)**

Cause For memory references (including IA-32), the data TLB entry needed by the data access is absent, and the hardware VHPT walker could not find the translation in the VHPT, or the hardware VHPT walker is not implemented on this processor.

Interruptions on this vector:

IR Data TLB fault  
Data TLB fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IHA – The virtual address of the hashed page table entry which corresponds to the reference that raised this fault.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The address of the data being referenced.

IIB0, IIB1 – If implemented, for Data TLB faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data TLB faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – If the interruption was due to a non-access operation then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0. For mandatory RSE fill or spill references, ISR.ed is always 0. The ISR.ni bit is 0 if PSR.ic was 1 when the interruption was taken, and is 1 if PSR.ic was in-flight. The ISR.code, ed, ei, ir, rs, sp and na bits are always 0 for IA-32 memory references. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
0								0								0								code{3:0}								
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	
0																						ed	ei	so	ni	ir	rs	sp	na	r	w	0

Notes The fault can only occur on an IA-32 or Itanium load, store, semaphore, or non-access operation when PSR.dt is 1, and the VHPT hardware walker is enabled for the referenced region. This fault can only occur on a mandatory RSE load/store operation if PSR.rt is 1, and the VHPT hardware walker is enabled for the referenced region. Refer to [“VHPT Environment”](#) on [page 2:67](#) for details on VHPT enabling.

The hardware VHPT walker may have failed due to an unimplemented page size, tag mismatch, illegal entry, or it may have terminated before reading the data. Software must be able to handle the case where the VHPT walker fails. The Data TLB fault is only taken if PSR.ic is 1 or in-flight, otherwise a Data Nested TLB fault is taken.

Name **Alternate Instruction TLB vector (0x0c00)**

Cause The instruction TLB entry needed by an instruction fetch (including IA-32) is absent, and the hardware VHPT walker was not enabled for this address.

Interruptions on this vector:

Alternate Instruction TLB fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the original instruction address. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – For Itanium memory references, the ISR.ei bits are set to indicate which instruction caused the exception and ISR.ni is set to 0 if PSR.ic was 1 when the interruption was taken, and set to 1 if PSR.ic was 0 or in-flight. For IA-32 memory references the ISR.ei and ni bits are 0. The defined ISR bits are specified below.

The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																			
0								0								0																																		
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																			
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

Notes This fault can only occur when the VHPT walker is disabled for the referenced region, and the fetched instruction bundle is to be executed. Refer to ["VHPT Environment"](#) on [page 2:67](#) for details on VHPT enabling.

Name **Alternate Data TLB vector (0x1000)**

Cause For memory references (including IA-32), the data TLB entry needed by data access is absent, and the hardware VHPT walker was not enabled for this address.

Interruptions on this vector:

IR Alternate Data TLB fault  
Alternate Data TLB fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The address of the data being referenced.

IIB0, IIB1 – If implemented, for Alternate Data TLB faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Alternate Data TLB faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – If the interruption was due to a non-access operation then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0. For mandatory RSE fill or spill references, ISR.ed is always 0. The ISR.ni bit is 0 if PSR.ic was 1 when the interruption was taken, and is 1 if PSR.ic was in-flight. For IA-32 memory references the ISR.code, ed, ei, ir, rs, sp and na bits are 0. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

Notes The fault can only occur on an IA-32 or Itanium load, store, semaphore, or non-access operation when PSR.dt is 1, and the VHPT hardware walker is disabled for the referenced region. This fault can only occur on a mandatory RSE load/store operation if PSR.rt is 1, and the VHPT hardware walker is disabled for the referenced region. The Alternate Data TLB fault is only taken if PSR.ic is 1 or in-flight, otherwise a Data Nested TLB fault is taken. Refer to [“VHPT Environment”](#) on [page 2:67](#) for details on VHPT enabling.

Name	<b>Data Nested TLB vector (0x1400)</b>
Cause	For memory references, the data TLB entry needed for a data reference is absent and PSR.ic is 0. Note: Data Nested TLB faults cannot occur during IA-32 instruction set execution, since PSR.ic must be 1.  Interruptions on this vector: <ul style="list-style-type: none"> <li>IR Data Nested TLB fault</li> <li>Data Nested TLB fault</li> </ul>
Parameters	IIP, IPSR, IIPA, IFS, ISR are <b>unchanged</b> from their previous values; they contain information relating to the original interruption.  ITIR – is <b>unchanged</b> from the previous value.  IFA – is <b>unchanged</b> from the previous value and contains the original address of the data being referenced.  IIB0, IIB1 – If implemented, the IIB registers are unchanged from their previous values. Please refer to <a href="#">Section 3.3.5.10, “Interrupt Instruction Bundle Registers (IIB0-1 – CR26, 27)”</a> on page 2:42 for details on the IIB registers.
Notes	This fault occurs when PSR.dt 1 and PSR.ic is 0 on a load, store, semaphore, and faulting non-access instructions. It also occurs when PSR.dt is 0 and PSR.ic is 0 for a regular_form probe instruction. Finally it can occur when PSR.rt is 1 and PSR.ic is 0 on a RSE mandatory load/store operation. Since the operating system is in control of the code executing at the time of the nested fault, it can by convention know which register contains the address that raised the nested event. As the PSR.ic bit is 0 on a nested fault, the IFA contains the original data address if the original interruption was caused by a data TLB fault. If the translation table entry required by the nested miss handler has not yet been allocated, then the address in the IFA will be passed to the operating system page fault handler. If the translation for the entry is available then the general register containing the nested fault address must be moved to the IFA prior to the insert. The ISR contains the ISR for the original faulting instruction, and not the ISR for the instruction that caused the nested fault.

Name **Instruction Key Miss vector (0x1800)**

Cause For instruction fetches (including IA-32), the PSR.it bit is 1, the PSR.pk bit is 1, and the access key from the TLB entry for the address of the executing instruction bundle does not match any of the valid protection keys.

Interruptions on this vector:

Instruction Key Miss fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the original instruction address. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. For IA-32 memory references the ISR.ei and ni bits are 0. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0							
0								0								0																						
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32							
0																				0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1



**Name**            **Data Key Miss vector (0x1c00)**

**Cause**            For memory references (including IA-32), the PSR.dt bit is 1, the PSR.pk bit is 1, and the access key from the TLB entry for the address referenced by a load, store, probe (regular\_form probe or probe.fault) or semaphore operation does not match any of the valid protection keys. The RSE may cause this fault if PSR.rt is 1, the PSR.pk bit is 1, and the access key from the TLB entry for the address referenced by an RSE mandatory load or store operation does not match any of the valid protection keys.

Interruptions on this vector:

IR Data Key Miss fault  
Data Key Miss fault

**Parameters**    IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – Faulting data address.

IIB0, IIB1 – If implemented, for Data Key Miss faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Key Miss faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – If the interruption was due to a non-access operation then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0. For mandatory RSE fill or spill references, ISR.ed is always 0. For IA-32 memory references, the ISR.code, ed, ei, ni, ir, rs, sp, and na bits are 0. The value for the ISR bits depend on the type of access performed and are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0																0						0					code{3:0}				
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

**Notes**            Probe (regular\_form probe or probe.fault) and the faulting variant of lfetch are the only non-access instructions that will cause a data key miss fault.

Name **Dirty-Bit vector (0x2000)**

Cause IA-32 or Itanium store or semaphore operations to a page with the dirty-bit (TLB.d) equal to 0 in the data TLB.

Interruptions on this vector:

Data Dirty Bit fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – Faulting data address.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE spill references, ISR.ed is always 0. For IA-32 memory references, ISR.ed, ei, ni, and rs are 0. If the interruption was due to a non-access operation then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	0	rs	0	na	r	1	0

Notes Dirty Bit fault can only occur in these situations:

- When PSR.dt is 1 on an IA-32 or Itanium store or semaphore operation
- When PSR.dt is 1 on a `probe.w.fault` or `probe.rw.fault`
- When PSR.rt is 1 on an RSE mandatory store operation

For `probe.w.fault` or `probe.rw.fault` the ISR.na bit is set, and the ISR.code field is written with a value of 5.

Only an IA-32 or Itanium semaphore, or `probe.rw.fault` operation would set ISR.r on a dirty bit fault.

Software is invoked to update the dirty bit in the data TLB entry and the Page table. The PSR.da bit can be used to suppress this fault for one executed instruction or one mandatory RSE store operation.

Name **Instruction Access-Bit vector (0x2400)**

Cause For instruction fetches (including IA-32), the access bit (TLB.a) in the TLB entry for this page is 0, and an instruction on the page is referenced.

Interruptions on this vector:

Instruction Access Bit fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. For IA-32 memory references the ISR.ei and ni bits are 0. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																			
0								0								0																																		
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																			
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

Notes The fault can only occur when PSR.it is 1 on an instruction reference (including IA-32). Software uses this fault for memory management page replacement algorithms. The PSR.ia bit can be used to suppress this fault for one executed instruction.

Name **Data Access-Bit vector (0x2800)**

Cause For data memory references (including IA-32), the access bit (TLB.a) in the TLB entry for this page is 0, and the page is referenced.

Interruptions on this vector:

IR Data Access Bit fault  
Data Access Bit fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – Faulting data address.

IIB0, IIB1 – If implemented, for Data Access Bit faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Access Bit faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0. For IA-32 memory references, ISR.code, ed, ei, ni, ir, rs, na and sp are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

Notes These faults can only occur in these situations:

- When PSR.dt is 1 on an IA-32 or Itanium load, store, or semaphore operation
- When PSR.dt is 1 on a `probe.fault`
- When PSR.dt is 1 on an `lfetch.fault`
- When PSR.rt is 1 on an RSE mandatory load/store operation

For `probe.fault` or `lfetch.fault` the ISR.na bit is set.

Software uses this fault for memory management page replacement algorithms. The PSR.da bit can be used to suppress this fault for one executed instruction or one mandatory RSE memory reference.

Name **Break Instruction vector (0x2c00)**

Cause An attempt is made to execute an Itanium `break` instruction.

Interruptions on this vector:

Break Instruction fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIM – Is updated with the break instruction immediate value.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, “Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0						
0								0								0																					
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32						
0																				0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Notes This fault cannot be raised by IA-32 instructions.

Name **External Interrupt vector (0x3000)**

Cause There are unmasked external interrupts pending from external devices, other processors, or internal processor events and:

- PSR.i is 1, while executing Itanium instructions
- PSR.i is 1 and (CFLAG.if is 0 or EFLAG.if is 1), while executing IA-32 instructions

IPSR.is indicates which instruction set was executing at the time of the interruption.

Interruptions on this vector:

External Interrupt

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IVR – Highest priority unmasked pending external interrupt vector number. If there are no unmasked pending interrupts the “spurious” interrupt vector (15) is reported.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction was to be executed when the external interrupt event was taken. The defined ISR bits are specified below. For external interrupts taken in the IA-32 instruction set, ISR.ei, ni and ir bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																		
0								0								0																																	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																		
0																					0	ei	0	ni	ir	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Notes: Software is expected to avoid situations which could cause ISR.ni to be 1.

Name **Virtual External Interrupt vector (0x3400)**

Cause The guest highest pending interrupt (GHPI) specified by the VMM is unmasked on the virtual processor.

IPSR.is indicates which instruction set was executing at the time of the interruption.

Interruptions on this vector:

Virtual External Interrupt

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction was to be executed when the external interrupt event was taken. The defined ISR bits are specified below. For external interrupts taken in the IA-32 instruction set, ISR.ei, ni and ir bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				0	ei	0	ni	ir	0	0	0	0	0	0	

Notes: Software is expected to avoid situations which could cause ISR.ni to be 1.

Name **Page Not Present vector (0x5000)**

Cause The bundle or IA-32 instruction being executed resides on a page for which the P-bit (TLB.p) in the instruction TLB entry is 0, or the data being referenced resides on a page for which the P-bit in the data TLB entry is 0.

Interruptions on this vector:

- IR Data Page Not Present fault
- Instruction Page Not Present fault
- Data Page Not Present fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IIB0, IIB1 – If implemented, for Data Page Not Present faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Page Not Present and Instruction Page Not Present faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

If the fault is due to a data page not present fault for both instruction and data original references:

- IFA – The virtual address of the data being referenced.
- ISR – If the interruption was due to a non-access operation then the ISR.code bits {3:0} are set to indicate the type of the non-access instruction; otherwise they are set to 0. The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0. For IA-32 memory references, ISR.code, ed, ei, ni, ir, rs, sp and na bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0																0						0						code{3:0}			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

If the fault is due to an instruction page not present fault:

- IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.
- ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below. For IA-32 memory references the ISR.ei and ni bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
0																0						0										
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	
0																					0	ei	0	ni	0	0	0	0	0	0	0	1

Notes This fault can only occur when PSR.it is 1 on an instruction reference, when PSR.dt is 1 on a load, store, semaphore, or non-access operation, or when PSR.rt is 1 on a RSE mandatory load/store operation.



**Name**            **Key Permission vector (0x5100)**

**Cause**            Data access (including IA-32): The PSR.dt bit is 1, the PSR.pk bit is 1 and read or write permission is disabled by the matching protection register on a load, store, or semaphore operation. The RSE may cause this fault if PSR.rt is 1, the PSR.pk bit is 1 and read or write permission is disabled by the matching protection register on an RSE mandatory load/store operation. Instruction access (including IA-32): The PSR.it bit is 1, the PSR.pk bit is 1 and execute permission is disabled by the matching protection register.

Interruptions on this vector:

- IR Data Key Permission fault
- Instruction Key Permission fault
- Data Key Permission fault

**Parameters**    IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IIB0, IIB1 – If implemented, for Data Key Permission faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Key Permission and Instruction Key Permission faults. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

If the fault is due to a data key permission fault:

- IFA – Faulting data address.
- ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0. For IA-32 memory references, the ISR.code, ed, ei, ni, ir, rs, sp bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0										
0																0																0						code{3:0}			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32										
0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0										

If the fault is due to an instruction key permission fault:

- IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.
- ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below. For IA-32 memory references, ISR.ei and ni are set to 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0						
0																0																0					
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32						
0																					0	ei	0	ni	0	0	0	0	0	0	0	1					

**Notes**            For `probe.fault` or `lfetch.fault` the ISR.na bit is set.

Name **Instruction Access Rights vector (0x5200)**

Cause For instruction fetches (including IA-32), the PSR.it bit is 1, and the access rights for this page do not allow execution or do not allow execution at the current privilege level.

Interruptions on this vector:

Instruction Access Rights fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below. For IA-32 memory references, ISR.ei and ni bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0									
0								0								0																								
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32									
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

Notes This fault does not occur if PSR.it is 0.

Name **Data Access Rights vector (0x5300)**

Cause For memory references (including IA-32), the PSR.dt bit is 1, and the access rights for this page do not allow read access or do not allow read access at the current privilege level for load and semaphore operations. The PSR.dt bit is 1, and the access rights for this page do not allow write access or do not allow write access at the current privilege level for store and semaphore operations.

The PSR.rt bit is 1, and the access rights for this page do not allow read access or do not allow read access at the current privilege level for the RSE mandatory load operation. The PSR.rt bit is 1, and the access rights for this page do not allow write access or do not allow write access at the current privilege level for the RSE mandatory store operation.

Interruptions on this vector:

IR Data Access Rights fault  
Data Access Rights fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

ITIR – The ITIR contains default translation information for the address contained in the IFA. The access key field within this register is set to the region id value from the referenced region register. The ITIR.ps field is set to the RR.ps field from the referenced region register. All other fields are set to 0.

IFA – Faulting data address.

IIB0, IIB1 – If implemented, for Data Access Rights faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Access Rights faults. Please refer to [Section 3.3.5.10, “Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0. For IA-32 memory references, ISR.code, ed, ei, ni, ir, rs, and sp bits are 0.

	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
	0										0										0							code{3:0}				
	63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
	0																					ed	ei	so	ni	ir	rs	sp	na	r	w	0

Notes For `probe.fault` or `lfetch.fault` the `ISR.na` bit is set.

Name **General Exception vector (0x5400)**

Cause An attempt is being made to execute an illegal operation, privileged instruction, access a privileged register, unimplemented field, unimplemented register, unimplemented address, or take an inter-instruction set branch when disabled.

Interruptions on this vector:

- IR Unimplemented Data Address fault
- Illegal Operation fault
- Illegal Dependency fault
- Privileged Operation fault
- Disabled Instruction Set Transition fault
- Reserved Register/Field fault
- Unimplemented Data Address fault
- Privileged Register fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP for the following faults:

- Illegal Operation fault
- Illegal Dependency fault
- Privileged Operation fault
- Disabled Instruction Set Transition fault
- Reserved Register/Field fault
- Unimplemented Data Address fault
- Privileged Register fault

The IIB registers are undefined for IR Unimplemented Data Address faults. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. For IA-32 instruction set faults, ISR.ei, ni, na, sp, rs, ir, ed bits are always 0.

- If the fault was caused by a non-access instruction, ISR.code{3:0} specifies which non-access instruction. See ["Non-access Instructions and Interruptions"](#) on [page 2:103](#).
- ISR.code{7:4} = 0: Illegal Operation fault. Cannot be raised by IA-32 instructions.
  - An attempt is being made to execute an illegal operation. Illegal operations include:
    - Attempts to execute instructions containing reserved major opcodes, reserved sub-opcodes, or reserved instruction fields, writing GR 0, FR 0 or FR 1, writing a read-only register, or accessing a reserved register.
    - Attempts to execute a reserved template encoding. An `rfi` to a reserved template encoding preserves IPSR.ri and will set ISR.ei to IPSR.ri.
    - Attempts to execute a bundle of template MLX when PSR.ri == 2. This can only be caused by doing an `rfi` with an improper setting of IPSR.ri. In this case, IPSR.ri and ISR.ei will both be 2.
    - Attempts to write outside the current register stack frame.
    - Attempts to specify the same GR, when the instruction has two GR targets (e.g., post-increment).

- If the instruction has two PR targets, and specifies the same PR for both, predicated-off unconditional compare, `fclass`, `tbit`, `tnat`, and `tf` instructions take this fault, even when their qualifying predicate is zero.
- Register bank conflict on a floating-point load pair instruction.
- An access to `BSPSTORE` or `RNAT` is performed with a non-zero `RSC.mode`, or a `loadrs` is performed with a non-zero `RSC.mode`.
- A `loadrs` is performed with a non-zero `CFM.sof` and a non-zero `RSC.loadrs`, or a `loadrs` causes more registers to be loaded from memory than can fit in the physical stacked register file.
- Attempts to predicate a `br.ia` instruction or to execute `br.ia` when `AR[BSPSTORE] != AR[BSP]`.
- Attempts to execute `epc` if `PFS.ppl` is less than `PSR.cpl`.
- Attempts to access interruption registers if `PSR.ic` is 1.
- Attempts to execute an `itc` or `itr` instruction if `PSR.ic` is 1.
- Attempts to allocate a stack frame larger than 96 registers, or with the rotating region larger than the stack frame, or with the size of locals larger than the stack frame, or specifying a qualifying predicate other than PR 0 on an `alloc` instruction.
- Attempts to execute instructions that are not supported by the processor.
- Attempts to execute a `ldfp` instruction with two odd-numbered physical FR targets or two even-numbered physical FR targets.
- Attempts to access an application register from the wrong unit type.
- Attempts to execute a `br.cloop`, `br.ctop`, `br.cexit`, `br.wtop`, or `br.wexit` other than in slot 2 of a bundle.
- Attempts to execute an `alloc`, `flushrs` or `loadrs` as other than the first instruction in an instruction group. (The result of such an attempt is undefined, and could result in an Illegal Operation fault, depending on the processor implementation. See [Section 3.5, “Undefined Behavior” on page 1:44](#) for details).
- Attempts to execute a `clrrrb`, `clrrrb.pr`, `cover`, `itc.d`, `itc.i`, `ptc.g` or `ptc.ga` instruction as other than the last instruction in an instruction group. (The result of such an attempt is undefined, and may possibly result in an Illegal Operation fault, depending on the processor See [Section 3.5, “Undefined Behavior” on page 1:44](#) for details).
- `ISR.code{7:4} = 1`: Privileged Operation fault. Cannot be raised by IA-32 instructions.
- `ISR.code{7:4} = 2`: Privileged Register fault. Cannot be raised by IA-32 instructions.
- `ISR.code{7:4} = 3`: Reserved Register/Field fault, Unimplemented Data Address fault or IR Unimplemented Data Address fault. Cannot be raised by IA-32 instructions. For Unimplemented Data Address fault:
  - If `ISR.rs = 0`: A data memory reference to an unimplemented address has occurred.
  - If `ISR.rs = 1`: A mandatory RSE reference to an unimplemented address has occurred.

For details, refer to [“Reserved and Ignored Registers and Fields” on page 1:23](#) and [“Unimplemented Address Bits” on page 2:73](#).

- $ISR.code\{7:4\} = 4$ : Disabled Instruction Set Transition fault. An instruction set transition was attempted while  $PSR.di$  was 1. This fault can be raised by either the Itanium `br.ia` instruction or the IA-32 `jmp` instruction.  $IPSR.is$  indicates the faulting instruction set.
- $ISR.code\{7:4\} = 8$ : Illegal Dependency fault. Cannot be raised by IA-32 instructions. The processor has detected a resource dependency violation.

If the fault is due to a Disabled ISA Transition fault, Illegal Dependency fault, Illegal Operation fault, Privileged Register fault or Reserved Register/Field fault:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0				
0								0								0							code{7:4}				0								
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32				
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0

Otherwise:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0				
0								0								0							code{7:4}				code{3:0}								
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32				
0																					0	ei	0	ni	ir	rs	0	na	r	w	0	0	0	0	0

Name **Disabled FP-Register vector (0x5500)**

Cause An attempt is made to reference a floating-point register set that is disabled.  
 When PSR.dfl is 1, execution of any IA-32 FP, SSE or MMX technology instructions raises a Disabled FP Register Low Fault (regardless of whether FR2 - FR31 are actually referenced).

When PSR.dfh is 1, execution of the first IA-32 instruction following a `br.ia` or `rfi` raises a Disabled FP Register High fault.

If concurrent IA-32 Disabled FP Register High and Low faults are generated, the Disabled FP Register High fault takes precedence and is reported in the ISR code, the Disabled FP Register Low fault is discarded and not reported in the ISR code.

Interruptions on this vector:

Disabled Floating-Point Register fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The defined ISR bits are specified below.

- `ISR.code{0}` = 1: FR2 - FR31 disabled and access attempted.
- `ISR.code{1}` = 1: FR32 - FR127 disabled and access attempted.

For IA-32 references, `ISR.ei`, `ni`, `sp`, `r`, and `w` bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0																0						0						code			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					0	ei	0	ni	0	0	sp	0	r	w	0

Name **NaT Consumption vector (0x5600)**

Cause A non-speculative operation (including IA-32) (e.g., load, store, control register access, instruction fetch etc.) read a NaT source register, NaTVal source register, or referenced a NaTPage.

Interruptions on this vector:

- IR Data NaT Page Consumption fault
- Instruction NaT Page Consumption fault
- Register NaT Consumption fault
- Data NaT Page Consumption fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to page 2:165 for a detailed description.

IIB0, IIB1 – If implemented, for Register NaT Consumption and Data NaT Page Consumption faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data NaT Page Consumption and Instruction NaT Page Consumption faults. Please refer to Section 3.3.5.10, “Interruptio

n Instruction Bundle Registers (IIB0-1 – CR26, 27)” on page 2:42 for details on the IIB registers. If the fault is due to a Data NaT Page Consumption fault or an IR Data NaT Page Consumption fault:

A non-speculative Itanium integer/FP instruction or instruction fetch or IA-32 data memory reference accessed a page with the NaTPage memory attribute.

- IFA – faulting data address.
- ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0. For the IA-32 instruction set, ISR.ed, ei, ni, ir, rs and na bits are 0. For probe.fault or lfetch.fault the ISR.na bit is set.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0		
0								0								0								2		code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32		
0																							0	ei	so	ni	ir	rs	0	na	r	w	0

If the fault is due to an Instruction NaT Page Consumption fault:

A non-speculative Itanium integer/FP instruction or instruction fetch accessed a page with the NaTPage memory attribute.

- IFA – The virtual address of the bundle or the 16 byte aligned IA-32 instruction address zero extended to 64-bits.
- ISR – The value for the ISR bits depend on the type of access performed and are specified below. For the IA-32 instruction set, ISR.ni and ei bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																								
0								0								0								2		0																													
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																								
0																							0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

If the fault is due to an Register NaT Consumption fault:

A non-speculative Itanium instruction reads a NaT’ed GR or an FR containing NaTVal. An IA-32 integer instruction reads a NaT’ed GR. For IA-32 instructions



behavior of NaT and NaTVal values is model specific, see [Section 6.2.4.3, "NaT/NaTVal Response for IA-32 Instructions"](#) on page 1:134 for details.

- ISR – The value for the ISR bits depend on the type of access performed and are specified below. For the IA-32 instruction set, ISR.ed, ei, ni, ir, rs, r, w, and na bits are 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0				1		code{3:0}									
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				0	ei	0	ni	0	0	0	na	r	w	0	

Name **Speculation vector (0x5700)**

Cause A `chk.a`, `chk.s`, or `fchkf` instruction needs to branch to recovery code, and the branching behavior is unimplemented by the processor. This fault cannot be raised by IA-32 instructions.

Interruptions on this vector:

Speculative Operation fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIM – contains the immediate value from the `chk.s`, `chk.a`, or `fchkf` instruction.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The `ISR.ei` bits are set to indicate which instruction caused the exception. The type of instruction which caused the fault is encoded in the lower four bits of the `ISR.code` field.

- If `ISR.code{3:0} = 0`: `chk.a` general register speculation fault.
- If `ISR.code{3:0} = 1`: `chk.s` general register speculation fault.
- If `ISR.code{3:0} = 2`: `chk.a` floating-point speculation fault.
- If `ISR.code{3:0} = 3`: `chk.s` floating-point speculation fault.
- If `ISR.code{3:0} = 4`: `fchkf` fault.

The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																			
0								0								0								code{3:0}																										
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																			
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Notes The Speculative Operation fault handler is required to perform the following steps:

1. Read the predicates and the IIM, IIP, IPSR, and ISR control registers, into scratch bank 0 general registers.
2. Copy the IIP value to IIPA.
3. Sign-extend the IIM value (from 21 bits to 64), shift it left by 4 bits, add it to the IIP value, and write this value back into IIP.
4. Set the `IPSR.ri` field to 0.
5. Check whether either `IPSR.tb` (Taken Branch trap) or `IPSR.ss` (Single Step enable) is 1. If not, emulation is complete, so restore the predicates and `rfi`. If so, then the check instruction would have taken one of these traps instead of branching to its target, so this handler needs to branch directly to the appropriate trap handler instead of performing the `rfi` (see steps 6 and 7).
6. If `IPSR.tb` was 1, then update `ISR.code` with its `tb` bit set to 1 and its `ss` bit also set to 1 if `IPSR.ss` was 1, and all other bits 0. Restore the predicates, execute a `srlz.d`, and branch to the taken branch vector (IVT offset 0x5f00).
7. If `IPSR.ss` was 1 (but not `IPSR.tb`), then update `ISR.code` with its `ss` bit set to 1, and all other bits 0. Restore the predicates, execute a `srlz.d`, and branch to the single step vector (IVT offset 0x6000).

The Speculative Operation fault handler does not need to check for unimplemented instruction addresses. They will be checked automatically by processor hardware when the handler executes its *rfi*. On processors which report unimplemented instruction addresses with an Unimplemented Instruction Address (UIA) trap, if an emulated check instruction targets an unimplemented address and also needs to take a Single Step trap or Taken Branch trap (or both), the UIA trap will not be raised until after the Single Step and/or Taken Branch trap has been handled, making it appear that the Unimplemented Instruction Address trap has the wrong priority. A Speculative Operation fault handler with this behavior is architecturally compliant. On processors which report unimplemented instruction addresses with an Unimplemented Instruction Address fault, the UIA fault will be taken at the target of the check rather than on the check instruction itself, so any Single Step trap and/or Taken Branch trap on the check will naturally become visible first.

Name **Debug vector (0x5900)**

Cause A debug fault has occurred. Either the instruction address matches the parameters set up in the instruction debug registers, or the data address of a load, store, semaphore, or mandatory RSE fill or spill matches the parameters set up in the data debug registers. All IA-32 instruction set debug events are delivered on the IA\_32\_Exception(Debug) vector; see [Chapter 9, "IA-32 Interruption Vector Descriptions."](#) IA-32 instructions can not raise this fault, IA-32 debug events are delivered on the IA\_32\_Exception(Debug) vector.

Interruptions on this vector:

- IR Data Debug fault
- Instruction Debug fault
- Data Debug fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, for Data Debug faults, the IIB registers contain the instruction bundle pointed to by IIP. The IIB registers are undefined for IR Data Debug and Instruction Debug faults. Please refer to [Section 3.3.5.10, "Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

If the fault is due to a data debug fault or an IR Data Debug fault:

- IFA – The address of the data being referenced.
- ISR – The value for the ISR bits depend on the type of access performed and are specified below. For mandatory RSE fill or spill references, ISR.ed is always 0.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								code{3:0}							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				ed	ei	0	ni	ir	rs	sp	na	r	w	0	

If the fault is due to an instruction debug fault:

- IFA – Faulting instruction fetch address.
- ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				0	ei	0	ni	0	0	0	0	0	0	0	1

Notes On an instruction reference this fault is suppressed if the PSR.db bit is 0 or if the PSR.id bit is 1. On a data reference this fault is suppressed if the PSR.db bit is 0 or if the PSR.dd bit is 1. The only non-access data operations which can cause a debug fault are the `probe.fault` and `lfetch.fault` instructions.

If unaligned accesses are being performed with debug faults enabled, this fault may be taken even though there is not a match for the address programmed in the breakpoint register. See [Section 7.1.2, "Debug Address Breakpoint Match Conditions"](#) on [page 2:154](#).

Name **Unaligned Reference vector (0x5a00)**

Cause If PSR.ac is 1, and the data address being referenced by an Itanium instruction is not aligned to the natural size of the load, store, or semaphore operation, or a data reference is made to a misaligned datum not supported by the implementation. See [“Memory Access Instructions” on page 1:57](#). For IA-32 data memory references, an IA\_32\_Exception(Alignment Check) fault is raised; see [Chapter 9, “IA-32 Interruption Vector Descriptions.”](#) IA-32 instructions can not raise this fault, IA-32 unaligned events are delivered on the IA\_32\_Exception(Alignment\_Check) vector.

If the data reference specified is both unaligned to the natural datum size and unsupported, then an Unaligned Data Reference fault is taken.

Interruptions on this vector:

Unaligned Data Reference fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IFA – The address of the data being referenced.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, “Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)” on page 2:42](#) for details on the IIB registers.

ISR – The value for the ISR bits depend on the type of access performed and are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	0	ni	0	0	sp	0	r	w	0

Name **Unsupported Data Reference vector (0x5b00)**

Cause An attempt was made to:

- Execute a `fetchadd`, `cmpxchg`, `xchg`, or unsupported `ld16`, `st16` or 10-byte memory reference (`ldfe` or `stfe`) instruction to a page that is neither cacheable with write-back write policy nor a NaTPage.
- Execute a `fetchadd` instruction to a page that is an uncacheable exported (UCE) page and the processor model does not support exporting of `fetchadd` instructions.

See “Effects of Memory Attributes on Memory Reference Instructions” on page 2:86 for details. IA-32 instructions can not raise this fault, IA-32 locked faults are delivered on the IA\_32\_Interrupt(Lock) vector.

If the data reference specified is both unaligned to the natural datum size and unsupported, then an Unaligned Data Reference fault is taken.

IA-32 data memory references that require an external atomic lock when `DCR.lc` is 1, raise an IA\_32\_Interrupt(Lock) fault; see Chapter 9, “IA-32 Interruption Vector Descriptions.”

Interruptions on this vector:

Unsupported Data Reference fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to page 2:165 for a detailed description.

IFA – The address of the data being referenced.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to Section 3.3.5.10, “Interruption Instruction Bundle Registers (IIB0-1 – CR26, 27)” on page 2:42 for details on the IIB registers.

ISR – The value for the ISR bits depend on the type of access performed and are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					ed	ei	0	ni	0	0	0	0	r	w	0

For `ldfe` and `stfe` instructions, the processor may optionally set both `ISR.r` and `ISR.w` to 1, although this is not recommended.

Name **Floating-point Fault vector (0x5c00)**

Cause A floating-point exception fault has occurred. IA-32 numeric instructions can not raise this fault, IA-32 floating point faults are delivered on the IA\_32\_Exception(Floating-Point) vector.

Interruptions on this vector:

Floating-Point Exception fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception.

ISR.code contains information about the FP exception fault. The ISR.code field has eight bits defined. See [Chapter 5](#) for details.

- ISR.code{0} = 1: IEEE V (invalid) exception (Normal or Parallel FP-HI)
- ISR.code{1} = 1: Denormal/Unnormal operand exception (Normal or Parallel FP-HI)
- ISR.code{2} = 1: IEEE Z (divide by zero) exception (Normal or Parallel FP-HI)
- ISR.code{3} = 1: Software assist (Normal or Parallel FP-HI)
- ISR.code{4} = 1: IEEE V (invalid) exception (Parallel FP-LO)
- ISR.code{5} = 1: Denormal/Unnormal operand exception (Parallel FP-LO)
- ISR.code{6} = 1: IEEE Z (divide by zero) exception (Parallel FP-LO)
- ISR.code{7} = 1: Software assist (Parallel FP-LO)

The defined ISR bits are specified below:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																		
0								0								0								code{7:0}																									
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																		
0																					0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Name **Floating-point Trap vector (0x5d00)**

Cause A floating-point exception trap has occurred. IA-32 numeric instructions can not raise this trap.

Interruptions on this vector:

Floating-Point Exception trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIPA. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception.

ISR.code contains information about the type of FP exception and IEEE information. The ISR code field contains a bit vector (see [Table 8-3 on page 2:170](#)) for all traps which occurred in the just-executed instruction. The defined ISR bits are specified below:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																	
0								0								0	fp trap code				0	0	0	ss	0	0	1																					
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																	
0																				0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



Name **Lower-Privilege Transfer Trap vector (0x5e00)**

Cause Two trapping conditions transfer control to this vector:

- An attempt is made to transfer control to an unimplemented address, resulting in either an Unimplemented Instruction Address trap or an Unimplemented Instruction Address fault. See “Unimplemented Address Bits” on page 2:73.
- The PSR.lp bit is 1, and a branch lowers the privilege level.

IA-32 instructions can not raise this trap.

Interruptions on this vector:

Unimplemented Instruction Address fault  
 Unimplemented Instruction Address trap  
 Lower-Privilege Transfer trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to page 2:165 for a detailed description.

**Note:** Please see “Interrupt Instruction Bundle Pointer (IIP – CR19)” on page 2:37 for a further clarification of the IIP value for an unimplemented instruction address trap.

IIB0, IIB1 – If implemented, for Lower-Privilege Transfer traps, the IIB registers contain the instruction bundle pointed to by IIPA. The IIB registers are undefined for Unimplemented Instruction Address faults and traps. Please refer to Section 3.3.5.10, “Interrupt Instruction Bundle Registers (IIB0-1 – CR26, 27)” on page 2:42 for details on the IIB registers.

ISR – For Unimplemented Instruction Address trap and Lower-Privilege Transfer trap, the ISR.ei bits are set to indicate which instruction caused the exception, and the ISR.code contains a bit vector (see Table 8-3 on page 2:170) for all traps which occurred in the just-executed instruction.

For Unimplemented Instruction Address fault ISR.fp\_trap\_code is set to 0.

The defined ISR bits are specified below.

If this vector was entered for an Unimplemented Instruction Address fault:

IFA – Faulting unimplemented instruction address

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
0								0								0	0								0	0	1	0	0	0	0	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	
0																				0	ri	0	ni	ir	0	0	0	0	0	0	0	1

If this vector was entered for an Unimplemented Instruction Address trap:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
0								0								0	fp trap code								0	0	1	ss	tb	lp	fp	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32	
0																				0	ei	0	ni	ir	0	0	0	0	0	0	0	0

If this vector was entered for a Lower-Privilege Transfer trap:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								0	0	0	ss	tb	1	0	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				0	ei	0	ni	ir	0	0	0	0	0	0	0

**Notes**

The Unimplemented Instruction Address trap can be the result of a taken branch, a taken *chk*, an *rfi*, or the execution of a slot 2 instruction in a bundle at the last implemented address. The lower privilege transfer trap is only taken on a branch demotion, and not an *rfi* return.

Processors may optionally report unimplemented instruction addresses with an Unimplemented Instruction Address fault on the fetch of the unimplemented address. To system software, this appears the same as if an Unimplemented Instruction Address trap had been taken, except that:

- any concurrent traps (Single Step, Taken Branch, Lower-Privilege Transfer, FP) will be taken first
- asynchronous interrupts (such as External interrupt) may be taken with IIP pointing to the unimplemented address before the Unimplemented Instruction Address fault is taken
- incomplete register stack frame interrupts may be taken with IIP pointing to the unimplemented address before the Unimplemented Instruction Address fault is taken
- *ISR.ei* will be equal to the value of *PSR.ri* at the time of the fault (and therefore will not indicate which instruction in the bundle pointed to by *IIPA* was responsible for the transition to an unimplemented address).

Name **Taken Branch Trap vector (0x5f00)**

Cause A taken branch was executed, and the PSR.tb bit is 1. IA-32 instructions can not raise this trap, IA-32 taken branch traps are delivered on the IA\_32\_Exception(Debug) vector.

The Taken Branch trap is not taken on an `rfi` instruction.

Interruptions on this vector:

Taken Branch trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

**Note:** Please see “[Interrupt Instruction Bundle Pointer \(IIP – CR19\)](#)” on [page 2:37](#) for a further clarification of the IIP value for an unimplemented instruction address trap or fault.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIPA. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The ISR.code contains a bit vector (see [Table 8-3 on page 2:170](#)) for all traps which occurred in the just-executed instruction. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								0								0								0	0	0	ss	1	0	0	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																				0	ei	0	ni	ir	0	0	0	0	0	0	0

Name **Single Step Trap vector (0x6000)**

Cause An instruction was successfully executed, and the PSR.ss bit is 1. For IA-32 instruction set, this condition is delivered on the IA\_32\_Exception(Debug) vector; see [Chapter 9, "IA-32 Interruption Vector Descriptions."](#) IA-32 instructions can not raise this trap, IA-32 single step events are delivered on the IA\_32\_Exception(Debug) vector.

The Single Step trap is not taken on an `rfi` instruction.

Interruptions on this vector:

Single Step trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIPA. Please refer to [Section 3.3.5.10, "Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception. The ISR.code contains a bit vector (see [Table 8-3 on page 2:170](#)) for all traps which occurred in the just-executed instruction. The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																		
0								0								0								0	0	0	1	0	0	0																			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																		
0																					0	ei	0	ni	ir	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Name **Virtualization vector (0x6100)**

Cause An attempt is made to execute an instruction which requires virtualization. This fault cannot be raised by IA-32 instructions.

Interruptions on this vector:

Virtualization fault

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers contain the instruction bundle pointed to by IIP. Please refer to [Section 3.3.5.10, "Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – The ISR.ei bits are set to indicate which instruction caused the exception.

The defined ISR bits are specified below.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0						
0								0								0								0													
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32						
0																				0	ei	0	ni	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Name **IA-32 Exception vector (0x6900)**

Cause A fault or trap was raised while executing from the IA-32 instruction set.

Interruptions on this vector:

- IA-32 Instruction Debug fault
- IA-32 Code Fetch fault
- IA-32 Instruction Length > 15 bytes fault
- IA-32 Device Not Available fault
- IA-32 FP Error fault
- IA-32 Segment Not Present fault
- IA-32 Stack Exception fault
- IA-32 General Protection fault
- IA-32 Divide by Zero fault
- IA-32 Alignment Check fault
- IA-32 Bound fault
- IA-32 INTO trap
- IA-32 Breakpoint (INT 3) trap
- IA-32 Data Breakpoint trap
- IA-32 Taken Branch trap
- IA-32 Single Step trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IFA – is undefined. The faulting IA-32 address is contained in IIPA.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, “Interrupt Instruction Bundle Registers \(IIB0-1 – CR26, 27\)”](#) on [page 2:42](#) for details on the IIB registers.

ISR – ISR.vector contains the IA-32 exception vector number. ISR.code contains the IA-32 error code for faults or a trap code listing concurrent trap events for traps.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								vector								error_code/trap_code															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					0	0	0	0	0	0	0	0	0	x	

Notes See [Chapter 9, “IA-32 Interruption Vector Descriptions”](#) for complete details on each IA-32 Exception and for error code and trap code definition.

Name **IA-32 Intercept vector (0x6a00)**

Cause An intercept fault or trap was raised while executing from the IA-32 instruction set. This vector handles all the IA-32 intercepts described in [Chapter 9, "IA-32 Interruption Vector Descriptions."](#)

Interruptions on this vector:

- IA-32 Invalid Opcode fault
- IA-32 Instruction Intercept fault
- IA-32 Locked Data Reference fault
- IA-32 System Flag Intercept trap
- IA-32 Gate Intercept trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIM – 64-bit information describing the cause of the intercept.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, "Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – ISR.vector contains a number specifying the type of intercept. ISR.code contains the IA-32 specific intercept information or a trap code listing concurrent trap events for traps.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								intercept_number								intercept_code/trap_code															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
0																					0	0	0	0	0	0	0	0	r	w	0

Notes See [Chapter 9, "IA-32 Interruption Vector Descriptions"](#) for complete details on each IA-32 Intercept and for the intercept code and trap code definition.

Name **IA-32 Interrupt vector (0x6b00)**

Cause An IA-32 software interrupt trap was executed. This vector handles all the IA-32 software interrupts described in [Chapter 9, "IA-32 Interruption Vector Descriptions."](#)

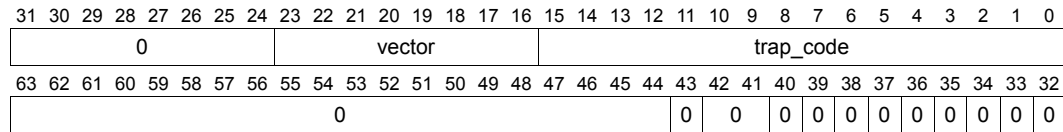
Interruptions on this vector:

IA-32 Software Interrupt (INT) trap

Parameters IIP, IPSR, IIPA, IFS – are defined; refer to [page 2:165](#) for a detailed description.

IIB0, IIB1 – If implemented, the IIB registers are undefined. Please refer to [Section 3.3.5.10, "Interruption Instruction Bundle Registers \(IIB0-1 – CR26, 27\)"](#) on [page 2:42](#) for details on the IIB registers.

ISR – ISR.vector contains the IA-32 defined interruption vector number. ISR.code contains a trap code listing concurrent trap events.



Notes See [Chapter 9, "IA-32 Interruption Vector Descriptions"](#) for complete details on this vector and the trap code definition.

§



This section gives detailed description of all possible IA-32 exceptions, interrupts and intercepts that can occur during IA-32 instruction set execution in the Itanium System Environment. Interruption resources not noted below are undefined after the interruption. For all cases where an interruption is taken out of the IA-32 instruction set, IPSR.is is set to 1.

## 9.1 IA-32 Trap Code

The following trap code is defined for concurrent traps reported during IA-32 instruction set execution. There is a bit for every possible concurrent trap condition.

**Figure 9-1. IA-32 Trap Code**

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0								b3	b2	b1	b0	ss	tb	0	

**Figure 9-2. IA-32 Trap Code**

Bit	Name	Description
2	tb	taken branch trap, set if an IA-32 branch is taken and branch traps are enabled (PSR.tb is 1).
3	ss	single step trap, set after the successful execution of every IA-32 instruction if PSR.ss or EFLAG.tf is 1.
4-7	b0 to b3	Data breakpoint trap due to a match with the corresponding Intel Itanium data breakpoint registers. Each bit indicates a match with the corresponding DBR registers; b0=DBR0/1, b1=DBR2/3, b2=DBR4/5, b3=DBR6/7. Zero, one or more bits may be set. These bits accumulate data breakpoint register matches that occurred during the duration of executing one IA-32 instruction. In order to be reported, the DBR register address and mask registers must precisely match the IA-32 data memory reference address, and the DBR read, write bits match the type of memory transaction, and the DBR privilege level mask match the value in PSR.cpl.

## 9.2 IA-32 Interruption Vector Definitions

Following are the definitions of IA-32 exceptions, interrupts and intercepts that can occur during IA-32 instruction set execution in the Itanium system environment.

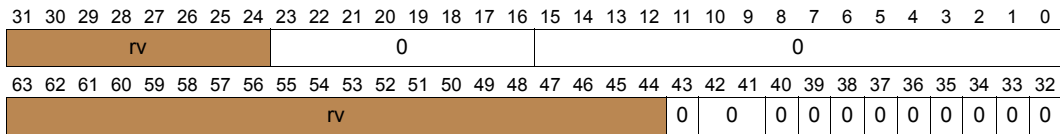
Name **IA\_32\_Exception (Divide) – Divide Fault**

Cause IA-32 IDIV or DIV instruction attempted a divide by zero operation. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 0.



Name **IA\_32\_Exception (Debug) – Code Breakpoint Fault**

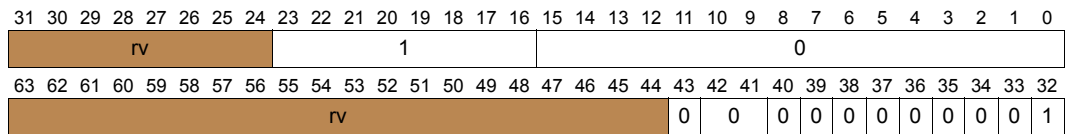
Cause The Itanium architecture debug facilities triggered an IA-32 code breakpoint fault on a IA-32 instruction fetch and PSR.id and EFLAG.rf are 0. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 1.

ISR.x – 1.



Name **IA\_32\_Exception (Debug) – Data Breakpoint, Single Step, Taken Branch Trap**

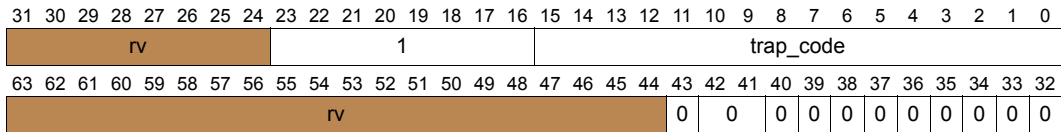
Cause The Itanium architecture debug facilities triggered an IA-32 data breakpoint, single-step or branch trap. In the Itanium System Environment, IA-32 Mov SS or Pop SS single step and data breakpoint traps are NOT deferred to the next instruction. Refer to the *Intel® 64 and IA-32 Architectures Software Developer’s Manual* for a complete definition of this trap.

Parameters IIPA – virtual address of the trapping IA-32 instruction (zero extended to 64-bits) if there was a taken branch trap. On `jmpbe` taken branch traps IIPA contains the address of the `jmpbe` instruction. For all other trap events, IIPA is undefined.

IIP – next Itanium instruction address or the virtual IA-32 instruction address zero extended to 64-bits.

ISR.vector – 1.

ISR.code – Trap Code, indicates Concurrent Single Step, Taken Branch, Data Breakpoint Trap events.



Name **IA\_32\_Exception (Break) – INT 3 Trap**

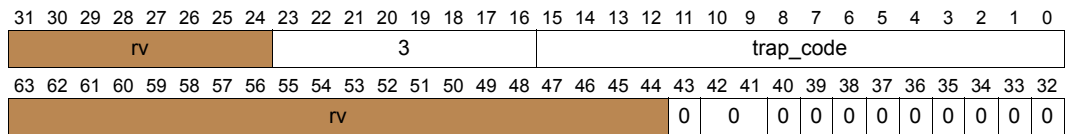
Cause IA-32 breakpoint instruction (INT 3) triggered a trap. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this trap.

Parameters IIPA – trapping virtual IA-32 instruction address zero extended to 64-bits.

IIP – next virtual IA-32 instruction address zero extended to 64-bits.

ISR.vector – 3.

ISR.code –Trap Code, indicates Concurrent Single Step condition.



Name **IA\_32\_Exception (Overflow) – Overflow Trap**

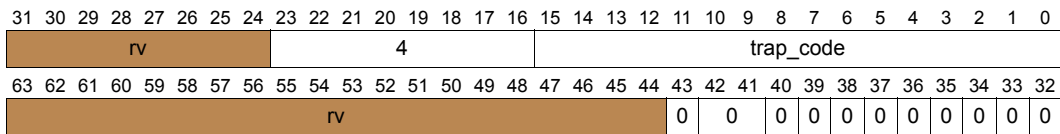
Cause IA-32 INTO instruction execution when EFLAG.of is set to one. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this trap.

Parameters IIPA – trapping virtual IA-32 instruction address zero extended to 64-bits.

IIP – next virtual IA-32 instruction address zero extended to 64-bits.

ISR.vector – 4.

ISR.code – Trap Code, indicates Concurrent Single Step.



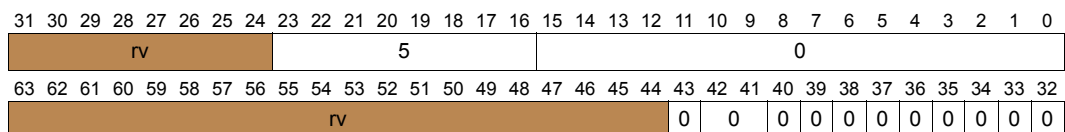
Name **IA\_32\_Exception (Bound) – Bounds Fault**

Cause Failed IA-32 Bound check instruction. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 5.



Name	<b>IA_32_Exception (InvalidOpcode) – Invalid Opcode Fault</b>
Cause	All IA-32 invalid opcode faults are delivered to the IA_32_Intercept(Instruction) handler, including IA-32 illegal, unimplemented opcodes, MMX technology and SSE instructions if CR0.EM is 1, and SSE instructions if CR4.fxsr is 0. All illegal IA-32 floating-point opcodes result in an IA_32_Intercept(Instruction) regardless of the state of CR0.em.



Name **IA\_32\_Exception (DNA) – Device Not Available Fault**

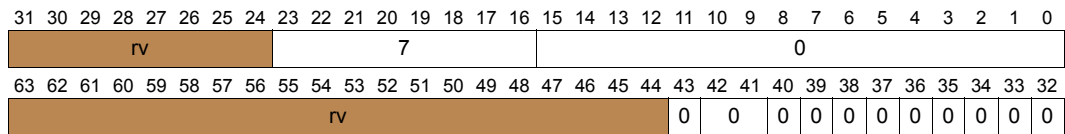
Cause The processor executed an IA-32 ESC or floating-point instruction with CR0.em is 1. Or an IA-32 WAIT, ESC, floating-point instruction, MMX technology or SSE instruction is executed and CR0.ts bit is 1.

Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 7.



Name        **Double Fault**

Cause        IA-32 Double Faults (IA-32 vector 8) are not generated by the processor in the Itanium System Environment.

Name        **Invalid TSS Fault**

Cause        IA-32 Invalid TSS Faults (IA-32 vector 10) are not generated in the Itanium System Environment.

Name **IA\_32\_Exception (NotPresent) – Segment Not Present Fault**

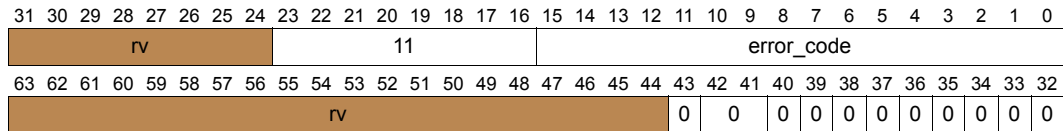
Cause Generated when the processor detects the Present-bit of the memory segment descriptor is zero during an IA-32 segment load or far control transfer instructions. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault and error codes.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 11.

ISR.code – IA-32 defined error code. See **Intel® 64 and IA-32 Architectures Software Developer’s Manual**.



Name **IA\_32\_Exception (StackFault) – Stack Fault**

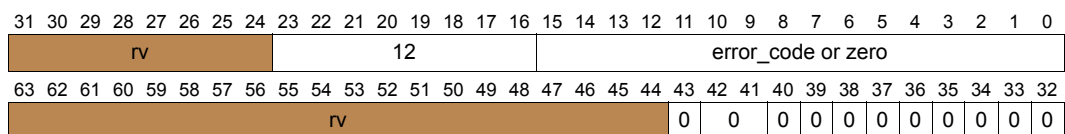
Cause IA-32 defined set of stack segment fault conditions detected during stack segment load operations or memory references relative to the stack segment, refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete list of all IA-32 faulting conditions. Stack faults can also be generated when the processor detects an inconsistent stack segment register descriptor value during an IA-32 stack reference instruction (e.g. PUSH, POP, CALL, RET,). See section “[Segment Descriptor and Environment Integrity](#)” for a list of possible inconsistent register descriptor conditions.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 12.

ISR.code – IA-32 defined ErrorCode. Zero if an inconsistent register descriptor is detected during a memory reference relative to the stack segment.



Name **IA\_32\_Exception (GPFault) – General Protection Fault**

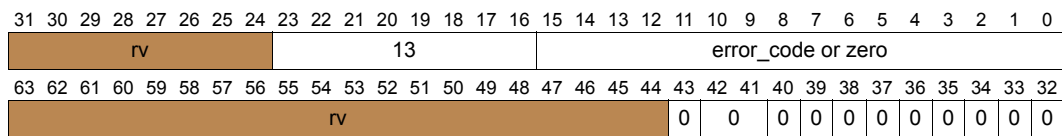
Cause IA-32 defined set of data and code segment fault conditions detected during data or code segment load operations or memory references relative to code or data segments, refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete list of all IA-32 General Protection Fault conditions. General Protection faults can also be generated when the processor detects an inconsistent code or data segment register descriptor value during an IA-32 code fetch or data memory reference. See section “[Segment Descriptor and Environment Integrity](#)” for a list of possible inconsistent register descriptor conditions.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 13.

ISR.code – IA-32 defined ErrorCode. Zero if an inconsistent register descriptor is detected during a memory reference relative to a code or data segment.



Name        **Page Fault**

Cause        IA-32 defined page faults (IA-32 vector 14) can not be generated in the Itanium System Environment.

Name **IA\_32\_Exception (FPError) – Pending Floating-point Error**

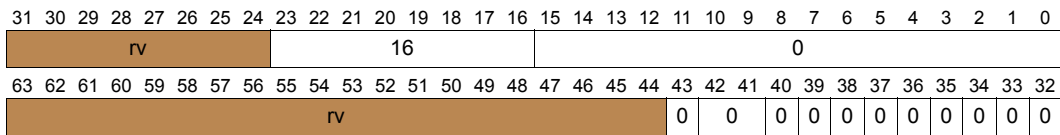
Cause An unmasked IA-32 floating-point exception is delivered on the next non-control IA-32 floating-point, MMX technology, WAIT, or `jmpbe` instruction trigger delivery of this exception. Floating-point errors are delivered regardless of the state of CR0.ne in the Itanium System Environment. IA-32 numeric exception delivery is not triggered by Itanium numeric exceptions or the execution of Itanium numeric instructions. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

FSR, FIR, FDR and FCR contain the IA-32 floating-point environment and exception information.

ISR.vector – 16.





Name **IA\_32\_Exception (AlignmentCheck) – Alignment Check Fault**

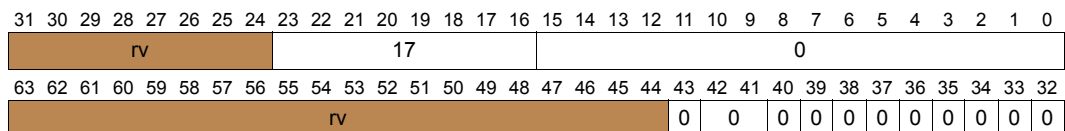
Cause An IA-32 instruction performed an unaligned data memory reference while PSR.ac is 1, or EFLAG.ac is 1 and CR0.am is 1 and the effective privilege level is 3. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

IFA – referenced virtual data address (byte granular) zero extended to 64-bits.

ISR.vector – 17.



Name        **Machine Check**

Cause        IA-32 Machine Check (IA-32 vector 18) is not generated in the Itanium System Environment.

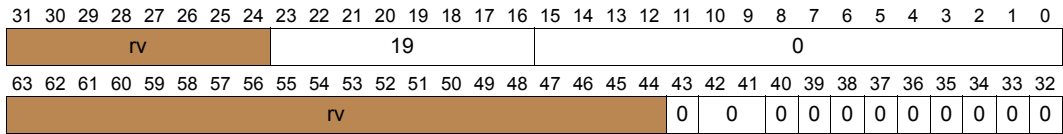
Name **IA\_32\_Exception (StreamingSIMD) – SSE Numeric Error Fault**

Cause An unmasked IA-32 SSE numeric error occurred. Numeric faults generated on SSE instructions are reported precisely on the faulting SSE instruction. SSE instructions do NOT trigger the report of any pending IA-32 floating-point exceptions. SSE instructions always ignore CR0.ne and the IGNNE pin. Refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for a complete definition of this fault.

Parameters IIP – virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

ISR.vector – 19.



Name **IA\_32\_Interrupt (Vector #N) – Software Trap**

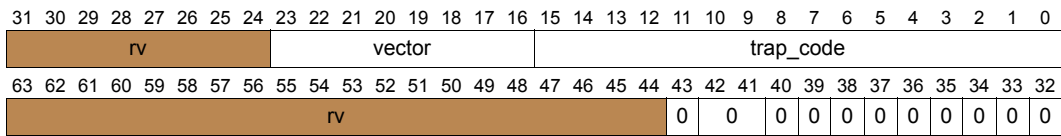
Cause The IA-32 INT n instruction forces an IA-32 interrupt trap. The IA-32 IDT is not consulted nor are any values pushed onto a memory stack.

Parameters IIPA – trapping virtual IA-32 instruction address (points to the INT instruction) zero extended to 64-bits.

IIP – next virtual IA-32 instruction address zero extended to 64-bits.

ISR.vector – vector number.

ISR.code – TrapCode, Indicates Concurrent Single Step Trap condition.



Name **IA\_32\_Intercept (Instruction) – Instruction Intercept Fault**

Cause Execution of unimplemented IA-32 opcodes, illegal opcodes or sensitive privileged IA-32 operating system instructions results in an instruction intercept. Intercepted opcodes include (but are not limited to); CLTS, HLT, INVD, INVLPG, IRET, LIDT, LGDT, LLDT, LMSW, LTR, MOV to CRs, MOV to/from DRs, RDMSR, RSM, SYSENTER, SYSEXIT, INT1, SIDT, SGDT, SLDT, SMSW, WBINVD, WRMSR, and all other unimplemented and illegal opcode patterns. If CR0.em is 1, execution of all IA-32 Intel MMX technology and IA-32 SSE instructions results in this intercept. If CR4.FXSR is 0, execution of all IA-32 SSE instructions results in this intercept. All illegal IA-32 floating-point opcodes result in an IA\_32\_Intercept(Instruction) regardless of the state of CR0.em. Intercepted opcodes are nullified and alter no architectural state.

Parameters IIP – Virtual IA-32 instruction address zero extended to 64-bits, points to the first byte of the intercepted IA-32 opcode (including prefixes).

IIPA – Virtual address of the faulting IA-32 instruction zero extended to 64-bits.

IIM – Opcode bytes, contains the first 8-bytes of the IA-32 instruction following all prefix bytes. All prefix bytes are decoded and presented as a bitmask in the Intercept Code along with the prefix length in bytes. Opcode bytes are loaded into IIM in the same format as encountered in memory and as defined in the **Intel® 64 and IA-32 Architectures Software Developer’s Manual**. The lowest memory address byte is placed in byte 0 of IIM, higher memory address bytes are placed in increasingly higher numbered bytes within IIM.

The 8-byte opcode loaded into IIM is stripped of the following prefixes; lock, repeat, address size, operand size, and segment override prefixes (opcode bytes 0xF3, 0xF2, 0xF0, 0x2E, 0x36, 0x3E, 0x26, 0x64, 0x65, 0x66, and 0x67). The 0x0F opcode series prefix is not stripped from the opcode bytes loaded into IIM. The opcode loaded into IIM includes all IA-32 opcode components, including 1 to 3 bytes of opcode, mod r/m bytes, sib bytes and any possible immediates and/or displacements.

If the opcode loaded in IIM is less than 8-bytes, the remainder higher order numbered bytes are set to 0. If the opcode is larger than 8-bytes, bytes after the 8th byte (following all stripped prefixes) are not reported. If required, emulation code must retrieve the extra opcode bytes by reading from the memory locations specified by IIP.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
byte3								byte2								byte1								byte0							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
byte7								byte6								byte5								byte4							

ISR.vector – 0, indicates instruction intercept.

ISR.code – Intercept Code indicates prefixes and prefix lengths.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0							
rv								0								intercept_code																						
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32							
rv																				0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 9-3 defines intercept codes for IA-32 instruction set intercepts. Intercept code fields are defined by Table 9-1 and Table 9-2 on page 2:234.

**Figure 9-3. IA-32 Intercept Code**

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
len				0	seg			sp	np	rp	lp	as	os	0	

**Table 9-1. Intercept Code Definition**

Field	Bits	Description
os	1	Operand Size – (OperandSize Prefix XOR CSD.d bit). When 1, indicates the effective operand size is 32-bits, when 0, 16-bits.
as	2	Address Size – (AddressSize Prefix XOR CSD.d bit). When 1, indicates the effective address size is 32-bits, when 0, 16-bits.
lp	3	Lock Prefix – If 1, indicates a lock prefix is present.
rp	4	REP or REPE/REPZ Prefix – If 1, indicates a REP/REPE/REPZ prefix is in effect.
np	5	REPNE/REPZ Prefix – If 1, indicates a REPNE/REPZ prefix is in effect.
sp	6	Segment Prefix – If 1, indicates a Segment Override prefix is present.
seg	7:9	Segment Value – Segment Prefix Override value, see <a href="#">Figure 9-2</a> for encodings. If there is no segment prefixes this field is undefined.
len	12:15	Length of Prefixes – Length of all prefix (in bytes) stripped from IIM. If there are no prefixes this field has a value of zero.

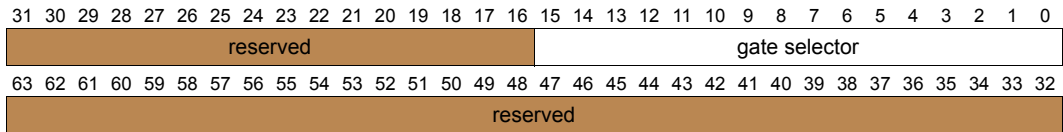
**Table 9-2. Segment Prefix Override Encodings**

Seg Value	Segment Prefix
0	ES Segment Override
1	CS Segment Override
2	SS Segment Override
3	DS Segment Override
4	FS Segment Override
5	GS Segment Override
6	reserved
7	reserved

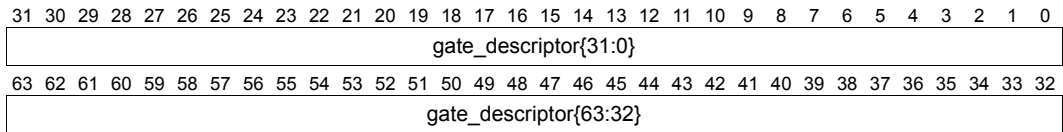
Name **IA\_32\_Intercept (Gate) – Gate Intercept Trap**

Cause If an IA-32 control transfer is initiated through a GDT/LDT descriptor that transfers control through a Call Gate, Task Gate or Task Segment this interception trap is generated.

Parameters IIPA – trapping virtual IA-32 instruction address zero extended to 64-bits.  
 IIP – next sequential virtual IA-32 instruction address zero extended to 64-bits.  
 IFA – Gate Selector. The gate selector is loaded in IFA{15:0}.



IIM – Gate, Task Gate or Task Segment Descriptor. The descriptor loaded in IIM adheres to the IA-32 GDT/LDT memory format, where byte 0 of the descriptor is in IIM{7:0}.



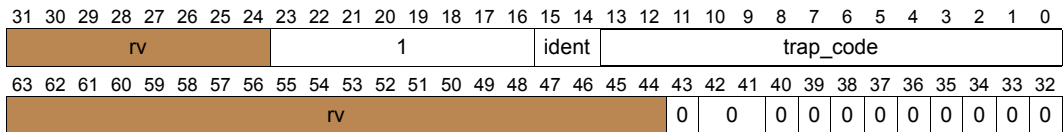
**Table 9-3. Gate Intercept Trap Code Identifier**

Instruction	ISR.code{15:14}
CALL	00
JMP	01

ISR.vector – 1, indicates gate interception.

ISR.code – TrapCode, Indicates Concurrent Data Debug, taken Branch, and Single Step Events.

ISR.code{15:14} – indicates whether CALL or JMP generated the trap. See [Table 9-3](#) for details.



Name **IA\_32\_Intercept (SystemFlag) – System Flag Trap**

Parameters System Flag Intercept Traps are generated for the following conditions:

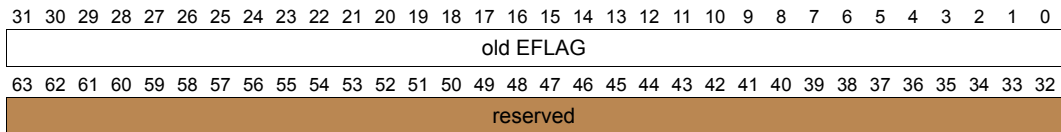
**CLI, STI, POPF, POPFD instructions.** If the EFLAG.if bit changes state and CFLG.ii is 1, or EFLAG.tf or EFLAG.ac change state, a System Flag intercept notification trap is delivered after the instruction completes. IIM contains the previous value of EFLAG before the trapping instruction executed. If IA-32 code does not have IOPL or CPL permission to modify the EFLAG bits, no intercept is generated. This intercept trap condition can be used to provide virtual interrupt services, and delay enabling of interrupts after the STI instruction.

**MOV SS, POP SS instructions.** After these instructions complete execution, a System Flag intercept notification trap is delivered. This intercept trap condition can be used to inhibit interrupts, and code breakpoints between Mov/Pop SS and the next instruction and to inhibit Single Step and Data Breakpoint traps on the Mov, or Pop SS instruction.

IIP – next virtual IA-32 instruction address zero extended to 64-bits.

IIPA – trapping virtual IA-32 instruction address zero extended to 64-bits.

IIM – contains the previous EFLAG value before the trapping instruction.



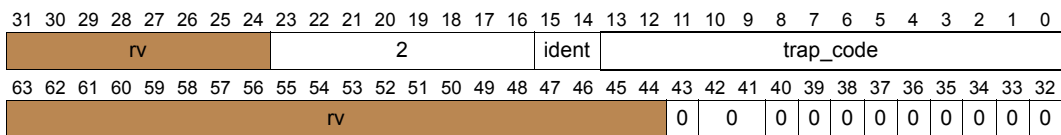
ISR.vector – 2.

ISR.code – Trap Code, indicates Concurrent Single Step Trap, Debug trap condition.

ISR.code{15:14} indicates which instruction generated the trap.

**Table 9-4. System Flag Intercept Instruction Trap Code Instruction Identifier**

Instruction	ISR.code{15:14}
CLI	00
STI	01
POPF, POPFD	10
MOV/POP SS	11





Name **IA\_32\_Intercept (Lock) – Locked Data Reference Fault**

Cause For IA-32 locked operations, if the DCR.lc bit is 1, and an atomic operation to made to non-write-back memory or to unaligned write-back memory that would result in a read-modify-write sequence being performed externally under an external bus lock, the processor raises a Locked Data Reference fault.

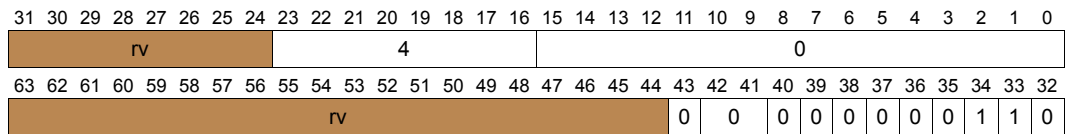
Parameters IIP – faulting virtual IA-32 instruction address zero extended to 64-bits.

IIPA – virtual address of the faulting IA-32 instruction zero extended to 64-bits.

IFA – faulting virtual data address (byte granular) zero extended to 64-bits.

ISR.vector – 4.

ISR.code – 0.



§



# Itanium® Architecture-based Operating System Interaction Model with IA-32 Applications

10

This section describes the IA-32 system execution model from the perspective of an Itanium architecture-based operating system interfacing with IA-32 code, while operating in the Itanium System Environment. The main features covered are:

- IA-32 system and control register behavior
- IA-32 virtual memory support
- IA-32 fault and trap handling
- IA-32 instruction behavior

## 10.1 Instruction Set Transitions

Instruction set transitions are defined in [Section 6.2.1, “Instruction Set Modes” on page 1:110](#). Operating systems can disable instruction set transitions (`jmp` and `br.ia`) by setting `PSR.di` to one. If `PSR.di` is one, execution of `jmp` or `br.ia` to IA-32 target results in a Disabled Instruction Set Transition Fault, and the operation is nullified.

The processor also transitions into an Itanium architecture-based operating system when IA-32 privileged system resources are accessed, on an interruption, or when the following conditions are detected:

- Instruction Interception – IA-32 system level privileged instructions are executed
- System Flag Interception – Various EFLAG system flags are modified, (e.g. AC, TF and IF-bits)
- Gate Interception – Control transfers are made through call gate, or transfers through a task switch (TSS segment or Task Gate).

All software interrupts, external interrupts, faults, traps and machine checks transition the processor to the Itanium instruction set, regardless of the state of `PSR.di`. IA-32 defined exceptions and software interrupts are delivered to Itanium architecture-based interruption handlers.

## 10.2 System Register Model

Registers are assigned the following conventions during transitions between IA-32 and Itanium instruction sets.

- **IA-32 State:** The register contains an IA-32 register during IA-32 instruction set execution. Expected IA-32 values should be loaded before switching to the IA-32 instruction set. After completion of IA-32 instructions, these registers contain the results of the execution of IA-32 instructions. These registers may contain any value during Itanium instruction execution according to Itanium software

conventions. Software should follow IA-32 and Itanium software calling conventions for these registers.

- **Shared:** Shared registers contain values that have similar functionality in either instruction set. For example, all Itanium control registers, debug registers are used for memory references (including IA-32). The stack pointer (ESP) and instruction pointer (IP) are also shared.
- **Unmodified:** These registers are not altered by IA-32 execution. Itanium architecture-based code can rely on these values not being modified during IA-32 instruction set execution. The register will have the same contents when entering the IA-32 instruction set and when exiting the IA-32 instruction set.
- **Undefined:** Registers marked as undefined may be used as scratch areas for execution of IA-32 instructions. Software can not rely on the value of these registers across an instruction set transition.

**Table 10-1. IA-32 System Register Mapping**

Intel® Itanium® Reg	IA-32 Reg	Convention	Size	Description
<b>Application Registers</b>				
EFLAG	EFLAG	IA-32 state	32	IA-32 System/Arithmetic flags, writes of some bits are conditioned by PSR.cpl and EFLAG.iopl.
CSD	CSD		64	IA-32 code segment (register format)
SSD	SSD		64	IA-32 stack segment (register format)
CFLG	CR0/CR4		64	IA-32 control flags, CR0=CFLG{31:0}, CR4=CFLG{63:32} <sup>a</sup> , writable at PSR.cpl=0 only.
<b>Kernel Registers</b>				
KR0	IOBASE <sup>b</sup>	IA-32 state	64	IA-32 virtual I/O port Base register
KR1	TSSD <sup>c</sup>			IA-32 TSS descriptor (register format)
KR2	CR3/CR2 <sup>d</sup>			IA-32 CR2=KR2{63:32}, CR3=KR2{31:0}
KR3-7		unmodified		Intel Itanium preserved registers
<b>Banked General Registers</b>				
GR16-31		unmodified		Preserved for operating system use
<b>Control Registers</b>				
DCR		unmodified, shared		Controls instruction set execution (including IA-32)
IFA, IIP, IPSR, ISR, IIM, IIPA, ITIR, IHA, IIB0-1, IFS, IVA		shared	64	Intel Itanium interruption registers may be overwritten on any TLB fault, interruption or exception encountered during IA-32 or Intel Itanium instruction set execution.
PTA		shared	64	Shared page table base for memory references (including IA-32)
ITM		shared		shared Intel Itanium interruption/timer resources

**Table 10-1. IA-32 System Register Mapping (Continued)**

Intel® Itanium® Reg	IA-32 Reg	Convention	Size	Description
LID, IVR, TPR, EOI, IRR0, IRR1, IRR2, IRR3, ITV, PMV, LRR0, LRR1, CMCV		shared	64	Intel Itanium external interrupt control registers are used to generate, prioritize and delivery external interrupts during IA-32 or Intel Itanium instruction set execution.
<b>Translation Resources</b>				
TRs		shared		All Intel Itanium virtual memory registers can be used for memory references (including IA-32).
TCs				
RRs				
PKRs				
Debug Registers				
IBRs	dr0-3, dr7	shared	64	Intel Itanium debug registers are used memory references (including IA-32).
DBRs	dr0-3, dr7			
<b>Performance Monitors</b>				
PMCs		shared	64	Intel Itanium performance monitors measure performance events (including IA-32).
PMDs		shared	64	reflect performance monitor results of execution (including IA-32)

- a. IA-32 MOV from CR0 and CR4 return the value in the CFLG register.
- b. The IOBase register is used by IN/OUT instructions. If IN/OUT operations are disabled via CFLG.io, this register can be used for other values.
- c. The TSSD registers are used by IN/OUT instructions for I/O permission via CFLG.io. If access to the TSS is disabled, these registers can be used for other values.
- d. The Mov from CR2,CR3 instructions return the value contained in KR2.

## 10.3 IA-32 System Segment Registers

System Descriptors are maintained in an unscrambled format shown in Figure 10-1 that differs from the IA-32 scrambled memory descriptor format. The unscrambled register format is designed to support fast conversion of IA-32 segmented 16/32-bit pointers into virtual addresses by Itanium architecture-based code. IA-32 segment register load instructions unscramble the GDT/LDT memory format into the descriptor register format on a segment register load. Itanium architecture-based software can also directly load descriptor registers provided they are properly unscrambled by software. When Itanium architecture-based software loads these registers, no data integrity checks are performed at that time if illegal values are loaded in any fields. For a complete definition of all bit fields and field semantics refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual**.

**Figure 10-1. IA-32 System Segment Register Descriptor Format (LDT, GDT, TSS)**

63	62	60	59	58	57	56	55	52	51	32	31	0
g	ig	p	dpl	s	stype	lim{19:0}			base{31:0}			

**Table 10-2. IA-32 System Segment Register Fields (LDT, GDT, TSS)**

Field	Bits	Description
base	31:0	Segment Base value. This value when zero extended to 64-bits, points to the start of the segment in the 64-bit virtual address space for IA-32 instruction set memory references. This value is ignored for Intel Itanium instruction set memory references.
lim	51:32	Segment Limit. Contains the maximum effective address value within the segment. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details and segment limit fault conditions.
stype	55:52	Segment Type identifier. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for encodings and definition.
s	56	Non System Segment. If 1, a data segment, if 0 a system segment.
dpl	58:57	Descriptor Privilege Level. The DPL is checked for memory access permission for IA-32 instruction set memory references.
p	59	Segment Present bit. If 0, and an IA-32 memory reference uses this segment an IA_Exception(GPFault) is generated.
ig	62:60	Ignored – For the LDT/GDT/TSS descriptors reads of this field return the last value written by Itanium architecture-based code. Reads of this field return zero if written by IA-32 descriptor loads. This field is ignored by the processor during IA-32 instruction set execution. This field may have a future use and should be set to zero by software.
g	63	Segment Limit Granularity. If 1, scales the segment limit by $\text{lim} = (\text{lim} \ll 12)   0\text{xFFF}$ for IA-32 instruction set memory references.

System segment selectors and descriptors for GDT and LDT are maintained in Itanium general registers to support segment register loads used extensively by segmented 16-bit code. On the transition into the IA-32 instruction set, GDT/LDT descriptor table must be initialized if IA-32 code will perform protected mode segment register loads or far control transfers.

Within the IA-32 System Environment, GDT and LDT are considered privileged operating system segmentation resources. However, in the Itanium System Environment, applications can transition between the IA-32 and Itanium instruction set and bypass IA-32 segmentation. Itanium user level instructions can also directly modify all selectors and descriptors including GDT and LDT. An operating system should either protect memory with virtual memory management mechanisms defined by the Itanium architecture or disabled application level instruction set transitions. Within the Itanium System Environment, GDT/LDT memory spaces must be mapped into user space, since supervisor overrides for accesses to GDT/LDT are disabled.

The TSSD descriptor points to the I/O Permission Bitmap. If CFLG.io is 1, IN, INS, OUT, and OUTS consult the TSSD I/O permission bitmap as defined in the *Intel® 64 and IA-32 Architectures Software Developer's Manual*. If CFLG.io is 0, the TSSD I/O permission bitmap is not checked. See [Section 10.7, "I/O Port Space Model"](#) for details on I/O port permission and for TLB-based access control. The TSSD register is not used within the Itanium System Environment to support task switches, or interlevel control transfers. If the TSSD is used for I/O Permissions, Itanium architecture-based operating system software must ensure that a valid 286 or 386 Task State Descriptor is loaded, otherwise IN/OUT operations to the TSSD I/O permission bitmap will result in undefined behavior.

The IDT descriptor is not supported or defined within the Itanium System Environment.

### 10.3.1 IA-32 Current Privilege Level

PSR.cpl is the current privilege level of the processor for instruction execution (including IA-32). PSR.cpl is used by the processor for all IA-32 descriptor segmentation and paging permission checks. PSR.cpl is a secured register. Typical IA-32 processors used SSD.dpl as the official privilege level of the processor. Since, SSD.dpl is not secured from user modification, processor implementations must base all privilege checks and state backups based on PSR.cpl.

### 10.3.2 IA-32 System EFLAG Register

The EFLAG (AR24) register is made of two major components, user arithmetic flags (CF, PF, AF, ZF, SF, OF, and ID) and system control flags (TF, IF, IOPL, NT, RF, VM, AC, VIF, VIP). None of the arithmetic or system flags affect Itanium instruction execution. The arithmetic flags are used by the IA-32 instruction set to reflect the status of IA-32 operations, control IA-32 string operations, and control branch conditions for IA-32 instructions. System flags are typically managed by an operating system and are used to control the overall operations of the processor. System flags are broken into two categories, system flags that control IA-32 instruction set execution behavior and virtualizable system flags. The NT system flag shown in bold font in Figure 10-2 is virtualized.

**Figure 10-2. IA-32 EFLAG Register**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
reserved (set to 0)										id	vip	vif	ac	vm	<b>rf</b>	0	nt	iopl	of	df	if	tf	sf	zf	0	af	0	pf	1	cf	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
reserved (set to 0)																															

System flags AC, TF, RF, VIF, VIP, IOPL and VM directly control the execution of IA-32 instructions. These bits do not control any Itanium instructions. See Table 10-3 for a complete definition these bits.

The NT bit does not directly control the execution of any IA-32 or Itanium instructions. All IA-32 instructions that modify this bit is intercepted (e.g. IRET, Task Switches)

See Table 10-3, "IA-32 EFLAG Field Definition" for the behavior on IA-32 and Itanium instruction reads/writes to this application register.

#### 10.3.2.1 Virtualized Interrupt Flag

To provide for virtualization of IA-32 code, the IF bit is virtualizable in the context of an operating system. Interrupts are enabled for IA-32 instructions, if (PSR.i and (~CFLG.if or EFLAG.if)) is true. For Itanium architecture-based code, interrupts are enabled if PSR.i is 1.

An optional System Flag intercept trap can be generated if CFLG.ii is 1, and the IF-flag changes state due to IA-32 code executing CLI, STI, or POPF. See Section 10.3.3.1, "IA-32 Control Registers" on page 2:246 for CFLG details. Using this model, virtualization code can set CFLG.if to 0 and CFLG.ii to 0, IA-32 instruction set modifications of EFLAG.if does not affect actual interrupt masking, therefore no notification events need be sent to virtualizing software. When virtualization code, detects and queues an external interrupt for delivery into a virtualized IA-32 operating

system/application, it can set CFLG.ii to 1 to force notification the next time the IF-bit changes state, indicating IA-32 code is either opening or closing the interrupt window. Setting CFLG.if to 1, allows for direct IA-32 control of interrupt masking.

Virtualization of the IF flag is independent of VME extensions. Both mechanisms can be used independently, see the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for the complete VME definition.

**Table 10-3. IA-32 EFLAG Field Definition**

EFLAG <sup>a</sup>	Bits	Description
EFLAG.cf	0	IA-32 Carry Flag. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
	1	Ignored – For IA-32 instructions, writes are ignored, reads return one. For Itanium instructions, the implementation can either ignore writes and return one on reads; or write the value, and return the last value written on reads.
	3,5,15	Ignored – For IA-32 instructions, writes are ignored, reads return zero. For Itanium instructions, the implementation can either ignore writes and return zero on reads, or write the value and return the last value written on reads.
EFLAG.pf	2	IA-32 Parity Flag. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
EFLAG.af	4	IA-32 Aux Flag. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
EFLAG.zf	6	IA-32 Zero Flag. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
EFLAG.sf	7	IA-32 Sign Flag. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
EFLAG.tf	8	IA-32 Trap Flag- In the Intel Itanium System Environment, IA-32 instruction single stepping is enabled when EFLAG.tf is 1 or PSR.ss is 1. EFLAG.tf does not control single stepping for Intel Itanium instruction set execution. When single stepping is enabled, the processor generates a IA_32_Exception(Debug) trap event after the successful execution of each IA-32 instruction. If EFLAG.tf is modified by the POPF or POPFD instruction an IA_32_Intercept(SystemFlag) trap is raised. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details on this bit.
EFLAG.if	9	IA-32 Interruption Flag. In the Intel Itanium System Environment, when PSR.i and (~CFLG.if or EFLAG.if) is 1, external interrupts are enabled during IA-32 instruction set execution, otherwise external interrupts are held pending. If CFLG.if is 1, modification of the EFLAG.if directly affects external interrupt enabling. If CFLG.if is 0, EFLAG.if does not affect interrupt enabling. The IF-bit does not affect external interrupt enabling for Intel Itanium instructions nor NMI interrupts. The IF bit can be modified by IA-32 and Itanium architecture-based code only when PSR.cpl is less than or equal to EFLAG.iopl. If PSR.cpl is greater than EFLAG.iopl, writes to the IF-bit are silently ignored. If CFLG.ii is 1, successful modification of the IF-bit by CLI, STI, or POPF results in an IA_32_Intercept(SystemFlag) trap, otherwise the IF-bit is modified without interception. Modification of this bit by Intel Itanium instructions does not result in an intercept. See the <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details on this bit.
EFLAG.df	10	IA-32 Direction Flag. See <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.
EFLAG.of	11	IA-32 Overflow Flag. See <b>Intel® 64 and IA-32 Architectures Software Developer’s Manual</b> for details.



**Table 10-3. IA-32 EFLAG Field Definition (Continued)**

EFLAG <sup>a</sup>	Bits	Description
EFLAG.iopl	13:12	IA-32 In/Out Privilege Level, controls accessibility by IA-32 IN/OUT instructions to the I/O port space and permission to modify the IF-bit for Intel Itanium and IA-32 instructions. If PSR.cpl > IOPL, permission is denied for IA-32 IN/OUT instructions, and modifications of EFLAG.if by either IA-32 or Intel Itanium instructions are ignored. IOPL can only be modified by IA-32 or Intel Itanium instructions executing at privilege level 0, otherwise modifications of this bit are silently ignored. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel<sup>®</sup> 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
EFLAG.nt	14	IA-32 Nested Task switch. In the IA-32 System Environment, indicates a nested task flag when chaining interrupted and called IA-32 tasks. IA-32 task switches are not directly supported in the Intel Itanium System Environment, since IRET, interruptions, calls, and jumps through task gates are always intercepted. EFLAG.nt can be modified by the POPF or POPFD instruction in both system environments. Modification of EFLAG.nt by POPF and POPFD does not result in a System Flag Intercept. See the <i>Intel<sup>®</sup> 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
EFLAG.rf	16	IA-32 Resume Flag. In the Intel Itanium System Environment, when EFLAG.rf or PSR.id is 1, code breakpoint faults are temporarily disabled for one IA-32 instruction, so that IA-32 instructions can be restarted after a code breakpoint fault without causing another code breakpoint fault. EFLAG.rf does not affect Intel Itanium Instruction Debug faults. After the successful execution of each IA-32 instruction, PSR.id and EFLAG.rf are cleared to zero. On entry into the IA-32 instruction set via <code>rfi</code> or <code>br.ia</code> , EFLAG.rf and PSR.id is not cleared until the successful completion of the first (target) IA-32 instruction. <code>jmpbe</code> clears the PSR.id and the EFLAG.rf bit. EFLAG.rf is set to 1 if a repeat string sequence (REP MOVSB, SCANS, CMPSB, LODSB, STOSB, INSB, OUTSB) takes an external interrupt, trap or fault before the final iteration. EFLAG.rf and PSR.id are set to 0 after the last iteration. For all other cases, external interrupts, faults, traps, and intercept conditions EFLAG.rf is unmodified. The RF-bit can be modified by Intel Itanium instructions running at any privilege level. IA-32 instructions cannot directly modify the RF-bit or PSR.id. Specifically, POPF cannot modify the RF-bit and execution of IRET is always intercepted in the Intel Itanium System Environment. See the <i>Intel<sup>®</sup> 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
EFLAG.vm	17	IA-32 Virtual Mode 86. When 1, IA-32 instructions execute in the VM86 environment. This bit can only be modified by IA-32 or Intel Itanium instructions executing at privilege ring 0, otherwise modifications of this bit by Intel Itanium or IA-32 instructions is silently ignored. Itanium architecture-based software is responsible for initializing the processor with the required VM86 register state before transferring to IA-32 VM86 environment. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel<sup>®</sup> 64 and IA-32 Architectures Software Developer's Manual</i> for complete details of the VM86 environment. Software must ensure the processor is in IA-32 Protected Mode when setting the VM bit.
EFLAG.ac	18	IA-32 Alignment Check. In the Intel Itanium System Environment, IA-32 instructions raise an IA_32_Exception(AlignmentCheck) fault if an unaligned reference is performed and PSR.ac is 1 or (CFLG.am is 1 and EFLAG.ac is 1 and memory is accessed at an effective privilege level of 3). Neither EFLAG.ac, CR0.am nor privilege level affect alignment check faults for Intel Itanium instructions. See <a href="#">Section 10.6.7, "Memory Alignment" on page 2:263</a> for details on alignment conditions. This bit can be modified by IA-32 and Intel Itanium instructions at any privilege level. Modification of this bit by the POPF instructions results in an IA_32_Intercept(SystemFlag) trap. See the <i>Intel<sup>®</sup> 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.

**Table 10-3. IA-32 EFLAG Field Definition (Continued)**

EFLAG <sup>a</sup>	Bits	Description
EFLAG.vif	19	IA-32 Virtual Interrupt Flag. See VME extensions in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details. Affects execution of POPF, PUSHF, CLI and STI. This bit is supported in both the IA-32 and Intel Itanium System Environments. A IA-32 Code Fetch fault (GPFault(0)) is generated on every IA-32 instruction (including the target of <i>rfi</i> and <i>br.ia</i> ), if the following condition is true: EFLAG.vip & EFLAG.vif & CFLG.pe & PSR.cpl=3 & (CFLG.pvi   (EFLAG.vm & CFLG.vme))
EFLAG.vip	20	IA-32 Virtual Interrupt Pending. See VME extensions in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for programming details. Affects execution of POPF, PUSHF, CLI and STI. This bit is supported in both the IA-32 and Intel Itanium System Environments.
EFLAG.id	21	IA-32 Identifier bit, can be written and read by IA-32 instructions, indicates IA-32 CPUID instruction is supported. This bit is supported in both the IA-32 and Intel Itanium System Environments.
	63:22	This field is reserved for IA-32 instructions – reads return zeros and non-zero writes causes IA_32_Exception (General Protection) faults. For Itanium instructions, the implementation can either raise Reserved Register/Field fault on non-zero writes and return zero on reads, or write the value (no Reserved Register/Field fault), and return the last value written on reads.

a. On entry into the IA-32 instruction set all bits may be read by subsequent IA-32 instructions, after exit from the IA-32 instruction set these bits represent the results of all prior IA-32 instructions. None of the EFLAG bits alter the behavior of Itanium instruction set execution.

### 10.3.3 IA-32 System Registers

IA-32 system registers such as CR3, CR2, debug registers, performance counters. IA-32 control registers do not affect execution of Itanium instructions. All IA-32 privileged instructions that access prior IA-32 system registers are intercepted.

#### 10.3.3.1 IA-32 Control Registers

IA-32 control registers CR0 and CR4 are mapped into the Itanium application register CFLG (AR27). IA-32 control bits, shown in [Figure 10-3](#), only control execution of the IA-32 instruction set. Additional CR0 bits are defined in CFLG to control virtualization of IA-32 code (namely the IO and IF bits) as shown in [Figure 10-3](#). CFLG is readable by Itanium architecture-based code at all privilege levels but can only be written at privilege level 0, otherwise a Privileged Register fault is generated. When Itanium architecture-based software loads this application register (AR24), a Reserved Register/Field fault will be raised if a non-zero value is written into bits listed as reserved.

**Figure 10-3. Control Flag Register (CFLG, AR27)**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0															
PG	CD	NW	ignored (set to 0)										AM	ig	WP	ignored (set to 0)										II	IF	IO	NE	ET	TS	EM	MP	PE												
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32															
reserved (set to 0)																						MM	EX	FX	SR	PCE	PG	EM	MCE	PAE	PSE	DE	TS	DP	PVI	VME										

- State in italics is virtualized. This state has no impact on a IA-32 or Itanium instruction set execution.
- State in bold only affects IA-32 instruction set execution, Itanium instruction execution is not affected.

Table 10-4 defines all IA-32 control register state and the behavior of each bit in the Itanium System Environment.

**Table 10-4. IA-32 Control Register Field Definition**

Field	Intel® Itanium® State	Bits	Description
CR0	CFLG{31:0}		CR0: IA-32 Mov to CR0 result in a instruction interception fault. Mov from CR0 returns the value contained in CFLG{31:0}. Modification of CFLG{31:0} by Intel Itanium instructions only alters the CR0 state, no side effects (such as TLB flushes) occur.
CR0.PE	CFLG.pe	0	Protected Mode Enable: This bit determines whether the processor operates in IA-32 Protected Mode or Real Mode. This bit affects only IA-32 instruction set execution, Intel Itanium operations are not affected by this bit. Modification of this bit by Itanium architecture-based code does have NOT any side effects such as flushing the TLBs. This bit is supported in both the IA-32 and Intel Itanium System Environments. See <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit and the Protected Mode environment.
CR0.MP	CFLG.mp	1	Monitor co-Processor: When CFLG.ts is 1 and CFLG.mp is 1, execution of IA-32 FWAIT/WAIT instructions results in an Device Not Available fault. This bit is ignored by Intel Itanium floating-point instructions. This bit is supported in both IA-32 and Intel Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
CR0.EM	CFLG.em	2	Emulation: When CFLG.em is set, execution of IA-32 ESC and floating-point instructions generates an IA_32_Exception(DNA) fault. When CFLG.em is 1, execution of IA-32 MMX technology or SSE instructions results in an IA_32_Intercept (Instruction) fault. This bit does not affect Intel Itanium floating-point instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
CR0.TS	CFLG.ts	3	Task Switched: When CFLG.ts is 1, execution of an IA-32 ESC, floating-point instruction, MMX technology or SSE instruction results in a IA_32_Exception(DNA) fault. When CFLG.ts is 1 and CFLG.mp is 1, execution of IA-32 FWAIT/WAIT instructions results in an IA_32_Exception(DNA) fault. This bit is ignored by Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
CR0.ET	CFLG.et	4	Extension Type: ET is ignored since i387 co-processor instructions are supported. This bit is reserved on all Pentium processors. Reads always return 1. This bit is supported in both the IA-32 and Intel Itanium System Environments.

**Table 10-4. IA-32 Control Register Field Definition (Continued)**

Field	Intel® Itanium® State	Bits	Description
CR0.NE	CFLG.ne	5	Numeric Error: Numeric errors are always enabled in the Intel Itanium System Environment. The NE bit and the IGNNE# pin are ignored by the processor and the FERR# pin is not asserted for any numeric errors on IA-32 or Intel Itanium floating-point instructions. In the IA-32 System Environment, this bit is supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> .
--	CFLG.io	6	I/O Enable: If CFLG.io is 1 and CPL>IOPL, IA-32 IN, INS, OUT, OUTS instructions consulted the TSS for I/O permission. If CFLG.io is 0 or CPL<=IOPL, permission is granted regardless of the state of the TSS I/O permission bitmap (the bitmap is not referenced). This bit always returns zero when read by the IA-32 Mov from CR0 instruction. This bit is not defined in the IA-32 System Environment.
--	CFLG.if	7	IF Enable: When CFLG.if is 1, EFLAG.if can be used to enable or disable external interrupts for IA-32 instructions. If CFLG.if is 0, EFLAG.if does not control external interrupt enabling. External interrupts are enabled for the IA-32 instruction set by if PSR.i and (~CFLG.if or EFLAG.if). This bit always returns zero when read by the IA-32 Mov from CR0 instruction. This bit is not defined in the IA-32 System Environment.
--	CFLG.ii	8	IF Intercept: When CFLG.ii is 1, successful modification of the EFLAG.if bit by IA-32 CLI, STI or POPF instructions result in a IA_32_Interrupt(SystemFlag) trap. This bit always returns zero when read by the IA-32 Mov from CR0 instruction. This bit is not defined in the IA-32 System Environment.
ignored		9:15, 17, 19:28	Ignored – This field is ignored by the processor during IA-32 instruction set execution. This field may have a future use and should be set to zero by IA-32 software. For Itanium instructions, the implementation can either ignore the writes and return zero on reads, or write the value and return the last value written on reads.
CR0.WP	CFLG.wp	16	Write Protect: This bit is ignored in the Itanium System Environment. In the IA-32 System Environment, WP controls supervisor write-protection for IA-32 paging. See <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.

**Table 10-4. IA-32 Control Register Field Definition (Continued)**

Field	Intel® Itanium® State	Bits	Description
CR0.AM	CFLG.am	18	Alignment Mask: For IA-32 instructions an IA_32_Exception(AlignmentCheck) fault is generated on a reference to an unaligned data memory operand if PSR.ac is 1 or (CFLG.am is 1 and EFLAG.ac is 1 and memory is accessed at an effective privilege level of 3). Neither EFLAG.ac, CR0.am nor privilege level affect alignment check faults for Itanium instructions. This bit is supported in both the IA-32 and Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on this bit.
CR0.NW	CFLG.nw	29	Not Write-through and Cache Disable: These bits are ignored in the Itanium System Environment. Cacheability is controlled virtual memory attributes. These bits are provided as storage for compatibility purposes.
CR0.CD	CFLG.cd	30	
CR0.PG	CFLG.pg	31	Paging Enable: In the Intel Itanium System Environment, this bit is ignored for IA-32 and Intel Itanium memory references. Virtual translations are enabled via PSR.it and PSR.dt. This bit is provided as storage for compatibility purposes. Modification of this bit by Itanium architecture-based code does NOT have any side effects such as flushing the TLBs. This bit is supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for the IA-32 System Environment.
CR2	KR2{63:32}		IA-32 Page Fault Virtual Address: IA-32 Mov to CR2 result in an interception fault. Mov from CR2 returns the value contained in KR2{63:32}. CR2 is replaced by IFA in the Intel Itanium System Environment.
CR3	KR2{31:0}		IA-32 Page Table Address: IA-32 Mov to CR3 result in an interception fault. Mov from CR3 return the value contained in KR2{31:0}. CR3 is replaced by PTA in the Intel Itanium System Environment. Modification of KR2{31:0} by Itanium architecture-based code does NOT have the side effect of flushing the TLBs.
CR3.PWT	KR4.pwt		Page Write-Through and Cache Disabled: In the Intel Itanium System Environment, these bits are ignored. This bit are provided as storage for compatibility purposes. These bits are supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for the IA-32 System Environment.
CR3.PCD	KR4.pcd		
CR4	CFLG{63:32}		CR4: A-32 Mov to CR4 result in an instruction interception fault. Mov from CR4 returns the value contained in CFLG{63:32}. Modification of CFLG{63:32} by Intel Itanium instructions only alters the register state, no side effects (such as TLB flushes) occur.

**Table 10-4. IA-32 Control Register Field Definition (Continued)**

Field	Intel® Itanium® State	Bits	Description
CR4.VME	CFLG.vme	32	IA-32 Virtual Machine Extension and Protected Mode Virtual Interrupt Enable: These bits control the VM86 VME extensions and Protected Mode Virtual Interrupt extensions defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for STI, CLI and PUSHF. These bits have no effect on Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments.
CR4.PVI	CFLG.pvi	33	
CR4.TSD	CFLG.tsd	34	Time Stamp Disable: IA-32 RDTSC user level reads of the Time Stamp Counter are enabled when CR4.tsd when zero. Otherwise execution of the RDTSC instruction results in a GPFault. CFLG.tsd is ignored by Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on these bits.
CR4.DE	CFLG.de	25	Debug Extensions: In the Intel Itanium System Environment, this bit is ignored by IA-32 or Intel Itanium references to the I/O port space. This bit is provided as storage for compatibility purposes. This bit is supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for the IA-32 System Environment.
CR4.PSE	CFLG.pse	36	Page Size Extensions: In the Intel Itanium System Environment, this bit is ignored by IA-32 or Intel Itanium references. In the IA-32 System Environment, this bit enables 4M-byte page extensions for IA-32 paging. Modification of this bit by Itanium architecture-based code does have any side effects such as flushing the TLBs.
CR4.PAE	CFLG.pae	37	Physical Address Extensions: In the IA-32 System Environment, this bit enables IA-32 Physical Address Extensions for IA-32 paging This bit is ignored in the Intel Itanium System Environment. Modification of this bit by Itanium architecture-based code does have any side effects such as flushing the TLBs.
CR4.MCE	CFLG.mce	38	Machine Check Enable: This bit is ignored in the Intel Itanium System Environment. This bit is provided as storage for compatibility purposes. This bit is supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for the IA-32 System Environment.
CR4.PGE	CFLG.pge	39	Paging Global Enable: This bit is ignored in the Intel Itanium System Environment. This bit is provided as storage for compatibility purposes. This bit is supported as defined in the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for the IA-32 System Environment, where this bit enables global pages for the IA-32 paging. Modification of this bit by Itanium architecture-based code does have any side effects such as flushing the TLBs.

**Table 10-4. IA-32 Control Register Field Definition (Continued)**

Field	Intel® Itanium® State	Bits	Description
CR4.PCE	CFLG.pce	40	Performance Counter Enable: IA-32 RDPMC user level reads of the performance counters are enabled when CR4.pce is 1. Otherwise execution of the RDPMC instruction results in a GPFault. CFLG.pce is ignored by Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on these bits.
CR4.FXSR	CFLG.FXSR	41	SSE FXSR Enable. When 1, enables the SSE register context. When 0, execution of all SSE instructions results in an IA_32_Interrupt(Instruction) fault. This bit does not control the behavior of Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on these bits.
CR4.MMXEX	CFLG.MMXEX	42	SSE Exception Enable: When 1, enables SSE unmasked exceptions. When 0, all SSE Exceptions are masked. This bit does not control the behavior of Intel Itanium instructions. This bit is supported in both the IA-32 and Intel Itanium System Environments. See the <i>Intel® 64 and IA-32 Architectures Software Developer's Manual</i> for details on these bits.
reserved		43:63	This field is reserved for IA-32 instructions – reads return zeros and non-zero writes causes IA_32_Exception (General Protection) faults. For Itanium instructions, the implementation can either raise Reserved Register/Field fault on non-zero writes and return zero on reads, or write the value (no Reserved Register/Field fault) and return the last value written on reads.

### 10.3.3.2 IA-32 Debug Registers

Within the Itanium System Environment, the IA-32 debug registers (DR0 - DR7) are superseded by the Itanium debug registers DBR0-7 and IBR0-7, see [Section 10.8.1, "Data Breakpoint Register Matching" on page 2:274](#) for details. Accesses to the IA-32 debug registers result in an interception fault.

The Itanium debug registers are designed to facilitate debugging of both IA-32 and Itanium architecture-based code. Specifically, instruction and data breakpoints can be programmed by loading 64-bit virtual addresses into IBR and DBR along with an address mask. Itanium defined single stepping mechanisms, and taken branch traps are also defined to trap on IA-32 instructions. See [Section 10.8.1, "Data Breakpoint Register Matching" on page 2:274](#) for details on IA-32 instruction set behavior with respect to the debug facilities defined by the Itanium architecture.



### **10.3.3.3 IA-32 Memory Type Range Registers (MTRRs)**

Within the Itanium System Environment, IA-32 MTRR registers are superseded by physical memory attributes supplied by the TLB, as defined in [Section 4.4.3, “Cacheability and Coherency Attribute”](#) on page 2:77. IA-32 instruction references to the MTRRs in the MSR register space results in an instruction intercept fault.

### **10.3.3.4 IA-32 Model Specific and Test Registers**

Within the Itanium System Environment, the IA-32 Model Specific Register space (MSRs) are superseded by the PAL firmware interface. Cache testing, initialization, processor configuration should be performed through the PAL interface. See [Section 11.10, “PAL Procedures”](#) on page 2:353 for a complete definition of the PAL functions and interfaces. Accesses to the IA-32 Model Specific Register space result in an instruction interception fault.

### **10.3.3.5 IA-32 Performance Monitor Registers**

Within the Itanium System Environment, the Itanium performance monitors are designed to measure IA-32 and Itanium instructions, and system performance through a unified performance monitoring facility. Itanium architecture-based code can program the performance monitors for IA-32 and/or Itanium events by configuring the PMC registers. Count values are accumulated in the PMD registers for both IA-32 and Itanium events. See implementation-specific documentation for the list of supported events and encodings.

IA-32 code can sample the performance counters by issuing the RDPMC instruction. RDPMC returns count values from the specified Itanium performance monitor. Operating systems can secure the monitors from being read by IA-32 code by setting PSR.sp to 1, or setting CR4.pce to 0, or setting the performance monitor’s pm-bit. Reads of a secured counter by RDPMC return a IA\_32\_Exception(GPFault(0)). IA-32 code cannot write or configure the performance monitors, all writes to the MSR register space are intercepted.

### **10.3.3.6 IA-32 Machine Check Registers**

Within the Itanium System Environment, IA-32 machine check registers are superseded by the Itanium machine check architecture. See [Section 11.3, “Machine Checks”](#) on page 2:296 for details. IA-32 accesses to the Pentium III Processor machine check registers results in an instruction intercept.

## **10.4 Register Context Switch Guidelines for IA-32 Code**

The following section gives operating system performance guidelines to minimize the amount of register context that must be saved and restored for IA-32 processes during a context switch.



## 10.4.1 Entering IA-32 Processes

High FP registers (FR32-127) – The processor requires access to all high FP registers during the execution of IA-32 instructions. It is recommended on entering an IA-32 process, that the OS save the high FP registers belonging to a prior context and then **enable** the high FP registers (PSR.dfh is 0). Otherwise, the processor will immediately raise a Disabled FP Register fault on the first IA-32 instruction executed in the IA-32 process. Performing the state save of the prior high FP register context during the context switch avoids the unnecessary generation of the Disabled FP Register fault.

Low FP registers (FR2-31) – The processor does not require access to the low FP registers unless executing IA-32 FP, MMX technology or SSE instructions. It is recommended on entry to an IA-32 process, that the OS **disable** the low FP registers by setting PSR.dfl to 1. PSR.dfl set to 0 indicates that there was a possibility that IA-32 FP, MMX technology or SSE instruction could execute and write FR8-31. If the low FP registers are enabled on entry to an IA-32 process (PSR.dfl is 0), all low FP registers will appear to be dirty on IA-32 process exit.

High Integer Registers (GR32-127) – Since the processor leaves all high registers in the register stack in an undefined state, these registers must be saved by the RSE before entering an IA-32 process.

Low Integer registers (GR1-31) – These registers must be explicitly saved before entering an IA-32 process.

## 10.4.2 Exiting IA-32 Processes

High FP registers (FR32-127) – PSR.mfh is unmodified when leaving the IA-32 instruction set. IA-32 instruction set execution leaves FR32-127 in an undefined state. Software can not rely on register values being preserved across an instruction set transition. These registers do NOT need to be preserved across a context switch.

Low FP registers (FR2-31) – PSR.mfl indicates there is a possibility that FR8-31 were modified by IA-32 FP, MMX technology, or SSE instruction. The modify bit is set by the processor when leaving the IA-32 instruction set, if PSR.dfl is 0, otherwise PSR.mfl is unmodified. During the state save of the outbound IA-32 process, it is recommended that the OS save FR2-31 if and only if the lower FP registers are marked as modified.

High Integer Registers (GR32-127) – Since the processor leaves all high registers undefined across an instruction set transition, these registers do NOT need to be preserved across an IA-32 context switch.

Low Integer registers (GR1-31) – These registers must be explicitly preserved across a context switch.

## 10.5 IA-32 Instruction Set Behavior Summary

Table 10-5 summarizes IA-32 instruction behavior within the Itanium System Environment. All IA-32 instructions are unchanged from the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** except where noted. IA-32 instructions can also generate additional Itanium register and memory faults as defined in

**Table 5-6.** Please refer to the **Intel® 64 and IA-32 Architectures Software Developer’s Manual** for the behavior of all IA-32 instructions in the IA-32 System Environment.

For all listed and unlisted IA-32 instructions in **Table 10-5** the following relationships hold:

- Writes of any IA-32 general purpose, floating-point or MMX technology or SSE registers by IA-32 instructions are reflected in the Itanium registers defined to hold that IA-32 state when the IA-32 instruction set completes execution.
- Reads of any IA-32 general purpose, floating-point or MMX technology or SSE registers by IA-32 instructions see the state of the Itanium registers defined to hold the IA-32 state after entering the IA-32 instruction set.
- IA-32 numeric instructions are controlled by and reflect their status in FCW, FSW, FTW, FCS, FIP, FOP, FDS and FEA. On exit from the IA-32 instruction set, Itanium registers defined to hold IA-32 state reflect the results of all IA-32 prior numeric instructions (FSR, FCR, FIR, FDR). Itanium numeric status and control resources defined to hold IA-32 state are honored by IA-32 numeric instructions when entering the IA-32 instruction set.

In **Table 10-5** *unchanged* indicates there is no change in behavior with respect to the IA-32 System Environment.

**Table 10-5. IA-32 Instruction Summary**

IA-32 Instruction	Intel® Itanium® System Environment	Comments
AAA, AAD, AAM, AAS	unchanged	
ADC, ADD, AND,		
ADDPS, ADDSS, ANDNPS, ANDPS		
ARPL		
BOUND		
BSF, BSR		
BSWAP		
BT, BTC, BTS, BTR		
CALL	near: no change far: no change gate more privilege: Gate Intercept gate same privilege: Gate Intercept task gate: Gate Intercept + additional taken branch trap	Intercept if through a call or task gate  If PSR.tb is 1, raise a taken branch trap.
CBW, CWDE, CDQ	unchanged	
CLC, CLD		
CLI	Optional System Flag Intercept	Intercept if EFLAG.if changes state and CFLG.ii is 1
CLTS	Instruction Intercept	IA-32 privileged instruction
CMC	unchanged	
CMOV		
CMP		
CMPPS, CMPSS, COMISS		
CMPS		

**Table 10-5. IA-32 Instruction Summary (Continued)**

IA-32 Instruction	Intel® Itanium® System Environment	Comments
CMPXCHG, 8B	Optional Lock Intercept	If Locks are disabled (DCR.Lc is 1) and a processor external lock transaction is required
CPUID	unchanged	
CWD, CDQ		
CVTPI2PS, CVTPS2PI, CVTSI2SS, CVTSS2SI, CVTTPS2PI, CVTTSS2SI		
DAA, DAS		
DEC		
DIV		
DIVPS, DIVSS		
ENTER		
EMMS		

**Table 10-5. IA-32 Instruction Summary (Continued)**

IA-32 Instruction	Intel® Itanium® System Environment	Comments
F2XM1	unchanged	IA-32 numeric instructions manipulate the IA-32 numeric register stack contained in f8-f15, status is reflected in FSR. Modification of the IA-32 numeric environment changes FIR, FDR FCR and FSR.
FABS		
FADD, FADDP, FIADD		
FBLD		
FBSTP		
FCHS		
FCLEX, FNCLEX		
FCMOV		
FCOM, FCOMPP		
FCOMI, FCOMIP		
FUCOMI, FUCOMIP		
FCOS		
FDECSTP		
FDIV, FDIVP, FIDIV		
FDIVR, FDIVRP, FDIVR		
FFREE		
FICOM, FICOMP		
FILD		
FINCSTP		
FINIT, FNINIT		
FIST, FISTP		
FLD		
FLD constant		
FLDCW		
FLDENV		
FMUL, FMULP, FIMUL		
FNOP		
FPATAN, FPTAN		
FPREM, FPREM1		
FRNDINT		
FRSTOR		
FSAVE, FNSAVE		
FSCALE		
FSIN, FSINCOS		
FSQRT		
FST, FSTP		
FSTCW, FNSTCW		
FSTENV, FNSTENV		
FSTSW, FNSTSW		
FSUB, FSUBP, FISUB		
FSUBR, FSUBRP, FISUBR		
FTST		
FUCOM, FUCOMP		
FWAIT		
FXAM		
FXCH		
EXTRACT		
FXRSTOR, FXSAVE		
FYL2X, FYL2XP1		

**Table 10-5. IA-32 Instruction Summary (Continued)**

IA-32 Instruction	Intel® Itanium® System Environment	Comments
HLT	Instruction Intercept	IA-32 privileged instruction
IDIV	unchanged	
IMUL		
IN, INS	unchanged + I/O ports are mapped virtually	If CFLG.io is 0, the TSS I/O permission bitmap is not consulted. Intel Itanium TLB faults control accessibility to I/O ports.
INC	unchanged	
INT 3, INTO	Mandatory Exception vector #	Delivered as an IA_32_Interrupt
INT n	Mandatory Interruption vector #	Delivered as an IA_32_Exception
INVD	Instruction Intercept	IA-32 privilege instruction
INVLPG		
IRET, IRETD	Real Mode: Instruction Intercept to VM86: Instruction Intercept from VM86: Instruction Intercept same privilege: Instruction Intercept less privilege: Instruction Intercept different task: Instruction Intercept	All forms of IRET result in an instruction intercept
Jcc	additional taken branch trap	If PSR.tb is 1, raise a taken branch trap.
JMP	near: no change far: no change gate task: Gate Intercept call gate: Gate Intercept additional taken branch trap	Intercept fault if through a call or task gate  If PSR.tb is 1, raise a taken branch trap.
JMPE		Jumps to the Intel Itanium instruction set
LAHF	unchanged	
LAR		
LDMXCSR		
LDS, LES, LFS, LGS, LSS		
LEA		
LEAVE		
LGDT, LLDT		
LIDT	Instruction Intercept	IA-32 privileged register resource
LMSW		
Lock prefix	Optional Lock Intercept	If Locks are disabled (DCR.lc is 1) and a processor external lock transaction is required
LODS	unchanged	
LOOP, LOOPcc	additional taken branch trap	If PSR.tb is 1, raise a taken branch trap.
LSL	unchanged	User level instruction
LTR	Instruction Intercept	IA-32 privileged register
MASKMOVQ	unchanged	
MAXPS, MAXSS, MINPS, MINSS		
MOV		
MOVNTPS, MOVNTQ		

**Table 10-5. IA-32 Instruction Summary (Continued)**

IA-32 Instruction	Intel® Itanium® System Environment	Comments		
MOV from CR	unchanged			
MOV to CR	Instruction Intercept	IA-32 privileged system registers		
MOV to/from DR				
Mov SS	System Flag Intercept Trap	System Flag Intercept Trap after instruction completes		
MOVAPS, MOVHPS, MOVLPS, MOVMSKPS, MOVSS, MOVUPS	unchanged			
MOVD, MOVQ				
MOVS				
MOVSX, MOVZX				
MUL				
MULPS, MULSS				
NEG				
NOP				
NOT				
OR				
ORPS				
OUT, OUTS			unchanged + I/O ports are mapped virtually	If CFLG.io is 0, the TSS I/O permission bitmap is not consulted. Intel Itanium TLB faults control accessibility to I/O ports.
PACKSS, PACKUS			unchanged	
PADD, PADD, PADDUS				
PAND, PANDN				
PCMPEQ, PCMPGT				
PEXTRW, PINSRW				
PMADD				
PMULHW, PMULLW, PMULHUW				
PMOVMASKB				
POP, POPA				
POP SS	System Flag Intercept	System Flag Intercept Trap after instruction completes		
POPF, POPFD	Optional System Flag Intercept	Intercept if EFLAG.if changes state and CFLG.ii is 1 Intercept if EFLAG.ac, or tf change state.		
POR	unchanged			
PREFETCH				
PSHUFW				
PSLL, PSRA, PSRL				
PSUB, PSUBS, PSUBUS				
PUNPCKH, PUNPCKL				
PXOR				
PUSH, PUSA				
PUSHF, PUSHFD	unchanged	Pushes value in EFLAG, no intercept		
RCL, RCR, ROL, ROR				
RCPPS, RSQRTPS				
RDMSR	Instruction Intercept	IA-32 privileged system register space		
RDTSC	Optional GPFault	No longer faults in VM86, GPFault if secured by PSR.si or CFLG.tsd.		
RDPMSR				
REP, REPcc prefix	unchanged			

**Table 10-5. IA-32 Instruction Summary (Continued)**

IA-32 Instruction	Intel® Itanium® System Environment	Comments
RET	near: no change far: no change less privilege: no change same privilege: no change + additional taken branch trap	If PSR.tb is 1, raise a taken branch trap.
RSM	Instruction Intercept	IA-32 privileged instruction
SAHF	unchanged	
SAL, SAR, SHL, SHR		
SBB		
SCAS		
SFENCE		
SETcc		
SGDT, SLDT	Instruction Intercept	IA-32 privileged instruction
SHLD, SHRD	unchanged	
SHUFPS, SQRTPS, SQRTSS		
SIDT	Instruction Intercept	IA-32 privileged instructions
SMSW		
STC, STD	unchanged	
STI	Optional System Flag Intercept	Intercept if EFLAG.if changes state and CFLG.ii is 1
STMXCSR	unchanged	
STOS		
STR	Instruction Intercept	IA-32 privileged instruction
SUB	unchanged	
SUBPS, SUBSS		
SYSENTER, SYSEXIT	Instruction Intercept	
TEST		
UCOMISS	unchanged	
UNPCKHPS, UNPCKLPS		
UD2	Instruction Intercept	Reserved undefined opcodes
VERR, VERW	unchanged	User level instruction
WAIT		
WBINVD	Instruction Intercept	IA-32 privileged instructions
WRMSR		
XADD	Optional Lock Intercept	If Locks are disabled (DCR.Ic is 1) and a processor external lock transaction is required than a Lock Intercept.
XCHG		
XLAT, XLATB	unchanged	
XOR		
XORPS		

## 10.6 System Memory Model

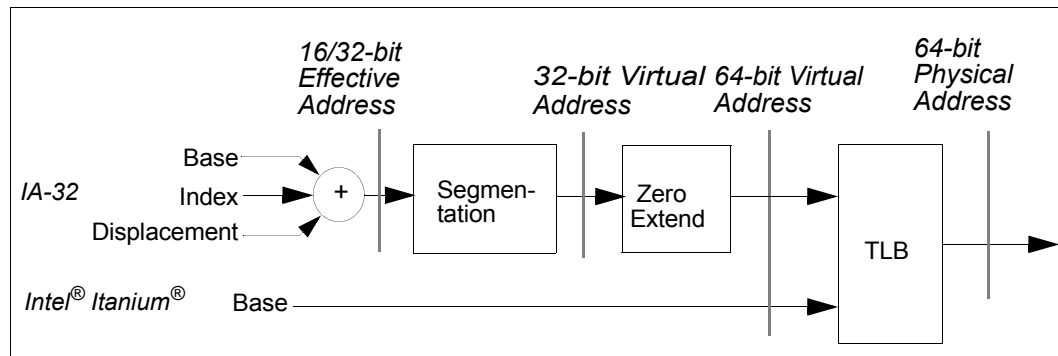
Within the Itanium System Environment, a unified memory model is presented to the programmer. Applications and the operating system see the same 64-bit virtual memory space and virtual addressing mechanisms regardless of the referencing instruction set. A virtual address points to the same physical storage location from both IA-32 and Itanium instruction sets.

Itanium architecture-based operating systems must not use IA-32 segmentation as a protected system resource. An Itanium architecture-based operating system must use virtual memory management defined by the Itanium architecture to secure IA-32 and Itanium architecture-based applications, memory and I/O devices. The Itanium architecture is defined to be *unsegmented architecture and all Itanium memory references bypass IA-32 segmentation and protection checks*. In addition, Itanium architecture-based user level code can directly modify IA-32 segment selector and descriptor values for all segments (including GDT and LDT). If operating systems do not rely on segmentation for protection, there are no security concerns for exposing IA-32 segment registers and descriptors to Itanium architecture-based user level applications

IA-32 instruction and data reference addresses are generated as 16/32-bit effective addresses as shown in [Figure 10-2](#). IA-32 segmentation is then applied to map 32-bit effective addresses into 32-bit virtual addresses, the processor then converts the address into a 64-bit virtual address by zero extension from 32 to 64-bits. Itanium instructions bypass all of these steps and directly generate addresses within the 64-bit virtual address space.

For both IA-32 and Itanium instruction set memory references, virtual memory management defined by the Itanium architecture is used to map a given virtual address into a physical address. Itanium architecture-based virtual memory management hardware does not distinguish between Itanium and IA-32 instruction set generated memory references during the conversion from a virtual to physical address.

**Figure 10-4. Virtual Memory Addressing**



### 10.6.1 Virtual Memory References

In the Itanium System Environment the following virtual memory options are available for supporting IA-32 and Itanium memory references.

- Software TLB fills (TLBs are enabled, but the VHPT is disabled).
- 8-byte short format VHPT, see [Section 4.1.5, "Virtual Hash Page Table \(VHPT\)"](#) on [page 2:61](#) for details.
- 32-byte long format VHPT.

Itanium virtual memory resources can be used by the operating system for all IA-32 memory references. These resources include virtual Region Registers (RR), Protection Key Registers (PKR), the Virtual Hash Page Table (VHPT), all supported range of page sizes, Translation Registers (ITR, DTR), the Translation Cache (ITC, DTC) and the complete set of Itanium virtual memory management faults defined in [Chapter 5](#).



## 10.6.2 IA-32 Virtual Memory References

By definition, IA-32 instruction and data memory references are confined to 32-bits of virtual addressing, the first 4 G-bytes of virtual region 0. However, IA-32 memory references can be mapped anywhere within the implemented physical address space by operating system code.

Virtual addresses are converted into physical addresses through the process defined in [Section 4.1, “Virtual Addressing” on page 2:45](#). IA-32 references use the Itanium TLB resources as follows.

- **Region Identifiers** – Operating systems can place IA-32 processes within virtual region 0, and use the entire  $2^{24}$  region identifier name space. By using region identifiers there is no requirement to flush IA-32 mappings on a context switch.
- **Protection Keys** – Operating systems can place mappings used by IA-32 processes within any number of protection domains. If PSR.pk is 1, all IA-32 references search the Protection Key Registers (PKR) for matching keys. If a key is not found, a Key Miss fault is generated. Otherwise, key read, write, execute permissions are verified.
- **TLB Access Bit** – If this bit is zero, an Access Bit fault is generated during Itanium or IA-32 instruction set memory references. Note: the processor does not automatically set the Access bit in the VHPT on every reference to the page. Access bit updates are controlled by the operating system.
- **TLB Dirty Bit** – If this bit is zero, a Dirty bit fault is generated during any Itanium or IA-32 instruction that stores to a dirty page. Note: the processor does not automatically set the Dirty bit in the VHPT on every write. Dirty bit updates are managed by the operating system.

## 10.6.3 IA-32 TLB Forward Progress Requirements

To ensure forward progress while executing IA-32 instructions, additional TLB resources and replacement policies must be defined over and above the definition given in [Section 4.1.1.2, “Translation Cache \(TC\)” on page 2:49](#). IA-32 instructions and data accesses may not be aligned resulting in a worst case scenario for two possible pages being referenced for every memory datum referenced during the execution of an IA-32 instruction. Furthermore, the worst case non-intercepted IA-32 opcode can reference up to 4 independent data pages.

The Translation Cache’s (TC) are required to have the following minimum set of resources to ensure forward progress. Given that software TLB fills can be used to insert entries into the TLB and a hardware page table walker is not necessarily used, the following requirements must be satisfied by the processor:

- Instruction Translation Cache – at least 1 way set associative with 2 sets, or 2 entries in a fully associative design. Replacement algorithms must not consistently displace the last 2 entries installed by software.
- Data Translation Cache – at least 4 way set associative with 2 sets, or 8 entries in a fully associative design. Replacement algorithms must not consistently displace the last 8 entries installed by software or the last 8 translations referenced by an IA-32 instruction.

- Unified Translation Cache – at least 5 way set associative with 2 sets, or 10 entries in a fully associative design. The processor must not consistently displace the last 10 entries installed or the last 10 translations referenced by an IA-32 instruction.

The processor must ensure that the minimum number of entries can co-exist in the TLB, and TC replacement algorithms allow software insertion of the required entries such that the required number of translations can be co-resident in the TLB.

The processor cannot ensure forward progress unless translations mapping the Itanium architecture-based TLB Miss handlers are statically mapped by the Instruction Translation Registers.

### 10.6.4 Multiprocessor TLB Coherency

Global TLB purges can not occur on another processor unless that processor is at an interruptible point. For IA-32 instruction set execution, interruptible points are defined as; 1) when the processor is between instructions (regardless of the state of PSR.i and EFLAG.if), and 2) each iteration of an IA-32 string instruction, regardless of the state of PSR.i and EFLAG.if

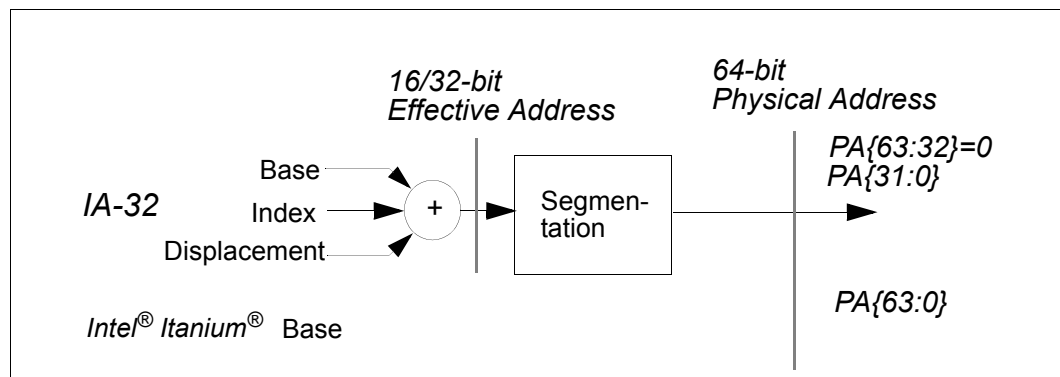
The processor may delay in its response and acknowledgment to a broadcast purge TC transaction until the processor executing an IA-32 instruction has reached a point (e.g. an IA-32 instruction boundary) where it is safe to process the purge TC request. The amount of the delay is implementation specific and can vary depending on the receiving processor and what instructions or operations are executing when it receives the purge request.

### 10.6.5 IA-32 Physical Memory References

When running IA-32 code, virtual addressing must be utilized by setting PSR.dt to 1 and PSR.it to 1, otherwise processor operation is undefined. Undefined behavior can include, but is not limited to: machine check abort on entry to IA-32 code, and unpredictable behavior for IA-32 self modifying code.

Operating systems must ensure PSR.dt and PSR.it are 1 before invoking IA-32 code. From a practical standpoint, the TLBs must be enabled so IA-32 code can access the virtual address space, and access memory areas other than WB (e.g. UC or the I/O port space).

**Figure 10-5. Physical Memory Addressing**



## 10.6.6 Supervisor Accesses

If the processor is operating in the Itanium System Environment, supervisor override is disabled, and LDT, GDT, TSS references are performed at the privilege level specified by PSR.cpl. Unaligned processor references to LDT, GDT, and TSS segments will never generate an EFLAG.ac enabled IA-32 Exception (AlignmentCheck) fault, even if PSR.cpl equals 3 and supervisor override is disabled.

Operating systems must ensure that the GDT/LDT are mapped to pages with user level read/write access.

Write permission is required if GDT, or LDT memory descriptor Access-bits are zero regardless of supervisor override conditions. If all GDT/LDT descriptor Access-bits are one, write permission can be removed. Otherwise, Access Rights, Key Miss or Key Miss faults can be generated during all segment descriptor referencing instructions.

If a fault is generated during a supervisory access, the ISR.so bit indicates that CPL is zero or a supervisor override condition was in effect (reference as made to GDT, LDT or TSS).

## 10.6.7 Memory Alignment

Depending on software conventions, memory structures may have different alignment or padding restrictions for the IA-32 and Itanium instruction sets. IA-32 and Itanium architecture-based software should use aligned memory operands as much as possible to avoid possible severe performance degradation associated with un-aligned values and extra over-head for unaligned data memory fault handlers.

The processor provides full functional support for all cases of un-aligned IA-32 data memory references. If PSR.ac is 1 or EFLAG.ac is 1 and CR0.am is 1 and the effective privilege level is 3, unaligned IA-32 memory references result in an IA-32 Exception (AlignmentCheck) fault. Unaligned processor references to LDT, GDT, and TSS segments will never generate an EFLAG.ac enabled IA-32 Exception (AlignmentCheck) fault, even if the effective privilege level is 3 and supervisor override is disabled.

Alignment conditions for Itanium memory references are not affected by the EFLAG.ac, CFLG.am bits.

If EFLAG.ac and CFLG.am are 1 and the reference is done at privilege level 3, IA-32 instruction set unaligned conditions are; 2-byte references not a 2-byte boundary, 4-byte references not on a 4-byte boundary, 8-byte references not on a 8-byte boundary, and 10-byte references not on a 8-byte boundary.

If PSR.ac is 1, IA-32 instruction set unaligned conditions are; 2-byte references not a 2-byte boundary, 4-byte references not on a 4-byte boundary, 8-byte references not on a 8-byte boundary, and 10-byte references not on a **16**-byte boundary.

The processor exhibits the following behavior when accesses are made to un-aligned data operands that span virtual page boundaries:

- IA-32 instruction set – If either page contains a fault, no memory location is modified. For reads, the destination register is not modified.
- Itanium instruction set – All page crossers result in an unaligned reference fault. Memory contents and register contents are not modified.

## 10.6.8 Atomic Operations

All Itanium load/stores and IA-32 non-locked memory references up to 64-bits that are aligned to their natural data boundaries are atomic.

Both IA-32 and Itanium atomic semaphore operations can be performed on the same shared memory location. The processor ensures IA-32 locked read-modify-write opcodes and Itanium semaphore operations are performed atomically even if the operations are initiated from the other instruction set by the same processors, or between multiple processors in an multiprocessing system.

There are some restrictions placed on Itanium atomic operations that may prevent Itanium architecture-based code from manipulating IA-32 semaphores in some rare cases:

- Unaligned Itanium semaphore operations result in an Unaligned Data Reference fault. Itanium architecture-based code manipulation of an IA-32 semaphore can only be performed if the IA-32 semaphore is aligned.
- Itanium semaphore operations to memory which is neither write-back cacheable nor a NaTPage result in an Unsupported Data Reference fault (regardless of the state of the DCR.lc). Itanium architecture-based code manipulation of an IA-32 semaphore can only be performed if the IA-32 semaphore is allocated in aligned write-back cacheable memory.

If an IA-32 locked atomic operation is defined as requiring a read-modify-write operation external to the processor under external bus lock and if DCR.lc is set to 1, an IA\_32\_Intercept(Lock) fault is generated. (IA-32 atomic memory references are defined to require an external bus lock for atomicity when the memory transaction is made to non-write-back memory or are unaligned across an implementation-specific non-supported alignment boundary.) If DCR.lc is set to 0, the processor may either execute the transaction as a series of non-atomic transactions or perform the transaction with an external bus lock, depending on the processor implementation. For processor implementations that do support external bus locks, software must ensure that the Bus Lock Mask bit is set to one, in order to ensure atomicity of these IA-32 operations when DCR.lc=0. The Bus Lock Mask bit is a feature controllable by the PAL\_BUS\_SET\_FEATURES procedure. (See [Table 11-63 on page 2:368](#) for more information).

If the processor supports external bus locks, unaligned IA-32 atomic references are supported, but their usage is strongly discouraged since they are typically performed outside the processor's cache which can severely degrade performance of the system. IA-32 external bus locks are not supported on all processor implementations.

For IA-32 semaphores, atomicity to uncached memory areas (UC) is platform specific, atomicity can only be ensured by the platform design and can not be enforced by the processor.

## 10.6.9 Multiprocessor Instruction Cache Coherency

The processor and platform ensure the processor's instruction cache is coherent for the following conditions:

- For all processors in the coherence domain, local and remote instruction cache coherency on all processors is enforced for any store generated by any processor running the IA-32 instruction set.
- For all processors in the coherence domain, instruction cache coherency on all processors is enforced for all coherent I/O traffic. (For non-coherent I/O, a processor may or may not see the results of an I/O operation.)
- For all processors in the coherence domain, instruction cache coherency is not enforced for stores generated by any processor running the Itanium instruction set. To ensure instruction cache coherency, Itanium architecture-based code must use the code sequence defined in [Section 4.4.6.2, “Memory Consistency”](#) on page 1:72.

**Table 10-6. Instruction Cache Coherency Rules**

Originating Instruction Set	Local Processor	External Processor	Coherent, I/O	Non-Coherent I/O
IA-32	Coherent	Coherent	Coherent	Maybe
Intel Itanium	May be Non-coherent	May be Non-coherent		Non-Coherent

### 10.6.10 IA-32 Memory Ordering

IA-32 memory ordering follows the Pentium III defined **processor ordered** model for cacheable and uncacheable memory. IA-32 *processor ordered* memory references are mapped onto the Itanium memory ordering model as follows:

- All IA-32 stores have **release** semantics. Except for IA-32 stores to write-coalescing memory that are unordered. Subsequent loads are allowed to bypass buffered local store data before it is globally visible. The amount of store buffering is model specific and can vary across processor generations.
- All IA-32 loads have **acquire** semantics. Some high performance processor implementations may speculatively issue *acquire* loads into the memory system for speculative memory types, if and only if the loads do not *appear* to pass other loads as observed by the program. If there is a coherency action that would result in the appearance to the program of a load bypassing other load, the processor will reissue the load operation(s) in program order.
- All IA-32 read-modify-write or locked instructions have memory **fence** semantics. All buffered stores are flushed.
- IA-32 IN, OUT and serializing operations (as defined in the *Intel® 64 and IA-32 Architectures Software Developer’s Manual*) have memory **fence** semantics. In addition, the processor will wait for completion (acceptance by the platform) of the IN or OUT before executing the next instruction. All buffered stores are flushed before the IN or OUT operation.
- IA-32 SFENCE has **release** semantics and will flush all buffered stores.

**Table 10-7. IA-32 Load/Store Sequentiality and Ordering**

IA-32 Memory Reference	Uncacheable	Write Coalescing	Cacheable
store	sequential release <sup>a</sup>	non-sequential unordered	non-sequential release <sup>b</sup>
load	sequential acquire <sup>a</sup>	non-sequential unordered	non-sequential acquire <sup>b</sup>

**Table 10-7. IA-32 Load/Store Sequentiality and Ordering (Continued)**

IA-32 Memory Reference	Uncacheable	Write Coalescing	Cacheable
locked or read-modify-write operation	sequential fence flush prior stores	non-sequential fence flush prior stores	non-sequential fence flush prior stores
IN, INS, OUT, OUTS	sequential fence flush prior stores	undefined	undefined
IA-32 Serialization	fence, flush prior stores		
SFENCE	release, flush prior stores		

- a. However, IA-32 loads/stores to uncacheable memory flush the write coalescing buffers.
- b. However, IA-32 load/stores to cacheable memory do not flush the write coalescing buffers.

Per [Table 10-7](#), IA-32 memory references can be expressed in terms of acquire, release, fence and sequential ordering rules defined by the Itanium architecture. IA-32 data memory references follow the same ordering relationships as defined for Itanium architecture-based code as defined in [Section 4.4.7, “Sequentiality Attribute and Ordering” on page 2:82](#). The following additional clarifications need to be made for IA-32 instruction set execution:

- IA-32 loads and instruction fetches to speculative memory can occur randomly. Read accesses to speculative memory can occur at arbitrary times even if the in-order execution of the program does not require a read or instruction fetch from a given memory location.
- IA-32 instruction fetches and loads to non-speculative memory occur in program order. IA-32 instruction cache line fetch accesses to uncached memory occur in the order specified by an in-order execution of the program. Note however that the same cache line may be fetched multiple times by the processor as multiple instructions within the cache line are executed. The size of a cache line and number of instruction fetches is model specific.
- IA-32 instruction fetches are not perceived as passing prior IA-32 stores. IA-32 stores into the IA-32 instruction stream are observed by the processor’s self modifying code logic. Speculative instruction fetches may be emitted by the processor before a store is seen to the instruction stream and then discarded. Self modifying code due to Itanium stores is not detected by the processor.
- IA-32 instruction fetches can pass prior loads or memory fence operations from the same processor. Data memory accesses, and memory fences are not ordered with respect to IA-32 instruction fetches.
- IA-32 instruction fetches can not pass any serializing instructions, including Itanium `sr1z.i` and IA-32 CPUID. For speculative memory types the processor may prefetch ahead of a serialization operation and then discard the prefetched instructions.
- IA-32 serializing operations wait for all previous stores and loads to complete, and for all prior stores buffered by the processor to become visible. IA-32 serializing instructions include CPUID.
- IA-32 OUT instructions may be buffered, however the processor will not start execution of the next IA-32 instruction until the OUT has completed (been accepted by the platform).
- The processor does not begin execution of the next IA-32 instruction until the IN or OUT has been completed (accepted) by the platform. This statement does not apply

for Itanium memory references to the I/O port space. The processor may issue instruction fetches and VHPT walks ahead of an IN or OUT.

- VHPT Walks are speculative and can occur at any time. VHPT walks can pass all prior IA-32 loads, stores, instruction fetches, I/O operations and serializing instructions.

### 10.6.10.1 Instruction Set Transitions

Instruction set transitions do not automatically fence memory data references. To ensure proper ordering software needs to take into account the following ordering rules.

#### 10.6.10.1.1 Transitions from Intel® Itanium® Instruction Set to IA-32 Instruction Set

- All data dependencies are honored, IA-32 loads see the results of all prior Itanium and IA-32 stores.
- IA-32 stores (*release*) can not pass any prior Itanium load or store.
- IA-32 loads (*acquire*) can pass prior Itanium unordered loads or any prior Itanium store to a different address. Itanium architecture-based software can prevent IA-32 loads from passing prior Itanium loads and stores by issuing an *acquire* operation (or *m.f*) before the instruction set transition.

#### 10.6.10.1.2 Transitions from IA-32 Instruction Set to Intel® Itanium® Instruction Set

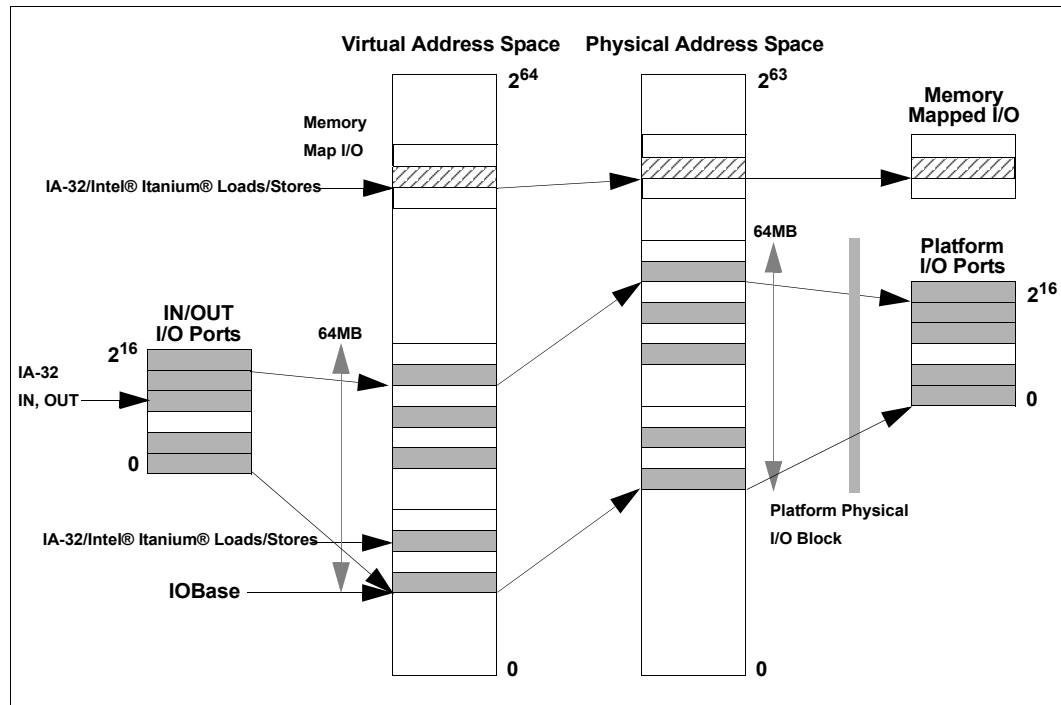
- All data dependencies are honored, Itanium loads see the results of all prior Itanium and IA-32 stores.
- Itanium stores or loads can not pass prior IA-32 loads (*acquire*).
- Itanium unordered stores or any Itanium load can pass prior IA-32 stores (*release*) to a different address. Itanium architecture-based software can prevent Itanium loads and stores from passing prior IA-32 stores by issuing a *release* operation (or *m.f*) after the instruction set transition.

## 10.7 I/O Port Space Model

A consistent unified addressing model is used for both IA-32 and Itanium references to the I/O port space. On prior IA-32 processors two I/O models exist; memory mapped I/O and the 64KB I/O port space. On processors based on the Itanium instruction set, the 64KB I/O port space defined by IA-32 processors is effectively mapped into the 64-bit virtual address space of the processor, producing a single memory mapped I/O model as shown in [Figure 10-1](#). This model allows Itanium normal load and store instructions to also access the I/O port space.

Itanium architecture-based operating system code can directly control IA-32 IN, OUT instruction and accessibility by IA-32 or Itanium load/store instructions to blocks of 4 virtual I/O ports using the TLBs. The entire range of virtual memory mechanisms defined by the Itanium architecture: access rights, dirty, access bits, protection keys, region identifiers can be used to control permission and addressability.

**Figure 10-1. I/O Port Space Model**



In the Itanium System Environment, the virtual location of the 64 MB I/O port space is determined by operating system. For IA-32 IN and OUT instructions, the operating system can specify the virtual base location via the I/O base register.

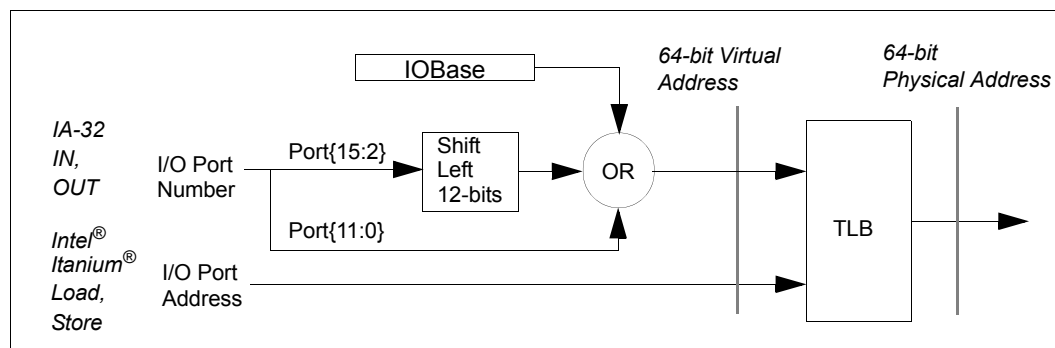
Any IA-32 or Itanium load or store within the virtual region mapped by the operating system to the platform's physical 64 MB I/O port block, directly accesses physical I/O devices within the I/O port space. The location of the 64 MB I/O port block within the 2<sup>63</sup> byte physical address space is determined by platform conventions, see [Section 10.7.2, "Physical I/O Port Addressing" on page 2:270](#) for details.

### 10.7.1 Virtual I/O Port Addressing

The IA-32 defined 64-KB I/O port space is expanded into 64 MB. This effectively places 4 I/O ports per each 4KB virtual and physical page. Since there are 4 ports per virtual page, the TLBs can be used port address translation, and permission checks as shown in [Figure 10-2](#).



**Figure 10-2. I/O Port Space Addressing**



For IA-32 IN and OUT instructions a port's virtual address is computed as:

```
port_virtual_address = IOBase | (port{15:2}<<12) | port{11:0}
```

This address computation places 4 ports on each 4K page and expands the space to 64MB, with the ports being at a relative offset specified by port{11:0} within each 4K-byte virtual page. IOBase is a kernel register (KR) maintained by the operating system that points to the base of the 64MB Virtual I/O port space. *The value in IOBase must be aligned on a 64MB boundary otherwise port address aliasing will occur and processor operation is undefined.*

For Itanium load and stores accesses to the I/O port space, a port's virtual address can be computed in the same manner, specifically.

```
port_virtual_address = IOBase | (port{15:2}<<12) | port{11:0}
```

In practice this address is a constant for any given physical I/O device.

**Note:** In the generation of the I/O port virtual address, software MUST ensure that port\_virtual\_address{11:2} are equal to port{11:2} bits. Otherwise, some processors implementations may place the port data on the wrong bytes of the processor's bus and the port will not be correctly accessed.

IA-32 IN and OUT instructions and Itanium or IA-32 load/store instructions can reference I/O ports in 1, 2, or 4-byte transactions. References to the legacy I/O port space cannot be performed with greater than 4 byte transactions due to bus limitations in most systems. Since an IA-32 IN/OUT instruction can access up to 4 bytes at port address 0xFFFF, the I/O port space effectively extends 3 bytes beyond the 64KB boundary. Operating systems can; 1) not map the excess 3 bytes, resulting in denial of permission for the excess 3 bytes, or 2) map via the TLB the excess 3 bytes back to port address 0 effectively wrapping the I/O port space at 64KB.

Operating system code can map each virtual I/O port space page anywhere within the physical address space using the Data Translation Registers or the Data Translation Cache. Large page translations can be used to reduce the number of mappings required in the TLB to map the I/O port space. For example, one 64MB translation is sufficient to map the entire expanded 64MB I/O port space. The **UC memory attribute** must be used for all I/O port space mappings to avoid speculative processor references to I/O devices, otherwise processor and platform operation is undefined.

**Operating System Warning:** Operating system code can not remap a given port to another port address within the I/O port space, such that `port_physical_address{21:12} != port_physical_address{11:2}`. Otherwise, based on the processor model, I/O port data may be placed on the wrong bytes of the processor's bus and the port will not be correctly accessed.

I/O port space breakpoints can be configured by loading the address and mask fields with the virtual address defined by the operating system to correspond to the I/O port space.

The processor (as defined in the next section) ensures that load, store references will not result in references to I/O devices for which permission was not granted.

All memory related faults defined in [Chapter 5, "Interruptions"](#) can be generated by IA-32 IN and OUT references to the I/O port space, including `IA_32_Exception(Debug)` traps for data address breakpoints and `IA_32_Exception(AlignmentCheck)` for unaligned references. (EFLAG.ac enabled unaligned port references are not detected by the processor). Itanium Data Breakpoint registers (DBRs) can be configured to generate debug traps for references into the I/O port space by either IA-32 IN/OUT instructions or by IA-32 or Itanium load/store instructions.

## 10.7.2 Physical I/O Port Addressing

Some processors implementations will provide an M/IO pin or bus indication by decoding physical addresses if references are within the 64MB physical I/O block. If so the 64MB I/O port space is compressed back to 64KB. Subsequent processor implementations may drop the M/IO pin (or bus indication) and rely on platform or chip-set decoding of a range of the 64MB physical address space.

Through the PAL firmware interface, the 64MB physical I/O block can be programmed to any arbitrary physical location. It is suggested that to be compatible with IA-32 based platforms, the platform physical location of the physical I/O block be programmed above 4G-bytes and above all useful DRAM, ROM and existing memory mapped I/O areas. See `PAL_PLATFORM_ADDR` on [page 2:442](#) for details.

Based on the platform design, some platforms can accept addresses for the expanded 64MB I/O block, while other platforms will require that the I/O port space be compressed back to 64KB by the processor. If the I/O port space needs to be compressed either the processor or platform (based on the implementation) will perform the following operation for all memory references within the physical I/O block.

```
IO_address{1:0} = PA{1:0}
IO_address{15:2} = PA{25:12} // byte strobes are generated
                        // from the lower I/O_address bits
```

The processor ensures that the bus byte strobes and bus port address are derived from `PA{25:12,1:0}`. Thus allowing an operating system to control security of each 4 ports via TLB management of `PA{25:12}`.

### 10.7.2.1 I/O Port Addressing Restrictions

For the 64MB physical I/O port block the following operations are undefined and may result in unpredictable processor operation; references larger than 4-bytes, instruction fetch references, references to any memory attribute other than UC, or semaphore references which require an atomic lock. To ensure I/O ports accesses are not granted for which the TLB has not been consulted, the processor ensures:

- All byte addresses within the same 4KB page alias to the 4 ports defined by the mapped physical I/O port page.
- All IA-32 and Itanium unaligned loads and stores that cross a 4-byte boundary to the processor's physical I/O port block are truncated. That is the bus transaction to the area before the 4-byte boundary is performed (the number of bytes emitted is model specific). No bus transaction is performed for the bytes that are beyond the 4-byte boundary. 4-byte crosser loads while return some undefined data. 4-byte crosser stores will not write all intended bytes.
- For IA-32 IN/OUT accesses that cross a 4-port boundary the processor will break the operation into smaller 1, 2, or 3 byte I/O port transactions within each 4KB virtual page.

### 10.7.3 IA-32 IN/OUT instructions

IA-32 I/O instructions (IN, OUT, INS, OUTS) defined in the **Intel® 64 and IA-32 Architectures Software Developer's Manual** are augmented as follows:

- I/O instructions first check for IOPL permission. If  $PSR.cpl \leq EFLAG.iopl$ , access permission is granted. Otherwise the TSS I/O permission bitmap may be consulted as defined below. If the Bitmap denies permission or is not consulted an `IA_32_Exception(GPFault)` is generated.
- If IOPL permission is denied and `CFLG.io` is 1, the TSS I/O permission bitmap is consulted for access permission. If the corresponding bit(s) for the I/O port(s) is 1, indicating permission is denied, a GPFault is generated. Otherwise access permission is granted. The TSS I/O permission bitmap provides 1 port permission control at the expense of additional processor data memory references. This mechanism can be used in the Itanium System Environment, but is not recommended since TLB access controls defined by the Itanium architecture are faster and provide a consistent control mechanism for both IA-32 and Itanium architecture-based code. Whereas, the TLB mechanism provides a control mechanism for both IA-32 and Itanium memory references.
- If `CFLG.io` is 0, the TSS I/O permission bitmap is not consulted and if the IOPL check failed an `IA_32_Exception(GPFault)` is generated. By setting `CFLG.io` to 0, operating system code can disable all processor references to the TSS. By setting  $IOPL < PSR.cpl$  and setting `CFLG.io` to 0, operating system code can block all user level execution of IA-32 I/O instructions, no TSS needs to be allocated or defined by the operating system.
- I/O port references generate a virtual port address relative to the IOBase register as defined in [Section 10.7.1, "Virtual I/O Port Addressing" on page 2:268](#).
- If data translations are enabled, the TLB is consulted for the required virtual to physical mapping. If the required mapping is not present a VHPT Translation, Data TLB Miss or Alternative Data TLB Miss fault is generated.
- If data translations are enabled, Access Rights, Permission Keys, Access, Dirty and Present bits are checked and faults generated.

- If data translations are disabled (PSR.dt is 0) or the referenced I/O port is mapped to an unimplemented virtual address (via the IOBase register), a GPFault is raised on the referencing IA-32 IN, OUT, INS, or OUTS instruction.
- Alignment and Data Address breakpoints are also checked and may result in an IA\_32\_Exception(AlignmentCheck) fault (if PSR.ac is 1) or IA\_32\_Exception(Debug) trap.
- If an IA-32 IN/OUT I/O port Accesses cross a 4-port boundary the processor will break the operation into smaller 1, 2, or 3 byte transactions.
- Assuming no faults, a physical transaction is emitted to the mapped or specified physical address.

The processor ensures that IA-32 IN, INS, OUT, OUTS references are fully ordered and will not allow prior or future data memory references to pass the I/O operation as defined in [Section 10.6.10, “IA-32 Memory Ordering” on page 2:265](#). The processor will wait for acceptance for IN and OUT operations before proceeding with subsequent externally visible bus transactions.

## 10.7.4 I/O Port Accesses by Loads and Stores

If an access is made to the I/O port block using IA-32 or Itanium loads and stores the following differences in behavior should be noted; EFLAG.iopl permission is not checked, TSS permission bitmap is not checked, and stores and loads do not honor IN and OUT memory ordering and acceptance semantics (the processor will not automatically wait for a store to be accepted by the platform).

Virtual addresses for the I/O port space should be computed as defined in [Section 10.7.1, “Virtual I/O Port Addressing” on page 2:268](#) If data translations are enabled, the TLB is consulted for mappings and permission, and the resulting mapped physical address used to address the physical I/O device.

If IA-32 ordering semantics are required to a particular I/O port device (or memory mapped I/O device), IA-32 or Itanium architecture-based software must enforce ordering to the I/O device. Software can either perform a memory ordering fence before and after the transaction, or use an load acquire or store release

To ensure the processor does not speculatively access an I/O device, all I/O devices must be mapped by the UC memory attribute.

If IA-32 acceptance semantics are required (i.e. additional data memory transactions are not initiated until the I/O transaction is completed), Itanium architecture-based code can issue a memory acceptance fence. Conversely, if certain I/O devices do not require IA-32 IN/OUT ordering or acceptance semantics, Itanium architecture-based code can relax ordering and acceptance requirements as shown below.

### OUT

```
[mf]//Fence prior memory references, if required

add port_addr = IO_Port_Base, Expanded_Port_Number
st.rel (port_addr), data
[mf.a] //Wait for platform acceptance, if required
[mf] //Fence future memory operations, if required
```

IN

```
[mf] //Fence prior memory references, if required
add port_addr = IO_Port_Base, Expanded_Port_Number
ld.acq data, (port_addr)
[mf.a] //Wait for platform acceptance, if required
[mf] //Fence future memory references, if required
```

## 10.8 Debug Model

The debug facilities defined by the Itanium architecture are designed to support debugging of both the Itanium and IA-32 instruction set. The following debug events can be triggered during IA-32 instruction set execution by Itanium debug resources.

- **Single Step trap** – When PSR.ss is 1 (or EFLAG.tf is 1), successful execution of each IA-32 instruction, results in an IA\_32\_Exception(Debug) trap. After the single step trap, IIP points to the next IA-32 instruction to be executed.
- **Breakpoint Instruction trap** – execution of INT 3 (breakpoint) instruction results in a IA\_32\_Exception(Debug) trap.
- **Instruction Debug fault** – When PSR.db is 1 and PSR.id is 0 and EFLAG.rf is 0, any IA-32 instruction fetch that matches the parameters specified by the IBR registers results in an IA\_32\_Exception(Debug) fault. After servicing a Debug fault, debuggers can set PSR.id (or EFLAG.rf for IA-32 instructions) before restarting the faulting instruction. If PSR.id is 1, Instruction Debug faults are temporarily disabled for one Itanium instruction. If PSR.id is 1 or EFLAG.rf is 1, Instruction Debug faults are temporarily disabled for one IA-32 instruction. The successful execution of an IA-32 instruction clears both PSR.id and EFLAG.rf bits. The successful execution of an Itanium instruction only clears PSR.id.
- **Data Debug traps** – When PSR.db is 1, any IA-32 data memory reference that matches the parameters specified by the DBR registers results in a IA\_32\_Exception(Debug) trap. IA-32 data debug events are traps, not faults as defined for Itanium instruction set data debug events. Trap behavior is required since any given IA-32 instruction can access several memory locations during its execution. The reported trap code returns the match status of the first four DBR registers that matched during the execution of the IA-32 instruction. Zero, one or DBR registers may be reported as matching.
- **Taken Branch trap** – When PSR.tb is 1, a IA\_32\_Exception(Debug) trap occurs on every IA-32 taken branch instruction (CALL, Jcc, JMP, RET, LOOP). After the trap, IIP points to the branch target.
- **Lower Privilege Transfer trap** – Does not occur during IA-32 instruction set execution.

For virtual memory accesses, breakpoint address registers contain the virtual addresses of the debug breakpoint. For physical accesses, the addresses in these registers are treated as a physical address. Software should be aware that debug registers configured to fault on virtual references, may also fault on a physical reference if translations are disabled. Likewise a debug register configured for physical references can fault on virtual references that match the debug breakpoint registers.

## 10.8.1 Data Breakpoint Register Matching

Each Itanium data breakpoint register has the following matching behavior for IA-32 instruction set data memory references:

- **DBR.addr** – IA-32 single or multi-byte data memory references that access ANY memory byte specified by the DBR address and mask fields results in a debug breakpoint trap regardless of datum size and alignment. The upper 32-bits of DBR.addr must be zero to detect IA-32 data memory references. Since IA-32 data breakpoints are traps, all processor implementations ensure data breakpoint traps are precise. Traps are only reported if any byte in the data memory reference ANDed with the DBR mask bitwise matches the DBR address field ANDed with the DBR mask. No spurious data breakpoint faults are generated for IA-32 data memory operands that are unaligned, nor are matches reported if no bytes of the operand lie within the address range specified by the DBR address and mask fields. Note, Itanium instruction set generated unaligned data memory references may result in spurious imprecise breakpoint faults.
- **DBR.mask** – by programming the mask a breakpoint range of 1, 2, 4, 8, or any power of 2 combination can be supported. Mask bits above bit 31 are checked by the processor during IA-32 data memory references
- **Trap code B bits** – are set indicating a match with the corresponding data breakpoint register DBR0-3. For IA-32 data debug traps, any number of B-bits can be set indicating a match.

The B-bits are only set and a data breakpoint trap generated if 1) the breakpoint register precisely matches the specified DBR address and mask, 2) it is enabled by the DBR read or write bits for the type of the memory transaction, 3) the DBR privilege field matches PSR.cpl, 4) PSR.db is 1, and 5) no other higher priority faults are taken.

I/O port space breakpoints can be configured by loading the address and mask fields with the virtual address defined by the operating system to correspond to the I/O port space.

## 10.8.2 Instruction Breakpoint Register Matching

Each Itanium instruction breakpoint register has the following matching behavior for IA-32 instruction set memory fetches:

- **IBR.addr** – an IBR register matches an IA-32 instruction fetch address, if the first byte of an IA-32 instruction address ANDed with the IBR mask bitwise matches the IBR address field ANDed with the IBR mask. Note that only the first byte is analyzed. The upper 32-bits of IBR.addr must be zero to detect IA-32 instruction fetches.
- **IBR.mask** – by programming the mask a breakpoint range of 1, 2, 4, 8, or any power of 2 combination can be supported. Mask bits above bit 31 are ignored during IA-32 instruction fetches.

The instruction breakpoint fault is generated if 1) the breakpoint register precisely matches the specified IBR address and mask, 2) it is enabled by the IBR execute bit, 3) the IBR privilege field matches PSR.cpl, 4) PSR.db is 1, 5) PSR.id is 0, and 6) no other higher priority faults are taken.

## 10.9 Interruption Model

Within the Itanium System Environment, all interruptions originating out of the IA-32 or Itanium instruction sets are delivered to Itanium architecture-based Interruption Handlers within the Itanium architecture-based operating system. Virtual memory management faults, machine checks, and external interrupts are always delivered to Itanium architecture-based interruption handlers regardless of the instruction set running at the time of the interruption. IA-32 exceptions, control transfers through gates, task switches, and accesses to sensitive IA-32 system resources are intercepted into Itanium architecture-based interruption handlers. Using these intercepts, Itanium architecture-based software can implement specific policies with regard to that resource. Policies may include virtualization, emulation of an IA-32 opcode or memory access, or various permission policies.

In general, if an interruption is independent of the executing instruction set (such as an external interruption or TLB fault) common Itanium architecture-based handlers are invoked. For classes of exceptions and intercept conditions that are specific to the IA-32 instruction set, three special Itanium architecture-based software handlers are invoked to deal with IA-32 instruction set interruptions. Table 10-8 shows the three interruption handlers defined to support IA-32 events. See Section 9.2, “IA-32 Interruption Vector Definitions” on page 2:213 for details on these interruption handlers.

**Table 10-8. IA-32 Interruption Vector Summary**

Handler	Description
IA_32_Intercept	Intercepted IA-32 instructions, I/O, system flag manipulation and gate transfers.
IA_32_Exception	IA-32 instruction set generated exceptions.
IA_32_Interrupt	IA-32 instruction set generated software interrupts

This grouping of interruption handlers simplifies software handlers such that they do not need to be concerned with behavior of both IA-32 and Itanium instruction sets.

Interruption registers (defined in Chapter 3) record the state of IA-32 execution at the point of interruption. For IA-32 exceptions, ISR contains IA-32 defined error codes and vector numbers as defined by the *Intel® 64 and IA-32 Architectures Software Developer’s Manual*. IA-32 instruction set related exceptions and software interruptions vector directly through the interruption mechanism defined by the Itanium architecture without consulting the IA-32 IDT or performing any memory stack pushes.

### 10.9.1 Interruption Summary

Table 10-9 summarizes the set of all IA-32 interruptions and how they are mapped to Itanium architecture-based interruption handlers within the Itanium System Environment. See Chapter 9 and Chapter 8 for a detailed definition of each interruption.

**Table 10-9. IA-32 Interruption Summary**

IA-32 Vector	Itanium® Architecture-based Interruption Handler	ISR Vector	ISR Code	Description
<b>IA-32 Defined Interruptions</b>				
0	IA_32_Exception (Divide)	0	0	IA-32 divide by zero fault.



**Table 10-9. IA-32 Interruption Summary (Continued)**

IA-32 Vector	Itanium® Architecture-based Interruption Handler	ISR Vector	ISR Code	Description
1	IA_32_Exception (Debug)	1	0	IA-32 instruction breakpoint fault.
1	IA_32_Exception (Debug)	1	TrapCode	IA-32 debug events. The Trap Code indicates concurrent taken branch, data breakpoint and single step trap conditions.
2	External Interrupt	0	0	NMI is delivered through the Intel Itanium External Interrupt vector.
3	IA_32_Exception(Break)	3	TrapCode	IA-32 INT 3 instruction.
4	IA_32_Exception(INTO)	4	TrapCode	IA-32 INTO detected overflow trap.
5	IA_32_Exception (Bound)	5	0	IA-32 BOUND range exceeded fault.
6	IA_32_Intercept(Inst)	0	InterceptCode	All IA-32 unimplemented and illegal opcodes.
7	IA_32_Exception(DNA)	7	0	IA-32 Device not available fault.
8	--	N/A		IA-32 Double fault can not be generated in the Intel Itanium System Environment, Intel reserved.
9	--	N/A		Intel reserved
10	--	N/A		IA-32 Invalid TSS fault can not generated in the Intel Itanium System Environment, Intel reserved,
11	IA_32_Exception(NotPresent)	11	ErrorCode <sup>a</sup>	IA-32 Segment Not present fault.
12	IA_32_Exception (Stack)	12	ErrorCode	IA-32 Stack Exception fault.
13	IA_32_Exception (GPFault)	13	ErrorCode	IA-32 General Protection fault.
14	Intel Itanium TLB faults	see Data TLB faults below		IA-32 Page fault can not be generated in the Intel Itanium System Environment, replaced by Intel Itanium TLB faults, Intel reserved,
15	--	N/A		Intel reserved.
16	IA_32_Exception (FPError)	16	0	IA-32 floating-point fault.
17	IA_32_Exception(AlignCheck)	17	0	IA-32 un-aligned data references.
18	Corrected MCHK	N/A		IA-32 Machine Check can not be generated in the Intel Itanium System Environment, replaced by the PAL Machine Check Architecture, Intel reserved.
19	IA_32_Exception (StreamSIMD)	19	0	IA-32 SSE Numeric Error fault.
20-31	--	N/A		Intel reserved.
0-255	External Interrupt	0	0	External interrupts are delivered through the Intel Itanium External Interrupt vector. Software must read the IVR register to determine the vector number.
0-255	IA_32_Interrupt (vector #)	Vect#	TrapCode	IA-32 Software Interrupt trap. ISR contains the vector number.
<b>IA-32 Interceptions</b>				



**Table 10-9. IA-32 Interruption Summary (Continued)**

IA-32 Vector	Itanium® Architecture-based Interruption Handler	ISR Vector	ISR Code	Description
	IA_32_Interrupt(Inst)	0	InterceptCode	Intercept for unimplemented, illegal or privileged IA-32 opcodes.
	IA_32_Interrupt(Gate)	1	TrapCode	Intercept for control transfers through a Call Gate, Task gate or Task Segment.
	IA_32_Interrupt(SystemFlag)	2	TrapCode	Intercept for modification of system flag values.
	IA_32_Interrupt(Lock)	4	0	IA-32 semaphore operation requires an external bus lock when DCR.Ic is 1.
		3,5-25 5	--	Intel reserved

a. The IA-32 Error Code is defined as a Selector Index, and TI, IDT and EXT bits are set based on the exception. See *Intel® 64 and IA-32 Architectures Software Developer's Manual* for the complete definition.

## 10.9.2 IA-32 Numeric Exception Model

IA-32 numeric instructions follow the IA-32 delayed floating-point exception model. Specifically IA-32 numeric exceptions are held pending until the next IA-32 numeric or MMX technology instruction as defined in the *Intel® 64 and IA-32 Architectures Software Developer's Manual*. Numeric faults generated on SSE instructions are reported precisely on the faulting SSE instruction. SSE instructions do NOT trigger the report of pending IA-32 numeric exceptions.

For voluntary transitions out of the IA-32 instruction, an implicit FWAIT operation is performed by the `jmpbe` instruction to ensure all pending numeric exceptions are reported. For involuntary transitions out of the IA-32 instruction set (external interruptions, TLB faults, exceptions, etc.) the processor does not perform a FWAIT operation. However, every IA-32 numeric instruction that generates a pending numeric exception loads the application registers FSR, FIR, and FDR with the IA-32 floating-point state on the instruction that generating the exception. This state contains information defined by the IA-32 FSTENV and FLDENV instructions. During a process context switch, the operating system must save and restore FSR, FIR, and FDR (effectively performing an FSTENV and FLDENV) to ensure numeric exceptions are correctly reported across a process switch.

## 10.10 Processor Bus Considerations for IA-32 Application Support

The section briefly discusses bus and platform considerations when supporting IA-32 applications in the Itanium System Environment.

Itanium architecture-based code does not assert the SPLCK and LOCK pins. The LOCK pin is used by IA-32 code to signal an external atomic bus transaction for which atomicity cannot be enforced within the processor's caches, whereas, SPLCK indicates if an unaligned external bus lock requires a split lock operation and hence several bus

transactions. For IA-32 code, if the platform does not support LOCK or SPLCK, the operating system must disable external bus lock transactions by setting DCR.lc to 1. When DCR.lc is 1, any IA-32 atomic reference not serviced internally in the processor's caches results in an IA\_32\_Intercept(Lock) fault. See [Section 3.3.4.1, "Default Control Register \(DCR – CR0\)" on page 2:31](#) for details. When DCR.lc is 0, operating system code is responsible for emulation of the IA-32 instruction and ensuring atomicity (if required).

The A20M and IGNE pins are ignored in the Itanium System Environment. FERR is not asserted in the Itanium System Environment.

In both IA-32 and Itanium System Environments, the M/IO pin (or an external bus indication) is asserted by any memory reference to the 64MB I/O port block range of the physical address space. See [Section 10.7, "I/O Port Space Model" on page 2:267](#) for details.

SMI and the SMM environment are not supported on processors based on the Itanium architecture. The PMI interrupt and PAL firmware environment replace them. See [Section 11.5, "Platform Management Interrupt \(PMI\)" on page 2:310](#) for details.

### 10.10.1 IA-32 Compatible Bus Transactions

Within the Itanium System Environment, the following bus transactions are initiated:

- INTA – Interrupt Acknowledge - emitted by the operating system (via a read to the INTA byte in the processor's Interrupt Block) to acquire the interrupt vector number from an external interrupt controller.
- HALT – Emitted when the processor has entered the halt state due to the operating system/platform firmware calling PAL\_HALT or PAL\_HALT\_LIGHT.
- SHUTDOWN – Emitted when the processor has entered the shutdown state. This can only be generated when the processor has entered into the IA-32 System Environment by calling PAL\_ENTER\_IA\_32\_ENV procedure call.
- STPACK – Stop Acknowledge. Emitted by calling an implementation-specific PAL firmware procedure. See the processor-specific firmware guide for more information.
- FLUSH – Emitted when the WBINVD or INVD instruction is executed when running in the IA-32 System Environment entered by calling PAL\_ENTER\_IA\_32\_ENV procedure call. Indicates that external caches (if any) should be invalidated.
- SYNC – Emitted when the WBINVD instruction is executed when running in the IA-32 System Environment entered by calling PAL\_ENTER\_IA\_32\_ENV procedure call. Indicates that external caches (if any) should copy all modified cache lines back to main memory.

This chapter defines the architectural requirements for the **Processor Abstraction Layer (PAL)** for all processors based on the Itanium architecture. It is intended for processor designers, firmware/BIOS designers, system designers, and writers of diagnostic and low level operating system software.

PAL is part of the Itanium processor architecture and its goal is to provide a consistent firmware interface to abstract processor implementation-specific features.

The objectives of this chapter are to define:

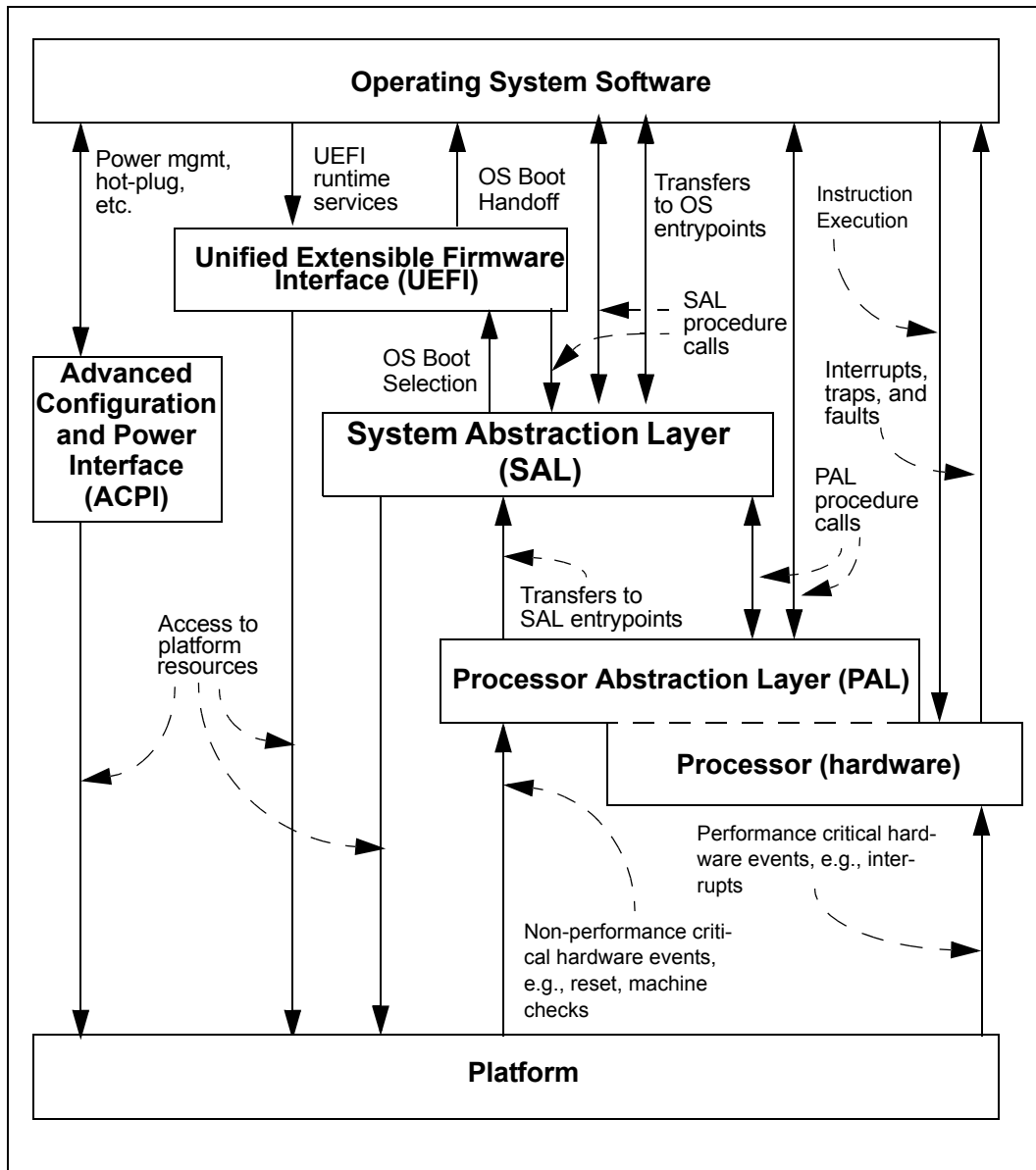
- The architectural behavior and interface requirements for processor testing, configuration and error recovery. This includes the hardware entrypoints into PAL and the PAL interfaces to platform firmware and system software.
- A set of boot and runtime PAL procedures to access processor implementation-specific hardware and to return information about processor implementation-dependent configuration.
- A computing environment for both PAL entrypoints and procedures such that:
  - Memory used by PAL procedures is allocated by the caller of PAL procedures.
  - PAL code runs little endian.
  - PAL interface is as endian neutral as possible.
  - PAL is Itanium architecture-based code.
  - PAL code runs at privilege level 0.
  - PAL procedures can be called without backing store, except where memory-based parameters are returned.
- The processor and platform hardware requirements for PAL. This includes minimizing PAL dependencies on platform hardware and clearly stating where those dependencies exist.
- A PAL interface and requirements to support firmware update and recovery.

## 11.1 Firmware Model

As shown in [Figure 11-1](#), Itanium architecture-based firmware consists of several major components: Processor Abstraction Layer (PAL), System Abstraction Layer (SAL), Unified Extensible Firmware Interface (UEFI) and Advanced Configuration and Power Interface (ACPI). PAL, SAL, UEFI and ACPI together provide processor and system initialization for an operating system boot. PAL and SAL provide machine check abort handling. PAL, SAL, UEFI and ACPI provide various run-time services for system functions which may vary across implementations. The interactions of the various services that PAL, SAL, UEFI and ACPI provide are illustrated in [Figure 11-1](#).

In the context of this model and throughout the rest of this chapter, the System Abstraction Layer (SAL) is a firmware layer which isolates operating system and other higher level software from implementation differences in the platform, while PAL is the firmware layer that abstracts the processor implementation.

**Figure 11-1. Firmware Model**



### 11.1.1 Processor Abstraction Layer (PAL) Overview

The purpose of the Processor Abstraction Layer, is to provide a firmware abstraction between the processor hardware implementation and system software and platform firmware, so as to maintain a single software interface for multiple implementations of the processor hardware. PAL is defined to be independent of the number of processors on a platform.

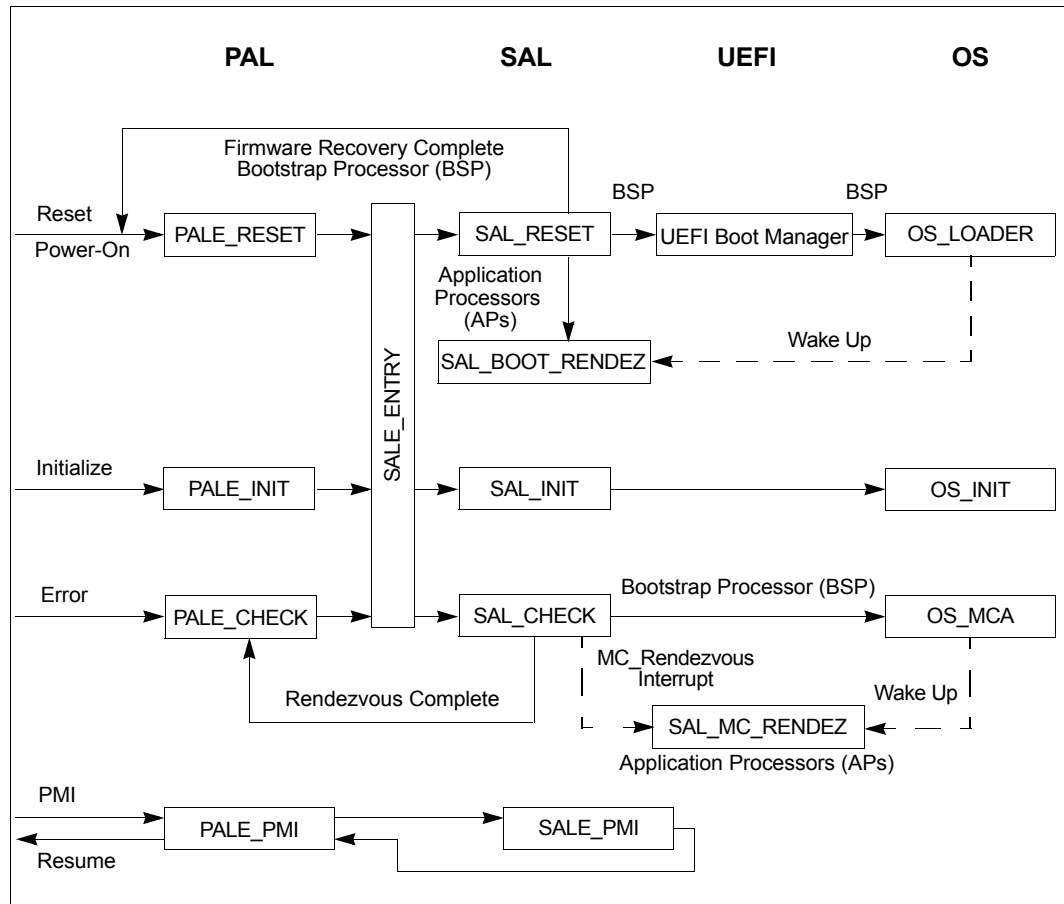
PAL encapsulates those processor functions that are likely to change on an implementation to implementation basis so that SAL firmware and operating system software can maintain a consistent view of the processor. These include non-performance critical functions dealing such as processor initialization, configuration and error handling.

PAL consists of two main components:

- Entrypoints, which are invoked directly by hardware events such as reset, init and machine checks. These interruption entrypoints perform functions such as processor initialization and error recovery.
- Procedures, which may be called by higher level firmware and software to obtain information about the identification, configuration, and capabilities of the processor implementation; to perform implementation-dependent functions such as cache initialization; or to allow software to interact with the hardware through such functions as power management or enabling/disabling processor features.

### 11.1.2 Firmware Entrypoints

Figure 11-2. Firmware Entrypoints Logical Model



### 11.1.3 PAL Entrypoints

The following hardware events can trigger the execution of a PAL entrypoint:

- Power-on/reset
- Hardware errors (both correctable and uncorrectable)
- Initialization event (via external interrupt bus message or processor pin)
- Platform management interrupt (via external interrupt bus message or processor pin)

These hardware events trigger the execution of one of the following PAL entrypoints (as shown in [Figure 11-2](#)):

- PALE\_RESET – Initializes and tests the processor following power-on or reset and then branches to SALE\_ENTRY to determine whether to perform firmware recovery update, or to boot the machine for OS use. See [Section 11.1.4, “SAL Entrypoints” on page 2:282](#).
- PALE\_CHECK – Determines if errors are processor related, saves processor related error information and corrects errors where possible (for example, by flushing a corrupted instruction cache line and marking the cache line as unusable). In all cases, PALE\_CHECK branches to SALE\_ENTRY to complete the error logging, correction, and reporting.
- PALE\_INIT – Saves the processor state, places the processor in a known state, and branches to SALE\_ENTRY. PALE\_INIT is entered as a response to an initialization event.
- PALE\_PMI – Saves the processor state and branches to SALE\_PMI. PALE\_PMI is entered as a response to a platform management interrupt.

### 11.1.4 SAL Entrypoints

There are two entrypoints from PAL into SAL:

- SALE\_ENTRY – PAL branches to this SAL entrypoint after a power-on, reset, machine check, or initialization event. If SALE\_ENTRY was invoked by a machine check or initialization event, SALE\_ENTRY branches to the appropriate routine:
  - SAL\_CHECK is invoked after a machine check.
  - SAL\_INIT is invoked after an initialization event.

If SALE\_ENTRY was invoked by a reset or power on, it checks to determine if a firmware recovery condition exists. If it does, SALE\_ENTRY performs the firmware update, then performs a RESET operation to invoke PAL\_RESET. If a recovery condition does not exist, SALE\_ENTRY returns to PAL\_RESET to complete processor self-test. PAL\_RESET then branches back to SALE\_ENTRY, which, in turn, branches to SAL\_RESET.

- SALE\_PMI – platform management interrupt. PALE\_PMI branches to this SAL entrypoint after saving processor state in response to the platform management interrupt.

### 11.1.5 OS Entrypoints

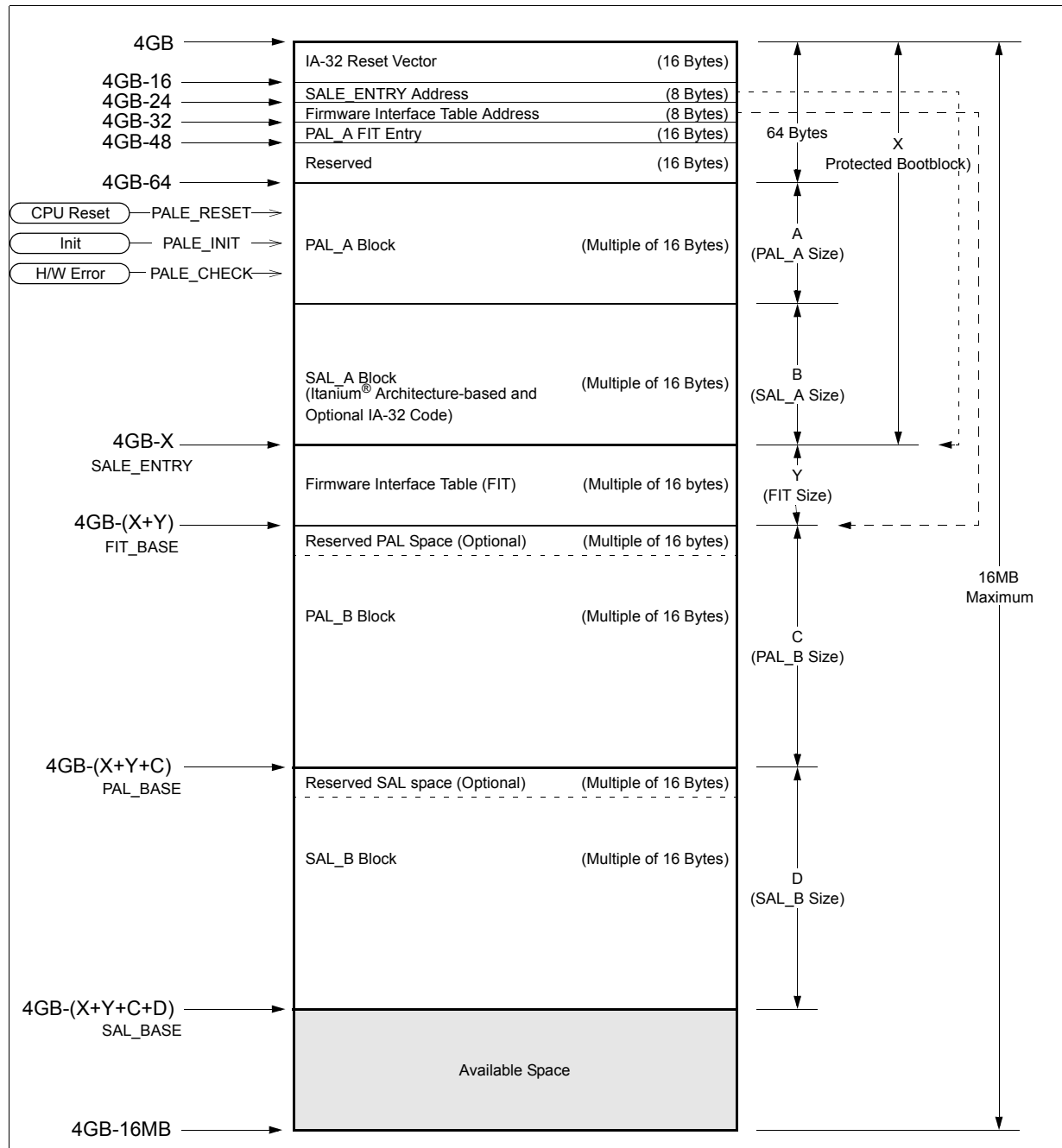
There are several entrypoints from SAL into an operating system (or equivalent software). Entrypoints from SAL into the operating system are expected to meet the following model:

- OS\_BOOT – Operating System Boot interface.
- OS\_MCA – Operating System Machine Check Abort Handler.
- OS\_INIT – Operating System Initialization Handler.
- OS\_RENDEZ – Operating System Multiprocessor Rendezvous interface.

### 11.1.6 Firmware Address Space

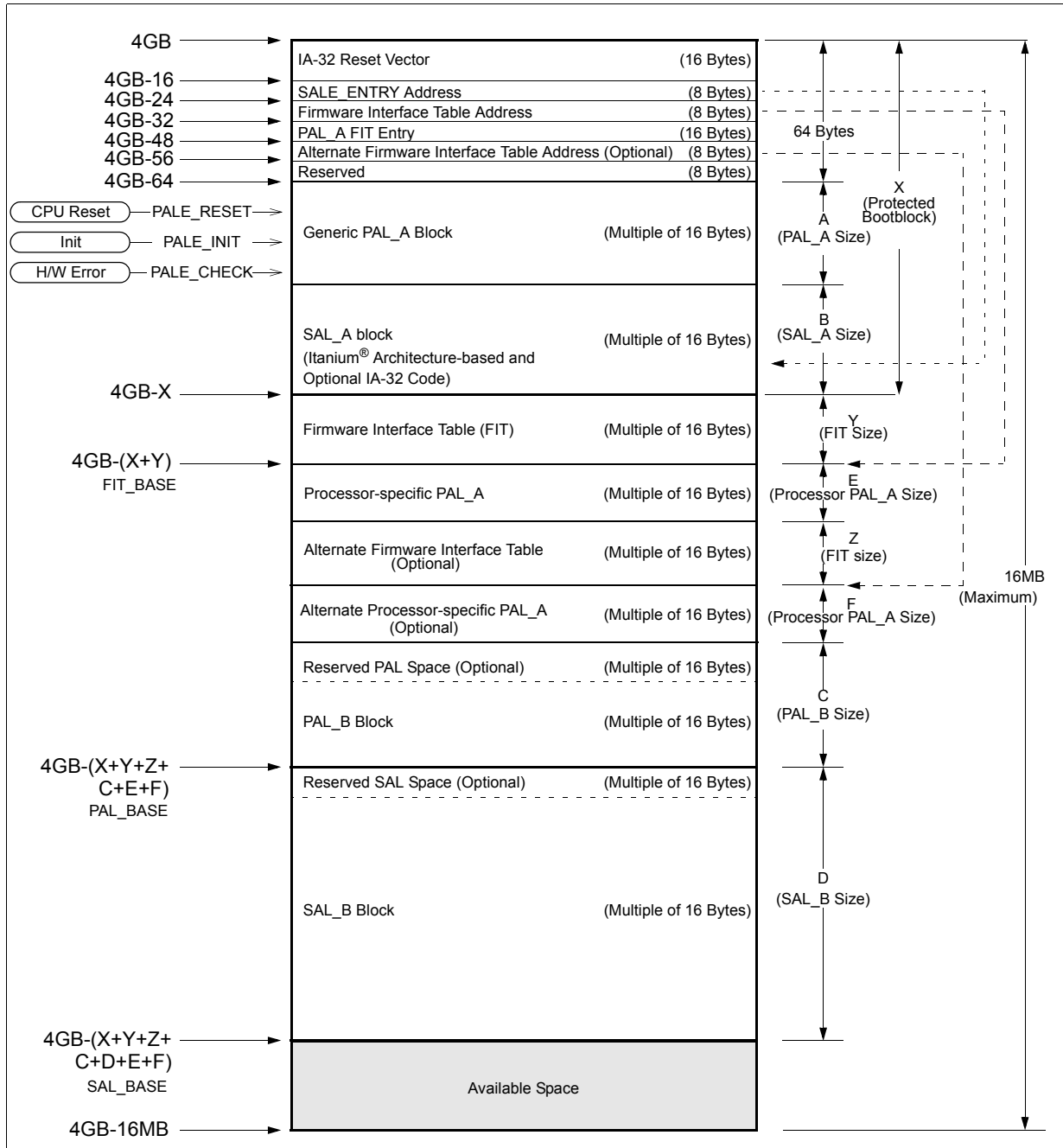
The firmware address space occupies the 16 MB region between 4 GB - 16 MB and 4 GB (addresses 0xFF00\_0000 through 0xFFFF\_FFFF). There are two primary layouts of this address space. The first version is shown in [Figure 11-3](#) and the second version is shown in [Figure 11-4](#). The first version has one PAL\_A component. This layout allows for robust recovery of PAL\_B and SAL\_B components. This layout is useful for cases where PAL\_A will not need to be upgraded. The second version splits the PAL\_A block into two components. The first component is referred to as the generic PAL\_A and the second component is the processor-specific PAL\_A. Splitting the PAL\_A up in this manner allows for a robust upgrade of the processor-specific PAL\_A firmware as well as the PAL\_B and SAL\_B components. This is very useful if a platform is designed to support multiple processor generations which would require a PAL\_A upgrade when the new processor generation is released. The generic PAL\_A which resides in the Protected Boot Block will work across processor generations for a given platform. The processor-specific PAL\_A resides outside the Protected Boot Block and works for a specific processor generation.

**Figure 11-3. Firmware Address Space**





**Figure 11-4. Firmware Address Space with Processor-specific PAL\_A Components**



The firmware address space is shared by SAL and PAL. Some of the SAL/PAL boundaries are implementation dependent. The address space contains the following regions and locations.

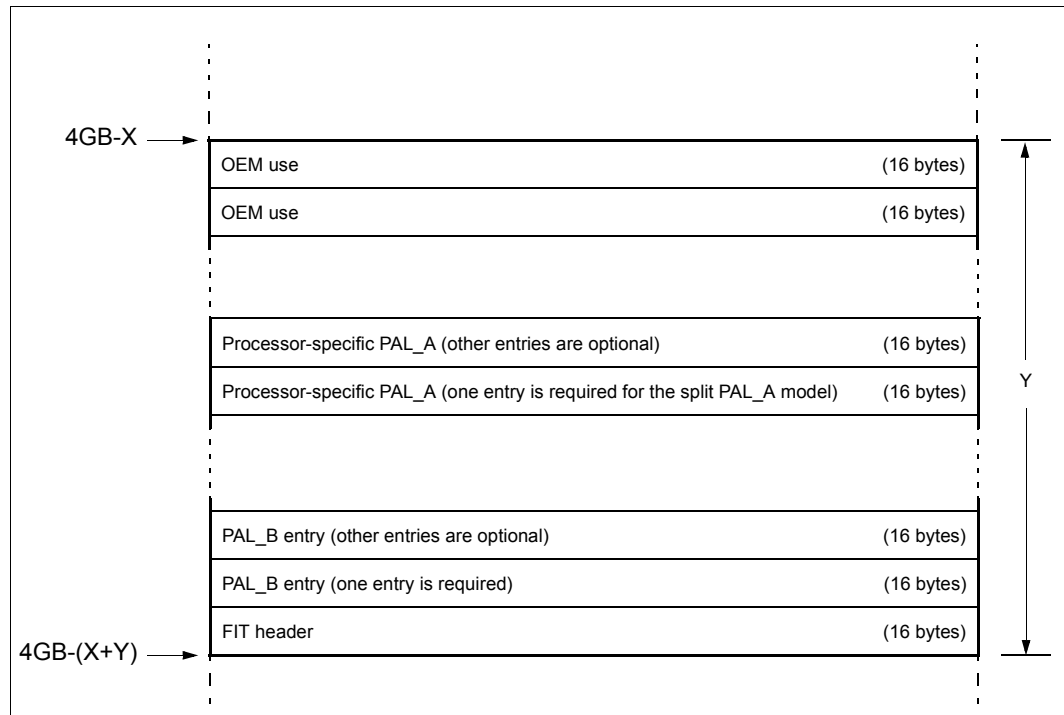
- The 16 bytes at 0xFFFF\_FFF0 (4GB-16) contain IA-32 Reset Code.
- The 8 bytes at 0xFFFF\_FFE8 (4GB-24) contain the physical address of the SALE\_ENTRY entrypoint.

- The 8 bytes at 0xFFFF\_FFE0 (4GB-32) contain the physical address of the Firmware Interface Table.
- The 16 bytes at 0xFFFF\_FF00 (4GB-48) contain the FIT entry for the PAL\_A (or generic PAL\_A in the split PAL\_A model) code provided by the processor vendor. The format of this FIT entry is described in [Figure 11-6](#).
- The 8 bytes at 0xFFFF\_FFC8 (4GB-56) contains the physical address of the alternate Firmware Interface Table. This pointer is optional and is only needed if the firmware contains an alternate FIT table. If no alternate FIT table is provided a value of 0x0 should be encoded in this entry.
- The 8 bytes at 0xFFFF\_FFC0 (4GB-64) are zero-filled and reserved for future use.
- PAL\_A code (also known as generic PAL\_A code in split PAL\_A model) resides below 0xFFFF\_FFC0. This area contains the hardware-triggered entrypoints PALE\_RESET, PALE\_INIT, and PALE\_CHECK. In the model where PAL\_A is not split, the PAL\_A code will perform any processor-specific initialization needed in order for SAL to perform a firmware recovery. In the split PAL\_A model, the generic PAL\_A will search the FIT table(s) to find the first compatible and error-free processor-specific PAL\_A code. It will then branch to this code to perform the processor-specific initialization needed in order for SAL to perform a firmware recovery. The PAL\_A code area is a multiple of 16 bytes in length.
- SAL\_A code occupies the region immediately below the PAL\_A code. This area contains the SALE\_ENTRY entrypoint as well as optional implementation-independent firmware update code. The SAL\_A code area is a multiple of 16 bytes in length.
- The collection of regions above from the beginning of the SAL\_A code to 4GB is called the Protected Bootblock. The size of the Protected Bootblock is SAL\_A size + PAL\_A size + 64.
- The Firmware Interface Table (FIT) comprises of 16-byte entries containing starting address and size information for the firmware components. The FIT is generated at build time, based on the size and location of the firmware components. Optionally, an alternate FIT may be included in the firmware. The alternate FIT will only be used if the primary FIT failed its checksum. In the split PAL\_A model, this allows the generic PAL\_A firmware to find the processor-specific PAL\_A component(s), even if the primary FIT is corrupt. This feature allows hand-off to the SAL recovery code, even if there is a primary FIT checksum failure.
- The processor-specific PAL\_A contains the code that is required to be run before handing off to SAL for a firmware recovery check. This component is only available on processors that support a split PAL\_A firmware model. One processor-specific PAL\_A is architecturally required in this model. The firmware may optionally contain two or more processor-specific PAL\_A components.
- The PAL\_B block is comprised of code that is not required to be executed for SAL to perform a firmware recovery update. The PAL\_B code area is a multiple of 16 bytes in length. The PAL\_B block must be aligned on a 32K byte boundary or a 64K byte boundary depending on the implementation. Processor specific documentation provides the requirement for alignment. An OEM can choose to have more than one PAL\_B block in the firmware image.
- The remainder of the firmware address space is occupied by SAL\_B code. SAL\_B may include IA-32 BIOS code. The location of the SAL\_B and IA-32 BIOS code within the firmware address space is implementation dependent.

At a minimum, all of the PAL firmware components, pointers at the top of the firmware address space, FIT tables and the portion of the SAL code that is executed at the RECOVERY CHECK hand-off must be accessible from the processor without any special system fabric initialization sequence. This implies that the system fabric is implicitly initialized at power on for accessing the portions of the firmware address space listed above or that the special hardware which contains the firmware code and data is implemented on the processor and not accessed across the system fabric. The entire firmware code and data area can also be implicitly initialized at power on from the processor as well, but the minimum set is listed above.

The Firmware Interface Table (FIT) contains starting addresses and sizes for the different firmware components. Because these code blocks may be compiled at different times and places, code in one block (such as PAL\_A) cannot branch to code in another block (such as PAL\_B) directly. The FIT allows code in one block to find entrypoints in another. Figure 11-5 below shows the FIT layout.

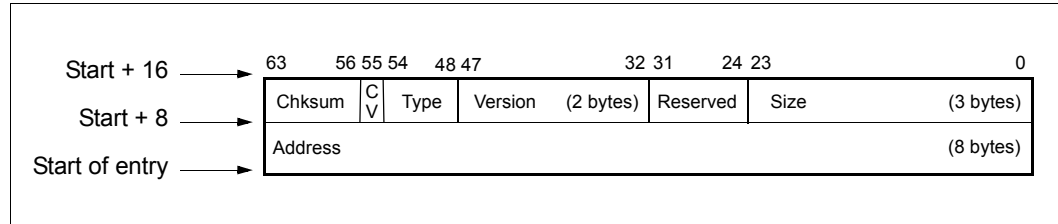
**Figure 11-5. Firmware Interface Table**



Each FIT entry contains information for the corresponding firmware component. The first entry contains size and checksum information for the FIT itself. The order of the following FIT entries must be arranged in ascending order by the type field, otherwise execution of firmware code will be unpredictable. Multiple FIT entries of the same type are allowed as shown in Figure 11-5.

When multiple entries of the same type exist for PAL components, PAL searches the FIT table in ascending order looking for the first entry that is compatible and error free for the processor it is currently executing on.

**Figure 11-6. Firmware Interface Table Entry**



- *Size* – A 3-byte field containing the size of the component in bytes divided by 16.
- *Reserved* – All fields listed as reserved must be zero filled.
- *Version* – A 2-byte field containing the component’s version number.
- *Type* – A 7-bit field containing the type code for the element. Types are defined in Table 11-1.

**Table 11-1. FIT Entry Types**

Type	Meaning
0x00	FIT Header
0x01	PAL_B (required)
0x02-0x0D	Reserved
0x0E	Processor-specific PAL_A
0x0F	PAL_A (also generic PAL_A) <sup>a</sup>
0x10-0x7E	OEM-defined
0x7F	Unused Entry

a. The PAL\_A FIT entry is located at 0xFFFF\_FFDO (4GB-48) and is not part of the actual FIT table.

OEMs may define unique types for one or more blocks of SAL\_B, IA-32 BIOS, etc., within the OEM-defined type range of 0x10 to 0x7E.

- *C\_V* – A 1-bit flag indicating whether the component has a valid checksum. If this field is zero, the value in the *Chksum* field is not valid.
- *Chksum* – A 1-byte field containing the component’s checksum. The modulo sum of all the bytes in the component and the value in this field (*Chksum*) must add up to zero. This field is only valid if the *C\_V* flag is non-zero. If the checksum option is selected for the FIT, in the FIT Header entry (FIT type 0), the modulo sum of all the bytes in the FIT table must add up to zero.

**Note:** The PAL\_A FIT entry is not part of the FIT table checksum.

- *Address* – An 8-byte field containing the base address of the component. For the FIT header, this field contains the ASCII value of “\_FIT\_<sp><sp><sp>” (<sp> represents the space character).

The FIT allows simpler firmware updates. Different components may be updated independently. This address layout can also support firmware images spanning multiple storage devices. FIT entries must be arranged in ascending order by the *type* field, otherwise execution of firmware code will be unpredictable.

## 11.2 PAL Power On/Reset

### 11.2.1 PALE\_RESET

The purpose of PALE\_RESET is to initialize and test the processor. Upon receipt of a power-on/reset event the processor begins executing code from the PALE\_RESET entrypoint in the firmware address space. PALE\_RESET initializes the processor and may perform a minimal processor self test. PAL may optionally perform authentication of the PAL firmware to ensure data integrity. If the authentication code runs cacheable by default, then a processor-specific mechanism will be provided to disable caching for diagnostic purposes.

PALE\_RESET then branches to SALE\_ENTRY to determine if a recovery condition exists, which would require an update of the firmware. If it does, SALE\_ENTRY performs the update and resets the system. If no firmware recovery is needed, SAL returns to PALE\_RESET to perform the processor self-tests and initialization. SAL can control the length and coverage of the PAL processor self-test by examining and modifying the self-test control word passed to SAL at the firmware recovery hand-off state. Please see Section 11.2.3, "PAL Self-test Control Word" for more information on the self-test control word.

The PAL processor self-tests are split into two phases. The first phase is written to test processor features that do not require external memory to be present to execute correctly. These tests are automatically run when SAL returns to PAL after the branch to SALE\_ENTRY for a firmware recovery check. This section is referred to as phase one of processor self-test and they are generally run early during the processor boot process. The second phase is written requiring that external memory is available to execute correctly. These tests are run when a call to the PAL procedure PAL\_TEST\_PROC is made with the correct parameters set up. These tests are referred to as phase two of processor self-test since they are usually run later in the processor boot process after external memory has been initialized on the platform.

PAL may execute IA-32 instructions to fully test and initialize the processor. This IA-32 code will not generate any special IA-32 bus transactions nor will it require any special platform features to correctly execute. PAL then branches to SALE\_ENTRY to conduct platform initialization and testing before loading the operating system software.

### 11.2.2 PALE\_RESET Exit State

- GRs: The contents of all general registers are undefined except the following:
  - GR20 (bank 1) contains the SALE\_ENTRY State Parameter as defined in [Figure 11-7](#). For the function field of the SALE\_ENTRY State Parameter, only the values 3, RECOVERY CHECK, for the first call to SALE\_ENTRY, and 0, RESET, for the second call to SALE\_ENTRY are valid.
  - GR32 contains 0 indicating that SALE\_ENTRY was entered from PALE\_RESET.
  - GR33 contains information about the geographically significant unique processor ID, and a mask that indicates which bits in the LID register (CR64) are read-only. Firmware should write the processor's local interrupt identifier in the programmable portion of the LID register. Writes to the read-only bits are ignored. See [Figure 11-8](#) for the definition of this parameter.

- GR34 contains the physical address for making a PAL procedure call. If the call is for RECOVERY CHECK, only the subset of PAL procedures needed for SALE\_ENTRY to perform firmware recovery will be available. These procedures are:
  - PAL\_FREQ\_RATIOS
  - PAL\_LOGICAL\_TO\_PHYSICAL
  - PAL\_PLATFORM\_ADDR
  - An implementation-specific PAL procedure for PAL authentication.
- GR35 contains the Self Test State Parameter as defined in [Figure 11-9](#).
- GR36 contains the PAL\_RESET return address for SALE\_ENTRY to return to if a recovery condition does not exist. When PAL\_RESET calls SALE\_ENTRY the second time to initialize the system for operating system use, this register will contain the physical address for making an implementation-specific PAL procedure call for PAL authentication.
 

**Note:** For all other PAL procedure calls, the physical address at GR34 should be used.
- GR37 contains the self-test control word as defined in [Figure 11-10](#). This control word is processor implementation-specific and informs SAL if self-test control is implemented and the number of controllable bits. If self-test control is implemented, PAL will read this value when SAL returns to PAL after firmware recovery check. If the self-test control is not supported, this register will be ignored when SAL returns to PAL after firmware recovery check.
- GR38 – Indicates if the PAL\_MEMORY\_BUFFER procedure is required to be called on this processor implementation for correct behavior. Also indicates the minimum buffer size required for the PAL\_MEMORY\_BUFFER procedure. [Table 11-2](#) defines the layout of this register.

**Table 11-2. GR38 Reset Layout**

Bit Field	Description
31:0	Unsigned integer denoting the minimum number of bytes required by the PAL_MEMORY_BUFFER procedure.
32:62	Reserved
63	Indicates if the PAL_MEMORY_BUFFER procedure is required by this processor implementation. A value of 1 indicates that it is required, a value of 0 indicates that it is not required.

- Banked GRs: All bank 0 general registers are undefined.
- FRs: The contents of all floating-point registers are undefined. The floating-point registers are enabled unless the *state* field of the Self Test State Parameter is FUNCTIONALLY RESTRICTED and the floating-point unit failed self test. Then, the floating-point registers are disabled. Refer to Section 11.2.2.3, “Definition of Self Test State Parameter” for the definition of FUNCTIONALLY RESTRICTED.
- Predicates: The contents of all predicate registers are undefined.
- BRs: The contents of all branch registers are undefined.
- ARs: The contents of all application registers are undefined except the following:
  - RSC: All fields in the register stack configuration register are 0, which places the RSE in enforced lazy mode.
- CFM: The CFM is set up so that all stacked registers are accessible, CFM.sof = 96 and all other CFM fields are 0.

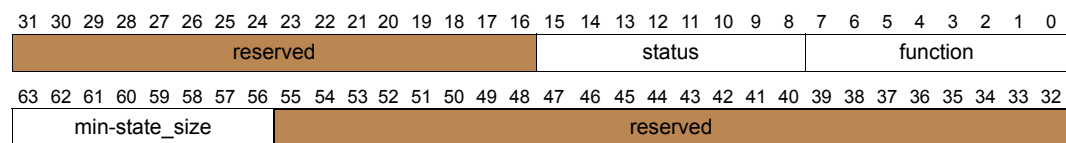
- PSR: PSR.bn is 1; PSR.df1 and PSR.dfh are 1 if the floating-point unit failed self test. All other PSR bits are 0. PSR.ic and PSR.i are zero to ensure external interrupts, NMI and PMI interrupts are disabled.
- CRs: The contents of all control registers are undefined except the following:
  - DCR: contains the value 0.
  - IVA: contains the physical address of an interruption vector table previously set up by PAL. SAL may choose to change this value. The IVA will be 0 when the SALE\_ENTRY State Parameter function is RECOVERY CHECK.
- RRs: The contents of all region registers are undefined.
- PKRs: The contents of all protection key registers are undefined.
- DBRs: The contents of all data breakpoint registers are undefined
- IBRs: The contents of all instruction breakpoint registers are undefined.
- PMCs: The contents of all performance monitor control registers are undefined.
- PMDs: The contents of all performance monitor data registers are undefined.
- Cache: The processor internal caches are enabled and invalidated. Unless directed otherwise by the self-test control word, phase one of the processor self-test verifies the caches themselves and the paths from the caches to the processor core. The path from external memory to the caches cannot be tested until phase two of the processor self-test.
 

**Note:** All cache contents will be invalidated when SAL returns to PAL after the RECOVERY\_CHECK hand-off. If the SAL uses the caches in their RECOVERY\_CHECK code, it is SAL's responsibility to write back any modified data in the caches before returning to PAL
- TLB: The TRs and TCs are initialized with all entries having been invalidated. The TLB is disabled because PSR.it=PSR.dt=PSR.rt=0. The TLBs cannot be fully tested until phase two of the processor self-test.

Prior to passing control to SALE\_ENTRY, PALE\_RESET must ensure that the processor Interrupt block pointer is set to point to address 0x0000\_0000\_FEE0\_0000.

### 11.2.2.1 Definition of SALE\_ENTRY State Parameter

**Figure 11-7. SALE\_ENTRY State Parameter**



- *function* – An 8-bit field indicating the reason for branching to SALE\_ENTRY.

**Table 11-3. function Field Values**

Function	Value	Description
RESET	0	System reset or power-on
MACHINE CHECK	1	Machine check event
INIT	2	Initialization event
RECOVERY CHECK	3	Check for recovery condition

All other values of *function* are reserved.

- *status* – A function-dependent 8-bit field indicating the firmware status on entry to SALE\_ENTRY. If the function value is RESET or RECOVERY\_CHECK, the *status* values are:

**Table 11-4. *status* Field Values**

Status	Value	Description
Normal	0	Normal reset.
FIT Header Failure	1	FIT header for FIT and alternate FIT (if supported) is incorrect
FIT Checksum Failure	2	FIT checksum for FIT and alternate FIT (if supported) is incorrect
PAL_B Checksum Failure	3	PAL_B checksum (for all compatible PAL_B's found) is incorrect
PAL_A Authentication Failure	4	PAL_A (generic in split model) failed authentication
PAL_B Authentication Failure	5	PAL_B (for all compatible PAL_B's found) failed authentication
PAL_B Not Found	6	FIT Entry for PAL_B missing from the FIT and alternate FIT (if supported)
Incompatible	7	No PAL_B was found in the FIT and alternate FIT (if supported) that is compatible with the processor stepping
32K Unaligned	8	No PAL_B was found in the FIT and alternate FIT (if supported) that was correctly aligned to a 32KB boundary
PAL_A_Spec Not Found / FIT Checksum Failure	9	No compatible processor-specific PAL_A was found in the FIT because of a FIT checksum failure and no compatible processor-specific PAL_A was found in the alternate FIT (if supported)
PAL_A_Spec Found / FIT Checksum Failure	10	A compatible processor-specific PAL_A was found in the alternate FIT. No compatible processor-specific PAL_A was found in the FIT due to a FIT checksum failure.
PAL_A_Spec Failure / Good PAL_A_Spec found in FIT	11	One or more compatible processor-specific PAL_A's found in the FIT failed its checksum or authentication. Another compatible processor-specific PAL_A was found in the FIT that passed its checksum and authentication.
PAL_A_Spec Auth Failure	12	No compatible processor-specific PAL_A's were found in the FIT or alternate FIT (if supported) that passed its checksum and authentication
PAL_A_Spec Auth Failure / Good PAL_A_Spec found in AF	13	One or more compatible processor-specific PAL_A's found in the FIT or alternate FIT (if supported) failed its checksum and authentication. Another compatible processor-specific PAL_A was found in the alternate FIT that passed its checksum and authentication.
PAL_A_Spec Not Found	14	No compatible processor-specific PAL_A was found in the FIT or alternate FIT (if supported)
PAL_A_Spec Not Found in FIT / Good PAL_A_Spec found in AF	15	No compatible processor-specific PAL_A was found in the FIT. A compatible processor-specific PAL_A was found in the alternate FIT.



**Table 11-4. status Field Values (Continued)**

Status	Value	Description
PAL_B Auth Failure / Good PAL_B found	16	One or more compatible PAL_B's failed authentication and checksum. Another compatible PAL_B was found that passed authentication and checksum.
64K Unaligned	17	No PAL_B was found in the FIT and alternate FIT (if supported) that was correctly aligned to a 64KB boundary.

All other values of *status* are reserved.

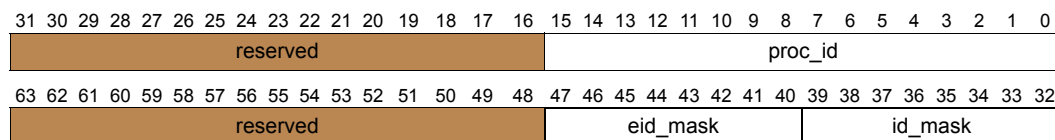
Definitions of *status* values for other values of *function* are listed in the machine check and init sections.

For the case of RECOVERY CHECK, authentication of PAL\_A and PAL\_B should be completed before call to SALE\_ENTRY.

- *min-state\_size* – An 8-bit field indicating the size in kilobytes (KB) of the min-state save area required for this implementation. A value of zero indicates a size of 4KB. A value greater than zero indicates the actual size in KB of the min-state save area required for this implementation. Values of 1-4 are reserved. For more information about the min-state save area, please refer to [Section 11.3.2.4, “Processor Min-state Save Area Layout”](#) on page 2:302.

### 11.2.2.2 Definition of Geographically Significant Processor Identifier Parameter

**Figure 11-8. Geographically Significant Processor Identifier**

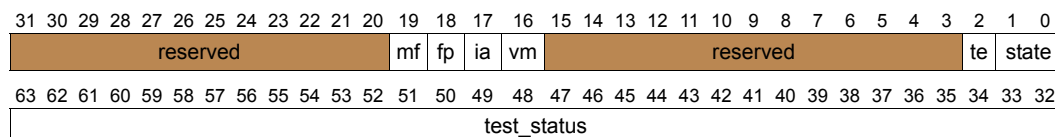


**Table 11-5. Geographically Significant Processor Identifier Fields**

Field	Bits	Description
proc_id	15:0	Geographically significant processor ID. The value returned in this field is the same as that returned by PAL_FIXED_ADDR.
Reserved	31:16	Reserved
id_mask	39:32	Mask indicating which bits in <i>id</i> are programmable: 0 = Programmable 1 = Read-only
eid_mask	47:40	Mask indicating which bits in <i>eid</i> are programmable: 0 = Programmable 1 = Read-only
Reserved	63:48	Reserved

### 11.2.2.3 Definition of Self Test State Parameter

**Figure 11-9. Self Test State Parameter**



- *state* – A 2-bit field indicating the state of the processor after self-test. If SAL directed PAL to skip some self-tests by modifying the self-test control word, failures related to these self-tests will not be reflected in this state.

**Table 11-6. *state* Field Values**

State	Value	Description
Catastrophic Failure	N/A	The processor is not capable of continuing. In this case it does not branch to SALE_ENTRY.
Healthy	00	No hardware failures have occurred in testing that would affect either the performance or functionality of the processor.
Performance Restricted	01	A hardware failure has occurred in testing that does not affect the functionality of the processor, but performance may be degraded.
Functionally Restricted	10	A hardware failure has occurred in testing that affects the functionality of the processor, but firmware code can still be run. The processor may also be performance restricted.

To further qualify FUNCTIONALLY RESTRICTED, the following requirements will be met:

- The processor has detected and isolated the failing component so that it will not be used.
- The processor must have at least one functioning memory unit, ALU, shifter, and branch unit.
- The floating-point unit may be disabled.
- The RSE is not required to work, but register renaming logic must work properly.
- The paths between the processor controlled caches and the register files have been shown to work. The path between the processor caches and memory cannot be validated until phase two of the processor self-test invoked by the PAL\_TEST\_PROC procedure.
- Loads and stores to firmware address space must work correctly.

Additional information about the failure can be obtained by examining the *test\_status* field of the *Self Test State Parameter*.

For the case of FUNCTIONALLY RESTRICTED, it is required that higher level firmware or OS not use failing functional units during their execution. PAL will not prevent failing functional units from being used.

- *te* – A 1-bit field indicating whether testing has occurred. If this field is zero, the processor has not been tested, and no other fields in the *Self Test State Parameter* are valid. The processor can be tested prior to entering SALE\_ENTRY for both RECOVERY CHECK and RESET functions.

If the *state* field indicates that the processor is functionally restricted, then the fields *vm*, *ia* & *fp* specify additional information about the functional failure.

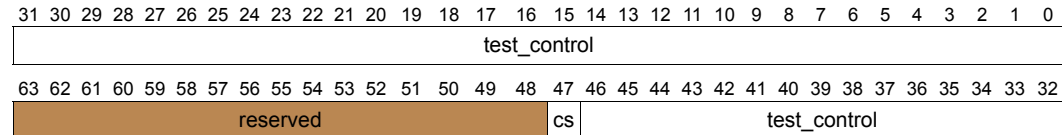
- *vm* – a 1-bit field, if set to 1, indicating that virtual memory features are not available
- *ia* – a 1-bit field, if set to 1, indicating that IA-32 execution is not available
- *fp* – a 1-bit field, if set to 1, indicating that floating-point unit is not available
- *mf* – a 1-bit field, if set to 1, indicating miscellaneous functional failure other than *vm*, *ia*, or *fp*. The *test\_status* field provides additional information about this failure on an implementation-specific basis.

- *test\_status* – An unsigned 32-bit-field providing additional information on test failures when the *state* field returns a value of PERFORMANCE RESTRICTED or FUNCTIONALLY RESTRICTED. The value returned is implementation dependent.

### 11.2.3 PAL Self-test Control Word

The PAL self-test control word is a 48-bit value. This bit field is defined in [Figure 11-10](#).

**Figure 11-10. Self-test Control Word**



- *test\_control* – This is an ordered implementation-specific control word that allows the user control over the length and runtime of the processor self-tests. This control word is ordered from the longest running tests up to the shortest running tests with bit 0 controlling the longest running test.

PAL may not implement all 47-bits of the *test\_control* word. PAL communicates if a bit provides control by placing a zero in that bit. If a bit provides no control, PAL will place a one in it.

PAL will have two sets of *test\_control* bits for the two phases of the processor self-test.

PAL provides information about implemented *test\_control* bits at the hand-off from PAL to SAL for the firmware recovery check. These *test\_control* bits provide control for phase one of processor self-test. It also provides this information via the PAL procedure call PAL\_TEST\_INFO for both the phase one and phase two processor tests depending on which information the caller is requesting.

PAL interprets these bits as input parameters on two occasions. The first time is when SAL passes control back to PAL after the firmware recovery check. The second time is when a call to PAL\_TEST\_PROC is made. When PAL interprets these bits it will only interpret implemented *test\_control* bits and will ignore the values located in the unimplemented *test\_control* bits.

PAL interprets the implemented bits such that if a bit contains a zero, this indicates to run the test. If a bit contains a one, this indicates to PAL to skip the test.

If the *cs* bit indicates that control is not available, the *test\_control* bits will be ignored or generate an illegal argument in procedure calls if the caller sets these bits.

- *cs* – Control Support: This bit defines if an implementation supports control of the PAL self-tests via the self-test control word. If this bit is 0, the implementation does not support control of the processor self-tests via the self-test control word. If this bit is 1, the implementation does support control of the processor self-tests via the self-test control word.

If control is not supported, GR37 will be ignored at the hand-off between SAL and PAL after the firmware recovery check and the PAL procedures related to the processor self-tests may return illegal arguments if a user tries to use the self-test control features.

## 11.3 Machine Checks

### 11.3.1 PALE\_CHECK

When a machine check abort (MCA) occurs, PALE\_CHECK is responsible for saving minimal processor state to a uncacheable platform-specific memory location previously registered with PAL via the PAL\_MC\_REGISTER\_MEM procedure. This platform location is called the Minimal State Save Area (min-state save area) and is described in [Section 11.3.2.4, "Processor Min-state Save Area Layout" on page 2:302](#). PALE\_CHECK is also responsible for correcting processor related errors whenever possible. PALE\_CHECK terminates either by returning to the interrupted context or by branching to SALE\_ENTRY, passing the state of the processor at the time of the error. The level of recovery provided by PALE\_CHECK is implementation dependent and is beyond the scope of this specification.

At the hand-off from PALE\_CHECK to SALE\_ENTRY, error information is passed in the Processor State Parameter described in [Section 11.3.2.1, "Processor State Parameter \(GR 18\)" on page 2:299](#). After exit from PALE\_CHECK, more detailed error information is available by calling the PAL\_MC\_ERROR\_INFO procedure. Information about implementation-dependent state is available by calling the PAL\_MC\_DYNAMIC\_STATE procedure. The interrupted process may be resumed by calling the PAL\_MC\_RESUME procedure. See [Section 11.3.3, "Returning to the Interrupted Process"](#) for more information on returning to the interrupted context and [Section 11.10, "PAL Procedures" on page 2:353](#) for detailed descriptions of all these procedure calls.

Code for handling machine checks must take into consideration the possibility that nested machine checks may occur. A nested machine check is a machine check that occurs while a previous machine check is being handled.

PALE\_CHECK is entered in the following conditions:

- When PSR.mc = 0 and an error occurs which results in a machine check, or
- When PSR.mc changes from 1 to 0 and there is a pending machine check from an earlier error.

PSR.mc is set to 1 by the hardware when PALE\_CHECK is entered. When PALE\_CHECK branches to SALE\_ENTRY, PSR.mc remains set (PSR.mc is restored to its original value if PALE\_CHECK terminates by returning to the interrupted context). SAL must not clear PSR.mc to 0 before all the information from the current machine check is logged. If SAL enables machine checks (by setting PSR.mc=0) during the SAL MCA handling, there is a potential for the error logs in the processor and the min-state save area to be overwritten by a subsequent MCA event.

The error information logged will reflect the state at the time the error occurred. State information from a different point in time will NOT be logged. If complete information is not available a code is logged which indicates that the information is not available.

- The processor state information used to resume a process for which an error has been corrected will reflect the state at the time the machine check interruption occurred and will be sufficient to resume the interrupted process.
- When a single error is signalled multiple times (for example, multiple operations to a single bad cache line), hardware and firmware will be able to perform the same logging and recovery as if the error had been signalled once.

For testing and configuration purposes, it may be necessary for software to intentionally generate a machine check. In this case PALE\_CHECK will log the error information, but not attempt recovery before branching to SALE\_ENTRY. To allow for this, the PAL\_MC\_EXPECTED procedure call is defined to indicate that PALE\_CHECK should not to attempt recovery.

### 11.3.1.1 Resources Required for Machine Check and Initialization Event Recovery

While the level of recovery from machine checks is implementation dependent, for each particular level of recovery there is a set of architecturally required resources. The following paragraphs define the required and optional resources needed to support firmware and software recovery of machine checks and initialization events.

- Minimal resources required to allow software recovery of machines checks when PSR.ic=1:
  - XR0 register: memory pointer to min-state save area previously registered with PAL via the PAL\_MC\_REGISTER\_MEM procedure. The layout of this memory area is described in [Section 11.3.2.4, "Processor Min-state Save Area Layout" on page 2:302](#).
  - Bank zero registers GR 24 through GR 31. These registers are not preserved across interruptions and may be used as scratch registers by machine check recovery code. See [Section 3.3.7, "Banked General Registers" on page 2:42](#) for the definition of the bank 0 registers.
- Additional resources required to allow software recovery of machine checks when PSR.ic=0. The presence of these resources is processor implementation specific. The PAL\_PROC\_GET\_FEATURES procedure described on [page 2:440](#) returns information on the existence of these optional resources.
  - XIP, XPSR, XFS: interruption resources implemented to store information about the IIP, IPSR and IFS when the machine check occurred. A model-specific version of the *rfi* instruction must also be implemented to restore the machine context from these resources.
  - XR1-XR3: scratch registers implemented to preserve bank 0 GR 24 through GR 31.

Each of the registers described above should be accessed only by PAL in order to support firmware and software recovery of machine checks.

### 11.3.2 PALE\_CHECK Exit State

The state of the processor on exiting PALE\_CHECK is listed below. For registers described as being saved to the min-state save area and available for use, the actual values in these registers are undefined unless specifically stated otherwise.

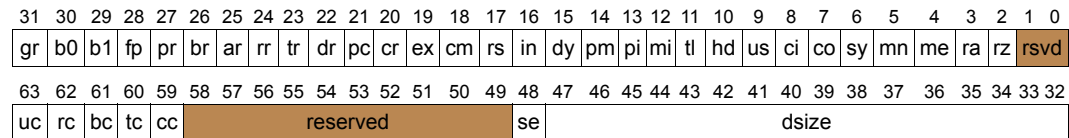
- GRs: The contents of all non-banked static registers (GR1-GR15), bank zero static registers and bank one static registers (GR16-31) at the time of the MCA have been saved in the min-state save area and are available for use.
  - If recovery is not supported when PSR.ic=0 then GR24 - GR31 (bank 0) are undefined and their contents have been lost. In this case, recovery is not possible. See [Section 11.3.1.1, "Resources Required for Machine Check and Initialization Event Recovery"](#) for details.

- GR16 through GR20 (bank 0) contain parameters which PALE\_CHECK passes to SALE\_ENTRY for diagnostic and recovery purposes:
  - GR16 contains the address to the first available location in the min-state save area for use by SAL. The address is 8-byte aligned.
  - GR17 contains the value of the min-state save area address stored in XR0.
  - GR18 contains the Processor State Parameter, as defined in [Figure 11-11](#).
  - GR19 contains the PALE\_CHECK return address for rendezvous, or 0 if no return is expected. (See Section 11.3.2.2, "Multiprocessor Rendezvous Requirements for Handling Machine Checks")
  - GR20 contains the SALE\_ENTRY State Parameter as defined in [Figure 11-4](#).
- FRs: The contents of all floating-point registers are unchanged from the time of the MCA.
- Predicates: All predicate registers have been saved in the min-state save area and are available for use.
- BRs: The contents of all branch registers are unchanged from the time of the MCA, except the following.
  - BR0 and BR1 have been saved to the min-state save area and are available for use. Either register may have been changed from the time of entry into PALE\_CHECK.
- ARs: The contents of all application registers are unchanged from the time of the MCA, except the RSE control register (RSC), the RSE backing store pointer (BSP), and the ITC and RUC counters. The RSC register is unchanged, except that the RSC.mode field will be set to 0 (enforced lazy mode) and the RSC register at the time of the MCA has been saved in the min-state save area. A cover instruction is executed in the PALE\_CHECK handler which allocates a new stack frame of zero size. BSP will be modified to point to a new location, since all the registers from the current frame at the time of interruption were added to the RSE dirty partition by the allocation of a new stack frame. The ITC register will not be directly modified by PAL, but will continue to count during the execution of the MCA handler. The RUC register will not be directly modified by PAL, but will continue to count during the execution of the MCA handler while the processor is active.
- CFM: The CFM register points to a zero-size current frame and all the rotating register bases are set to zero. The CFM register at the time of the MCA has been saved to the min-state save area in either the IFS or XFS slot depending on the implementation.
- RSE: Is in enforced lazy mode, and stacked registers are unchanged from the time of the MCA.
- PSR: PSR.mc is 1; PSR.mfl, PSR.mfh, and PSR.pk are unchanged; all other bits are 0. The PSR at the time of the MCA is saved in the min-state save area.
- CRs: The contents of all control registers are unchanged from the time of the MCA with the exception of interruption resources, which are described below.
- RRs: The contents of all region registers are unchanged from the time of the MCA.
- PKRs: The contents of all protection key registers are unchanged from the time of the MCA.
- DBR/IBRs: The contents of all breakpoint registers are unchanged from the time of the MCA.
- PMCs/PMDs: The contents of the PMC registers are unchanged from the time of the MCA. The contents of the PMD registers are not modified by PAL code, but may be modified if events it is monitoring are encountered.

- Cache: The processor internal cache is enabled and is unchanged from the time of the MCA except for any lines that were invalidated to correct the error.
- TLB: The TCs may be initialized and the TRs are unchanged from the time of the MCA.
- Interruption Resources:
  - IRR: PALE\_CHECK may not change the IRR, but interrupts may have arrived asynchronously, changing the contents of the IRRs.
  - The contents of IIP, IPSR and IFS at the time of the MCA are saved to the min-state save area and are available for use.

### 11.3.2.1 Processor State Parameter (GR 18)

**Figure 11-11. Processor State Parameter**



The term “valid” in [Table 11-7](#) indicates that the registers are either unchanged from the time of interruption or that the values have been preserved in the min-state save area.

**Table 11-7. Processor State Parameter Fields**

Field	Bits	Description
rsvd	1:0	Reserved
rz	2	The attempted processor rendezvous was successful if set to 1.
ra	3	A processor rendezvous was attempted if set to 1.
me	4	Distinct multiple errors have occurred, not multiple occurrences of a single error. Software recovery may be possible if error information has not been lost.
mn	5	Min-state save area has been registered with PAL if set to 1.
sy	6	Storage integrity synchronized. A value of 1 indicates that all loads and stores prior to the instruction on which the machine check occurred completed successfully, and that no loads or stores beyond that point occurred. See <a href="#">Table 11-8</a> .
co	7	Continuable. A value of 1 indicates that all in-flight operations from the processor where the machine check occurred were either completed successfully (such as a load), were tagged with an error indication (such as a poisoned store), or were suppressed and will be re-issued if the current instruction stream is restarted. This bit can only be set if the architectural state saved on a machine check is all valid. If this bit is set, then <i>us</i> must be cleared to 0, and <i>ci</i> must be set to 1. See <a href="#">Table 11-8</a> .
ci	8	Machine check is isolated. A value of 1 indicates that the error has been isolated by the system, it may or may not be recoverable. If 0, the hardware was unable to isolate the error within the CPU and memory hierarchy. The error may have propagated off the system (to persistent storage or the network). If <i>ci</i> = 0 then <i>us</i> will be set to 1, and <i>co</i> and <i>sy</i> are cleared to 0. See <a href="#">Table 11-8</a> .
us	9	Uncontained storage damage. A value of 1 indicates the error is contained within the CPU and memory hierarchy, but that some memory locations may be corrupt. If <i>us</i> is set to 1, then <i>co</i> and <i>sy</i> will always be cleared to 0. See <a href="#">Table 11-8</a> .
hd	10	Hardware damage. A value of 1 indicates that as a result of the machine check some non essential hardware is no longer available causing this processor to execute with degraded performance (no functionality has been lost).

**Table 11-7. Processor State Parameter Fields (Continued)**

Field	Bits	Description
tl	11	Trap lost. A value of 1 indicates the machine check occurred after an instruction was executed but before a trap that resulted from the instruction execution could be generated.
mi	12	More information. A value of 1 indicates that more error information about the machine check event is available by making the PAL_MC_ERROR_INFO procedure call.
pi	13	Precise instruction pointer. A value of 1 indicates that the machine logged the instruction pointer to the bundle responsible for generating the machine check.
pm	14	Precise min-state save area. A value of 1 indicates that the min-state save area contains the state of the machine for the instruction responsible for generating the machine check. When this bit is set, the <i>pi</i> bit will always be set as well.
dy	15	Processor Dynamic State is valid. (1=valid, 0=not valid) See the PAL_MC_DYNAMIC_STATE procedure call for more information.
in	16	Interruption caused by INIT. (0=machine check, 1=INIT)
rs	17	The RSE is valid. (1=valid, 0=not valid)
cm	18	The machine check has been corrected. (1=corrected, 0=not corrected)
ex	19	A machine check was expected. (1=expected, 0=not expected)
cr	20	Control registers are valid. (1=valid, 0=not valid)
pc	21	Performance counters are valid. (1=valid, 0=not valid)
dr	22	Debug registers are valid. (1=valid, 0=not valid)
tr	23	Translation registers are valid. (1=valid, 0=not valid)
rr	24	Region registers are valid. (1=valid, 0=not valid)
ar	25	Application registers are valid. (1=valid, 0=not valid)
br	26	Branch registers are valid. (1=valid, 0=not valid)
pr	27	Predicate registers are valid. (1=valid, 0=not valid)
fp	28	Floating-point registers are valid. (1=valid, 0=not valid)
b1	29	Preserved bank one general registers are valid. (1=valid, 0=not valid)
b0	30	Preserved bank zero general registers are valid. (1=valid, 0=not valid)
gr	31	General registers are valid. (1=valid, 0=not valid) (does not include banked registers)
dsize	47:32	Size in bytes of Processor Dynamic State returned by PAL_MC_DYNAMIC_STATE.
se	48	Shared Error. Machine check corresponds to structure shared by multiple logical processors.
rsvd	58:49	Reserved
cc	59	Cache check. A value of 1 indicates that a cache related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information. This bit must not be set for non-cacheable transaction errors.
tc	60	TLB check. A value of 1 indicates that a TLB related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
bc	61	Bus check. A value of 1 indicates that a bus related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
rc	62	Register file check. A value of 1 indicates that a register file related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
uc	63	Uarch check. A value of 1 indicates that a micro-architectural related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.



### 11.3.2.1.1 Using Processor State Parameter to Determine if Software Recovery of a Machine Check is Possible

The *us*, *ci*, *co*, and *sy* bits in the Processor State Parameter are valid only if the error has not been previously corrected in hardware or firmware (*cm* bit is 0). Even then, only the bit combinations shown in Table 11-8 are valid. If the multiple error bit is set (*me*=1) both the *co* and *sy* bits must be 0. The *us* and *ci* bits will be set according to the worst case of the errors that occurred.

**Table 11-8. Software Recovery Bits in Processor State Parameter**

cm	us	ci	co	sy	Description
1	x	x	x	x	The machine check is corrected. The <i>us</i> , <i>ci</i> , <i>co</i> , and <i>sy</i> bits are not valid.
0	1	0	0	0	The error was not isolated. Software must reset system. Data on disk may be corrupt.
0	1	1	0	0	The error was isolated but not contained. Corrupt data was not written to I/O, but may remain in the CPU or memory untagged. Software must reset system.
0	0	1	0	0	The error was isolated and contained, but is not continuable. The current instruction stream cannot be restarted without loss of information. Partial recovery may be possible.
0	0	1	1	0	The error was isolated, contained, and is continuable. If software can correct the error the current instruction stream can be continued with no loss of information.
0	0	1	1	1	The error was isolated, contained, and is continuable. The instruction pointer points to the instruction where the error occurred. If software can correct the error the current instruction stream can be continued with no loss of information.

### 11.3.2.2 Multiprocessor Rendezvous Requirements for Handling Machine Checks

When PALE\_CHECK has determined that an error has occurred which could cause a multiprocessor system to lose error containment, it must rendezvous the other processors in the system before proceeding with further processing of the machine check. This is accomplished by branching to SALE\_ENTRY with a non-zero return vector address in GR19. It is then the responsibility of SAL to rendezvous the other processors and return to PALE\_CHECK through the address in GR19. If the rendezvous was successful GR19 must be set to 0 before return.

At the time PALE\_CHECK makes the rendezvous call to SALE\_ENTRY, the processor state is exactly the same as defined in See "PALE\_CHECK Exit State" on page 2:297. with the following requirement on the use of registers by SAL:

Any processor state not listed below must be either unchanged or restored by SAL before returning to PALE\_CHECK.

- SAL will preserve the values in GR4-GR7 and GR17-GR18.
- SAL will return to PALE\_CHECK via the address in GR19.
- SAL will set up GR19 to indicate the success of the rendezvous before returning to PAL.
  - GR19 is zero to indicate the rendezvous was successful.
  - GR19 is non zero to indicate that the rendezvous was unsuccessful.
- All other non-banked (GR1-3, GR8-15), bank 0 GRs (GR20-GR31) and BR0 are undefined and available for use by SAL.

After return from the SAL rendezvous call, PALE\_CHECK will complete processing the machine check if the rendezvous was successful and then branch to SALE\_ENTRY with GR19 set to zero. The processor state when transferring to SAL is as defined in [Section 11.3.2, “PALE\\_CHECK Exit State” on page 2:297](#). If the rendezvous failed PALE\_CHECK will simply construct the Processor State Parameter and branch to SALE\_ENTRY.

Any further discussion of multiprocessor rendezvous, including platform requirements and implications, is beyond the scope of this specification. See the relevant SAL/Error handling documents for further information.

### 11.3.2.3 Unconsumed Data-Poisoning Event Handling

If, during the transfer/access of information between levels of the cache/memory hierarchy, there is data found to have an uncorrectable error and is marked poison, error reporting events may be raised. If such an error event is sent to a processor that doesn't consume the corrupted data, then the error is termed an **unconsumed data-poisoning event**.

Unconsumed data-poisoning events are by default reported as a CMC and can optionally be promoted to an MCA via bit 53 of *feature\_set* 0 of PAL\_PROC\_SET\_FEATURES. When they are signaled as a CMC the PSP.cm is set to 1 to indicate that the error has been corrected (in the sense that the line has been marked poison, preventing any silent data corruption).

If bit 53 is 1, unconsumed data-poisoning events are reported as MCAs. To immediately report unconsumed data-poisoning events as **uncorrected errors** (in the sense that the data in question has been lost), the caller can set bit 53 to 1. PSP settings for a data-poisoning event with bit 53 equal to 1 are given in the table below. See also [Table 11-8](#).

**Table 11-9. PSP Bit Settings for Unconsumed Data-poisoning Events on MCA**

cm	us	ci	co	sy
0	0	1	1	0

When promotion is enabled (bit 53 is 1), and a continuable data-poisoning event is indicated (i.e., the PSP bits are set as in the above table, and either *cache\_check.dp*, *bus\_check.dp* or both are 1), and if no other MCAs occur at the same time (i.e., no other errors are indicated in the error information from PAL\_MC\_ERROR\_INFO), the interrupted process is always continuable. Promotion to MCA with bit 53 allows the OS to take proactive measures to recover from the poisoned data, but this is not required for the interrupted process to be continuable.

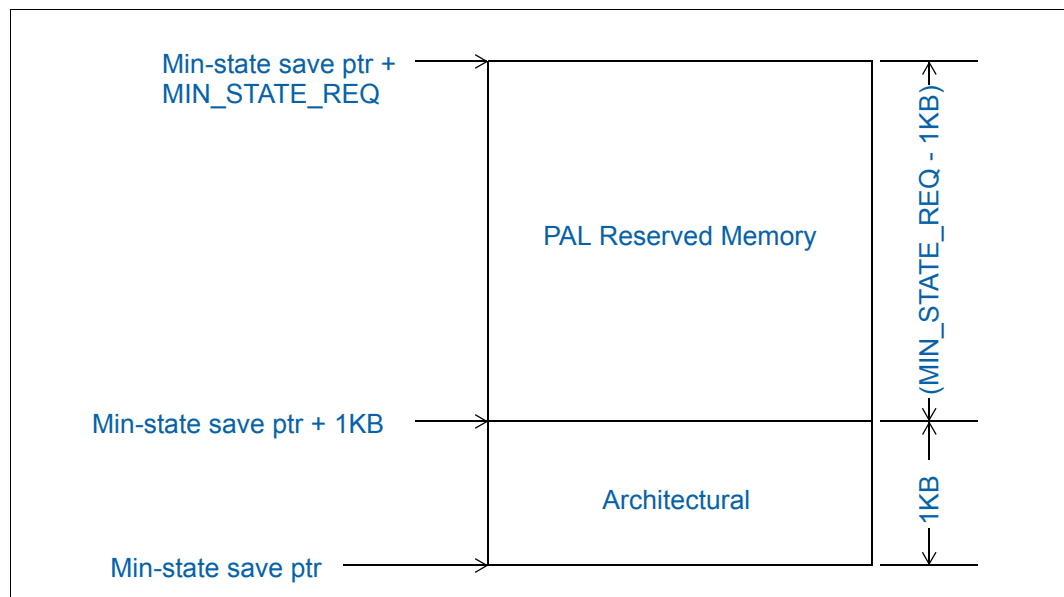
### 11.3.2.4 Processor Min-state Save Area Layout

The processor min-state save area is minimally 4KB in size, but an implementation may require larger sizes. The reset hand-off state indicates if a size greater than 4KB is required and also provides the required size. Please refer to [Section 11.2.2.1, “Definition of SALE\\_ENTRY State Parameter” on page 2:291](#) for more information on the reset hand-off state. The required size is referred to as MIN\_STATE\_REQ. The min-state save area is required to be in an uncacheable region. The first 1KB of this

area is architectural state needed by the PAL code to resume during MCA and INIT events (architected min-state save area + reserved). The remaining space in the buffer is a scratch space reserved exclusively for PAL use, therefore SAL and OS must not use this area. The layout of the processor min-state save area is shown in [Figure 11-1](#).

The processor min-state save area is 4KB in size and must be in an uncacheable region. The first 1KB of this area is architectural state needed by the PAL code to resume during MCA and INIT events (architected min-state save area + reserved). The remaining 3KB is a scratch buffer reserved exclusively for PAL use, therefore SAL and OS must not use this area. The layout of the processor min-state save area is shown in [Figure 11-1](#).

**Figure 11-1. Processor Min-state Save Area Layout**



The layout for the processors portion of the architectural 1KB processor min-state save area is shown in [Figure 11-2](#). When SAL registers the area with PAL, it passes in a pointer to offset zero of the area. When PALE\_CHECK is entered as a result of a machine check, it fills in processor state, processes the machine check, and branches to SALE\_ENTRY with a pointer to the first available memory location that SAL can use in GR16. SAL may allocate a variable sized area above the address passed in GR16 up to the 1KB architectural limit, but this is internal to SAL and not known to PAL.

The base address of the min-state save area must minimally be aligned to a 512-byte boundary, but larger alignments are allowed. All saves and restores to and from the min-state save area are made using 8-byte wide load and store instructions. If the processor min-state save area is not registered via the PAL\_MC\_REGISTER\_MEM procedure prior to the machine check, software recovery is not possible.



The NaT bits stored in the first entry of the min-state save area have the following layout.

**Figure 11-3. NaT Bits for Saved GRs**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
NaT bits for Bank 0 GR16 to GR31																NaT bits for GR15 to GR1											UD				
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
Undefined (not used)																NaT bits for Bank 1 GR16 to GR31															

**Table 11-10. NaT Bits for Saved GRs**

Bits	Description
0	Undefined (not used)
15:1	NaT bits for GR15 to GR1. Bit 1 represents GR1 and subsequent bits follow the ascending pattern.
31:16	NaT bits for Bank 0 GR16 to GR31. Bit 16 represents Bank 0 GR16 and subsequent bits follow the ascending pattern.
47:32	NaT bits for Bank 1 GR16 to GR31. Bit 32 represents Bank 1 GR16 and subsequent bits follow the ascending pattern.
63:48	Undefined (not used)

The value passed in GR16 to SAL may point beyond the defined processor state shown in Figure 11-2. PAL may use this area for implementation-dependent processor state that needs to be saved and restored.

### 11.3.2.5 Definition of SALE\_ENTRY State Parameter

**Figure 11-4. SALE\_ENTRY State Parameter**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
reserved																								function							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
reserved																															

- *function* – An 8-bit field indicating the reason for branching to SALE\_ENTRY.

**Table 11-11. function Field Values**

Function	Value	Description
RESET	0	System reset or power-on
MACHINE CHECK	1	Machine check event
INIT	2	Initialization event
RECOVERY CHECK	3	Check for recovery condition in SAL

All other values of *function* are reserved.

### 11.3.3 Returning to the Interrupted Process

The PAL\_MC\_RESUME procedure is defined to return to the interrupted context after handling a machine check or initialization event. See page 2:436 for a description of the PAL\_MC\_RESUME procedure. If software attempts to return to the interrupted context without using this procedure, processor behavior is undefined.

There are certain error cases that may require returning to a new context in order to recover from the machine check. If this occurs a new context can be returned to via the PAL\_MC\_RESUME procedure with the *new\_context* flag set. The caller needs to set up the new processor min-state save area as shown in [Figure 11-2](#) for all the listed register states. If the caller wants to return to a context where PSR.ic is zero (i.e., an interruption handler) the IIP, IPSR and IFS values in the min-state save area must be set up with the first level return values. These are the values for the IP, PSR and CFM of the interruption handler it wishes to return to. The XIP, XPSR, XFS values in the min-state save area must be set up with the second level return values. These are the IP, PSR and CFM values for where the interruption handler will return to. If the caller wants to return to a context where PSR.ic is one, it must set up the IIP, IPSR, IFS and the XIP, XPSR, and XFS both to contain the new instruction pointer, PSR value, and CFM values.

When returning to a new context, the memory area from BR1 up to the 1KB architectural limit is ignored by the PAL\_MC\_RESUME procedure. The software constructing the new context min-state save area does not have to worry filling in this memory area with any values. When a new context is returned to, the state originally saved in the min-state save area (old context) shall be discarded and never used again.

In order to return to the interrupted context without loss of any architectural state, the caller must restore all register state that is not stored in the processors min-state save area before making the PAL\_MC\_RESUME procedure call. Since BR0 and BR1 are the only two branch registers saved in the min-state save area, the caller must only use these two branch registers when making the PAL\_MC\_RESUME procedure call.

## 11.4 PAL Initialization Events

### 11.4.1 PALE\_INIT

PALE\_INIT is entered when an initialization event (INIT) occurs, as a result of the assertion on an INIT signal to the processor or an INIT interruption occurring. If PSR.mc = 1, the initialization event is held pending until PSR.mc becomes 0. The purpose of PALE\_INIT is to save the architecturally defined processor state to the Minimal State Save Area (min-state save area) and to branch to SALE\_ENTRY. The code sequence interrupted by the initialization event can be restarted via PAL\_MC\_RESUME if PSR.ic = 1. The code sequence interrupted by the initialization event can be restarted if PSR.ic = 0 and the processor has implemented the optional recovery resources described in [Section 11.3.1.1, “Resources Required for Machine Check and Initialization Event Recovery”](#) on page 2:297. If PSR.ic = 0 and the optional recovery resources have not been implemented, then the initialization event is not recoverable.

### 11.4.2 PALE\_INIT Exit State

The state of the processor on exiting PALE\_INIT is listed below. For registers described as being saved to the min-state save area and available for use, the actual values in these registers are undefined unless specifically stated otherwise.

- GRs: The contents of all non-banked static registers (GR1-GR15), bank zero static registers and bank one static registers (GR16-31) at the time of the INIT have been saved in the min-state save area and are available for use.

- If recovery is not supported when PSR.ic=0 then GR24 - GR31 (bank 0) are undefined and their contents have been lost. In this case, recovery is not possible. See Section 11.3.1.1, “Resources Required for Machine Check and Initialization Event Recovery” for details.
- GR16 through GR20 (bank 0) contain parameters which PALE\_INIT passes to SALE\_ENTRY for diagnostic and recovery purposes:
  - GR16 contains the address to the first available location in the min-state save area for use by SAL. The address is 8-byte aligned.
  - GR17 contains the value of the min-state save area address stored in XR0.
  - GR18 contains the Processor State Parameter, as defined in [Figure 11-5 on page 2:308](#).
  - GR19 contains the PALE\_INIT return address for rendezvous, or 0 if no return is expected. (See Section 11.3.2.2, “Multiprocessor Rendezvous Requirements for Handling Machine Checks”)
  - GR20 contains the SALE\_ENTRY state as defined in [Figure 11-4](#).
- FRs: The contents of all floating-point registers are unchanged from the time of the INIT.
- Predicates: All predicate registers have been saved in the min-state save area and are available for use.
- BRs: The contents of all branch registers are unchanged from the time of the INIT except the following:
  - BR0 and BR1 have been saved to the min-state save area and are available for use. Either register may have been changed from the time of entry into PALE\_CHECK.
- ARs: The contents of all application registers are unchanged from the time of the INIT, except the RSE control register (RSC), the RSE backing store pointer (BSP), and the ITC and RUC counters. The RSC register is unchanged, except that the RSC.mode field will be set to 0 (enforced lazy mode) and the RSC register at the time of the INIT has been saved in the min-state save area. A cover instruction is executed in the PALE\_INIT handler which allocates a new stack frame of zero size. BSP will be modified to point to a new location, since all the registers from the current frame at the time of interruption were added to the RSE dirty partition by the allocation of a new stack frame. The ITC register will not be directly modified by PAL, but will continue to count during the execution of the INIT handler. The RUC register will not be directly modified by PAL, but will continue to count during the execution of the INIT handler while the processor is active.
- CFM: The CFM register points to a zero-size current frame and all the rotating register bases are set to zero. The CFM register at the time of the INIT has been saved to the min-state save area in either the IFS or XFS slot depending on the implementation.
- RSE: The RSE is in enforced lazy mode, and all stacked registers are unchanged from the time of the INIT.
- PSR: PSR.mc is 1; PSR.mfl, PSR.mfh, and PSR.pk are unchanged; all other bits are 0. The PSR at the time of the INIT is saved in the min-state save area.
- CRs: The contents of all control registers are unchanged from the time of the INIT with the exception of the interruption resources, which are described below.
- RRs: The contents of all region registers are unchanged from the time of the INIT.
- PKRs: The contents of all protection key registers are unchanged from the time of the INIT.

- DBR/IBRs: The contents of all breakpoint registers are unchanged from the time of the INIT.
- PMCs/PMDs: The contents of the PMC registers are unchanged from the time of the INIT. The contents of the PMD registers are not modified by PAL code, but may be modified if events it is monitoring are encountered.
- Cache: The contents of the caches are unchanged from the time of the INIT.
- TLB: The TCs may be initialized and the TRs are unchanged from the time of the INIT.
- Interruption Resources:
  - IRR: PALE\_INIT may not change the IRR, but interrupts may have arrived asynchronously, changing the contents of the IRRs.
  - The contents of IIP, IPSR and IFS at the time of INIT are saved to the min-state save area and are available for use.

### 11.4.2.1 Processor State Parameter (GR18)

**Figure 11-5. Processor State Parameter**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
gr	b0	b1	fp	pr	br	ar	rr	tr	dr	pc	cr	ex	cm	rs	in	dy	pm	pi	mi	tl	hd	us	ci	co	sy	mn	me	ra	rz	rsvd	
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
uc	rc	bc	tc	cc	reserved										se	dsize															

The term “valid” in [Table 11-7](#) indicates that the registers are either unchanged from the time of interruption or that the values have been preserved in the min-state save area.

**Table 11-12. Processor State Parameter Fields**

Field	Bits	INIT value	Description
rsvd	1:0		Reserved
rz	2	x <sup>a</sup>	The attempted processor rendezvous was successful if set to 1.
ra	3	x <sup>a</sup>	A processor rendezvous was attempted if set to 1.
me	4	0	Distinct multiple errors have occurred, not multiple occurrences of a single error. Software recovery may be possible if error information has not been lost.
mn	5	x <sup>a</sup>	Min-state save area has been registered with PAL if set to 1.
sy	6	0	Storage integrity synchronized. A value of 1 indicates that all loads and stores prior to the instruction on which the machine check occurred completed successfully, and that no loads or stores beyond that point occurred. See <a href="#">Table 11-8</a> .
co	7	1	Continuable. A value of 1 indicates that all in-flight operations from the processor where the machine check occurred were either completed successfully (such as a load), were tagged with an error indication (such as a poisoned store), or were suppressed and will be re-issued if the current instruction stream is restarted. This bit can only be set if the architectural state saved on a machine check is all valid. If this bit is set, then <i>us</i> must be cleared to 0, and <i>ci</i> must be set to 1. See <a href="#">Table 11-8</a> .
ci	8	1	Machine check is isolated. A value of 1 indicates that the error has been isolated by the system, it may or may not be recoverable. If 0, the hardware was unable to isolate the error within the CPU and memory hierarchy. The error may have propagated off the system (to persistent storage or the network). If <i>ci</i> = 0 then <i>us</i> will be set to 1, and <i>co</i> and <i>sy</i> are cleared to 0. See <a href="#">Table 11-8</a> .



**Table 11-12. Processor State Parameter Fields (Continued)**

Field	Bits	INIT value	Description
us	9	0	Uncontained storage damage. A value of 1 indicates the error is contained within the CPU and memory hierarchy, but that some memory locations may be corrupt. If <i>us</i> is set to 1, then <i>co</i> and <i>sy</i> will always be cleared to 0. See <a href="#">Table 11-8</a> .
hd	10	0	Hardware damage. A value of 1 indicates that as a result of the machine check some non essential hardware is no longer available causing this processor to execute with degraded performance (no functionality has been lost).
tl	11	0	Trap lost. A value of 1 indicates the machine check occurred after an instruction was executed but before a trap that resulted from the instruction execution could be generated.
mi	12	0	More information. A value of 1 indicates that more error information about the machine check event is available by making the PAL_MC_ERROR_INFO procedure call.
pi	13	0	Precise instruction pointer. A value of 1 indicates that the machine logged the instruction pointer to the bundle responsible for generating the machine check.
pm	14	0	Precise min-state save area. A value of 1 indicates that the min-state save area contains the state of the machine for the instruction responsible for generating the machine check. When this bit is set, the <i>pi</i> bit will always be set as well.
dy	15	x <sup>a</sup>	Processor Dynamic State is valid. (1=valid, 0=not valid) See the PAL_MC_DYNAMIC_STATE procedure call for more information.
in	16	1	Interruption caused by INIT. (0=machine check, 1=INIT)
rs	17	x <sup>a</sup>	The RSE is valid. (1=valid, 0=not valid)
cm	18	0	The machine check has been corrected. (1=corrected, 0=not corrected)
ex	19	0	A machine check was expected. (1=expected, 0=not expected)
cr	20	x <sup>a</sup>	Control registers are valid. (1=valid, 0=not valid)
pc	21	x <sup>a</sup>	Performance counters are valid. (1=valid, 0=not valid)
dr	22	x <sup>a</sup>	Debug registers are valid. (1=valid, 0=not valid)
tr	23	x <sup>a</sup>	Translation registers are valid. (1=valid, 0=not valid)
rr	24	x <sup>a</sup>	Region registers are valid. (1=valid, 0=not valid)
ar	25	x <sup>a</sup>	Application registers are valid. (1=valid, 0=not valid)
br	26	x <sup>a</sup>	Branch registers are valid. (1=valid, 0=not valid)
pr	27	x <sup>a</sup>	Predicate registers are valid. (1=valid, 0=not valid)
fp	28	x <sup>a</sup>	Floating-point registers are valid. (1=valid, 0=not valid)
b1	29	x <sup>a</sup>	Preserved bank one general registers are valid. (1=valid, 0=not valid)
b0	30	x <sup>a</sup>	Preserved bank zero general registers are valid. (1=valid, 0=not valid)
gr	31	x <sup>a</sup>	General registers are valid. (1=valid, 0=not valid) (does not include banked registers)
dsiz	47:32	x <sup>a</sup>	Size in bytes of Processor Dynamic State returned by PAL_MC_DYNAMIC_STATE.
se	48	0	Shared Error. Machine check corresponds to structure shared by multiple logical processors.
rsvd	58:49		Reserved
cc	59	0	Cache check. A value of 1 indicates that a cache related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
tc	60	0	TLB check. A value of 1 indicates that a TLB related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
bc	61	0	Bus check. A value of 1 indicates that a bus related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.

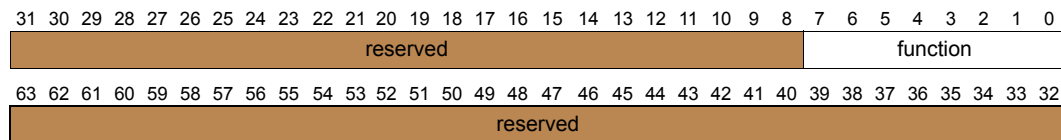
**Table 11-12. Processor State Parameter Fields (Continued)**

Field	Bits	INIT value	Description
rc	62	0	Register file check. A value of 1 indicates that a register file related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.
uc	63	0	Uarch check. A value of 1 indicates that a micro-architectural related machine check occurred. See the PAL_MC_ERROR_INFO procedure call for more information.

a. The values of the fields marked with x are set by the PAL INIT handler based on the INIT handling.

### 11.4.2.2 Definition of SALE\_ENTRY State Parameter

**Figure 11-6. SALE\_ENTRY State Parameter**



- *function* – An 8-bit field indicating the reason for branching to SALE\_ENTRY.

**Table 11-13. function Field Values**

Function	Value	Description
RESET	0	System reset or power-on
MACHINE CHECK	1	Machine check event
INIT	2	Initialization event
RECOVERY CHECK	3	Check for recovery condition in SAL

All other values of *function* are reserved.

## 11.5 Platform Management Interrupt (PMI)

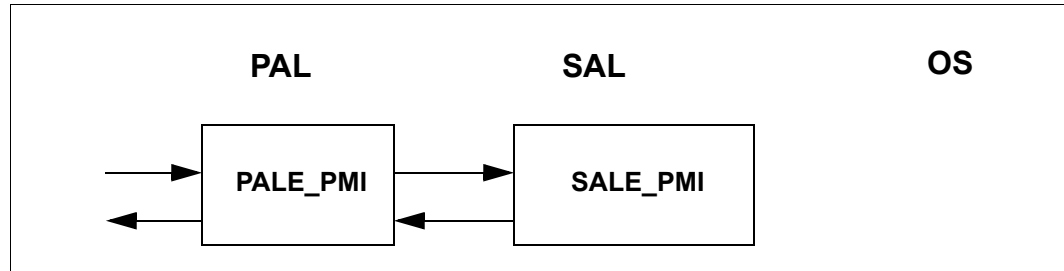
### 11.5.1 PMI Overview

PMI is an asynchronous interrupt that encapsulates a collection of platform-specific interrupts. Platform Management Interrupts occur during instruction processing, causing the flow of control to be passed to the PAL PMI handler. In the process, state is saved in the interruption registers (IIP, IPSR) by the processor hardware and the processor starts executing instructions at the PALE\_PMI entrypoint. The PAL code will save some additional state in the bank 0 registers. The PAL will either handle the PMI if it is PAL related PMI or transition to the SAL PMI code if it is a SAL related PMI. Upon completion of processing, the SAL PMI code returns to PAL PMI code to restore the interrupted processor state and to resume execution at the interrupted instruction.

As shown in [Figure 11-7](#), PMI code consists of two major components, namely the PAL PMI handler which handles all processor-specific processing, and the SAL PMI handler which handles all platform-related processing. The location of the PALE\_PMI and SALE\_PMI handlers are programmable. The location of the PALE\_PMI handler can be programmed by the PAL\_COPY\_PAL procedure described on [page 2:389](#). The SALE\_PMI handler can be programmed by the PAL\_PMI\_ENTRYPOINT procedure described on [page 2:443](#). If a PMI is taken very early in the boot sequence before PAL has a chance

to register its PALE\_PMI endpoint, processor operation is undefined. If a SAL related PMI is seen before the SAL PMI handler is registered, the PAL PMI code will just return to the interrupted context

**Figure 11-7. PMI Entrypoints**



The hardware events that can cause the PMI request are referred to as PMI events. PMI events are asynchronous interrupts higher priority than all external interrupts and are only maskable when the system software is processing very critical tasks with PSR.ic=0. When PSR.ic is 1, PMI events are unmasked. PSR.i has no effect on PMI events. All PMI events are internally latched into an array of implementation-specific latches in the processor. The PAL PMI handler reads the latches to determine what PMI vector requests are pending and dispatches them in priority order. Table 11-14 lists the PMI events and their priority.

**Table 11-14. PMI Events and Priorities**

PMI Events	Priority
PMI message for PAL (vectors 4-15)	High
PMI message for SAL (vectors 1-3)	
PMI pin <sup>a</sup> (vector 0)	Low

a. PMI pin is not required to be present on all systems.


PMI messages can be delivered by an external interrupt controller, or as an inter-processor interrupt using delivery mode 010. Table 11-15 shows the PMI message vector assignments. Vectors 4-15 are reserved for PAL, and within these PAL vectors, a higher vector number has higher priority. Vectors 1-3 are available for SAL to use, and within these SAL vectors, a higher vector number has higher priority. A PMI pin event, when the PMI pin<sup>1</sup> is present, is indicated by vector 0. The PMI vector number is passed to the SAL PMI handler in GR 24.

**Table 11-15. PMI Message Vector Assignments**

Priority		Vector	Description
Low	SAL Vectors	0	PMI pin
↓		1	Available for SAL firmware
		2	
		3	
High			

1. PMI pin is not required to be present. Software can query the presence of PMI pin via the PAL\_PROC\_GET\_FEATURES procedure call.

**Table 11-15. PMI Message Vector Assignments**

Priority		Vector	Description	
Low  High	PAL Reserved	4	PAL Reserved	
		5		
		6		
		7		
		8		
		9		
		10		
		11		
		12		
		13		IA-32 Machine Check Rendezvous
		14		PAL Reserved
		15		

### 11.5.2 PALE\_PMI Exit State

The state of the processor on exiting PALE\_PMI is:

- GRs: The contents of non-banked general registers are unchanged from the time of the interruption.
  - Bank 1 GRs: The contents of all bank one general registers are unchanged from the time of the interruption.
  - Bank 0:GR16-23: The contents of these bank zero general registers are unchanged from the time of the interruption.
  - Bank 0:GR24-31: contain parameters which PALE\_PMI passes to SALE\_PMI:
    - GR24 contains the value decoded as follows:
      - Bits 7-0: PMI Vector Number
      - Bit 63-8: Reserved
    - GR25 contains the value of the min-state save area address stored in XR0.
    - GR26 contains the value of saved RSC. The contents of this register shall be preserved by SAL PMI handler.
    - GR27 contains the value of saved B0. The contents of this register shall be preserved by SAL PMI handler.
    - GR28 contains the value of saved B1. The contents of this register shall be preserved by SAL PMI handler.
    - GR29 contains the value of the saved predicate registers. The contents of this register shall be preserved by SAL PMI handler
    - GR30-31 are scratch registers available for use.
- FRs: The contents of all floating-point registers are unchanged from the time of the interruption.
- Predicates: The contents of all predicate registers are undefined and available for use.
- BRs: The contents of all branch registers are unchanged, except the following which contain the defined state.
  - BR1 is undefined and available for use.

- BR0 PAL PMI return address.
- ARs: The contents of all application registers are unchanged from the time of the interruption, except the RSE control register (RSC) and the ITC and RUC counters. The RSC.mode field will be set to 0 (enforced lazy mode) while the other fields in the RSC are unchanged. The ITC register will not be directly modified by PAL, but will continue to count during the execution of the PMI handler. The RUC register will not be directly modified by PAL, but will continue to count during the execution of the PMI handler while the processor is active.
- CFM: The contents of the CFM register is unchanged from the time of the interruption.
- RSE: Is in enforced lazy mode, and stacked registers are unchanged from the time of the interruption.
- PSR: PSR.mc, PSR.mfl, PSR.mfh, and PSR.pk are unchanged; all other bits are 0.
- CRs: The contents of all control registers are unchanged from the time of the interruption with the exception of interruption resources, which are described below.
- RRs: The contents of all region registers are unchanged from the time of the interruption.
- PKRs: The contents of all protection key registers are unchanged from the time of the interruption.
- DBR/IBRs: The contents of all breakpoint registers are unchanged from the time of the interruption.
- PMCs/PMDs: The contents of the PMC registers are unchanged from the time of the PMI. The contents of the PMD registers are not modified by PAL code, but may be modified if events it is monitoring are encountered
- Cache: The processor internal cache is not specifically modified by the PMI handler but may be modified due to normal cache activity of running the handler code.
- TLB: The TCs are not modified by the PALE\_PMI handler and the TRs are unchanged from the time of the interruption.
- Interruption Resources:
  - IRRs: The contents of IRRs are unchanged from the time of the interruption.
  - IIP and IPSR contain the value of IP and PSR. The IFS.v bit is reset to 0.

### 11.5.3 Resume from the PMI Handler

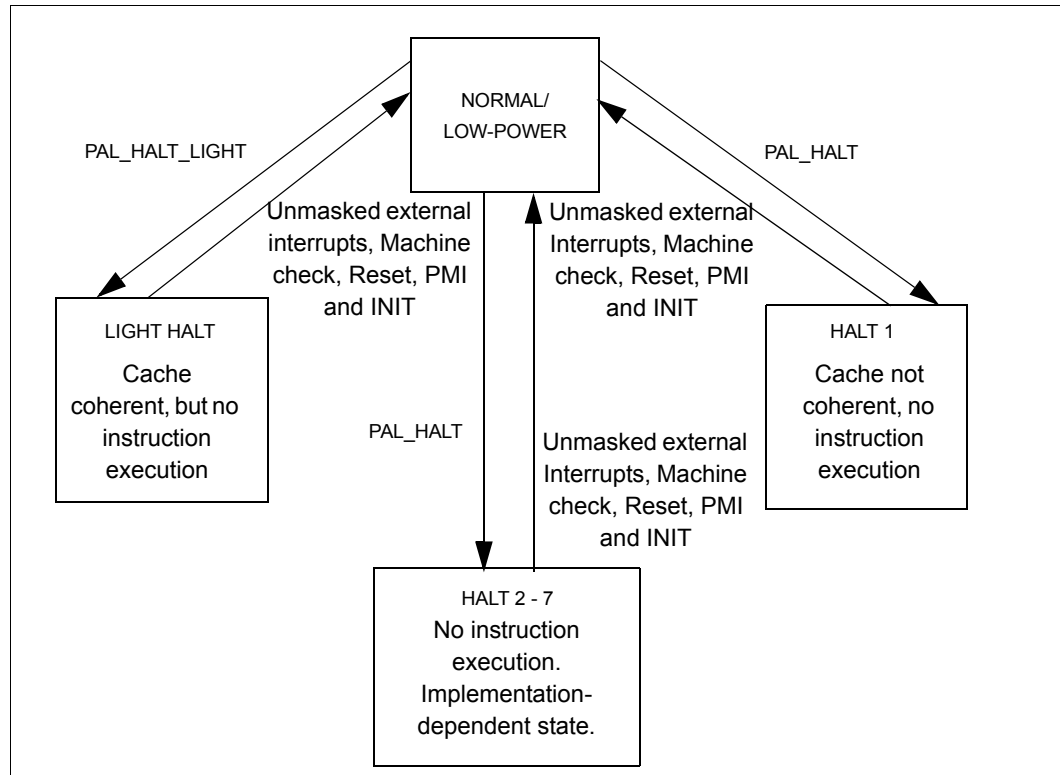
To return to the instruction that was interrupted by the PMI event, SAL PMI must branch to the PAL PMI target address in BR0. All register contents must be preserved as specified in [Section 11.5.2, "PALE\\_PMI Exit State" on page 2:312](#).

## 11.6 Power Management

This section describes the architecturally supported set of required and optional power states that may be implemented to reduce power consumption in implementations where this is a design goal. In addition, the PAL interfaces required to manage these states are described.

Figure 11-8 shows state transitions for the various power states and the software interfaces required for the transitions.

**Figure 11-8. Power States**



- **NORMAL** – The normal, fully functional, highest power state.
- **LOW-POWER** – An implementation may choose to dynamically reduce power via microarchitectural low power techniques. The operation of interrupts, snoops, etc., in low-power mode will be identical to those in normal-power mode. This dynamic power reduction is optional for an implementation to support. The PAL procedures `PAL_PROC_GET_FEATURES` and `PAL_PROC_SET_FEATURES` returns whether an implementation supports dynamic power reduction. If an implementation supports dynamic power reduction then this procedure will allow the caller to enable or disable this feature.

The following software controllable low power states may be provided. They are described below.

- **LIGHT\_HALT** – Entered by calling `PAL_HALT_LIGHT`. This state reduces power by stopping instruction execution, but maintains cache and TLB coherence in response to external requests. The processor transitions from this state to the **NORMAL** state in response to any unmasked external interrupt (including NMI), machine check, reset, PMI or INIT. An unmasked external interrupt is defined to be an interrupt that is permitted to interrupt the processor based on the current setting of the `TPR.mic` and `TPR.mmi` fields. This state is a required state.
- **HALT 1** – Entered by calling `PAL_HALT` with a power state argument equal to one. This implementation-dependent low-power state will maintain the processor caches but will ignore any coherency bus traffic. This state is optional for a processor to

implement. It is the responsibility of the caller to ensure cache coherency in this state.

- HALT 2 - 7 – These are optional implementation-dependent states entered by calling PAL\_HALT with a power state argument in the range of 2-7. Before making this procedure call, the operating system software should first ascertain that the states are implemented by calling PAL\_HALT\_INFO. The information returned from the PAL\_HALT\_INFO procedure will also specify the coherency of caches and TLBs for each of these low-power states.

The interval timer within the processor will function at a constant frequency in all the power states as long as the input clock to the processor is maintained. If all logical processors on the physical processor are in a halt state, the resource utilization counter for the last logical processor to enter a halt state will function at a constant frequency as long as the input clock to the processor is maintained. However, the performance monitor event that counts the number of processor clock cycles will only increment in either the NORMAL or LOW-POWER state.

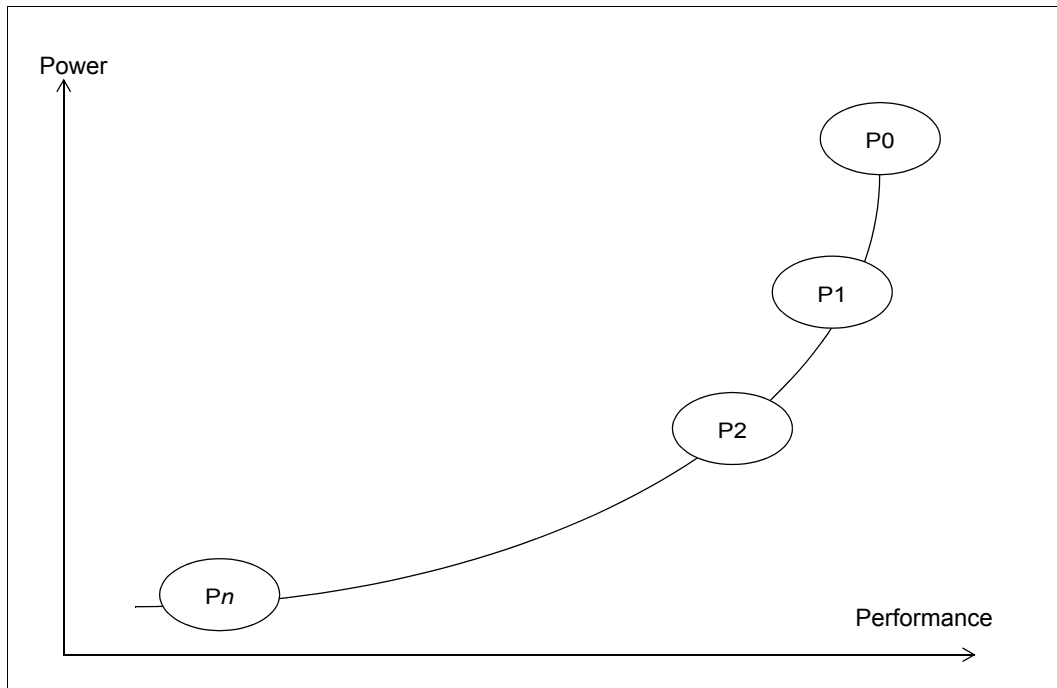
The PAL procedure PAL\_HALT\_INFO returns information about the power states implemented in a particular processor. This information allows the caller to decide which low power states are implemented and which ones to call based on the callers requirements.

### **11.6.1 Power/Performance States (P-states)**

This section describes the power/performance states (hence to be referred to as P-states) supported by the Itanium architecture. P-states enable the caller to adjust the power/performance characteristics of the processor in response to changing workload requirements. This allows for implementation of a processor-level power management policy which is driven by system demand and response time requirements.

The P-states are defined within the context of the active/executing processor state. At the highest performing P-state (referred to as the P0 state), the processor uses its maximum performance capability and may consume maximum power. In the next P-state (P1), the processor performance capability is limited below the maximum performance, and it consumes less than the maximum power. Successive P-states continue to have reduced performance capabilities and reduced power consumption. The Itanium architecture supports a maximum of 16 P-states, with the highest numbered P-state that is available on an implementation providing the least possible performance capability and minimal power consumption while remaining in a non-HALT state.

**Figure 11-9. Power and Performance Characteristics for P-states**

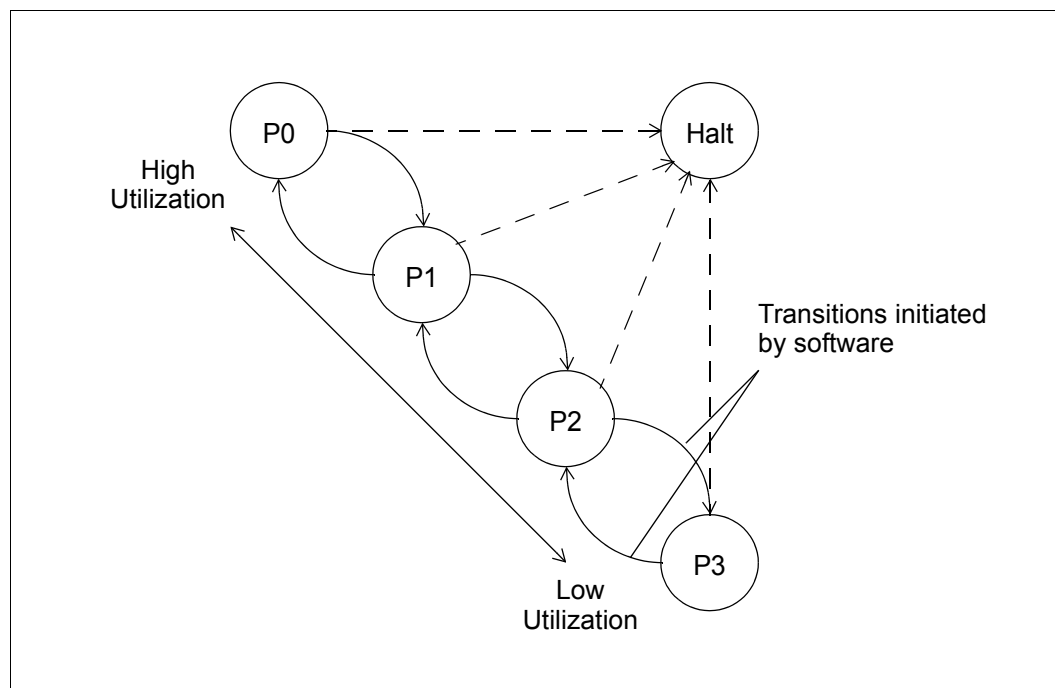


P-states can be utilized by software to implement a demand-based dynamic power management policy where it would continuously try to adapt the processor performance to the current workload characteristics. This allows software to achieve power savings at the system level, while allowing it to quickly respond to changing workload requirements.

The example in [Figure 11-10](#) assumes four P-states (P0, P1, P2 and P3), and a software policy that transitions between the states depending on the current system utilization. During times of high utilization, the software migrates the processor towards lower-numbered P-states, which increases processor performance and increases the dissipated power. When system utilization is low, the software policy migrates the processor towards higher-numbered P-states, thereby reducing the processor performance and reducing dissipated power. The figure also shows the HALT state, which the software can transition to at any time from a given P-state.



**Figure 11-10. Example of a P-state Transition Policy**



### 11.6.1.1 Power Dependency Domains

The concept of P-states applies to each logical processor, and this gives software the required granularity to individually control the power/performance characteristics for each available thread of execution in the system. In the most simplistic case, the processor package has only one thread of execution, and this allows software to apply the same P-state policy at the package-level as well as at the logical processor level. However, with implementations that support multithreading and multiple cores, a single package can have multiple logical processors (threads of execution). These may have P-state dependencies among them, which may not allow for individual P-state control flexibility at the software level. For example, these logical processors may be sharing the same clock and power delivery network. In such circumstances, software would need to know which logical processors have dependencies and what the nature of the dependencies is, so that appropriate coordination techniques can be applied. To allow the architecture definition to comprehend multi-threaded/multi-core designs, we define the concept of dependency domain and coordination mechanisms.

A **dependency domain** is comprised of logical processors that share a common set of implementation-dependant domain parameters that affect power consumption and performance for all logical processors in that domain. As an example, a processor package comprised of two cores controlled by the same clock and power distribution network are part of the same dependency domain, since changing either the operating frequency or voltage will affect power consumption and performance for both cores. Alternatively, if these two cores on the processor package had independent distribution networks for clocks and power, then a change in the parameters for one core would not have any effect on the other core, and in that case, the cores would not belong to the same dependency domain. Software can utilize P-states to effect changes in the domain

parameters. Each P-state maps to a set of values for the domain parameters, and hence a P-state transition results in a change in the underlying power/performance characteristics for the logical processor.

The Itanium architecture supports different types of dependency domains, which enables software to have different degrees of control for P-state changes affecting logical processors in the domain.

A **software-coordinated dependency domain (SCDD)** relies on the software to coordinate P-state changes among the processors in that dependency domain. Software will have knowledge about logical processors belonging to that domain, and will decide when it is appropriate to request the P-state transition. The software policy has to be aware that a P-state change on any logical processor will change the P-state for all logical processors in that domain. As an example, let us assume that the SCDD consisted of two cores with the same clock and power distribution networks and the intent of the software policy was to lower power/performance only when the workload utilization was low on both cores. Software could then monitor utilization on both cores, and when both cores were under-utilized (i.e., were running at a higher performance P-state than required by the current system demand), it could migrate one of the cores to a lower performance P-state. This transition would simultaneously reduce performance and power dissipation for both cores, and would result in both cores operating at the same lower P-state.

A **hardware-coordinated dependency domain (HCDD)** relies on hardware-based mechanisms to synchronize P-state changes. Software can make independent P-state change requests on individual processors, recognizing that hardware is responsible for the required coordination with other processors in the same HCDD. Hardware-based coordination mechanisms would be implemented to allow for changes to the logical processor's power and performance local parameters (which are implementation-dependant), in addition to the existing domain parameters. Hardware would use a combination of changes to both of these parameters to satisfy the software-initiated P-state change request. This type of coordination mechanism is effective when it is desired to have individual control over all logical processors, and when the hardware has local parameters for power/performance at the logical processor level. The local parameters allow for fine-grained control (affecting only the logical processor power/performance), whereas the domain parameters allow for coarse-grained control (affecting all logical processors). Domain parameters are set by hardware according to the highest requested power/performance level (i.e., the lowest numbered P-state) of the logical processors in the power domain. As an example, let us assume that the HCDD consisted of two cores with the same clock and power distribution networks, and that there were also some other techniques to affect power and performance which were local to each logical processor. Let us also assume that software has initially set both cores to the P0 state. When software initiates a P-state transition to P1 (which is a lower power/performance level) on the first core, hardware would use only the local parameters to carry out the request, and the domain parameters would remain at P0. Suppose software on the second core then initiates a P-state transition to P3. Hardware would then set the local parameters for the second core to reflect this request, undo the changes to the local parameters for the first core plus initiate changes to the domain parameters to transition the domain to the P1 state (the highest requested power/performance level of the two cores).

A **hardware-independent dependency domain (HIDD)** is a self-contained domain that typically means that every logical processor is the only logical processor in that domain, and its domain parameters are individually controllable. Since there are no dependencies with any other logical processors, there is no P-state coordination needed for such domains. Software can make P-state change requests independently on that logical processor.

### 11.6.1.2 Platform Power-Cap and P-states

Some processor implementations include mechanisms which allow the platform hardware and firmware to temporarily decrease the operating frequency of logical processors, to implement fast-response power capping. This is referred to as a **Platform Power-Cap**. In such implementations, the P-state requested by software is not changed by the platform power-cap. Software is able to change its P-state request during platform power-caps; when the platform power-cap is removed, the processor operating frequency returns to the frequency determined by software's most recent P-state settings.

Platform power caps are meant to have a very short duration and very low duty cycle so they do not significantly affect software methods for managing power through P-states. Platform power-caps do not affect the instantaneous operating P-state observed by software, but do affect the weighted-average performance index reported to software by PAL, so that software may take into account any small effects. (See the `PAL_GET_PSTATE` procedure for details.)

### 11.6.1.3 PAL Interfaces for P-states

The PAL procedure `PAL_PROC_GET_FEATURES` returns whether an implementation supports P-states. If an implementation supports P-states then the `PAL_PROC_SET_FEATURE` procedure will allow the caller to enable or disable this feature.

The Itanium architecture provides three PAL procedures to enable P-state functionality.

**PAL\_PSTATE\_INFO:** This procedure returns information about the P-states implemented on a particular processor. For details on the information returned by this procedure, please refer to the procedure description on [page 2:396](#). The Itanium architecture supports a maximum of 16 P-states.

**PAL\_SET\_PSTATE:** This procedure allows the caller to request the transition of the processor to a new P-state. The procedure can either return with transition success (request was accepted) or transition failure (request was not accepted) depending on hardware capabilities, implementation-specific event conditions, and the spacing between successive `PAL_SET_PSTATE` procedure calls.

If hardware has the ability to either preempt a previous in-progress P-state transition, or to queue successive P-state requests while the first request is in transition, then the implementation has a pre-emptive policy for P-state request handling. The architecture also allows for a non-preemptive policy for P-state request handling, whereby a new `PAL_SET_PSTATE` request is not accepted if a previous P-state transition is already in progress. The `PAL_SET_PSTATE` procedure returns different status values corresponding to the accepted and not accepted cases for P-state requests. If the transition is not accepted, no P-state transition is initiated by the `PAL_SET_PSTATE`

procedure, and the caller is expected to make another PAL\_SET\_PSTATE request to transition to the desired P-state. The *transition\_latency\_2* field in the *pstate\_buffer* returned by PAL\_PSTATE\_INFO indicates the time interval the caller needs to wait to have a reasonable chance of success when initiating another PAL\_SET\_PSTATE call.

Implementation-specific event conditions may prevent a PAL\_SET\_PSTATE request from being accepted (e.g., due to a thermal protection mechanism), in which case the PAL procedure returns a status of *transition failure*. Such events are expected to be rare and to happen only in abnormal situations.

It should be noted that platform power-caps do not cause a PAL\_SET\_PSTATE request to fail. The requested P-state is registered with PAL, and the procedure returns a status of *transition success*.

SCDD: If the logical processor belongs to a software-coordinated dependency domain, the PAL\_SET\_PSTATE procedure will change the domain parameters resulting in a transition to the requested P-state for all logical processors in that domain.

HCDD: If the logical processor belongs to a hardware-coordinated dependency domain, the PAL\_SET\_PSTATE procedure will attempt to change the power/performance characteristics for that logical processor. Since the power/performance characteristics for the domain depend on the P-state settings of the other logical processors in the domain, a PAL\_SET\_PSTATE call on one logical processor may result in either partial or complete transition to the requested P-state. In case of partial transition (see [Figure 11-11, "Computation of performance\\_index" on page 2:321](#) for an example, where the logical processor transitions from state P0 to state P3 in partial increments), the logical processor may attempt to perform changes at a later time to the local parameters and/or domain parameters to transition to the originally requested P-state based on P-state transition requests on other logical processors. Software can also approximate the behavior of a SCDD by forcing P-state transitions. See the description of the PAL\_SET\_PSTATE procedure for more details.

HIDD: If the logical processor belongs to a hardware-independent dependency domain, the PAL\_SET\_PSTATE procedure will attempt to change the domain parameters, which will transition the logical processor in that domain to the requested P-state.

**PAL\_GET\_PSTATE:** This procedure returns the performance index of the logical processor, relative to the highest available P-state (P0). A value of 100 in P0 represents the minimum processor performance in the P0 state. For example, if the value returned by the procedure is 80, this indicates that the performance of the logical processor over the last time period was 20% lower than the minimum P0 performance. For processors that support variable P-states, it is possible for a processor to report a number greater than 100, representing that the processor is running at a performance level greater than the minimum P0 performance. For example, if the value returned by the processor is 120, it indicates that the performance of the logical processor over the last time period was 20% higher than the minimum P0 performance. The performance index is measured over the time interval since the last PAL\_GET\_PSTATE call with a type operand of 1. If the processor supports variable P-state performance then the PAL\_PROC\_SET\_FEATURE procedure can be used to enable or disable this feature. Software may choose, on each invocation of the PAL\_GET\_PSTATE procedure, whether to reset the internal performance measurement logic; resetting the measurement logic

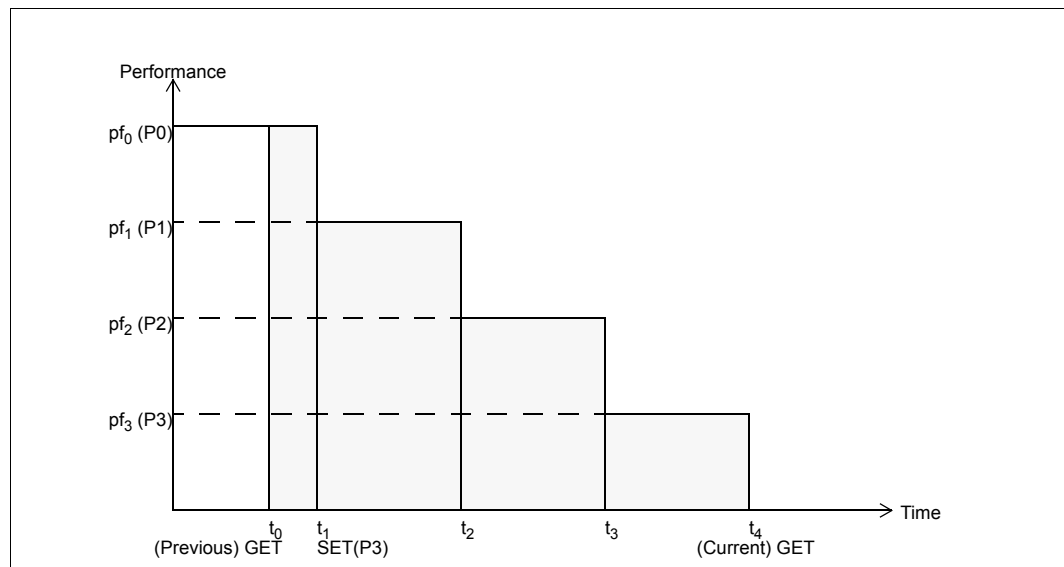
initiates a new *performance\_index* count, which is reported when the next PAL\_GET\_PSTATE procedure call is made. A call to PAL\_GET\_PSTATE with a *type* operand of 1 resets the performance measurement logic.

SCDD: If the logical processor belongs to a software-coordinated dependency domain, the performance index returned (for either *type*=0 or 3) corresponds to the target P-state requested by the most recent successful PAL\_SET\_PSTATE procedure call. No weighted average (*type*=1 or 2) is computed by PAL; calling PAL\_GET\_PSTATE with *type*=1 or 2 on a SCDD logical processor is undefined.

HCDD: If the logical processor belongs to a hardware-coordinated dependency domain, the performance index returned (*type*=1 or 2) will be a weighted-average sum of the *performance\_index* values corresponding to the different P-states that the logical processor was operating in since performance measurement was last reset. Note that this return value may not necessarily correspond to the performance index of the target P-state requested by the most recent PAL\_SET\_PSTATE procedure call. For example, let's assume that the previous PAL\_GET\_PSTATE procedure was called at time  $t_0$ , when the processor was operating in state P0. The previous PAL\_SET\_PSTATE procedure requested a transition from P0 to P3. The transition happened over a period of time, such that the logical processor went through states P1 at time  $t_1$ , P2 at time  $t_2$  and P3 at time  $t_3$ , and was in state P3 at time  $t_4$  when the current PAL\_GET\_PSTATE procedure was called. The *performance\_index* returned is calculated as:

$$\begin{aligned}
 \text{performance\_index} = & \\
 & ((\text{time spent in P0 after the previous PAL\_GET\_PSTATE}) * (\text{performance\_index for P0}) + \\
 & (\text{time spent in P1}) * (\text{performance\_index for P1}) + \\
 & (\text{time spent in P2}) * (\text{performance\_index for P2}) + \\
 & (\text{time spent in P3 up to the current PAL\_GET\_PSTATE}) * (\text{performance\_index for P3})) / \\
 & (\text{time interval between previous and current PAL\_GET\_PSTATE}) = \\
 & \frac{(t_1 - t_0) \times pf_0 + (t_2 - t_1) \times pf_1 + (t_3 - t_2) \times pf_2 + (t_4 - t_3) \times pf_3}{t_4 - t_0}
 \end{aligned}$$

**Figure 11-11. Computation of *performance\_index***



As seen above, for a HCDD, the PAL\_GET\_PSTATE procedure allows the caller to get feedback on the dynamic performance of the processor over a software-controlled time period. The caller can use this information to get better system utilization over a subsequent time period by changing the P-state in correlation with the current workload demand. The caller can also use PAL\_GET\_PSTATE to see the most recent P-state set for this logical processor (*type=0*) and the instantaneous current P-state that the domain parameters are set to (*type=3*). Platform power-caps do not affect either of these return values.

HIDD: If the logical processor belongs to a hardware-independent dependency domain, a weighted-average performance index can be returned by PAL\_GET\_PSTATE (*type=1* or *2*). Since software could calculate the performance index based on P-states it set, the weighted-average performance index is only of value when factoring in the effect of platform power-caps.

Note that P-state transitions typically do not happen instantaneously. An implementation-specific amount of time is required for a given transition to complete. The computation of the weighted-average *performance\_index* may not take into account the fact that transitions of power/performance are gradual, but may be done as though they were instantaneous at the point when the transition starts. The expectation is that any errors in computing the *performance\_index* due to non-instantaneous transitions to higher and lower P-states will tend to cancel out, and to the extent that they do not, will be insignificant.

#### 11.6.1.4 Variable P-state Performance

Some processors support variable P-state performance which allows the frequency to vary within a given P-state in order to achieve the maximum performance for that P-state's power budget. The PAL\_PROC\_GET\_FEATURES procedure indicates whether the processor supports variable P-state performance (see "[PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)](#)" on page 2:446 for details).

Since the frequency within a P-state can vary, the performance index calculation is slightly different when a processor supports variable P-state performance. Frequencies for a given P-state are represented by an index value  $F_{x,y}$ . The value  $x$  is the P-state number and  $y$  represents a frequency point in the range from 0 to  $N$ . A value of 0 represents the minimum frequency index value for the given P-state. For example:

$F_{0,0}$  to  $F_{0,N}$  – Frequency index values for the P0 state  
 $F_{1,0}$  to  $F_{1,N}$  – Frequency index values for the P1 state  
 ...etc.

$F_{0,0}$  is the minimum frequency index for the P0 state and its value is 100.  $F_{0,1}$  represents a higher frequency point for P0 and will have a value greater than 100. For example, if  $F_{0,1}$  frequency is 5% greater than  $F_{0,0}$  it would have a value of 105.

The *performance\_index* equation for P0 is calculated as follows:

$$\frac{((F_{0,0} * \text{time spent in } F_{0,0}) + (F_{0,1} * \text{time spent in } F_{0,1}) + \dots + (F_{0,N} * \text{time spent in } F_{0,N}))}{(\text{Total Time spent in } P_0)}$$

For example, let's say the minimum frequency of P0 is 1 GHz and the maximum frequency of P0 is 1.5 GHz. If we are at 1 GHz for a time period of 4, 1.25 GHz for a time period of 16 and 1.5 GHz for a time period of 20, the average performance index is:

$$((100 * 4) + (125 * 16) + (150 * 20)) / (5 + 15 + 20) = 135$$

The *performance\_index* equation for other P-states can be calculated in a similar manner using their respective frequency index values.

The total *performance\_index* equation for a processor with four P-states (P0, P1, P2, P3) would be:

$$\begin{aligned} & ((F_{0,0} * \text{time spent in } F_{0,0}) + (F_{0,1} * \text{time spent in } F_{0,1}) + \dots (F_{0,N} * \text{time spent in } F_{0,N}) + \\ & (F_{1,0} * \text{time spent in } F_{1,0}) + (F_{1,1} * \text{time spent in } F_{1,1}) + \dots (F_{1,N} * \text{time spent in } F_{1,N}) + \\ & (F_{2,0} * \text{time spent in } F_{2,0}) + (F_{2,1} * \text{time spent in } F_{2,1}) + \dots (F_{2,N} * \text{time spent in } F_{2,N}) + \\ & (F_{3,0} * \text{time spent in } F_{3,0}) + (F_{3,1} * \text{time spent in } F_{3,1}) + \dots (F_{3,N} * \text{time spent in } F_{3,N})) / \\ & (\text{Total Time}) \end{aligned}$$

#### 11.6.1.5 Interaction of P-states with HALT State

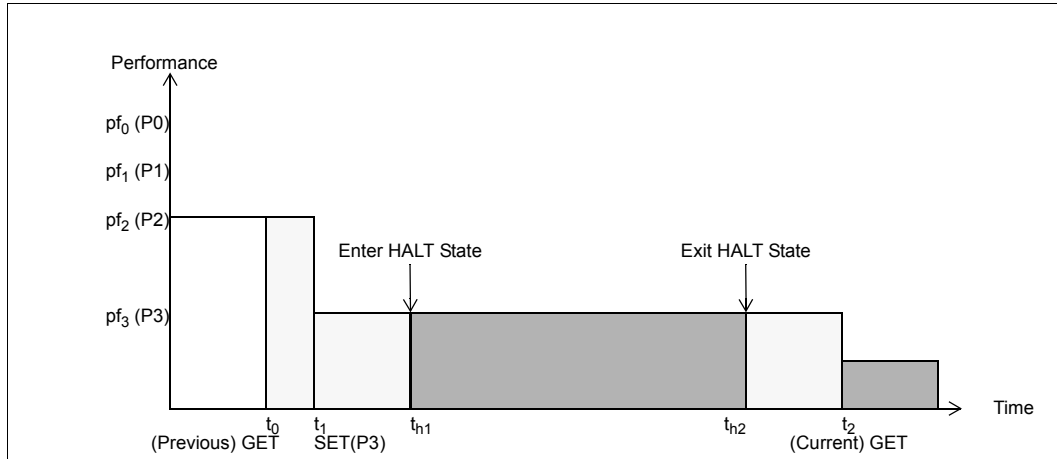
It is possible for a logical processor to enter and exit a HALT state between two consecutive calls to PAL\_GET\_PSTATE. Since the logical processor is not executing any instructions while in the HALT state, the performance index contribution during this period is essentially 0, and will not be accounted for in the *performance\_index* value returned when the next PAL\_GET\_PSTATE procedure call is made.

For example, let us assume that the previous PAL\_GET\_PSTATE procedure was called at time  $t_0$ , when the processor was operating in state P2. The previous PAL\_SET\_PSTATE procedure initiated a transition from P2 to P3 at time  $t_1$ . The processor entered HALT state at time  $t_{h1}$ , and exited the HALT state at time  $t_{h2}$ , and was in state P3 at time  $t_2$  when the current PAL\_GET\_PSTATE procedure was called. The *performance\_index* returned is calculated as:

$$\begin{aligned} \text{performance\_index} = & ((\text{time in P2 after the previous PAL\_GET\_PSTATE}) * (\text{performance\_index for P2}) + \\ & (\text{time in P3 before entering HALT state}) * (\text{performance\_index for P3}) + \\ & (\text{time in P3 after exiting HALT up to current PAL\_GET\_PSTATE})) * (\text{performance\_index for P3}) / \\ & (\text{time interval between previous and current GET, excluding time spent in HALT}) = \end{aligned}$$

$$\frac{(t_1 - t_0) \times pf_2 + (t_{h1} - t_1) \times pf_3 + (t_2 - t_{h2}) \times pf_3}{(t_2 - t_0) - (t_{h2} - t_{h1})}$$

**Figure 11-12. Interaction of P-states with HALT State**



As shown above, the value returned for *performance\_index* does not account for the performance during the time spent by the logical processor in the HALT state. This provides for better accuracy in the value reported for *performance\_index*, allowing the caller to make optimal adjustments to the system utilization even in scenarios where we have interactions between P-states and HALT state.

## 11.7 PAL Virtualization Support

This section describes the PAL architectural support for Itanium processor virtualization.

On processors in the Itanium Processor Family that support processor virtualization, the PAL virtualization support described in this document will be available. Itanium processor virtualization support can be determined by calling `PAL_PROC_GET_FEATURES`.

The virtualization support in PAL presents an implementation-independent interface to enable the VMM to implement software policies to manage/support virtualization of Itanium processors.

The PAL extensions for virtualization consist of three main components:

1. A set of procedures to support virtualization operations. These procedures allow the VMM to configure logical processors for virtualization operations and suspend/resume virtual processors on logical processors. Details for this component are described in [Section 11.10, "PAL Procedures" on page 2:353](#).
2. A set of services to provide low-latency, low-overhead support for performance-critical VMM operations. Details for this component are described in [Section 11.11, "PAL Virtualization Services" on page 2:486](#).
3. A PAL intercept interface to allow PAL to deliver virtualization events to the VMM in a low-latency, low-overhead manner. This PAL-to-VMM interface also allows PAL to provide optimizations for VMM operations. Details for this component are described in [Section 11.7.3, "PAL Intercepts in Virtual Environment" on page 2:332](#).



The VMM is responsible for managing the set of available system resources (CPU, memory, peripherals) and implement policies to virtualize these resources. In order to support virtual processor operations, the VMM will create a **virtual environment** and associate logical processors with the virtual environment. A virtual environment consists of one or more logical processors plus the memory resource allocated by the VMM during PAL\_VP\_INIT\_ENV.

The VMM creates a virtual environment by calling PAL\_VP\_ENV\_INFO to obtain the memory requirement for creating a virtual environment, and then by calling PAL\_VP\_INIT\_ENV on each logical processor that is to be part of the virtual environment. After a virtual environment is created, the VMM can create and initialize virtual processors to run in the environment by calling PAL\_VP\_CREATE.

The state of a virtual processor belonging to a virtual environment can be restored/saved on a logical processor in the environment by calling PAL\_VP\_RESTORE or PAL\_VP\_SAVE respectively. The VMM starts virtual processor operations on a logical processor by invoking either PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER.

The VMM can add/remove a logical processor from a virtual environment at any time by calling PAL\_VP\_INIT\_ENV or PAL\_VP\_EXIT\_ENV respectively.

### 11.7.1 Virtual Processor Descriptor (VPD)

The Virtual Processor Descriptor (VPD) represents the abstraction of processor resources of a single virtual processor. The VPD consists of per-virtual-processor control information together with performance-critical architectural state. The VPD is 64K in size and the base must be 32K aligned. [Table 11-16](#) shows the fields and layout of the VPD. The values in the VPD can be stored in little or big endian format, depending on the setting of *be* field setting in “[config\\_options – Global Configuration Options](#)” during PAL\_VP\_INIT\_ENV call. See “[PAL\\_VP\\_INIT\\_ENV – PAL Initialize Virtual Environment \(268\)](#)” on [page 2:478](#) for details. The VPD is divided into two classes – the first class stores control information and the second class stores the performance-critical architectural state of the virtual processor.

The VMM must keep the virtual processor state in the VPD for a particular state entry either: always, or only when one or more particular accelerations is enabled, as described in the Class columns of [Table 11-16](#), [Table 11-17](#) and [Table 11-18](#). See [Section 11.7.4.2, “Virtualization Accelerations”](#) on [page 2:337](#) for details.

**Note:** Not all architectural state of the virtual processor is included in the VPD. The VMM is responsible for setting up all the required virtual processor state in the architectural registers as well as in the VPD prior to resuming virtual processor execution. See [Table 11-122, “Virtual Processor Settings in Architectural Resources for PAL\\_VPS\\_RESUME\\_NORMAL and PAL\\_VPS\\_RESUME\\_HANDLER”](#) on [page 2:489](#) and [Table 11-123, “Processor Status Register Settings for Virtual Processor Execution”](#) on [page 2:490](#) for details.

**Table 11-16. Virtual Processor Descriptor (VPD)**

Name	Entries	Offset	Description	Class
vac	1	0	Virtualization Acceleration Control – these control bits enable virtualization acceleration of a particular resource or instruction. See <a href="#">Section 11.7.1.1, “Virtualization Controls”</a> on <a href="#">page 2:329</a> for details.	Control [always]
vdc	1	8	Virtualization Disable Control – these control bits disable the virtualization of a particular resource or instruction. See <a href="#">Section 11.7.1.1, “Virtualization Controls”</a> on <a href="#">page 2:329</a> for details.	Control [always]
virt_env_vaddr	1	16	PAL Virtual Environment Buffer Address – this field stores the host virtual address of the virtual environment which the virtual processor belongs to. The value in this field must be the same as the <i>vbase_addr</i> field during PAL_VP_INIT_ENV call.	Control [always]
Reserved	29	24	Reserved Area – Reserved for future expansion.	Reserved
vhpi	1	256	Virtual Highest Priority Pending Interrupt – Specifies the current highest priority pending interrupt for the virtual processor. See <a href="#">Table 11-124, “vhpi – Virtual Highest Priority Pending Interrupt”</a> on <a href="#">page 2:495</a> for details.	Control [a_int]
Reserved	95	264	Reserved Area – Reserved for future expansion.	Reserved
vgr[16-31]	16	1024	Virtual General Registers – Represent the bank 1 general registers 16-31 of the virtual processor. When the virtual processor is running and <i>vpsr.bn</i> is 1, the values in these entries are undefined.	Architectural State [a_bsw]
vbgr[16-31]	16	1152	Virtual Banked General Registers – Represent the bank 0 general registers 16-31 of the virtual processor. When the virtual processor is running and <i>vpsr.bn</i> is 0, the values in these entries are undefined.	Architectural State [a_bsw]
vnat	1	1280	Virtual General Register NaTs – Bits 0-15 represent the NaT values corresponding to <i>vgr16-31</i> , where the NaT bit for <i>vgr16</i> is in bit 0. Bits 16-63 are don't cares.	Architectural State [a_bsw]
vbnat	1	1288	Virtual Banked Register NaTs – Bits 16-31 represent the NaT values corresponding to <i>vbgr16-31</i> , where the NaT bit for <i>vbgr16</i> is in bit 16. Bits 0-15 and 32-63 are don't cares.	Architectural State [a_bsw]
vcpuid[0-4]	5	1296	Virtual CPUID Registers – Represent <i>cpuid</i> registers 0-4 of the virtual processor. NOTE: If <i>a_tf</i> is disabled or not supported, <i>vcpuid[0-1]</i> and <i>vcpuid[4]{63:32}</i> must contain the same values as the corresponding values of the logical processor on which this virtual processor is running. If <i>a_tf</i> is enabled, The VMM may maintain a different <i>VCPUID[4]{63:32}</i> value from the <i>CPUID[4]{63:32}</i> value of the logical processor on which the virtual processor is running.	Architectural State [a_from_cpuid, a_tf <sup>a</sup> ]

**Table 11-16. Virtual Processor Descriptor (VPD) (Continued)**

Name	Entries	Offset	Description	Class
Reserved	11	1336	Reserved Area – Reserved for future expansion.	Reserved
vpsr	1	1424	Virtual Processor Status Register – Represents the Processor Status Register of the virtual processor.	Architectural State See <a href="#">Table 11-17</a> for details.
vpr	1	1432	Virtual Predicate Registers – Represents the Predicate Registers of the virtual processor. The bit positions in vpr correspond to predicate registers in the same manner as with the mov predicates instruction. The contents in this field are undefined except at virtualization intercept handoff. The VMM can not rely on the contents in this field to be preserved when the virtual processor is running.	Architectural State [always]
Reserved	76	1440	Reserved Area – Reserved for future expansion. This area may also be used by PAL to hold additional machine-specific processor state.	Reserved
vcr[0-127]	128	2048	Virtual Control Registers – Represent the control registers of the virtual processor. For the reserved control registers, the corresponding VPD entries are reserved.	Architectural State See <a href="#">Table 11-18</a> for details.
Reserved	128	3072	Reserved Area – Reserved for future expansion. This area may also be used by PAL to hold additional machine-specific processor state	Reserved
Reserved	3456	4096	Reserved Area – Reserved for future expansion. This area may also be used by PAL to hold additional machine-specific processor state	Reserved
vmm_avail	128	31744	Available for VMM use. This area is ignored by the processor and PAL.	Ignored
Reserved	4096	32768	Reserved Area – Reserved for future expansion. This area may also be used by PAL to hold additional machine-specific processor state	Reserved

a. The a\_tf acceleration only requires vcpuid[4] be kept in the VPD.

[Table 11-17](#) provides details on which vpsr bits are required to be store in the VPD for different accelerations. Two bits, vpsr.ic and vpsr.si are always required to be in the VPD. The remaining vpsr bits are only required to be stored in the VPD if certain virtualization accelerations are enabled. Even though some fields are not required to be stored in the VPD, the VMM is free to store the entire vpsr in the VPD.

**Table 11-17. Virtual Processor Descriptor (VPD) – VPSR**

Field	Bits	Class
User Mask = PSR{5:0}		
rv	0	Reserved
be	1	No accelerations require these fields. <sup>a</sup>
up	2	
ac	3	
mfl	4	
mfh	5	
System Mask = PSR{23:0}		
ic	13	Always
i	14	a_int, a_from_psr
pk	15	a_from_psr
rv	12:6, 16	Reserved
dt	17	a_from_psr
dfi	18	
dfh	19	
sp	20	
pp	21	
di	22	
si	23	Always
PSR.I = PSR{31:0}		
db	24	a_from_psr
lp	25	
tb	26	
rt	27	
rv	31:28	Reserved
PSR{63:0}		
cpl	33:32	No accelerations require these fields.
is	34	
mc	35	a_from_psr
it	36	
id	37	No accelerations require these fields.
da	38	
dd	39	
ss	40	
ri	42:41	
ed	43	
bn	44	a_bsw
ia	45	No accelerations require these fields.
vm	46	
rv	63:47	Reserved

a. The user mask is not virtualized. See [Section 11.7.4.2.4, “MOV-from-PSR Optimization”](#) on page 2:341 and [Section 11.7.4.2.10, “Interrupt Collection and User Mask Optimization”](#) on page 2:345 for further details.

**Table 11-18. Virtual Processor Descriptor (VPD) – VCR[0-127]**

Register	Name	Class
VCR0-15		No accelerations require these virtual control registers.
VCR16	VIPSR	a_from_int_cr, a_to_int_cr
VCR17	VISR	
VCR18		No accelerations require this virtual control register.
VCR19	VIIP	a_from_int_cr, a_to_int_cr
VCR20	VIFA	Always
VCR21	VITIR	Always
VCR22	VIIPA	a_from_int_cr, a_to_int_cr
VCR23	VIFS	a_cover, a_from_int_cr, a_to_int_cr
VCR24	VIIM	a_from_int_cr, a_to_int_cr
VCR25	VIHA	
VCR26	VIIB0	
VCR27	VIIB1	
VCR28-65		No accelerations require these virtual control registers.
VCR66	VTPR	a_int
VCR67-127		No accelerations require these virtual control registers.

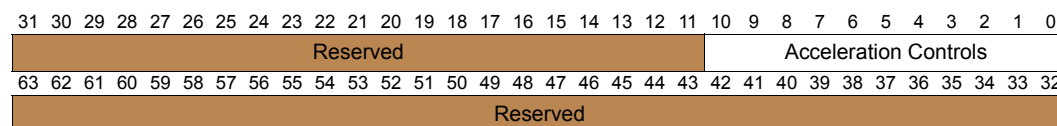
### 11.7.1.1 Virtualization Controls

The Virtualization Acceleration Control (*vac*) and Virtualization Disable Control (*vd*) fields in the VPD contain configuration control bits which define the set of events that will cause an intercept from PAL to the VMM. The virtualization controls are divided into two categories:

1. Virtualization Acceleration Control – these control bits enable virtualization optimization support of a particular resource or instruction. [Figure 11-13](#) and [Table 11-19](#) describe these control bits.
2. Virtualization Disable Control – these control bits disable the virtualization of a particular resource or instruction. [Figure 11-14](#) and [Table 11-20](#) describe these control bits.

The *vac* and *vd* settings are specified by the VMM during virtual processor initialization when the PAL\_VP\_CREATE procedure is called, and cannot be changed until the virtual processor is terminated by PAL\_VP\_TERMINATE.

**Figure 11-13. Virtualization Acceleration Control (*vac*)**



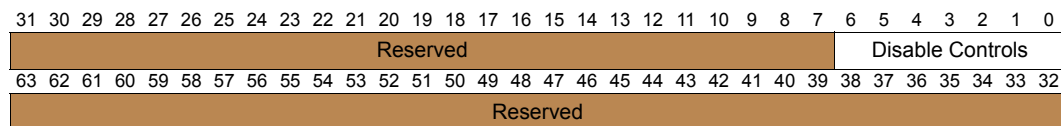
**Table 11-19. Virtualization Acceleration Control (*vac*) Fields**

Field	Bit	Description
a_int	0	Enable the virtual external interrupt optimization. See <a href="#">Section 11.7.4.2.1, “Virtual External Interrupt Optimization”</a> on page 2:338 for details.
a_from_int_cr	1	Enable the interruption control register (CR16-27) read optimization. See <a href="#">Section 11.7.4.2.2, “Interruption Control Register Read Optimization”</a> on page 2:340 for details.

**Table 11-19. Virtualization Acceleration Control (vac) Fields (Continued)**

Field	Bit	Description
a_to_int_cr	2	Enable the interruption control register (CR16-27) write optimization. See Section 11.7.4.2.3, "Interruption Control Register Write Optimization" on page 2:341 for details.
a_from_psr	3	Enable the processor status register read optimization. See Section 11.7.4.2.4, "MOV-from-PSR Optimization" on page 2:341 for details.
a_from_cpuid	4	Enable the CPUID read optimization. See Section 11.7.4.2.5, "MOV-from-CPUID Optimization" on page 2:342 for details.
a_cover	5	Enable the <code>cover</code> instruction optimization. See Section 11.7.4.2.6, "Cover Optimization" on page 2:343 for details.
a_bsw	6	Enable the <code>bsw</code> instruction optimization. See Section 11.7.4.2.7, "Bank Switch Optimization" on page 2:343 for details.
a_all_probes	7	Enable virtualization of probe instructions. See Section 11.7.4.2.8, "Probe Instruction Virtualization" on page 2:344 for details.
a_select_probes	8	
a_tf	9	Enable the test feature optimization. See Section 11.7.4.2.9, "Test Feature Optimization" on page 2:345 for details.
a_ic_um	10	Enable the interruption collection and user mask optimization. See Section 11.7.4.2.10, "Interruption Collection and User Mask Optimization" on page 2:345 for details.
Reserved	63:11	Reserved

**Figure 11-14. Virtualization Disable Control (vdc)**



**Table 11-20. Virtualization Disable Control (vdc) Fields**

Field	Bits	Description
d_vmsw	0	Disable <code>vmsw</code> instruction – If 1, disables <code>vmsw</code> instruction on the logical processor. Execution of the <code>vmsw</code> instruction, independent of the state of <code>PSR.vm</code> , will cause a virtualization intercept.
d_extint	1	Disable external interrupt control register virtualization – If 1, accesses (reads/writes) of the external interrupt control registers (CR65-71) are not virtualized. Code running with <code>PSR.vm==1</code> can read and write the external interrupt control registers of the logical processor directly and without handling off to the VMM. See Section 11.7.4.3.2, "Disable External Interrupt Control Register Virtualization" on page 2:347 for details.
d_ibr_dbr	2	Disable breakpoint register virtualization – If 1, accesses (reads/writes) of the data and instruction breakpoint registers (IBR/DBR) are not virtualized. Code running with <code>PSR.vm==1</code> can read and write the data/instruction breakpoint registers of the logical processor directly and without handling off to the VMM. If 0, accesses of the breakpoint registers with <code>PSR.vm==1</code> result in virtualization intercepts.
d_pmc	3	Disable PMC virtualization – If 1, accesses (reads/writes) of the performance monitor configuration registers (PMCs) are not virtualized. Code running with <code>PSR.vm==1</code> can read and write the performance monitor configuration registers of the logical processor directly and without handling off to the VMM. If 0, accesses of the performance counter configuration registers with <code>PSR.vm==1</code> result in virtualization intercepts.

**Table 11-20. Virtualization Disable Control (vdc) Fields (Continued)**

Field	Bits	Description
d_to_pmd	4	Disable PMD write virtualization – If 1, writes to the performance monitor data registers (PMDs) are not virtualized. Code running with PSR.vm==1 can write the performance monitor data registers of the logical processor directly and without handling off to the VMM. If 0, writes of the performance counter data registers with PSR.vm==1 result in virtualization intercepts.
d_itm	5	Disable ITM virtualization – If 1, writes to the Interval Timer Match (ITM) register are not virtualized. Code running with PSR.vm==1 can write the ITM register of the logical processor directly and without handling off to the VMM. If 0, writes of the ITM register with PSR.vm==1 result in virtualization intercepts.
d_psr_i	6	Disable PSR.i virtualization – If 1, accesses (reads/writes) to the interrupt bit in processor state register (PSR.i) are not virtualized. Code running with PSR.vm==1 can read and write only the interrupt bit via the <code>ssm</code> and <code>rsm</code> instructions directly without handling off to the VMM. Attempts to modify other PSR bits in addition to the interrupt bit via the <code>ssm</code> and <code>rsm</code> instructions will result in virtualization intercepts. Attempts to modify the interrupt bit with the <code>mov psr.l</code> instruction will continue to result in virtualization intercepts. If 0, accesses to the PSR.i bit with PSR.vm==1 result in virtualization intercepts.
Reserved	63:7	Reserved

## 11.7.2 Interruption Handling in a Virtual Environment

For logical processors which have been added to a virtual environment through `PAL_VP_INIT_ENV`, all IVA-based interruptions continue to be delivered to the **host IVT** independent of the state of `PSR.vm` at the time of interruption. All IVA-based interruptions are serviced by the host IVT pointed to by the IVA (CR2) control register on the logical processor.

IVA-based interruptions that do not represent virtualization events will be delivered to the **guest IVT** by the VMM. The guest IVT is specified by the VIVA control register in the VPD of the virtual processor.

For IVA-based interruption handling during virtual processor operations, PAL provides maximum flexibility to the VMM by supporting **per-virtual-processor host IVTs**. This allows the VMM to provide a different host IVT with optimizations specific to a particular guest operating system on the virtual processor. The VMM can also choose to provide the same IVT for some or all of the virtual processors in a virtual environment.

Hence, at any time in a virtual environment, the IVA (CR2) control register of the logical processor will be pointing to either:

- The per-virtual-processor host IVT
- The generic host IVT not specific to any virtual processor

The per-virtual-processor host IVT for each virtual processor is setup by PAL when the virtual processor is first created (`PAL_VP_CREATE`) or registered (`PAL_VP_REGISTER`) in the virtual environment. The VMM passes a pointer to the host IVT specific to the virtual processor as an incoming parameter to the `PAL_VP_CREATE` or `PAL_VP_REGISTER` procedures. The per-virtual-processor host IVT is setup to perform long branches to the corresponding vector of the IVT specified in the incoming parameter for all IVA-based

interruptions except the Virtualization vector. Virtualization vector will be delivered as virtualization intercept in the per-virtual-processor host IVT. See [Section 11.7.3, “PAL Intercepts in Virtual Environment”](#) on page 2:332 for details on PAL intercepts.

In the virtual environment, the IVA (CR2) control register will be set by PAL virtualization-related procedures and services as summarized in [Table 11-21](#).

**Table 11-21. IVA Settings after PAL Virtualization-related Procedures and Services**

<b>PAL Virtualization-related Procedure / Service</b>	<b>Description</b>
PAL_VP_CREATE	These procedures do not change the IVA control register.
PAL_VP_ENV_INFO	
PAL_VP_EXIT_ENV	This procedure sets the IVA control register to point to the IVT specified by the caller.
PAL_VM_INIT_ENV	These procedures do not change the IVA control register.
PAL_VP_REGISTER	
PAL_VP_RESTORE / PAL_VPS_RESTORE	This procedure / service sets the IVA control register to point to the per-virtual-processor host IVT.
PAL_VP_SAVE / PAL_VPS_SAVE	This procedure / service does not change the IVA control register.
PAL_VP_TERMINATE	This procedure sets the IVA control register to point to the IVT specified by the caller.

After successful execution of PAL\_VP\_RESTORE procedure or PAL\_VPS\_RESTORE service, the IVA control register on the logical processor is set to point to the per-virtual-processor host IVT. After successful completion of PAL\_VP\_RESTORE procedure, the VMM must not change the IVA control register on the logical processor until after the next invocation of PAL\_VP\_SAVE or PAL\_VPS\_SAVE.

On IVA-based interruptions when a virtual processor is running (after PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER), the IVA control register on the logical processor is unchanged and will continue to point to the per-virtual-processor host IVT. On resume execution to the same virtual processor through PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER PAL services, the VMM must ensure the IVA control register on the logical processor is set to point to the per-virtual-processor host IVT at the time of interruption.<sup>1</sup>

### 11.7.3 PAL Intercepts in Virtual Environment

When the IVA control register on the logical processor is set to point to the per-virtual-processor host IVT, virtualization intercepts will be raised at the Virtualization vector or at an optional virtualization intercept handler specified by the VMM. By default, virtualization intercepts are delivered to the Virtualization vector of the IVT specified by the VMM during PAL\_VP\_CREATE / PAL\_VP\_REGISTER. If the VMM specified the optional virtualization intercept handler, all virtualization intercepts are delivered to that handler (instead of the Virtualization vector.)

---

1. In other words, the VMM is allowed to change to another IVT after IVA-based interruptions happening during virtual processor execution. The VMM must ensure the per-virtual processor IVT is restored before resuming to the same virtual processor through PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER.



Section 11.7.3.1, “PAL Virtualization Intercept Handoff State” on page 2:333 describes the handoff state of the PAL intercepts. For all interruption vectors other than Virtualization vector, the architectural state at the corresponding IVA-based interruption vector is the same as defined in Chapter 8, “Interruption Vector Descriptions” in Volume 2.

### 11.7.3.1 PAL Virtualization Intercept Handoff State

The state of the logical processor at virtualization intercept handoff is:

- GRs:
  - Non-banked GRs: The contents of non-banked general registers are preserved from the time of the interruption.
  - Bank 1 GRs: The contents of all bank one general registers are preserved from the time of the interruption.
  - Bank 0: GR16-23: The contents of these bank zero general registers are preserved from the time of the interruption.
  - Bank 0: GR24-31: Scratch, contains parameters/state for VMM:
    - GR24 indicates the cause of virtualization intercept. See Table 11-22, “PAL Virtualization Intercept Handoff Cause (GR24)” for details. This field is not provided to the VMM if the value of the *cause* field in the *config\_options* parameter passed to PAL\_VP\_INIT\_ENV is 0. If the value of the *cause* field in the *config\_options* parameter passed to PAL\_VP\_INIT\_ENV is 0, the value of GR24 on virtualization intercept handoff is undefined.
    - GR25 contains the 41-bit opcode in little endian format and the type of the instruction which caused the fault, excluding the qualifying predicate (qp) field. See Figure 11-15, “PAL Virtualization Intercept Handoff Opcode (GR25),” on page 2:335 for details.
    - GR26-31 are available for the VMM to use.
- FRs: The contents of all floating-point registers are preserved from the time of the interruption.
- Predicates: The contents of all predicate registers are undefined and available for use. The original contents are saved in the VPD.
- BRs: The contents of all branch registers are preserved from the time of the interruption.
- ARs: The contents of all application registers are preserved from the time of the interruption, except the ITC and RUC counters. The ITC register will not be directly modified by PAL, but will continue to count during the execution of the virtualization intercept handler. The RUC register will not be directly modified by PAL, but will continue to count during the execution of the virtualization intercept handler while the processor is active.
- CFM: The contents of the CFM register is preserved from the time of the interruption.
- RSE: All RSE state is preserved from the time of the interruption.
- PSR: PSR fields are set according to the “Interruption State” column in Table 3-2, “Processor Status Register Fields” on page 2:24. PSR.up and pp are set to 0 when *fr\_pmc* field in *config\_options* parameter during PAL\_VP\_INIT\_ENV is 1.
- CRs: The contents of all control registers are preserved from the time of the interruption with the exception of resources described below.

- IRRs: The contents of IRRs are not changed by PAL. Incoming interruptions may change the contents.
- IFS: IFS is unchanged from the time of the interruption.
- IIP: Contains the value of IP at the time of the interruption.
- IPSR: Contains the value of PSR at the time of the interruption.
- RRs: The contents of all region registers are preserved from the time of the interruption.
- PKRs: The contents of all protection key registers are preserved from the time of the interruption.
- DBRs/IBRs: The contents of all breakpoint registers are preserved from the time of the interruption.
- PMCs/PMDs: The contents of the PMC registers are preserved from the time of the virtualization intercept. The contents of the PMD registers are not modified by PAL code, but may be modified if events being monitored are encountered. The performance counters will be frozen if specified by the VMM through a parameter of PAL\_VP\_INIT\_ENV procedure.
- Cache: The processor internal cache is not specifically modified by PAL handler but may be modified due to normal cache activity of running the handler code.
- TLB: The TRs are unchanged from the time of the interruption.

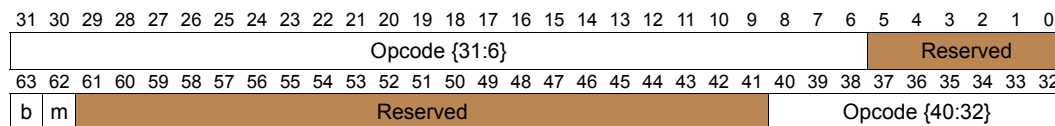
**Table 11-22. PAL Virtualization Intercept Handoff Cause (GR24)**

Value	Cause	Description
1	toAR	Due to MOV-to-AR instruction.
2	toARimm	Due to MOV-to-AR-imm instruction.
3	fromAR	Due to MOV-from-AR instruction.
4	toCR	Due to MOV-to-CR instruction.
5	fromCR	Due to MOV-from-CR instruction.
6	toPSR	Due to MOV-to-PSR instruction.
7	fromPSR	Due to MOV-from-PSR instruction.
8	itc_d	Due to <code>itc.d</code> instruction.
9	itc_i	Due to <code>itc.i</code> instruction.
10	toRR	Due to MOV-to-RR instruction.
11	toDBR	Due to MOV-to-DBR instruction.
12	toIBR	Due to MOV-to-IBR instruction.
13	toPKR	Due to MOV-to-PKR instruction.
14	toPMC	Due to MOV-to-PMC instruction.
15	toPMD	Due to MOV-to-PMD instruction.
16	itr_d	Due to <code>itr.d</code> instruction.
17	itr_i	Due to <code>itr.i</code> instruction.
18	fromRR	Due to MOV-from-RR instruction.
19	fromDBR	Due to MOV-from-DBR instruction.
20	fromIBR	Due to MOV-from-IBR instruction.
21	fromPKR	Due to MOV-from-PKR instruction.
22	fromPMC	Due to MOV-from-PMC instruction.
23	fromCPUID	Due to MOV-from-CPUID instruction.
24	ssm	Due to <code>ssm</code> instruction.
25	rsm	Due to <code>rsm</code> instruction.
26	ptc_l	Due to <code>ptc.l</code> instruction.

**Table 11-22. PAL Virtualization Intercept Handoff Cause (GR24) (Continued)**

Value	Cause	Description
27	ptc_g	Due to <code>ptc.g</code> instruction.
28	ptc_ga	Due to <code>ptc.ga</code> instruction.
29	ptr_d	Due to <code>ptr.d</code> instruction.
30	ptr_i	Due to <code>ptr.i</code> instruction.
31	thash	Due to <code>thash</code> instruction.
32	ttag	Due to <code>ttag</code> instruction.
33	tpa	Due to <code>tpa</code> instruction.
34	tak	Due to <code>tak</code> instruction.
35	ptc_e	Due to <code>ptc.e</code> instruction.
36	cover	Due to <code>cover</code> instruction.
37	rfi	Due to <code>rfi</code> instruction.
38	bsw_0	Due to <code>bsw.0</code> instruction.
39	bsw_1	Due to <code>bsw.1</code> instruction.
40	vmsw	Due to <code>vmsw</code> instruction.
41	probe	Due to <code>probe</code> instruction.
All other values	Reserved	Reserved for future expansion.

**Figure 11-15. PAL Virtualization Intercept Handoff Opcode (GR25)**



## 11.7.4 Virtualization Optimizations

After the `PAL_VP_INIT_ENV` procedure is called, execution of the virtualized instructions listed in [Table 3-10, “Virtualized Instructions”](#) on page 2:44 with `PSR.vm==1` results in virtualization intercepts to the VMM. Virtualization optimizations reduce overall virtualization overhead by allowing these instructions to execute, with `PSR.vm==1`, without causing intercepts to the VMM. There are two types of virtualization optimizations – global and local. Local virtualization optimizations are further divided into virtualization accelerations and virtualization disables.

Global virtualization optimizations are specified during initialization of the virtual environment (i.e., during `PAL_VP_INIT_ENV`). The specified optimizations are applicable to all the virtual processors running in the virtual environment. See [Section 11.7.4.1, “Global Virtualization Optimizations”](#) for details on the global virtualization optimizations supported in the architecture.

Local virtualization optimizations are specified during the creation of the virtual processor (i.e., during `PAL_VP_CREATE`). The optimization settings were specified in the VPD and hence local to each virtual processor. The VMM can specify different local optimization settings for different virtual processors. The two classes of local virtualization optimizations are:

- Virtualization accelerations – Virtualization accelerations optimize the execution of virtualized instructions by supporting fast access to the virtual instance of the

resource and perform the virtualized operations based on the virtual instance of the resource without handing off to the VMM. [Section 11.7.4.2, “Virtualization Accelerations” on page 2:337](#) describes the supported Virtualization accelerations in the architecture.

- Virtualization disables – Virtualization disables optimize the execution of virtualized instructions by disabling virtualization of a particular resource or instruction. Accesses to the virtualization-disabled resources or executions of virtualization-disabled instructions, even with `PSR.vm=1`, will not cause intercepts to the VMM. [Section 11.7.4.3, “Virtualization Disables” on page 2:346](#) describes the supported Virtualization disables in the architecture.

### 11.7.4.1 Global Virtualization Optimizations

[Table 11-23](#) summarizes the global virtualization optimizations supported in Itanium architecture.

**Table 11-23.Global Virtualization Optimizations Summary**

Optimization	<i>config_options</i> <sup>a</sup>	Description
Virtualization Opcode Optimization	opcode	<a href="#">Section 11.7.4.1.1</a>
Virtualization Cause Optimization	cause	<a href="#">Section 11.7.4.1.2</a>
Guest MOV-from-AR.ITC Optimization	gitc	<a href="#">Section 11.7.4.1.3</a>

a. *config\_options* is a parameter for the PAL\_VP\_INIT\_ENV procedure. See [“PAL\\_VP\\_INIT\\_ENV – PAL Initialize Virtual Environment \(268\)” on page 2:478](#) for details.

Certain global virtualization optimizations have VPD synchronization requirements. Please refer to the corresponding section of each global virtualization optimizations for more details on these requirements.

#### 11.7.4.1.1 Virtualization Opcode Optimization

Virtualization opcode optimization is always enabled. Opcode information is provided to the VMM during PAL intercepts in the virtual environment. In some processor implementations, the opcode provided may not be guaranteed to be the opcode that triggered the intercept; virtual machine monitors can determine whether this is guaranteed from the *vp\_env\_info* return value of PAL\_VP\_ENV\_INFO.

[Table 11-24](#) and [Table 11-16, “Virtual Processor Descriptor \(VPD\)” on page 2:326](#) shows the synchronization requirements and the VPD states that will be accessed for this optimization.

**Table 11-24.Synchronization Requirements for Virtualization Opcode Optimization**

VPD Resource	Synchronization Required
vpsr.ic	Write
vpsr.si	Write
vifa	Write
vitir	Write

### 11.7.4.1.2 Virtualization Cause Optimization

Virtualization cause optimization is enabled by the *cause* bit in the *config\_options* parameter of *PAL\_VP\_INIT\_ENV*. When enabled, the causes of virtualization intercepts will be provided to the VMM during PAL intercept handoffs within the virtual environment. When disabled, no cause information will be provided during PAL intercept handoffs.

This optimization requires no special synchronization.

### 11.7.4.1.3 Guest MOV-from-AR.ITC Optimization

Guest MOV-from-AR.ITC optimization allows software running with *PSR.vm*=1 to execute MOV-from-AR.ITC instructions without any intercepts to the VMM. The value returned will be the sum of the value in the interval timer counter register (ITC) and interval timer offset register (ITO), unless a fault condition is detected (see [Table 11-25, “Behavior of Guest MOV-from-AR.ITC Instruction in Virtual Environment”](#) for details). The VMM is responsible for programming the ITO register to provide the desired return value for guest execution with *PSR.vm* = 1 of the MOV-from-ITC instruction when this optimization is enabled.

This optimization is enabled by the *gipc* bit in the *config\_options* parameter of *PAL\_VP\_INIT\_ENV*. The behavior of the guest MOV-from-AR.ITC instruction is affected by the settings of *psr.ic* and *vpsr.ic* as well, as shown in [Table 11-25](#).

This optimization requires no special synchronization.

This optimization is not supported on all processor implementations. Software can call *PAL\_VP\_ENV\_INFO* to determine the availability of this feature.

**Table 11-25. Behavior of Guest MOV-from-AR.ITC Instruction in Virtual Environment**

<i>gipc</i> <sup>a</sup>	<i>psr.si</i>	<i>vpsr.si</i>	MOV-from-AR.ITC when <i>PSR.vm</i> =1
0	0	0	No virtualization intercept – guest reads AR.ITC
	0	1	Invalid setting – behavior is undefined.
	1	0	Virtualization intercept
	1	1	If <i>vpsr.cpl</i> is not zero: Privileged Register fault If <i>vpsr.cpl</i> is zero: Virtualization intercept
1	0	0	No virtualization intercept – guest reads the sum of ITC and ITO
	0	1	If <i>vpsr.cpl</i> is not zero: Privileged Register fault If <i>vpsr.cpl</i> is zero: No Virtualization intercept – guest reads the sum of ITC and ITO
	1	0	Virtualization intercept.
	1	1	If <i>vpsr.cpl</i> is not zero: Privileged Register fault If <i>vpsr.cpl</i> is zero: Virtualization intercept

a. *gipc*=0: Optimization disabled; *gipc*=1: Optimization enabled.

## 11.7.4.2 Virtualization Accelerations

[Table 11-26](#) summarizes the virtualization accelerations supported in Itanium architecture.

**Table 11-26. Virtualization Accelerations Summary**

Optimization	Virtualization Acceleration Control (vac) <sup>a</sup>	Description
Virtual External Interrupt Optimization	a_int	Section 11.7.4.2.1
Interrupt Control Register Read Optimization	a_from_int_cr	Section 11.7.4.2.2
Interrupt Control Register Write Optimization	a_to_int_cr	Section 11.7.4.2.3
MOV-from-PSR Optimization	a_from_psr	Section 11.7.4.2.4
MOV-from-CPUID Optimization	a_from_cpuid	Section 11.7.4.2.5
Cover Optimization	a_cover	Section 11.7.4.2.6
Bank Switch Optimization	a_bsw	Section 11.7.4.2.7
Virtualize all <code>probe</code> instructions	a_all_probes	Section 11.7.4.2.8
Virtualize selected <code>probe</code> instructions	a_select_probes	
Test Feature Optimization	a_tf	Section 11.7.4.2.9
Interrupt Collection and User Mask Optimization	a_ic_um	Section 11.7.4.2.10

a. The Virtualization Acceleration Control (vac) field resides in the Virtual Processor Descriptor (VPD), see Section 11.7.1, “Virtual Processor Descriptor (VPD)” on page 2:325 for details.

For each of the accelerations, certain virtual processor control and architectural state is managed directly by hardware/firmware, and hence must be maintained in the VPD, and synchronization is required when the VMM reads or writes this state in the VPD. Some entries must be maintained in the VPD independent of any accelerations. (These are marked as [always].) See Table 11-16 for details on which VPD state is used with each of the accelerations. See Section 11.11, “PAL Virtualization Services” on page 2:486 for a description of the synchronization services.

#### 11.7.4.2.1 Virtual External Interrupt Optimization

The virtual external interrupt optimization allows the VMM to specify the virtual highest priority pending interrupt so that a virtual external interrupt is raised on changes of `vtpr` or `vpsr.i` only when that the virtual highest priority pending interrupt is unmasked. For details on virtual external interrupts, see “Virtual External Interrupt vector (0x3400)” on page 2:187.

The virtual external interrupt optimization is enabled by the `a_int` bit in the Virtualization Acceleration Control (vac) field in the VPD. When this optimization is enabled, the VMM specifies the virtual highest priority pending interrupt (vhpi) through the `PAL_VPS_SET_PENDING_INTERRUPT` service described in Section 11.11.2, “PAL Virtualization Service Specifications” on page 2:488. If this optimization is disabled, processor behavior is undefined if `PAL_VPS_SET_PENDING_INTERRUPT` is invoked.

When this optimization is enabled, execution of `rsm` and `ssm` instructions<sup>1</sup>, with `PSR.vm==1`, which modify only `vpsr.i` will not intercept to the VMM and `vpsr.i` is updated with the new value, unless a fault condition is detected (see Table 11-29 for details).

1. The execution of `rsm` and `ssm` instructions with `PSR.vm==1` is affected by both the virtual external interrupt optimization (`a_int`) and the interruption collection and user mask optimization (`a_ic_um`). Software can enable or disable both optimizations together, or enable each optimization independently. Section 11.7.4.4.1, “Virtual External Interrupt Optimization and Interruption Collection and User Mask Optimization” on page 2:349 describes the behavior when both optimizations are enabled.

When this optimization is enabled, execution of `rsm` and `ssm` instructions<sup>1</sup>, with `PSR.vm==1` and system mask equal to zero (0x0), will not intercept to the VMM unless a fault condition is detected (see [Table 11-29](#) for details).

A virtual external interrupt is raised if the virtual highest priority pending interrupt (vhpi) is unmasked by the new `vpsr.i` and `vtpr`. If the virtual highest priority pending interrupt (vhpi) is still masked by the new `vpsr.i` or `vtpr`, no virtual external interrupt will be raised. Note that execution of MOV-to-PSR instructions with `PSR.vm==1` always results in a virtualization intercept no matter which PSR bits are modified.

When this optimization is enabled, execution of `rsm` and `ssm` instructions<sup>1</sup>, with `PSR.vm==1`, which modify any bits in addition to `vpsr.i` result in a virtualization intercepts. No virtual external interrupts are raised and the VMM is responsible for delivering a virtual external interrupt if the virtual highest priority pending interrupt (vhpi) is unmasked.

When this optimization is enabled, execution of a MOV-from-CR instruction, with `PSR.vm==1`, targeting `vtpr` reads the most recent value, unless a fault condition is detected (see [Table 11-29](#) for details).

When this optimization is enabled, on execution of MOV-to-TPR instructions with `PSR.vm==1`, `vtpr` will be updated with the new value without handling off to the VMM, unless a fault condition is detected (see [Table 11-29](#) for details). A virtual external interrupt is raised if the virtual highest priority pending interrupt (vhpi) is unmasked by the new `vpsr.i` and `vtpr`. No virtual external interrupt is raised if the virtual highest priority pending interrupt is still masked by `vpsr.i` or `vtpr`.

When this optimization is enabled, after completion of an instruction with `PSR.vm==1` which modifies `vtpr` or `vpsr.i` (if the instruction completes without an intercept), a determination is made as to whether the new state unmask the virtual highest priority pending interrupt. If it does, then a virtual external interrupt will be raised and the VMM will be entered on the Virtual External Interrupt vector. See [Table 11-27](#) for details on the detection of virtual external interrupts.

**Table 11-27. Detection of Virtual External Interrupts**

Condition	Virtual External Interrupt
<code>vhpi &lt;= (!vpsr.i &lt;&lt; 5   vtpr.mmi &lt;&lt;4   vtpr.mic)</code>	No – virtual highest priority pending interrupt is still masked.
<code>vhpi &gt; (!vpsr.i &lt;&lt; 5   vtpr.mmi &lt;&lt;4   vtpr.mic)</code>	Yes – virtual highest priority pending interrupt is unmasked.

Synchronization is required when this optimization is enabled, see [Table 11-28](#) for details.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-28. Synchronization Requirements for Virtual External Interrupt Optimization**

VPD Resource	Synchronization Required
<code>vtpr</code>	Read, Write
<code>vpsr.i</code>	Read, Write
<code>vhpi</code>	Write

**Table 11-29. Interruptions when Virtual External Interrupt Optimization is Enabled**

Instructions	Interruptions
<code>rsm, ssm</code>	When the virtual external interrupt optimization is enabled, execution of <code>rsm</code> and <code>ssm</code> instructions with <code>PSR.vm==1</code> which modify only <code>vpsr.i</code> , may raise the following faults: <ul style="list-style-type: none"> <li>Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> </ul>
MOV-from-TPR	When the virtual external interrupt optimization is enabled, execution of MOV-from-CR instruction targeting <code>vtpr</code> with <code>PSR.vm==1</code> , may raise the following faults: <ul style="list-style-type: none"> <li>Illegal Operation fault – if the target operand specifies GR 0 or an out-of-frame stacked register</li> <li>Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> </ul>
MOV-to-TPR	When the virtual external interrupt optimization is enabled, execution of MOV-to-CR instruction targeting <code>vtpr</code> with <code>PSR.vm==1</code> , may raise the following faults: <ul style="list-style-type: none"> <li>Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> <li>Register NaT Consumption fault – if the NaT bit in the source register is one</li> <li>Reserved Register/Field fault – if the reserved field in the <code>vtpr</code> is being written with a non-zero value</li> </ul>

**Note:** This field cannot be enabled together with `d_extint` or `d_psr_i` virtualization disables. If this control is enabled together with any one of described disables, an error will be returned during `PAL_VP_CREATE` and `PAL_VP_REGISTER`. See [Section 11.7.4.4, “Virtualization Optimization Combinations”](#) on page 2:349 for details.

#### 11.7.4.2.2 Interruption Control Register Read Optimization

The interruption control register read optimization is enabled by the `a_from_int_cr` bit in the Virtualization Acceleration Control (`vac`) field in the VPD. When this optimization is enabled, and `vpsr.ic` is 0, software running with `PSR.vm==1` will be able to read the virtual interruption control registers (`vipsr`, `visr`, `viip`, `vifa`, `vitir`, `viipa`, `vifs`, `viim`, `viha`, `viib0-1`) without any intercepts to the VMM, unless a fault condition is detected (see [Table 11-31](#) for details).

If this optimization is disabled, a read of the interruption CRs with `PSR.vm==1` results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see [Table 11-30](#) for details.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16, “Virtual Processor Descriptor \(VPD\)”](#) on page 2:326.

**Table 11-30. Synchronization Requirements for Interruption Control Register Read Optimization**

VPD Resource	Synchronization Required
<code>vipsr</code> , <code>visr</code> , <code>viip</code> , <code>vifa</code> , <code>vitir</code> , <code>viipa</code> , <code>vifs</code> , <code>viim</code> , <code>viha</code> , <code>viib0-1</code>	Write



**Table 11-31. Interruptions when Interruption Control Register Read Optimization is Enabled**

Instructions	Interruptions
Move from interruption control registers	When the interruption control register read optimization is enabled, reads of interruption control registers with PSR.vm==1, may raise the following faults: <ul style="list-style-type: none"> <li>Illegal Operation fault – if vpsr.ic is not zero or the target operand specifies GR 0 or an out-of-frame stacked register</li> <li>Privileged Operation fault – if vpsr.cpl is not zero</li> </ul>

#### 11.7.4.2.3 Interruption Control Register Write Optimization

The interruption control register write optimization is enabled by the `a_to_int_cr` bit in the Virtualization Acceleration Control (`vac`) field in the VPD. When this optimization is enabled, and `vpsr.ic` is 0, software running with `PSR.vm==1` will be able to write the virtual interruption control registers (`vipsr`, `visr`, `viip`, `vifa`, `vitir`, `viipa`, `vifs`, `viim`, `viha`, `viib0-1`) without any intercepts to the VMM, unless a fault condition is detected (see [Table 11-33](#) for details).

If this optimization is disabled, a write of the interruption control registers with `PSR.vm==1` results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see [Table 11-32](#) for details.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-32. Synchronization Requirements for Interruption Control Register Write Optimization**

VPD Resource	Synchronization Required
<code>vipsr</code> , <code>visr</code> , <code>viip</code> , <code>vifa</code> , <code>vitir</code> , <code>viipa</code> , <code>vifs</code> , <code>viim</code> , <code>viha</code> , <code>viib0-1</code>	Read

**Table 11-33. Interruptions when Interruption Control Register Write Optimization is Enabled**

Instructions	Interruptions
Move to interruption control registers	When the interruption control register write optimization is enabled, writes to interruption control registers with <code>PSR.vm==1</code> , may raise the following faults: <ul style="list-style-type: none"> <li>Illegal Operation fault – if <code>vpsr.ic</code> is not zero</li> <li>Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> <li>Register NaT Consumption fault – if the NaT bit of the source operand is one</li> <li>Reserved Register/Field fault – if any reserved field in the specified control register is written with a non-zero value</li> <li>Unimplemented Data Address fault – if writing to <code>vifa</code> and an unimplemented virtual address is specified</li> </ul>

#### 11.7.4.2.4 MOV-from-PSR Optimization

The MOV-from-PSR optimization is enabled by the `a_from_psr` bit in the Virtualization Acceleration Control (`vac`) field in the VPD. When this optimization is enabled, software running with `PSR.vm==1` will be able to execute MOV-from-PSR instructions to read

the virtual processor status register without any intercepts to the VMM; and the last value written to the vpsr will be returned, unless a fault condition is detected (see [Table 11-35](#) for details). The value returned for the fml, mfh, ac, up and be bits are simply the values of those bits in the PSR of the logical processor, since those bits are not virtualized.

If this optimization is disabled, execution of a MOV-from-PSR instruction with PSR.vm==1 results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see [Table 11-34](#) for details.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-34. Synchronization Requirements for MOV-from-PSR Optimization**

VPD Resource	Synchronization Required
vpsr{36:35, 31:6} See <a href="#">Table 11-17</a> , “Virtual Processor Descriptor (VPD) – VPSR” on page 2:328 for details.	Write

**Table 11-35. Interruptions when MOV-from-PSR Optimization is Enabled**

Instructions	Interruptions
MOV-from-PSR	When the MOV-from-PSR optimization is enabled, MOV-from-PSR instructions with PSR.vm==1, may raise the following faults: <ul style="list-style-type: none"> <li>Illegal Operation fault – if the target operand specifies GR 0 or an out-of-frame stacked register</li> <li>Privileged Operation fault – if vpsr.cpl is not zero</li> </ul>

**Note:** This field cannot be enabled together with the d\_psr\_i virtualization disable control (vdc) described in [Section 11.7.4.3.7](#), “Disable PSR Interrupt-bit Virtualization” on page 2:348. If this control is enabled together with the d\_psr\_i control, an error will be returned during PAL\_VP\_CREATE and PAL\_VP\_REGISTER. See [Section 11.7.4.4](#), “Virtualization Optimization Combinations” on page 2:349 for details.

#### 11.7.4.2.5 MOV-from-CPUID Optimization

The MOV-from-CPUID optimization is enabled by the a\_from\_cpuid bit in the Virtualization Acceleration Control (vac) field in the VPD. When this optimization is enabled, software running with PSR.vm==1 will be able to execute MOV-from-CPUID instruction to read the virtual CPUID registers without any intercepts to the VMM; and the corresponding VCPUID value from the VPD will be returned, unless a fault condition is detected (see [Table 11-37](#) for details).

If this optimization is disabled, execution of a MOV-from-CPUID instruction with PSR.vm==1 results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see [Table 11-36](#) for details.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-36. Synchronization Requirements for MOV-from-CPUID Optimization**

VPD Resource	Synchronization Required
vcpuid0-4	Write

**Table 11-37. Interruptions when MOV-from-CPUID Optimization is Enabled**

Instructions	Interruptions
MOV-from-CPUID	<p>When the MOV-from-CPUID optimization is enabled, MOV-from-CPUID instructions with PSR.vm==1, may raise the following faults:</p> <ul style="list-style-type: none"> <li>• Illegal Operation fault – if the target operand specifies GR 0 or an out-of-frame stacked register</li> <li>• Register NaT Consumption fault – if the NaT bit in the target register is one</li> <li>• Reserved Register/Field fault – if a reserved CPUID register is being read</li> </ul>

#### 11.7.4.2.6 Cover Optimization

The cover optimization is enabled by the `a_cover` bit in the Virtualization Acceleration Control (`vac`) field in the VPD. When this optimization is enabled, software running with PSR.vm==1 will be able to execute `cover` instructions without any intercepts to the VMM, unless a fault condition is detected (see Table 11-39 for details). The `cover` instruction will execute and `vcr.ifs` will be updated if `vpsr.ic` is 0.

If this optimization is disabled, execution of the `cover` instruction with PSR.vm==1 results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see Table 11-38 for details.

When this optimization is enabled, certain VPD state is accessed, as described in Table 11-16, “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-38. Synchronization Requirements for Cover Optimization**

VPD Resource	Synchronization Required
vifs	Read, Write

**Table 11-39. Interruptions when Cover Optimization is Enabled**

Instructions	Interruptions
<code>cover</code>	<p>When the cover optimization is enabled, <code>cover</code> instructions with PSR.vm==1, may raise the following faults:</p> <ul style="list-style-type: none"> <li>• Illegal Operation fault – if the instruction is not the last instruction in an instruction group</li> </ul>

#### 11.7.4.2.7 Bank Switch Optimization

The bank switch optimization is enabled by the `a_bsw` bit in the Virtualization Acceleration Control (`vac`) field in the VPD. When this optimization is enabled, execution of the `bsw` instruction with PSR.vm==1 spills the currently active banked registers and the corresponding NaT bits to the VPD, and loads the other banked registers and the

corresponding NaT bits from the VPD. `vpsr.bn` is updated to reflect the new register bank without any intercepts to the VMM, unless a fault condition is detected (see [Table 11-46](#) for details).

If this optimization is disabled, execution of the `bsw` instruction with `PSR.vm==1` results in a virtualization intercept.

Synchronization is required when this optimization is enabled, see [Table 11-40](#) for details.

**Table 11-40. Synchronization Requirements for Bank Switch Optimization**

VPD Resource	Synchronization Required
<code>vpsr.bn</code>	Read, Write

**Table 11-41. Interruptions when Bank Switch Optimization is Enabled**

Instructions	Interruptions
<code>bsw</code>	<p>When the bank switch optimization is enabled, <code>bsw</code> instructions with <code>PSR.vm==1</code>, may raise the following faults:</p> <ul style="list-style-type: none"> <li>• Illegal Operation fault – if the instruction is not the last instruction in an instruction group</li> <li>• Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> </ul>

**Note:** This field cannot be enabled together with the `d_psr_i` virtualization disable control (`vdc`) described in [Section 11.7.4.3.7, “Disable PSR Interrupt-bit Virtualization” on page 2:348](#). If this control is enabled together with the `d_psr_i` control, an error will be returned during `PAL_VP_CREATE` and `PAL_VP_REGISTER`. See [Section 11.7.4.4, “Virtualization Optimization Combinations” on page 2:349](#) for details.

#### 11.7.4.2.8 Probe Instruction Virtualization

The probe instruction virtualization is controlled by the `a_all_probes` and `a_select_probes` bits in the Virtualization Acceleration Control (`vac`) field in the VPD.

When the `a_all_probes` bit is set to 1, all `probe` instructions running at all privilege levels with `PSR.vm==1` will result in virtualization intercepts.

When the `a_select_probes` bit is set to 1, the following `probe` instructions will raise virtualization intercepts when executed with `PSR.vm==1` at the most privileged level (`VPSR.cpl==0`):

- `probe` instructions in immediate-form, with immediate field equal to privilege level 0
- All `probe` instructions in register-form

Please refer to the instruction description page for the `probe` instruction for details on the usage of immediate-form and register-form of the instruction.

**Note:** Software cannot enable both `a_all_probes` and `a_select_probes` bits together - an error will be returned during `PAL_VP_CREATE` and `PAL_VP_REGISTER`.

The virtualization of `probe` instructions is not supported on all processor implementations. Software can call `PAL_VP_ENV_INFO` to determine the availability of this feature.

There is no synchronization requirement for the virtualization of `probe` instructions.

#### 11.7.4.2.9 Test Feature Optimization

The test feature optimization is enabled by the `a_tf` bit in the Virtualization Acceleration Control (`vac`) field in the VPD.

When this optimization is enabled, test feature (`tf`) instructions running with `PSR.vm==1` will test the `VCPUID[4]` register in the VPD. The VMM may maintain a different `VCPUID[4]{63:32}` value from the `CPUID[4]{63:32}` value of the logical processor on which the virtual processor is running.

If the VMM indicates to a guest that an instruction is not supported by clearing the corresponding bit in `VCPUID[63:32]`, then guest execution of that instruction, when `a_tf` is enabled, will behave the same as it would in implementations that do not implement that instruction. See [Table 11-42](#) for more information.

**Table 11-42. Impact of clearing VCPUID bits with the `a_tf` optimization**

VCPUID[4] bit	Instructions affected	Behavior when <code>vcpuid[4]</code> bit is 0
32	<code>clz</code>	Illegal Operation fault
33	<code>mpy4</code> <code>mpyshl4</code>	Illegal Operation fault Illegal Operation fault

If this optimization is disabled or not supported, execution of the test feature (`tf`) instruction with `PSR.vm==1` will test the `CPUID[4]` register. The VMM must maintain the same `VCPUID[4]{63:32}` value as the `CPUID[4]{63:32}` value of the logical processor on which the virtual processor is running.

Synchronization is required when this optimization is enabled; see [Table 11-43](#) for details.

This optimization is not supported on all processor implementations. Software can call `PAL_VP_ENV_INFO` to determine the availability of this feature.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-43. Synchronization Requirements for Test Feature Optimization**

VPD Resource	Synchronization Required
<code>vcpuid[4]{63:32}</code>	Write

#### 11.7.4.2.10 Interruption Collection and User Mask Optimization

The interruption collection and user mask optimization is enabled by the `a_ic_um` bit in the Virtualization Acceleration Control (`vac`) field in the VPD.

When this optimization is enabled and `PSR.vm==1`, execution of `rsm` and `ssm` instructions<sup>1</sup> with a mask targeting no fields other than the `ic` and user mask fields will not intercept to the VMM, unless a fault condition is detected (see [Table 11-45](#) for details). The `ic` field in `vpsr` and user mask bits in PSR targeted by the mask will be updated with the new value.

When this optimization is enabled, execution of `rsm` and `ssm` instructions, with `PSR.vm==1` and system mask equal to zero (0x0), will not intercept to the VMM unless a fault condition is detected (see [Table 11-45](#) for details).

When `PSR.vm==1`, execution of `rsm` and `ssm` instructions<sup>1</sup>, which modify any bits other than `vpsr.ic` and user mask fields will result in virtualization intercepts independent of whether this optimization is enabled or not.

Synchronization is required when this optimization is enabled; see [Table 11-44](#) for details.

This optimization is not supported on all processor implementations. Software can call `PAL_VP_ENV_INFO` to determine the availability of this feature.

When this optimization is enabled, certain VPD state is accessed, as described in [Table 11-16](#), “Virtual Processor Descriptor (VPD)” on page 2:326.

**Table 11-44. Synchronization Requirements for Interrupt Collection and User Mask Optimization**

VPD Resource	Synchronization Required
<code>vpsr.ic</code>	Read, Write

**Table 11-45. Interruptions when Interrupt Collection and User Mask Optimization is Enabled**

Instructions	Interruptions
<code>rsm</code> , <code>ssm</code>	When the interruption collection and user mask optimization is enabled, execution of <code>rsm</code> and <code>ssm</code> instructions with <code>PSR.vm==1</code> which modify <code>vpsr.ic</code> and any user mask fields, may raise the following faults: <ul style="list-style-type: none"> <li>Privileged Operation fault – if <code>vpsr.cpl</code> is not zero</li> </ul>

### 11.7.4.3 Virtualization Disables

[Table 11-26](#) summarizes the virtualization disables supported in Itanium architecture.

**Table 11-46. Virtualization Disables Summary**

Disable	Virtualization Disable Control ( <i>vdc</i> ) <sup>a</sup>	Description
Disable VMSW Instruction	<code>d_vmsw</code>	<a href="#">Section 11.7.4.3.1</a>
Disable External Interrupt Control Register Virtualization	<code>d_extint</code>	<a href="#">Section 11.7.4.3.2</a>
Disable Breakpoint Register Virtualization	<code>d_ibr_dbr</code>	<a href="#">Section 11.7.4.3.3</a>
Disable PMC Virtualization	<code>d_pmc</code>	<a href="#">Section 11.7.4.3.4</a>
Disable MOV-to-PMD Virtualization	<code>d_to_pmd</code>	<a href="#">Section 11.7.4.3.5</a>

1. The execution of `rsm` and `ssm` instructions with `PSR.vm==1` is affected by both the virtual external interrupt optimization (`a_int`) and the interruption collection and user mask optimization (`a_ic_um`). Software can enable or disable both optimizations together, or enable each optimization independently. [Section 11.7.4.4.1](#), “Virtual External Interrupt Optimization and Interruption Collection and User Mask Optimization” on page 2:349 describes the behavior when both optimizations are enabled.

**Table 11-46. Virtualization Disables Summary (Continued)**

Disable	Virtualization Disable Control (vdc) <sup>a</sup>	Description
Disable ITM Virtualization	d_itm	<a href="#">Section 11.7.4.3.6</a>
Disable PSR Interrupt-bit Virtualization	d_psr_i	<a href="#">Section 11.7.4.3.7</a>

a. The Virtualization Disable Control (vdc) field resides in the Virtual Processor Descriptor (VPD), see [Section 11.7.1, “Virtual Processor Descriptor \(VPD\)” on page 2:325](#) for details.

#### 11.7.4.3.1 Disable VMSW Instruction

The VMSW instruction disable is controlled by the d\_vmsw bit in the Virtualization Disable Control (vdc) field in the VPD. When this control is set to 1, the vmsw instruction is disabled on the logical processor. Execution of the vmsw instruction, independent of the state of PSR.vm, results in a virtualization intercept.

If this control is set to 0, the vmsw instruction can be executed by both the VMM and guest without virtualization intercepts, if PSR.it is 1 and the vmsw instruction is executed on a page with access rights of 7.

#### 11.7.4.3.2 Disable External Interrupt Control Register Virtualization

The external interrupt control register virtualization disable is controlled by the d\_extint bit in the Virtualization Disable Control (vdc) field in the VPD. When this control is set to 1, the external interrupt control registers (CR65-71) are not virtualized, and code running with PSR.vm==1 can read and write these resources directly without any intercepts to the VMM.

If this control is set to 0, accesses (reads/writes) to the external interruption control registers with PSR.vm==1 result in virtualization intercepts.

**Note:** This field cannot be enabled together with the a\_int virtualization acceleration control (vac) described in [Section 11.7.4.2.1, “Virtual External Interrupt Optimization” on page 2:338](#). If this control is enabled together with the a\_int control, an error will be returned during PAL\_VP\_CREATE and PAL\_VP\_REGISTER. See [Section 11.7.4.4, “Virtualization Optimization Combinations” on page 2:349](#) for details.

#### 11.7.4.3.3 Disable Breakpoint Register Virtualization

The breakpoint register virtualization disable is controlled by the d\_ibr\_dbr bit in the Virtualization Disable Control (vdc) field in the VPD. When this control is set to 1, accesses (reads/writes) to the data and instruction breakpoint registers (DBR/IBR) are not virtualized, and code running with PSR.vm==1 can read and write these resources directly without any intercepts to the VMM.

If this control is set to 0, accesses (reads/writes) to the breakpoint registers with PSR.vm==1 result in virtualization intercepts.

#### 11.7.4.3.4 Disable PMC Virtualization

The PMC virtualization disable is controlled by the `d_pmc` bit in the Virtualization Disable Control (`vd_c`) field in the VPD. When this control is set to 1, accesses (reads/writes) to the performance monitor configuration registers (PMCs) are not virtualized, and code running with `PSR.vm==1` can read and write these resources directly without any intercepts to the VMM.

If this control is set to 0, accesses (reads/writes) to the performance counter configuration registers with `PSR.vm==1` result in virtualization intercepts.

#### 11.7.4.3.5 Disable MOV-to-PMD Virtualization

The MOV-to-PMD<sup>1</sup> virtualization disable is controlled by the `d_to_pmd` bit in the Virtualization Disable Control (`vd_c`) field in the VPD. When this control is set to 1, writes to the performance monitor data registers (PMDs) are not virtualized, and code running with `PSR.vm==1` can write these resources directly without any intercepts to the VMM.

If this control is set to 0, writes to the performance monitor data registers with `PSR.vm==1` result in virtualization intercepts.

#### 11.7.4.3.6 Disable ITM Virtualization

The ITM virtualization disable is controlled by the `d_itm` bit in the Virtualization Disable Control (`vd_c`) field in the VPD. When this control is set to 1, writes to the Interval Timer Match (ITM) register are not virtualized, and code running with `PSR.vm==1` can write this resource directly without any intercepts to the VMM.

If this control is set to 0, writes to the ITM register with `PSR.vm==1` result in virtualization intercepts.

#### 11.7.4.3.7 Disable PSR Interrupt-bit Virtualization

The PSR interrupt-bit virtualization disable is controlled by the `d_psr_i` bit in the Virtualization Disable Control (`vd_c`) field in the VPD. When this control is set to 1, accesses (reads/writes) to the interrupt bit in processor state register (`PSR.i`) are not virtualized. Code running with `PSR.vm==1` can read and write to `PSR.i` through `ssm` and `rsm` instructions without any intercepts to the VMM. Attempts to modify other PSR bits in addition to the interrupt bit via the `ssm` and `rsm` instructions will result in virtualization intercepts.

This control has no effect on `mov psr.l` instructions; attempts to modify the interrupt bit with the `mov psr.l` instruction result in virtualization intercepts.

**Note:** This field cannot be enabled together with `a_int`, `a_from_psr` or `a_bsw` virtualization accelerations. If this control is enabled together with any one of described accelerations, an error will be returned during `PAL_VP_CREATE` and `PAL_VP_REGISTER`. See [Section 11.7.4.4, "Virtualization Optimization Combinations"](#) on page 2:349 for details.

- 
1. The MOV-from-PMD instruction is not virtualized. Hence there is no need to provide optimizations for the MOV-from-PMD instruction.



#### 11.7.4.4 Virtualization Optimization Combinations

Table 11-47 describes the supported combinations of virtualization accelerations and disables.

**Table 11-47. Supported Virtualization Optimization Combinations**

	d_vmsw	d_extint	d_ibr_dbr	d_pmc	d_to_pmd	d_itm	d_psr_i
a_int	o <sup>a</sup>	x <sup>b</sup>	o	o	o	o	x
a_from_int_cr	o	o	o	o	o	o	o
a_to_int_cr	o	o	o	o	o	o	o
a_from_psr	o	o	o	o	o	o	x
a_from_cpuid	o	o	o	o	o	o	o
a_cover	o	o	o	o	o	o	o
a_bsw	o	o	o	o	o	o	x
a_all_probes	o	o	o	o	o	o	o
a_select_probes	o	o	o	o	o	o	o
a_tf	o	o	o	o	o	o	o
a_ic_um	o	o	o	o	o	o	o

a. “o” indicates the corresponding virtualization acceleration and disable can be enabled together.

b. “x” indicates the corresponding virtualization acceleration and disable cannot be enabled together.

##### 11.7.4.4.1 Virtual External Interrupt Optimization and Interruption Collection and User Mask Optimization

The execution of `rsm` and `ssm` instructions with `PSR.vm==1` is affected by both of these optimizations:

- Virtual External Interrupt Optimization (`a_int`), described in [Section 11.7.4.2.1, “Virtual External Interrupt Optimization”](#), and
- Interruption Collection and User Mask Optimization (`a_ic_um`), described in [Section 11.7.4.2.10, “Interruption Collection and User Mask Optimization”](#).

Software can enable or disable both optimizations together, or enable each optimization independently.

When both optimizations are enabled and `PSR.vm==1`, `rsm` and `ssm` instructions with a mask targeting any fields in `i`, `ic` and user mask will not be intercepted to the VMM, unless a fault condition is detected. The `i` and `ic` fields in `vpsr` and user mask in `PSR` will be updated with the new value.

When `PSR.vm==1`, `rsm` and `ssm` instructions with a mask targeting any fields other than `i`, `ic` and user mask fields will result in virtualization intercepts independent of whether these two optimizations are enabled or not.

#### 11.7.4.5 Virtualization Synchronizations

When certain virtualization accelerations described in [Section 11.7.4.2, “Virtualization Accelerations” on page 2:337](#) are enabled, processor implementations can provide implementation-specific control resources to enhance the performance of virtual processors. Two PAL services are provided to synchronize the implementation-specific control resources and the resources in the VPD. There are two types of synchronizations:

1. **Read synchronization** – When a specific acceleration is enabled, after interruptions and intercepts that occur when PSR.vm was 1, the VMM must invoke PAL\_VPS\_SYNC\_READ to synchronize the related resources before reading their values from the VPD.
2. **Write synchronization** – When a specific acceleration is enabled, the VMM must invoke PAL\_VPS\_SYNC\_WRITE to synchronize the related resources after modifying their values in the VPD and before resuming the virtual processor.

For details on PAL\_VPS\_SYNC\_READ and PAL\_VPS\_SYNC\_WRITE, see [Section 11.11.2, “PAL Virtualization Service Specifications” on page 2:488](#).

Read and/or write synchronizations are required only if the specific acceleration is enabled. For the resources that require synchronizations if the acceleration is enabled, failure to perform the proper synchronizations will result in undefined processor behavior<sup>1</sup>.

The synchronization requirements of the related resources for each acceleration are described in the corresponding sections for each acceleration in [Section 11.7.4.2, “Virtualization Accelerations” on page 2:337](#).

No synchronization is required for any of the virtualization disables.

## 11.8 PAL Glossary

### Corrected Error

All errors of this type are corrected by the platform or processor in either hardware or firmware. This severity is for logging purposes only. There is no architectural damage caused by the detecting and reporting functions. Corrected errors require no operating system intervention to correct the error.

### Corrected Machine Check (CMC)

A corrected machine check is a machine check that has been successfully corrected by hardware and/or firmware. Information about the cause of the error is recorded, and an interrupt is set to allow the Operating System software to examine and diagnose the error. Return is controlled to the program executing at the time of the error.

### Entrypoint

A firmware entrypoint is a piece of code which is triggered by a hardware event, usually the assertion of a processor pin, or the receipt of an interruption. If return to the caller is done, it is through the RFI instruction. The currently defined PAL entrypoints are PALE\_RESET, PALE\_INIT, PALE\_PMI, and PALE\_CHECK.

### Fatal Error

An uncorrected error which can corrupt state, and the state information is not known. These type of errors cannot be corrected by the hardware, firmware, or the operating system. The integrity of the system, including the IO devices is not guaranteed and may require I/O device initialization and a system reboot to continue. Fatal errors may or may not be contained within the processor or memory hierarchy.

---

1. Virtual machine monitors must perform all the required synchronizations specified. Virtual machine monitors not conforming to this specification are not guaranteed to work on all processor implementations.

**Machine Check (MC)**

A machine check is a hardware event that indicates that a hardware error or architectural violation has occurred that threatens to damage the architectural state of the machine, possibly causing data corruption. The occurrence of the error triggers the execution of firmware code which records information about the error, and attempts to recover when possible.

**OLR**

On line replacement. The replacement of a computer component while the system is up and running without requiring a reboot.

**PAL Intercepts**

Interfaces where PAL transfers control to the VMM on virtualization events (execution of virtualized instructions/operations with `PSR.vm=1`). For details see [Section 11.7.3, "PAL Intercepts in Virtual Environment"](#) on page 2:332.

**Power-on**

The reset event that occurs when the power input to the processor is applied and the reset input to the processor is asserted.

**Preserved**

When applied to an entrypoint, preserved means that the value contained in a register at exit from the entrypoint code is the same as the value at the time of the hardware event that caused the entrypoint to be invoked. When applied to a procedure, preserved means that the value contained in a register at exit from the procedure is the same as the value at entry to the procedure. The value may have been changed and restored before exit.

**Processor Abstraction Layer (PAL)**

PAL is firmware that abstracts processor implementation differences and provides a consistent interface to higher level firmware and software. PAL has no knowledge of platform implementation details.

**Procedure**

A firmware procedure is a piece of code which is called from other firmware or software, and which uses the return mechanism of the Itanium Runtime Calling Conventions to return to its caller.

**Recoverable Error**

An uncorrected error which can corrupt state, but the state information is known. Recoverable errors cannot be corrected by either the hardware or firmware. This type of error requires operating system analysis and a corrective action to recover. System operation/state may be impacted.

**Reserved**

When applied to a data variable, it means that the variable must not be used to convey information. All software passing the variable must place a value of zero in the variable. The occurrence of a non-zero value may cause undefined results.

When applied to a value or range of values, any values not defined in the range and specified as reserved must not be used. The occurrence of a reserved value may cause undefined results.

**Reset**

The reset event that occurs when the reset input to the processor is asserted.

**Scratch**

When applied to either an entrypoint or procedure, scratch means that the contents of the register has no meaning and need not be preserved. Further the register is available for the storage of local variables. Unless otherwise noted, the register should not be relied upon to contain any particular value after exit.

**Stacked Calling Convention**

The firmware calling convention which adheres fully to the Itanium Runtime Calling Conventions. To use this calling convention, the RSE must be working and usable.

**Static Calling Convention**

The firmware calling convention which adheres to the Itanium Runtime Calling Conventions for the static general registers, branch registers, predicate registers, but for which all other registers are unused except for the RSE control registers. The RSE is placed in enforced lazy mode, and the stacked general registers or memory are not referenced.

**System Abstraction Layer (SAL)**

SAL is firmware that abstracts platform implementation differences for higher level software. SAL has no knowledge of processor implementation details.

**Unchanged**

When applied to an entrypoint, unchanged means that the register referenced has not been changed from the time of the hardware event that caused the entrypoint to be invoked until it exited to higher level firmware or software. When applied to a procedure, unchanged means that the register referenced has not been changed from procedure entry until procedure exit. In all cases, the value at exit is the same as the value at entry or the occurrence of the hardware event.

**Virtual Machine Monitor (VMM)**

The VMM is the system software which implements software policies to manage/support virtualization of processor and platform resources.

**Virtual Processor Descriptor (VPD)**

Represents the abstraction of the processor resources of a single virtual processor. The VPD consists of per-virtual-processor control information together with performance-critical architectural state. See [Section 11.7.1, "Virtual Processor Descriptor \(VPD\)" on page 2:325](#) for details.

**Virtual Processor State**

A memory data structure which represents the architectural state of a virtual processor. Part of the virtual processor state is located in the Virtual Processor Descriptor (VPD), and the rest is located in memory data structures maintained by the virtual machine monitor.

## 11.9 PAL Code Memory Accesses and Restrictions

PAL issues load and store operations to memory in the following cases with the following memory attributes:

- During machine check/INIT handling to the min-state save area memory region registered with PAL using the UC memory attribute.

- During the execution of PAL procedures to the memory buffer allocated by the caller of the procedure using the memory attribute of the address passed by the caller.
- PAL may also issue loads from the architected firmware address space and loads/stores from the registered min-state save area whenever it is executing a PAL procedure or handling PAL-based interruptions (reset, MCA, INIT and PMI). PAL code may use either the UC or WBL memory attribute when accessing these areas.

PAL code will not send IPIs that require any special support from the platform.

## 11.10 PAL Procedures

PAL procedures may be called by higher-level firmware and software to obtain information about the identification, configuration, and capabilities of the processor implementation, or to perform implementation-dependent functions such as cache initialization. These procedures access processor implementation-dependent hardware to return information that characterizes and identifies the processor or implements a defined function on that particular processor.

PAL procedures are implemented by a combination of firmware code and hardware. The PAL procedures are defined to be relocatable from the firmware address space. Higher level firmware and software must perform this relocation during the reset flow. The PAL procedures may be called both before and after this relocation occurs, but performance will usually be better after the relocation. In order to ensure no problems occur due to the relocation of the PAL procedures, these procedures are written to be position independent. All references to constant data done by the procedures is done in an IP relative way.

PAL procedures are provided to return information or allow configuration of the following processor features:

- Cache and memory features supported by the processor
- Processor identification, features, and configuration
- Machine Check Abort handling
- Power state information and management
- Processor self test
- Firmware utilities

PAL procedures are implemented as a single high level procedure, named `PAL_PROC`, whose first argument is an index which specifies which PAL procedure is being called. Indices are assigned depending on the nature of the PAL procedure being referenced, according to [Table 11-48](#).

**Table 11-48. PAL Procedure Index Assignment**

Index	Description
0	Reserved
1 - 255	Architected procedures; static register calling conventions
256 - 511	Architected procedures; stacked register calling conventions
512 - 767	Implementation-specific procedures; static registers calling conventions
768 - 1023	Implementation-specific procedures; stacked register calling conventions
1024 +	Reserved

The assignment of indices for all architected procedures is controlled by this document. The assignment of indices for implementation-specific procedures is controlled by the specific processor for which the procedures are implemented. No implementation-specific procedure calls are required for the correct operation of a processor. No SAL or operating system code should ever have to call an implementation-specific procedure call for normal activity. They are reserved for diagnostic and bring-up software and the results of such calls may be unpredictable.

Architected procedures may be designated as required or optional. If a procedure is designated as optional, a unique return code will be returned to indicate the procedure is not present in this PAL implementation. It is the caller's responsibility to check for this return code after calling any optional PAL procedure

In addition to the calling conventions described below, PAL procedure calls may be made in physical mode (PSR.it=0, PSR.rt=0, and PSR.dt=0) or virtual mode (PSR.it=1, PSR.rt=1, and PSR.dt=1). All PAL procedures may be called in physical mode. Only those procedures specified later in this chapter may be called in virtual mode. PAL procedures written to support virtual mode, and the caller of PAL procedures written in virtual mode must obey the restrictions documented in this chapter, otherwise the results of such procedure calls may be unpredictable.

### 11.10.1 PAL Procedure Summary

The following tables summarize the PAL procedures by application area. Included are the name of the procedure, the index of the procedure, the class of the procedure (whether required or optional), the calling convention used for the procedure (static or stacked), and whether the procedure can be called in physical mode only, virtual mode only, or both physical and virtual modes.

On processor implementations with multiple logical processors in a physical processor package, calling a certain PAL procedures may affect resources shared by the logical processors. In the following tables, procedures that may affect resources on multiple processors are marked next to the corresponding procedure names; procedures that are not marked have no effects on other logical processors.

**Table 11-49. PAL Cache and Memory Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_CACHE_FLUSH <sup>a</sup>	1	Req.	Static	Both	No	Flush the instruction or data caches.
PAL_CACHE_INFO	2	Req.	Static	Both	No	Return detailed instruction or data cache information.
PAL_CACHE_INIT <sup>a</sup>	3	Req.	Static	Phys.	No	Initialize the instruction or data caches.

**Table 11-49.PAL Cache and Memory Procedures (Continued)**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_CACHE_PROT_INFO	38	Req.	Static	Both	No	Return instruction or data cache protection information.
PAL_CACHE_SHARED_INFO	43	Opt.	Static	Both	No	Returns information on which logical processors share caches.
PAL_CACHE_SUMMARY	4	Req.	Static	Both	No	Return a summary of the cache hierarchy.
PAL_MEM_ATTRIB	5	Req.	Static	Both	No	Return a list of supported memory attributes.
PAL_PREFETCH_VISIBILITY	41	Req.	Static	Both	No	Used in architected sequence to transition pages from a cacheable, speculative attribute to an uncacheable attribute. See <a href="#">Section 4.4.11.2, “Physical Addressing Attribute Transition – Disabling Prefetch/Speculation and Removing Cacheability”</a> on page 2:90.
PAL_PTCE_INFO	6	Req.	Static	Both	No	Return information needed for <code>ptc.e</code> instruction to purge entire TC.
PAL_VM_INFO	7	Req.	Static	Both	No	Return detailed information about virtual memory features supported in the processor.
PAL_VM_PAGE_SIZE	34	Req.	Static	Both	No	Return virtual memory TC and hardware walker page sizes supported in the processor.
PAL_VM_SUMMARY	8	Req.	Static	Both	No	Return summary information about virtual memory features supported in the processor.
PAL_VM_TR_READ	261	Req.	Stacked	Phys.	No	Read contents of a translation register.

a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-50.PAL Processor Identification, Features, and Configuration Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_BRAND_INFO	274	Opt.	Stacked	Both	No	Provides processor branding information.
PAL_BUS_GET_FEATURES	9	Req.	Static	Phys.	No	Return configurable processor bus interface features and their current settings.
PAL_BUS_SET_FEATURES <sup>a</sup>	10	Req.	Static	Phys.	No	Enable or disable configurable features in processor bus interface.
PAL_DEBUG_INFO	11	Req.	Static	Both	No	Return the number of instruction and data breakpoint registers.
PAL_FIXED_ADDR	12	Req.	Static	Both	No	Return the fixed component of a processor’s directed address.
PAL_FREQ_BASE	13	Opt.	Static	Both	No	Return the frequency of the output clock for use by the platform, if generated by the processor.
PAL_FREQ_RATIOS	14	Req.	Static	Both	No	Return ratio of processor, bus, and interval time counter to processor input clock or output clock for platform use, if generated by the processor.
PAL_GET_HW_POLICY	48	Opt.	Static	Both	Dep.	Get current hardware resource sharing policy.
PAL_LOGICAL_TO_PHYSICAL	42	Opt.	Static	Both	No	Return information on which logical processors map to a physical processor package.
PAL_PERF_MON_INFO	15	Req.	Static	Both	No	Return the number and type of performance monitors.
PAL_PLATFORM_ADDR <sup>a</sup>	16	Req.	Static	Both	No	Specify processor interrupt block address and I/O port space address.
PAL_PROC_GET_FEATURES	17	Req.	Static	Phys.	No	Return configurable processor features and their current setting.

**Table 11-50. PAL Processor Identification, Features, and Configuration Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_PROC_SET_FEATURES <sup>a</sup>	18	Req.	Static	Phys.	No	Enable or disable configurable processor features.
PAL_REGISTER_INFO	39	Req.	Static	Both	No	Return AR and CR register information.
PAL_RSE_INFO	19	Req.	Static	Both	No	Return RSE information.
PAL_SET_HW_POLICY <sup>a</sup>	49	Opt.	Static	Both	Dep.	Set current hardware resource sharing policy.
PAL_VERSION	20	Req.	Static	Both	No	Return version of PAL code.

a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-51. PAL Machine Check Handling Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_MC_CLEAR_LOG <sup>a</sup>	21	Req.	Static	Both	No	Clear all error information from processor error logging registers.
PAL_MC_DRAIN	22	Req.	Static	Both	No	Ensure that all operations that could cause an MCA have completed.
PAL_MC_DYNAMIC_STATE	24	Opt.	Static	Both	No	Return Processor Dynamic State for logging by SAL.
PAL_MC_ERROR_INFO	25	Req.	Static	Both	No	Return Processor Machine Check Information and Processor Static State for logging by SAL.
PAL_MC_ERROR_INJECT <sup>a</sup>	276	Opt.	Stacked	Both	Dep.	Injects the requested processor error or returns information on the supported injection capabilities for this particular processor implementation.
PAL_MC_EXPECTED	23	Req.	Static	Phys.	No	Set/Reset Expected Machine Check Indicator.
PAL_MC_HW_TRACKING	51	Opt.	Static	Both	Dep.	Query which hardware structures are performing hardware status tracking
PAL_MC_REGISTER_MEM	27	Req.	Static	Phys.	No	Register min-state save area with PAL for machine checks and inits.
PAL_MC_RESUME	26	Req.	Static	Phys.	No	Restore minimal architected state and return to interrupted process.

a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-52. PAL Power Information and Management Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_GET_PSTATE	262	Opt.	Stacked	Both	Dep.	Returns information on the performance index of the processor.
PAL_HALT	28	Opt.	Static	Phys	No	Enter the low-power HALT state or an implementation-dependent low-power state.
PAL_HALT_INFO	257	Req.	Stacked	Both	No	Return the low power capabilities of the processor.
PAL_HALT_LIGHT	29	Req.	Static	Both	No	Enter the low power LIGHT HALT state.
PAL_PSTATE_INFO	44	Opt.	Static	Both	No	Returns information about the P-states supported by the processor.
PAL_SET_PSTATE <sup>a</sup>	263	Opt.	Stacked	Both	Dep.	Request processor to enter power/performance state.
PAL_SHUTDOWN	45	Opt.	Static	Phys	Dep.	Puts the processor in a low power state which can be exited only by a reset event.



a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-53. PAL Processor Self Test Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_CACHE_LINE_INIT <sup>a</sup>	31	Req.	Static	Phys.	No	Initialize tags and data of a cache line for processor testing.
PAL_CACHE_READ	259	Opt.	Stacked	Phys.	No	Read tag and data of a cache line for diagnostic testing.
PAL_CACHE_WRITE <sup>a</sup>	260	Opt.	Stacked	Phys.	No	Write tag and data of a cache for diagnostic testing.
PAL_TEST_INFO	37	Req.	Static	Phys.	No	Returns alignment and size requirements needed for the memory buffer passed to the PAL_TEST_PROC procedure as well as information on self-test control words for the processor self tests.
PAL_TEST_PROC <sup>a</sup>	258	Req.	Stacked	Phys.	No	Perform late processor self test.

a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-54. PAL Support Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_COPY_INFO	30	Req.	Static	Phys.	No	Return information needed to relocate PAL procedures and PAL PMI code to memory.
PAL_COPY_PAL	256	Req.	Stacked	Phys.	No	Relocate PAL procedures and PAL PMI code to memory.
PAL_MEMORY_BUFFER <sup>a</sup>	277	Opt.	Stacked	Phys.	No	Provides cacheable memory to PAL for exclusive use during runtime.
PAL_PMI_ENTRYPOINT <sup>a</sup>	32	Req.	Static	Phys.	No	Register PMI memory entrypoints with processor.

a. Calling this procedure may affect resources on multiple processors. Please refer to implementation-specific reference manuals for details.

**Table 11-55. PAL Virtualization Support Procedures**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_VP_CREATE	265	Opt.	Stacked	Virt.	Dep.	Initializes a new VPD for the operation of a new virtual processor in the virtual environment.
PAL_VP_ENV_INFO	266	Opt.	Stacked	Virt.	Dep.	Returns the parameters needed to enter a virtual environment.
PAL_VP_EXIT_ENV	267	Opt.	Stacked	Virt.	Dep.	Allows a logical processor to exit a virtual environment.
PAL_VP_INFO	50	Opt.	Static	Phys.	No	Returns information about virtual processor features.
PAL_VP_INIT_ENV	268	Opt.	Stacked	Virt.	Dep.	Allows a logical processor to enter a virtual environment.
PAL_VP_REGISTER	269	Opt.	Stacked	Virt.	Dep.	Register a different host IVT for the virtual processor.
PAL_VP_RESTORE	270	Opt.	Stacked	Virt.	Dep.	Restore virtual processor state on the logical processor.

**Table 11-55. PAL Virtualization Support Procedures (Continued)**

Procedure	Idx	Class	Conv.	Mode	Buffer	Description
PAL_VP_SAVE	271	Opt.	Stacked	Virt.	Dep.	Save virtual processor state on the logical processor.
PAL_VP_TERMINATE	272	Opt.	Stacked	Virt.	Dep.	Terminates operation for the specified virtual processor.

## 11.10.2 PAL Calling Conventions

The following general rules govern the definition of the PAL procedure calling conventions.

### 11.10.2.1 Overview of Calling Conventions

There are two calling conventions supported for PAL procedures: static registers only and stacked registers. Any single PAL procedure will support only one of the two calling conventions. In addition, PAL procedure may be called in either physical mode (PSR.it=0, PSR.rt=0, and PSR.dt=0) or virtual mode (PSR.it=1, PSR.rt=1, and PSR.dt=1).

#### 11.10.2.1.1 Static Registers Only

This calling convention is intended for boot time usage before main memory may be available or error recovery situations, where memory or the RSE may not be reliable. All parameters are passed in the lower 32 static general registers. The stacked registers will not be used within the procedure. No memory arguments may be passed as parameters to or from PAL procedures written using the static register calling convention. To avoid RSE activity, static register PAL procedures must be called with the br.cond instruction, not the br.call instruction. Please refer to [Table 11-59](#) for a detailed list of the general register usage for static registers only calling convention.

#### 11.10.2.1.2 Stacked Registers

This calling convention is intended for usage after memory has been made available, and for procedures which require memory pointers as arguments. The stacked registers are also used for parameter passing and local variable allocation. This convention conforms to the *Itanium Software Conventions and Runtime Architecture Guide*. Thus, procedures using the stacked register calling convention can be written in the C language. There are two exceptions to the runtime conventions.

1. The first argument to the procedure must also be copied to GR28 prior to making the procedure call. This allows procedures written using both static and stacked register calling conventions to call a single PAL\_PROC entrypoint. This should be accomplished by having the stacked register procedures call a stub module which copies GR32 to GR28, then call PAL\_PROC. It is the responsibility of the caller to provide this stub. Please refer to [Table 11-60](#) for a detailed list of the general register usage for the stacked register calling convention.
2. Floating point registers 32-127 are preserved (rather than scratch, as in the normal Itanium Software Conventions), except on the PAL\_TEST\_PROC procedure. This allows OSs to avoid having to save and restore these registers around a stacked-convention PAL procedure call.

### 11.10.2.1.3 Making PAL Procedure Calls in Physical or Virtual Mode

PAL procedure calls which are made in physical mode must obey the calling conventions described in this chapter, but there are no additional restrictions beyond those noted above. PAL procedure calls made in virtual mode must have the region occupied by PAL\_PROC virtually mapped with an ITR. It needs to map this same area with either a DTR or DTC using the same translation as the ITR. In addition, it must also provide a DTR or DTC mapping for any memory buffer pointers passed as arguments to a procedure. It is the responsibility of the caller to provide these mappings.

If the caller chooses to map the PAL\_PROC area or any memory pointers with a DTC it must call the procedure with PSR.ic = 1 to handle any TLB faults that could occur. The PAL\_PROC code needs to be mapped with an ITR. Unpredictable results may occur if it is mapped with an ITC register.

### 11.10.2.1.4 Dependence on the PAL Memory Buffer

The PAL\_MEMORY\_BUFFER procedure must be called to establish a PAL memory buffer before calling certain PAL procedures that are dependent on the buffer.

## 11.10.2.2 Processor State

The PAL procedures are only available to the code running at privilege level 0. They must run in physical mode (unless specified as callable in virtual mode). PAL procedures are not interruptible by external interrupt or NMI, since PSR.i must be 0 when the PAL procedure is called. PAL procedures are not interruptible by PMI events, if PSR.ic is 0. If PSR.ic is 1, PAL procedures can be interrupted by PMI events. PAL procedures can be interrupted by machine checks and initialization events.

Generally PAL procedures will not enable interruptions not already enabled by the caller. Any PAL call that might cause interruptions (besides data TLB faults, see Section 11.10.2.1.3, "Making PAL Procedure Calls in Physical or Virtual Mode"), must install an IVA handler to handle them. PAL\_TEST\_PROC may generate any interruptions it needs to test the processor.

Table 11-56 defines the requirements for the PSR at entry to and at exit from a PAL procedure call. The operating system must follow the state requirements for PSR shown below. PAL procedure calls made by SAL may impose additional requirements. PAL\_TEST\_PROC may change PSR bits shown as unchanged in order to test the processor. These bits will be preserved in this case. PSR bits are described in increasing bit number order. Any PSR bit numbers not specified are reserved and unchanged.

**Table 11-56. State Requirements for PSR**

PSR Bit	Description	Entry	Exit	Class
be	big-endian memory access enable	0	0	preserved
up	user performance monitor enable	C	C	unchanged
ac	alignment check	C	C	preserved
mfl	floating-point registers f2-f31 written	C	C	preserved
mfh	floating-point registers f32-f127 written	C	C	preserved
ic	interruption state collection enable	0	0	unchanged
		1	1	preserved
i	interrupt enable	0	0	unchanged

**Table 11-56. State Requirements for PSR (Continued)**

PSR Bit	Description	Entry	Exit	Class
pk	protection key validation enable	C	C	unchanged
dt	data address translation enable <sup>a</sup>	0	0	unchanged
		1	1	preserved
dfi	disabled FP register f2 to f31	0	0	unchanged
dfh	disabled FP register f32 to f127 <sup>b</sup>	0	0	unchanged
		1	1	unchanged
sp	secure performance monitors	C	C	unchanged
pp	privileged performance monitor enable	C	C	unchanged
di	disable ISA transition	C	C	preserved
si	secure interval timer	C	C	unchanged
db	debug breakpoint fault enable	0	0	unchanged
lp	lower-privilege transfer trap enable	0	0	unchanged
tb	taken branch trap enable	0	0	unchanged
rt	register stack translation enable <sup>a</sup>	0	0	unchanged
		1	1	preserved
cpl	current privilege level	0	0	unchanged
is	instruction set	0	0	preserved
mc	machine check abort mask <sup>c</sup>	0	0	preserved
		1	1	unchanged
it	instruction address translation enable <sup>a</sup>	0	0	unchanged
		1	1	preserved
id	instruction debug fault disable	0	0	unchanged
da	data access and dirty-bit fault disable	0	0	unchanged
dd	data debug fault disable	0	0	unchanged
ss	single step trap enable	0	0	unchanged
ri	restart instruction	0	0	preserved
ed	exception deferral	0	0	preserved
bn	register bank	1	1	preserved
ia	instruction access-bit fault disable	0	0	unchanged
vm	processor virtualization	0	0	unchanged

- a. PAL procedures which are called in physical mode must remain in physical mode for the duration of the call. PAL procedures which are called in virtual mode, may perform specific actions in physical mode, but must return to the same virtual mode state before returning from the call.
- b. PAL\_TEST\_PROC and an implementation-specific authentication procedure call need to be called with PSR.dfh equal to 0. If they are not they will return invalid argument. All other PAL procedure calls may be called with PSR.dfh equal to 0 or 1.
- c. Most PAL runtime procedures should be called with PSR.mc = 0 except for code flow involved in handling machine checks.

**11.10.2.2.1 Definition of Terms**

The terms used in the definition of the requirements have the following meaning:

**Table 11-57. Definition of Terms**

Term	Description
entry	Start of the first instruction of the PAL procedure.
exit	Start of the first instruction after return to caller's code.

**Table 11-57. Definition of Terms**

Term	Description
0	Must be zero at entry to the procedure or on exit from the procedure. If the value at entry is not zero, the procedure may return an illegal argument or execute in an undefined manner.
1	Must be one at entry to the procedure or on exit from the procedure. If the value at entry is not one, the procedure may return an illegal argument or execute in an undefined manner.
reserved	When any input parameter is listed as reserved, this value must be zero. If an input value has a range of values, any values outside the range, listed as reserved, must not be used. For either case, the PAL procedure may return an illegal argument or execute in an undefined manner.
C	The state of bits marked with C are defined by the caller. If the value at exit is also C, it must be the same as the value at entry.
unchanged	The PAL procedure must not change these values from their entry values during execution of the procedure.
scratch	The PAL procedure may modify these values as necessary during execution of the procedure. The caller cannot rely on these values.
preserved	The PAL procedure may modify these values as necessary during execution of the procedure. However, they will be restored to their entry values prior to exit from the procedure.

#### 11.10.2.2.2 System Registers

The PAL\_TEST\_PROC procedure may change system registers marked as unchanged in order to fully test the processor. When this is done, the values of the system registers will be preserved.

**Table 11-58. System Register Conventions**

Name	Description	Class
DCR	Default Control Register	preserved
ITM	Interval Timer Match Register	unchanged
IVA	Interrupt Vector Address	preserved <sup>a</sup>
PTA	Page Table Address	preserved
GPTA	Guest Page Table Address	preserved
IPSR	Interrupt Processor Status Register	scratch
ISR	Interrupt Status Register	scratch
IIP	Interrupt Instruction Bundle Pointer	scratch
IFA	Interrupt Faulting Address	scratch
ITIR	Interrupt TLB Insertion Register	scratch
IIPA	Interrupt Instruction Previous Address	scratch
IFS	Interrupt Function State	scratch
IIM	Interrupt Immediate Register	scratch
IHA	Interrupt Hash Address	scratch
IIB0-1	Interrupt Instruction Bundle Registers	scratch
LID	Local Interrupt ID	unchanged
IVR	Interrupt Vector Register (read only)	unchanged
TPR	Task Priority Register	unchanged
EOI	End Of Interrupt	unchanged
IRR0-IRR3	Interrupt Request Registers 0-3 (read only)	unchanged
ITV	Interval Timer Vector	unchanged
PMV	Performance Monitoring Vector	unchanged

**Table 11-58. System Register Conventions (Continued)**

Name	Description	Class
CMCV	Corrected Machine Check Vector	unchanged
LRR0-LRR1	Local Redirection Registers 0-1	unchanged
RR	Region Registers	preserved
PKR	Protection Key Registers	preserved
TR	Translation Registers	unchanged <sup>b</sup>
TC	Translation Cache	scratch
IBR/DBR	Break Point Registers	preserved <sup>c</sup>
PMC	Performance Monitor Control Registers	preserved
PMD	Performance Monitor Data Registers	unchanged <sup>d</sup>

- a. On some implementations, PAL virtualization support procedures may program IVA to a different value. Refer to the description of the PAL virtualization procedures for details.
- b. If an implementation provides a means to read TRs for PAL, this should be preserved.
- c. The PAL\_MC\_ERROR\_INJECT may modify these registers if the caller is using the triggering capability. Refer to “PAL\_MC\_ERROR\_INJECT – Inject Processor Error (276)” on page 2:421 for more information.
- d. No PAL procedure writes to the PMD. Depending on the PMC, the PMD may be kept counting performance monitor events during a procedure call. The exception is PAL\_TEST\_PROC, which tests the performance counters.

### 11.10.2.2.3 General Registers

PAL will use one of two general register calling conventions described in [Section 11.10.2.1, “Overview of Calling Conventions” on page 2:358](#), depending on the availability of memory and the stacked registers at the time of the call. The following tables describe the contents of the general registers.

**Table 11-59. General Registers – Static Calling Convention**

Register	Conventions
GR0	always 0
GR1	preserved
GR2 - GR3	scratch, used with 22 bit immediate add
GR4 - GR7	preserved
GR8 - GR11	scratch, procedure return value
GR12	preserved
GR13	unchanged
GR14 - GR27	scratch
GR28 - GR31	input arguments, scratch (PAL index must be passed in GR28)
Bank 0 Registers (GR16 - GR23)	preserved
Bank 0 Registers (GR 24 - GR31)	scratch
GR32 - GR127	unchanged

**Table 11-60. General Registers – Stacked Calling Conventions**

Register	Conventions
GR0	always 0
GR1	preserved
GR2 - GR3	scratch, used with 22 bit immediate add
GR4 - GR7	preserved

**Table 11-60. General Registers – Stacked Calling Conventions (Continued)**

Register	Conventions
GR8 - GR11	scratch, procedure return value
GR12	special, stack pointer (sp)
GR13	special, thread pointer (tp)
GR14 - GR27	scratch
GR28	input argument, scratch (PAL Index must be passed in GR28)
GR29-GR31	scratch
Bank 0 Registers (GR16 - GR23)	preserved
Bank 0 Registers (GR 24 - GR31)	scratch
GR32 - GR127	stacked registers; in0 - in95: input arguments (PAL index must be in0) loc0 - loc95: local variables out0 - out95: output arguments

The caller must initialize SP for physical and virtual procedure calls only prior to calling a PAL procedure. A minimum 8 KB of room must be available for the stack space of the PAL procedure. The caller to a PAL procedure should set up the RSE backing store before making any procedure calls using the stacked calling conventions. The backing store memory should have a minimum of 8 KB of room for RSE spills.

PAL shall be called with PSR.bn=1. The GR specifications for GR16 through GR31 are for bank one. The bank zero register requirements are specified as a separate line item.

#### 11.10.2.2.4 Floating-point Registers

Floating point registers 32-127 are preserved. PAL must either not use these, or must save and restore them, except on the PAL\_TEST\_PROC procedure, which may overwrite these registers without preserving them. The remainder of the floating-point register conventions are the same as those of the *Itanium Software Conventions and Runtime Architecture Guide*.

#### 11.10.2.2.5 Predicate Registers

The conventions for the predicate registers follow the *Itanium Software Conventions and Runtime Architecture Guide*.

#### 11.10.2.2.6 Branch Registers

The conventions for the branch registers follow the *Itanium Software Conventions and Runtime Architecture Guide*.

#### 11.10.2.2.7 Application Registers

**Table 11-61. Application Register Conventions**

Register	Description	Class
KR0-7	Kernel Registers	unchanged
RSC	Register Stack Configuration Register	unchanged
BSP	Backing Store Pointer (read only)	unchanged <sup>a</sup>

**Table 11-61. Application Register Conventions**

Register	Description	Class
BSPSTORE	Backing Store Pointer for Memory Stores	unchanged <sup>a</sup>
RNAT	RSE NaT Collection Register	unchanged <sup>a</sup>
FCR	IA-32 Floating-point Control Registers	preserved
EFLAG	IA-32 EFLAG register	preserved
CSD	IA-32 Code Segment Descriptor	preserved
SSD	IA-32 Stack Segment Descriptor	preserved
CFLAG	IA-32 Combined CR0 and CR4 Register	preserved
FSR	IA-32 Floating-point Status Register	preserved
FIR	IA-32 Floating-point Instruction Register	preserved
FDR	IA-32 Floating-point Data Register	preserved
CCV	Compare and Exchange Compare Value Register	scratch
UNAT	User NaT Collection Register	according to GR class
FPSR	Floating-point Status Register	preserved
ITC	Interval Time Counter	unchanged <sup>b</sup>
RUC	Resource Utilization Counter	unchanged <sup>c</sup>
PFS	Previous Function State	preserved
LC	Loop Counter Register	preserved
EC	Epilog Counter Register	preserved

- a. BSP, BSPSTORE, and RNAT may not be changed by PAL, but the value at exit may be different due to RSE activity. PAL\_TEST\_PROC is an exception to this rule, and the RSE contents may not be relied on after making this procedure call.
- b. No PAL procedure writes to the ITC. The value at exit is the value at entry plus the elapsed time of the procedure call.
- c. No PAL procedure writes to the RUC. The value at exit is the value at entry plus the number of cycles provided to the processor during the procedure call.

PAL procedures that use the static calling conventions do not use stacked registers or the RSE. Therefore RSE internal state and the CFM are unchanged by these procedures.

### 11.10.2.3 Return Buffers

Any addresses passed to PAL procedures as buffers for return parameters must be 8-byte aligned. Unaligned addresses may cause undefined results.

### 11.10.2.4 Invalid Arguments

The PAL procedure calling conventions specify rules that must be followed. These rules specify certain PSR values, they specify that reserved fields and arguments must be zero filled and specify that values not defined in a range and defined as reserved must not be used.

If the caller of a PAL procedure does not follow these rules, an invalid argument return value may be returned or undefined results may occur during the execution of the procedure. If the caller passes in a PAL procedure index value that is not defined, PAL will return an Unimplemented procedure (-1) status to the caller.



### **11.10.3 PAL Procedure Specifications**

The following pages provide detailed interface specifications for each of the PAL procedures defined in this document. Included in the specification are the input parameters, the output parameters, and any required behavior.

## PAL\_BRAND\_INFO – Provides Processor Branding Information (274)

**Purpose:** Provides processor branding information.

**Calling Conv:** Stacked Registers

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_BRAND_INFO within the list of PAL procedures.
info_request	Unsigned 64-bit integer specifying the information that is being requested. (See <a href="#">Table 11-62</a> )
address	Unsigned 64-bit integer specifying the address of the 128-byte block to which the processor brand string shall be written.
Reserved	0

Return Value	Description
status	Return status of the PAL_BRAND_INFO procedure.
brand_info	Brand information returned. The format of this value is dependent on the input values passed.
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-6	Input argument is not implemented
-9	Call requires PAL memory buffer

**Description:** PAL\_BRAND\_INFO procedure calls are used to ascertain the processor branding information.

The *info\_request* input argument for PAL\_BRAND\_INFO describes which processor branding information is being requested. The *info\_request* values are split into two categories: architected and implementation-specific. The architected *info\_request* have values from 0-15. The implementation-specific *info\_request* have values 16 and above. The architected *info\_request* are described in this document. The implementation-specific *info\_request* are described in processor-specific documentation.

This call returns the processor brand information as requested with the *info\_request* argument. [Table 11-62](#) describes the values.

**Table 11-62. Processor Brand Information Requested**

Value	Description
0	The ASCII brand identification string will be copied to the address specified in the address input argument. The processor brand identification string is defined to be a maximum of 128 characters long; 127 bytes will contain characters and the 128th byte is defined to be NULL (0). A processor may return less than the 127 ASCII characters as long as the string is null terminated. The string length will be placed in the <i>brand_info</i> return argument.
All Other Values	Reserved

This procedure will return an invalid argument if an unsupported *info\_request* argument is passed as an input or a -6 if the requested information was not available on the current processor.

## PAL\_BUS\_GET\_FEATURES – Get Processor Bus Dependent Configuration Features (9)

**Purpose:** Provides information about configurable processor bus features.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_BUS_GET_FEATURES within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_BUS_GET_FEATURES procedure.
	features_avail	64-bit vector of features implemented. See <a href="#">Table 11-63</a> . (1=implemented, 0=not implemented)
	feature_status	64-bit vector of current feature settings. See <a href="#">Table 11-63</a> .
	feature_control	64-bit vector of features controllable by software. (1=controllable, 0= not controllable)

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** [Table 11-63](#) defines the set of possible bus interface features and their bit position in the return vector. Different busses will implement similar features in different ways. For example, data error detection may be implemented by ECC or parity. In other cases, certain features may be tied together. In this case, enabling any one feature in a group will enable all features in the group, and similarly, disabling any one feature in a group will disable all features. Caller algorithms should be written to obtain desired results in these instances. Only those configuration features for which a 1 is returned in *feature\_control* can be changed via PAL\_BUS\_SET\_FEATURES.

For all values in [Table 11-63](#), the *Class* field indicates whether a feature is required to be available (Req.) or is optional (Opt.). The *Control* field indicates which features are required to be controllable. These features will either be controllable through this PAL call or through other hardware means like forcing bus pins to a certain value during processor reset. The *control* field applies only when the feature is available. PALE\_CHECK and PALE\_INIT should not modify these features. An operating system should not modify bus features without detailed information about the platform it is running on.

**Table 11-63. Processor Bus Features**

Bits	Class	Control	Description
63	Opt.	Req.	Disable Bus Data Error Checking. When 0, bus data errors are detected and single bit errors are corrected. When 1, no error detection or correction is done.
62	Opt.	Req.	Disable Bus Address Error Signalling. When 0, bus address errors are signalled on the bus. When 1, no bus errors are signalled on the bus. If Disable Bus Address Error Checking is 1, this bit is ignored.
61	Opt.	Req.	Disable Bus Address Error Checking. When 0, bus errors are detected, single bit errors are corrected., and a CMCI or MCA is generated internally to the processor. When 1, no bus address errors are detected or corrected.
60	Opt.	Req.	Disable Bus Initialization Event Signaling. When 0, bus protocol errors (BINIT#) are signaled by the processor on the bus. When 1, bus protocol errors (BINIT#) are not signaled on the bus. If Disable Bus Initialization Event Checking is 1, this bit is ignored.
59	Opt.	Req.	Disable Bus Initialization Event Checking. When 0, bus protocol errors (BINIT#) are detected and sampled and an MCA is generated internally to the processor. When 1, the processor will ignore bus protocol error conditions (BINIT#).
58	Opt.	Req.	Disable Bus Requester Bus Error Signalling. When 0, BERR# is signalled if a bus error is detected. When 1, bus errors are not signalled on the bus.
57	Opt.	Req.	Disable Bus Requester Internal Error Signalling. When 0, BERR# is signalled when internal processor requestor initiated bus errors are detected. When 1, internal requester bus errors are not signalled on the bus.
56	Opt.	Req.	Disable Bus Error Checking. When 0, the processor takes an MCA if BERR# is asserted. When 1, the processor ignores the BERR# signal.
55	Opt.	Req.	Disable Response Error Checking. When 0, the processor asserts BINIT# if it detects a parity error on the signals which identify the transactions to which this is a response. When 1, the processor ignores parity on these signals.
54	Opt.	Req.	Disable Transaction Queuing. When 0, the in-order transaction queue is limited only by the number of hardware entries. When 1, the processor's in-order transactions queue is limited to one entry.
53	Opt.	Req.	Enable a bus cache line replacement transaction when a cache line in the exclusive state is replaced from the highest level processor cache and is not present in the lower level processor caches. When 0, no bus cache line replacement transaction will be seen on the bus. When 1, bus cache line replacement transactions will be seen on the bus when the above condition is detected.
52	Opt.	Req.	Enable a bus cache line replacement transaction when a cache line in the shared or exclusive state is replaced from the highest level processor cache and is not present in the lower level processor caches. When 0, no bus cache line replacement transaction will be seen on the bus. When 1, bus cache line replacement transactions will be seen on the bus when the above condition is detected.
51:32	N/A	N/A	Reserved
31	Opt.	Opt.	Enable Half transfer rate. When 0, the data bus is configured at the 2x data transfer rate. When 1, the data bus is configured at the 1x data transfer rate,
30	Opt.	Req.	Disable Bus Lock Mask. When 0, the processor executes locked transactions atomically. When 1, the processor masks the bus lock signal and executes locked transactions as a non-atomic series of transactions.
29	Opt.	Req.	Request Bus Parking. When 0, the processor will deassert bus request when finished with each transaction. When 1, the processor will continue to assert bus request after it has finished, if it was the last agent to own the bus and if there are no other pending requests.
28:0	N/A	N/A	Reserved

## PAL\_BUS\_SET\_FEATURES – Set Processor Bus Dependent Configuration Features (10)

**Purpose:** Enables/disables specific processor bus features.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_BUS_SET_FEATURES within the list of PAL procedures.
feature_select	64-bit vector denoting desired state of each feature (1=select, 0=non-select).
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_BUS_SET_FEATURES procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Can not complete call without error

**Description:** PAL\_BUS\_GET\_FEATURES should be called to ascertain the implemented processor bus configuration features, their current setting, and whether they are software controllable, before calling PAL\_BUS\_SET\_FEATURES. The list of possible processor features is defined in [Table 11-63](#). Attempting to enable or disable any feature that cannot be changed will be ignored.

## PAL\_CACHE\_FLUSH – Flush Data or Instruction Caches (1)

**Purpose:** Flushes the processor instruction or data caches.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_CACHE_FLUSH within the list of PAL procedures.
	cache_type	Unsigned 64-bit integer indicating which cache to flush. See <a href="#">Table 11-64</a> .
	operation	Formatted bit vector indicating the operation of this call. See <a href="#">Figure 11-1</a> .
	progress_indicator	Unsigned 64-bit integer specifying the starting position of the flush operation.

Returns:	Return Value	Description
	status	Return status of the PAL_CACHE_FLUSH procedure.
	vector	Unsigned 64-bit integer specifying the vector number of the pending interrupt.
	progress_indicator	Unsigned 64-bit integer specifying the starting position of the flush operation.
	Reserved	0

Status:	Status Value	Description
	2	Call completed without error, but a PMI was taken during the execution of this procedure.
	1	Call has not completed flushing due to a pending interrupt
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** Flushes the instruction or data caches controlled by the processor as specified by the *cache\_type* parameter. Encoding for the *cache\_type* parameter follows:

**Table 11-64. *cache\_type* Encoding**

Value	Description
1	Flush all caches containing instructions.
2	Flush all caches containing data.
3	Flush all caches (instruction and data).
4	Make local instruction caches coherent with the data caches.

All other values of *cache\_type* are reserved. If the cache is unified, both instruction and data lines are flushed, regardless of the value of *cache\_type*.

Flushing all caches containing instructions, causes the instruction and unified caches to be flushed. Flushing all caches containing data, causes all data and unified caches to be flushed. Flushing all caches causes all data, instruction, and unified caches to be flushed.

When the caller specifies to make local instruction caches coherent with the data caches, this procedure will ensure that the instruction caches on the processor that this procedure call was made, will see the effects of stores to instruction code performed by this processor. This procedure is not required to ensure coherency of instruction caches on other processors in the system when this input argument is used. Refer to [Section 4.4.3, “Cacheability and Coherency Attribute” on page 2:77](#) for more information on stores and their coherency requirements with local instruction caches.

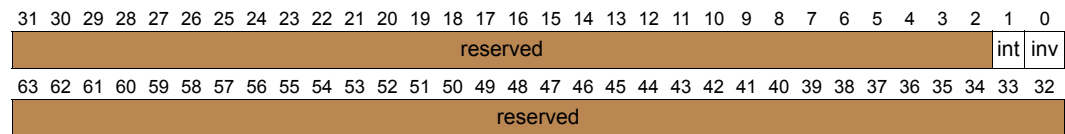
The effects of flushing data and unified caches is broadcast throughout the coherence domain. The effects of flushing instruction caches may or may not be broadcast

throughout the coherence domain. The procedure will perform the necessary serialization and synchronization as required by the architecture.

This call does not ensure that data in the processors coalescing buffers are flushed to memory. See Section 4.4.5, “Coalescing Attribute” on page 2:78 on how to flush the coalescing buffers.

The *operation* parameter controls how this call will operate. The *operation* parameter has the following format:

**Figure 11-1. operation Parameter Layout**



- inv* – 1 bit field indicating whether to invalidate clean lines in the cache.

If this bit is 0, all modified cache lines for the specified *cache\_type* are copied back to memory. Optimally, modified and non-modified cache lines are left valid in the specified cache in a clean (non-modified) state. However, based on the processor implementation, cache lines in the specified cache may alternatively be invalidated. If this bit is 1, all modified cache lines for the specified *cache\_type* are flushed by copying the cache line to memory. All cache lines in the specified cache are then invalidated.

If *cache\_type* is equal to 4 (make local instruction caches coherent with the data caches) the *inv* bit will be ignored.

Table 11-65 will clarify the effects of the *inv* bit. The modified state represents a cache line that contains modified data. The clean state represents a cache line that contains no modified data.
- int* – 1 bit field indicating if the processor will periodically poll for external interrupts while flushing the specified *cache\_type(s)*.

If this bit is a 0, unmasked external interrupts will not be polled. The processor will ignore all pending unmasked external interrupts until all cache lines in the specified *cache\_type(s)* are flushed. Depending on the size of the processor’s caches, bus bandwidth and implementation characteristics, flushing the caches can take a long period of time, possibly delaying interrupt response times and potentially causing I/O devices to fail.

If this bit is a 1, external interrupts will be polled periodically and will exit the procedure if one is seen. If an unmasked external interrupt becomes pending, this procedure will return and allow the caller to service the interrupt before all cache lines in the specified *cache\_type(s)* are flushed.

**Table 11-65. Cache Line State when *inv* = 0**

Old State	New State	Comments
Invalid	Invalid	
Clean	Clean <sup>a</sup>	
Modified	Clean <sup>a</sup>	Modified data is copied back to memory

a. Based on the processor implementation the cache line can be invalidated or left in a model-specific clean state

**Table 11-66. Cache Line State when *inv* = 1**

Old State	New State	Comments
Invalid	Invalid	Modified data is copied back to memory.
Clean	Invalid	
Modified	Invalid	

The *progress\_indicator* is an unsigned 64-bit integer specifying the starting position of the flush operation. Values in this parameter are model specific and will vary across processor implementations.

The first time this procedure is called, the *progress\_indicator* must be set to zero. If this procedure exits due to an external interrupt and this procedure is then again called to resume flushing, the *progress\_indicator* must be set to the value previously returned by PAL\_CACHE\_FLUSH. Software must program no value other than zero or the value previously returned by PAL\_CACHE\_FLUSH otherwise behavior is undefined.

This procedure makes one flush pass through all caches specified by *cache\_type* and all sets and associativities within those caches. The specified *cache\_type(s)* are ensured to be flushed only of cache lines resident in the caches prior to PAL\_CACHE\_FLUSH initially being called with the *progress\_indicator* set to 0.

This procedure ensures that prefetches initiated prior to making this call with *progress\_indicator* set to 0 are flushed based on the *cache\_type* argument passed.

- If *cache\_type* specifies to flush all instruction caches then the call ensures all prior instruction prefetches are flushed.
- If *cache\_type* specifies to flush all data caches then the call ensures all prior data prefetches are flushed.
- If *cache\_type* specifies to flush all caches then the call ensures all prior instruction and data prefetches are flushed from the caches.
- If *cache\_type* specifies to make local instruction caches coherent with the data caches, then the call will ensure all prior instruction prefetches are flushed.

Due to the following conditions, software cannot assume that when this procedure completes the entire flush pass that the specified *cache\_type(s)* are empty of all clean and/or modified cache lines.

- After an interruption, the flush pass resumes at the interruption point (specified by *progress\_indicator*). Due to execution of the interrupt handlers during the flush pass, the specified caches may contain new and possibly modified cache lines in sections of the caches already flushed. The caller specifies if this procedure should poll for interrupts via the *int* bit of the *operation* parameter.
- Prior prefetches initiated before this procedure is called are disabled and flushed from the cache as described above. However, if a speculative translation exists in either the ITLB or DTLB, speculative instruction or data prefetch operation could immediately reload a non-modified cache line after it was flushed. To ensure prefetches do not occur, software must remove all speculative translation before



calling this routine. Alternatively, software can disable the TLBs by setting PSR.it, PSR.dt, and PSR.rt to 0.

- The specified caches may also contain PAL firmware code cache entries required to flush the cache.
- The specified caches may contain PAL and SAL PMI code if this call was made with PSR.ic = 1 and a PMI interrupt is seen during the execution of the call.
- The specified caches may contain SAL or OS machine check or INIT code if these handlers run in a cacheable mode and a machine check or INIT event is seen.
- In a processor that contains multiple logical processors, the specified caches may contain new and possibly modified cache lines in sections of the cache already flushed due to execution of instructions on other logical processors that share the specified caches. Information about how caches are shared among logical processors is described in the PAL\_CACHE\_SHARED\_INFO procedure on [page 2:382](#). Information about logical processors on the same physical processor package are described in the PAL\_LOGICAL\_TO\_PHYSICAL procedure on [page 2:404](#).

This procedure does ensure that all cache lines resident in the specified *cache\_type(s)* prior to this procedure being initially called are flushed regardless of intervening external interrupts. It also ensures that prefetches initiated prior to the initial call to this procedure that affect the caches specified in *cache\_type*, as described above, are flushed regardless of intervening external interrupts.

To ensure forward progress, PAL\_CACHE\_FLUSH advances through the cache flush sequence at least by one cache line before sampling for pending external interrupts. The amount of flushing that occurs before interrupts are polled will vary across implementations.

PAL\_CACHE\_FLUSH will return the following values to indicate to the caller the status of the call.

- *status* – When the call returns a 1, it indicates that the call did not have any errors but is returning due to a pending unmasked external interrupt. To continue flushing the caches, the caller must call PAL\_CACHE\_FLUSH again with the value returned in the *progress\_indicator* return value.

When the call returns a 0, it indicates that the call completed without any errors. All cache lines that were present in the cache (when the most recent call to PAL\_CACHE\_FLUSH with a *progress\_indicator* of zero) are flushed and possibly invalidated. All intermediate calls must have used the proper *progress\_indicator*, otherwise behavior is undefined.

When the call returns a 2, it indicates that the call completed without any errors but that a PMI was taken during the execution of this call. This indicates to the caller that all cache lines that were present in the cache (when the most recent call to PAL\_CACHE\_FLUSH with a *progress\_indicator* of zero) are flushed but that code and data related to handling PMIs may be present in the cache.

- *vector* – If the return status is 1 and this procedure exited due to a pending unmasked external interrupt, this field returns the interrupt vector number. The external interrupt will have been removed. The interrupt is considered to be “in-service” and software must service this interrupt for the specified vector and then issue EOI. If the return status is not 1, the values returned is undefined.
- *progress\_indicator* – When the return status is 1, specifies the current position in the flush pass. The value returned is model specific and will vary across processor implementations. If the return status is not 1, the value returned is undefined.

## PAL\_CACHE\_INFO – Get Detailed Cache Information (2)

**Purpose:** Returns information about a particular processor instruction or data cache at a specified level in the cache hierarchy.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_CACHE_INFO within the list of PAL procedures.
	cache_level	Unsigned 64-bit integer specifying the level in the cache hierarchy for which information is requested. This value must be between 0 and one less than the value returned in the <i>cache_levels</i> return value from PAL_CACHE_SUMMARY.
	cache_type	Unsigned 64-bit integer with a value of 1 for instruction cache and 2 for data or unified cache. All other values are reserved.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_CACHE_INFO procedure.
	config_info_1	The format of <i>config_info_1</i> is shown in Figure 11-2.
	config_info_2	The format of <i>config_info_2</i> is shown in Figure 11-3.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This call describes in detail the characteristics of a given processor controlled cache in the cache hierarchy. It returns information in the *config\_info\_1* and *config\_info\_2* returns parameters.

The *config\_info\_1* return value has the following structure:

**Figure 11-2. config\_info\_1 Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
stride								line_size								associativity								reserved		at		u			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
load_hints								store_hints								load_latency								store_latency							

- *u* – Bit that is 1 if the cache is unified and 0 if the cache is split.
- *at* - 2-bit field denoting cache memory attributes, as follows:

**Table 11-67. Cache Memory Attributes**

Value	Description
0	Write through cache
1	Write back cache
2-3	Reserved

- *associativity* – Unsigned 8-bit integer denoting the associativity of the cache. A value of 0 indicates a fully associative cache. A value of 1 indicates a direct mapped cache.
- *line\_size* – Unsigned 8-bit integer denoting the binary logarithm (log2) of the minimum write back size of a flush operation to memory or the line size of the

cache if the cache contents never get flushed to memory (for example an instruction cache).

- *stride* – Unsigned 8-bit integer denoting the binary log of the most effective stride in bytes for flushing the cache.
- *store\_latency* – Unsigned 8-bit integer denoting the number of cycles after issue until the value is written into the cache. If the cache cannot accept a store (like an instruction cache) the value returned will be 256 (0xff).
- *load\_latency* – Unsigned 8-bit integer denoting the number of processor cycles after issue until the value may be used if it is found in the cache.
- *store\_hints* – 8-bit vector denoting hints implemented by the processor store instruction. For instruction caches this bit vector will be zero indicating no store hints are supported.

**Table 11-68. Cache Store Hints**

Bits	Description
0	Temporal, level 1
2:1	Reserved
3	Non-temporal, all levels
7:4	Reserved

- *load\_hints* – 8-bit vector denoting hints implemented by the processor load instruction.

**Table 11-69. Cache Load Hints**

Bits	Hint
0	Temporal, level 1
1	Non-temporal, level 1
2	Reserved
3	Non-temporal, all levels
7:4	Reserved

The *config\_info\_2* return value has the following structure:

**Figure 11-3. *config\_info\_2* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
cache_size																															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
reserved								tag_ms_bit								tag_ls_bit								alias_boundary							

- *cache\_size* – Unsigned 32-bit integer denoting the size of the cache in bytes.
- *alias\_boundary* – Unsigned 8-bit integer indicating the binary log of the minimum number of bytes which must separate aliased addresses in order to obtain the highest performance.
- *tag\_ls\_bit* – Unsigned 8-bit integer denoting the least-significant address bit of the tag.
- *tag\_ms\_bit* – Unsigned 8-bit integer denoting the most-significant address bit of the tag.

## PAL\_CACHE\_INIT – Initialize Caches (3)

**Purpose:** Initializes the processor controlled caches.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_CACHE_INIT within the list of PAL procedures.
level	Unsigned 64-bit integer containing the level of cache to initialize. If the cache level can be initialized independently, only that level will be initialized. Otherwise implementation-dependent side-effects will occur.
cache_type	Unsigned 64-bit integer with a value of 1 to initialize the instruction cache, 2 to initialize the data cache, or 3 to initialize both. All other values are reserved.
restrict	Unsigned 64-bit integer with a value of 0 or 1. All other values are reserved. If <i>restrict</i> is 1 and initializing the specified level and <i>cache_type</i> of the cache would cause side-effects, PAL_CACHE_INIT will return -4 instead of initializing the cache.

Return Value	Description
status	Return status of the PAL_CACHE_INIT procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error
-4	Call could not initialize the specified level and <i>cache_type</i> of the cache without side-effects and <i>restrict</i> was 1.

**Description:** Initializes one or all the processor's caches. The effect of this procedure is to initialize the caches without causing writebacks. This procedure cannot be used where coherency is required because any data in the caches will be lost.

The *level* argument must either be -1, indicating all cache levels, or a non-negative number indicating the specific level to initialize. In the latter case, *level* must be in the range from 0 up to one less than the *cache\_levels* return value from PAL\_CACHE\_SUMMARY:

**Table 11-70. PAL\_CACHE\_INIT level Argument Values**

Value	Description
-1	Initializes all cache levels ( <i>cache_type</i> and <i>restrict</i> are ignored)
0 to N	Initialize only the specified cache level.

The *restrict* argument specifies how to handle potential side-effects, where:

**Table 11-71. PAL\_CACHE\_INIT restrict Argument Values**

Value	Description
0	No restriction: initialize the specified level and <i>cache_type</i> of the cache, even if doing so will cause side effects in other caches.
1	Restrict initialization to the specified level and <i>cache_type</i> without side effects to other cache levels. If this cannot be done, return -4.

All other values of *restrict* are reserved.

## PAL\_CACHE\_LINE\_INIT – Initialize a Data Cache Line (31)

**Purpose:** Initializes the tags and data of a data or unified cache line of a processor controlled cache to known values without the availability of backing memory.

**Calling Conv:** Static

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_CACHE_LINE_INIT within the list of PAL procedures.
	address	Unsigned 64-bit integer value denoting the physical address from which the physical page number is to be generated. The address must be an implemented physical address, bit 63 must be zero.
	data_value	64-bit data value which is used to initialize the cache line.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_CACHE_LINE_INIT procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Can not complete call without error

**Description:** A line in the data or unified cache is initialized to the values passed in the arguments of this procedure. The physical page number of the line is derived from the *address* value passed. The tags of the line are set to Private, Dirty, and Valid. The cache line is initialized using *data\_value* repeated until it fills the line. This procedure replicates *data\_value* to a size equal to the largest line size in the processor-controlled cache hierarchy.

This procedure call cannot be used where coherency is required.

## PAL\_CACHE\_PROT\_INFO – Get Detailed Cache Protection Information (38)

**Purpose:** Returns protection information about a particular processor instruction or data cache at a specified level in the cache hierarchy.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_CACHE_PROT_INFO within the list of PAL procedures.
	cache_level	Unsigned 64-bit integer specifying the level in the cache hierarchy for which information is requested. This value must be between 0 and one less than the value returned in the <i>cache_levels</i> return value from PAL_CACHE_SUMMARY.
	cache_type	Unsigned 64-bit integer with a value of 1 for instruction cache and 2 for data or unified cache. All other values are reserved.
	Reserved	0

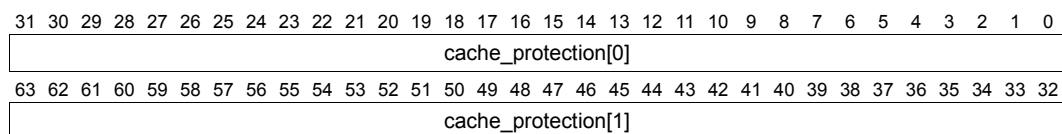
Returns:	Return Value	Description
	status	Return status of the PAL_CACHE_PROT_INFO procedure.
	config_info_1	The format of <i>config_info_1</i> is shown in <a href="#">Figure 11-4</a> .
	config_info_2	The format of <i>config_info_2</i> is shown in <a href="#">Figure 11-5</a> .
	config_info_3	The format of <i>config_info_3</i> is shown in <a href="#">Figure 11-6</a> .

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** PAL\_CACHE\_PROT\_INFO returns information about the data and tag protection method for the specified cache. The three returns compose a six-element array of 32-bit protection information structures.

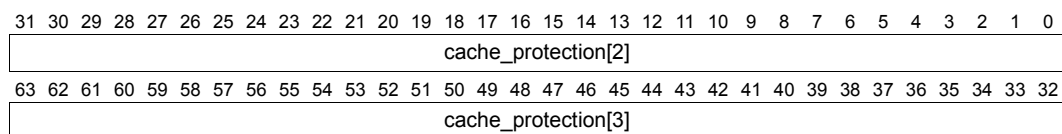
The *config\_info\_1* return value has the following structure:

**Figure 11-4. *config\_info\_1* Return Value**



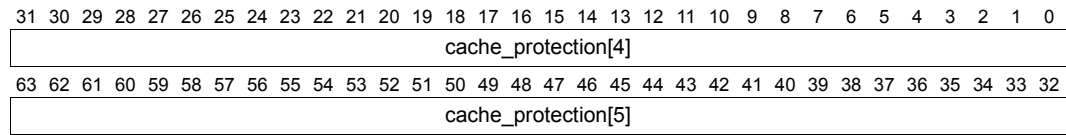
The *config\_info\_2* return value has the following structure:

**Figure 11-5. *config\_info\_2* Return Value**



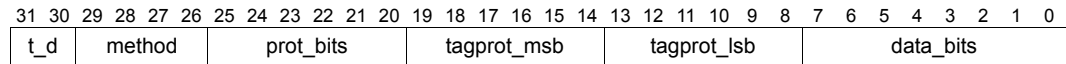
The *config\_info\_3* return value has the following structure:

**Figure 11-6. config\_info\_3 Return Value**



Each *cache\_protection* element has the following structure:

**Figure 11-7. cache\_protection Fields**



- *data\_bits* – Unsigned 8-bit integer denoting the number of data bits that each unit of protection covers. For example, if the cache hardware generates 8 bits of ECC per 64 bits of data, *data\_bits* would be 64. This field is only valid if *t\_d* is 0, 2, or 3.
- *tagprot\_lsb* – Unsigned 6-bit integer denoting the least-significant tag address bit that this protection method covers. This field is only valid if *t\_d* is 1, 2, or 3.
- *tagprot\_msb* – Unsigned 6-bit integer denoting the most-significant tag address bit that this protection method covers. This field is only valid if *t\_d* is 1, 2, or 3.
- *prot\_bits* – Unsigned 6-bit integer denoting the number of protection bits generated for the field specified by the *t\_d* element.
- *method* – Unsigned 4-bit integer denoting the protection method, where:

**Table 11-72. method Values**

Value	Description
0	no ECC or parity protection
1	odd parity protection
2	even parity protection
3	ECC protection

All other values of *method* are reserved.

- *t\_d* – 2-bit field denoting whether this protection method applies to the tag, the data, or both. If over both, the tag and data unit could be concatenated with the tag either to the left (more significant) or to the right (less significant) than a unit of data. For the values of 2 and 3, the difference of *tagprot\_msb* and *tagprot\_lsb* indicates the number of tag bits that are protected with the data bits. The data bits are described by the *data\_bits* field below. This field is encoded as follows:

**Table 11-73. t\_d Values**

Value	Description
0	Data protection
1	Tag protection
2	Tag+data protection (tag is more significant)
3	Data+tag protection (data is more significant)

When obtaining cache information via this call, information for the data cache should be obtained first, then if the *u* bit of the *config\_info\_1* parameter is not set, obtain the information for the instruction cache.

## PAL\_CACHE\_READ – Read Values from the Processor Cache (259)

**Purpose:** Reads the data and tag of a processor-controlled cache line for diagnostic testing.

**Calling Conv:** Stacked Registers

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_CACHE_READ within the list of PAL procedures.
line_id	8-byte formatted value describing where in the cache to read the data.
address	64-bit 8-byte aligned physical address from which to read the data. The address must be an implemented physical address on the processor model with bit 63 set to zero.
Reserved	0

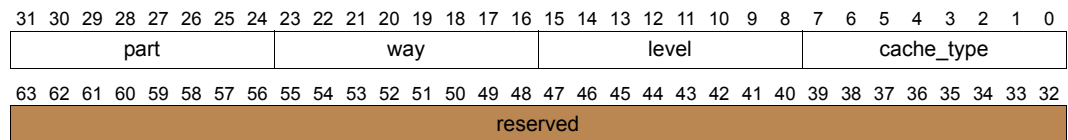
Return Value	Description
status	Return status of the PAL_CACHE_READ procedure.
data	Right-justified value returned from the cache line.
length	The number of bits returned in <i>data</i> .
mesi	The status of the cache line.

Status Value	Description
1	The word at <i>address</i> was found in the cache, but the line was invalid.
0	Call completed without error.
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error.
-5	The word at <i>address</i> was not found in the cache.
-7	The operation requested is not supported for this <i>cache_type</i> and <i>level</i> .

**Description:** A value is read from the specified cache line, if present. This procedure allows reading cache data, tag, protection, or status bits.

The *line\_id* argument is an 8-byte quantity in the following format:

**Figure 11-8. Layout of *line\_id* Return Value**



- *cache\_type* – Unsigned 8-bit integer denoting whether to read from instruction (1) or data/unified (2) cache. All other values are reserved.
- *level* – Unsigned 8-bit integer specifying which cache within the cache hierarchy to read. This value must be in the range from 0 up to one less than the *cache\_levels* return value from PAL\_CACHE\_SUMMARY.
- *way* – Unsigned 8-bit integer denoting within which cache way to read. If the cache is direct-mapped this argument is ignored.
- *part* – Unsigned 8-bit integer denoting which portion of the specified cache line to read:



**Table 11-74. *part* Input Values**

Value	Description
0	data
1	tag
2	data protection bits
3	tag protection bits
4	combined protection bits for data and tags <sup>a</sup>

a. Note that for this *part* no data is returned. Only protection bits are returned.

All other values of *part* are reserved.

The *data* return value contains the value read from the cache. Its contents are interpreted according to the *line\_id.part* field as follows:

**Table 11-75. *part* Input Values and corresponding *data* Return Values**

Part	Data
0	64-bit data.
1	right-justified tag of the specified line.
2	right-justified protection bits corresponding to the 64 bits of data at <i>address</i> . If the cache uses less than 64-bits of data to generate protection, <i>data</i> will contain more than one value. For example if a cache generates parity for every 8-bits of data, this return value would contain 8 parity values. The PAL_CACHE_PROT_INFO call returns information on how a cache generates protection information in order to decode this return value. If a cache uses greater than 64-bits of data to generate protection, <i>data</i> will contain the value to use for the portion of the cache line indicated by <i>address</i> .
3	right-justified protection bits for the cache line tag.
4	right-justified protection bits for the cache line tag and 64 bits of data at <i>address</i> .

The *length* return value contains the number of valid bits returned in *data*.

The *mesi* return value contains the status bits of the cache line. Values are defined as follows:

**Table 11-76. *mesi* Return Values**

Value	Description
0	invalid
1	shared
2	exclusive
3	modified

All other values of *mesi* are reserved.

To guarantee correct behavior for this procedure, it is required that there shall be no RSE activity that may cause cache side effects.

## PAL\_CACHE\_SHARED\_INFO – Get Information on Caches Shared by Logical Processors (43)

**Purpose:** Returns information on caches shared between logical processors.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

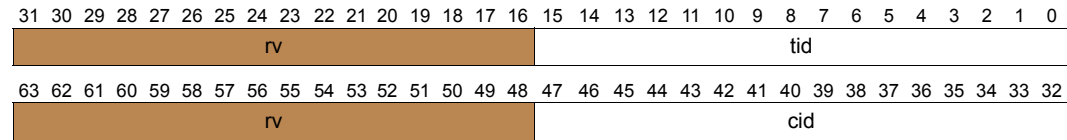
Argument	Description
index	Index of PAL_CACHE_SHARED_INFO within the list of PAL procedures.
cache_level	Unsigned 64-bit integer specifying the level in the cache hierarchy for which information is requested. This value must be between 0 and one less than the value returned in the cache_levels return value from PAL_CACHE_SUMMARY.
cache_type	Unsigned 64-bit integer with a value of 1 for instruction cache and 2 for data or unified cache. All other values are reserved.
proc_number	Unsigned 64-bit integer that specifies for which logical processor information is being requested. This input argument must be zero for the first call to this procedure and can be a maximum value of one less than the number of logical processors sharing this cache, which is returned by the <i>num_shared</i> return value.

Return Value	Description
status	Return status of the PAL_CACHE_SHARED_INFO procedure.
num_shared	Unsigned integer that returns the number of logical processors that share the processor cache level and type, for which information was requested.
proc_n_cache_info1	The format of <i>proc_n_cache_info1</i> is shown in <a href="#">Figure 11-9</a> .
proc_n_cache_info2	The format of <i>proc_n_cache_info2</i> is shown in <a href="#">Figure 11-10</a> .

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

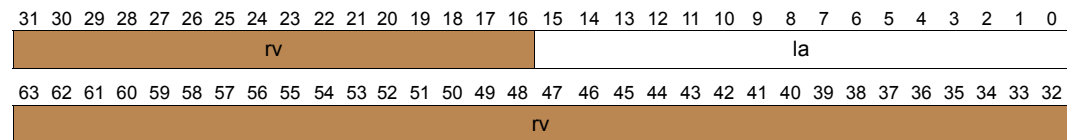
**Description:** This procedure will return information about how the processor caches are shared among logical processors (See [“PAL\\_LOGICAL\\_TO\\_PHYSICAL – Get Information on Logical to Physical Processor Mappings \(42\)”](#) on page 2:404 for a definition of a logical processor). If the caller is only interested in how many logical processors are sharing a particular cache level, this procedure will only need to be called once. If the caller is interested in identifying which logical processors are sharing the processor caches, this procedure will need to be called a number of times equal to the value returned in *num\_shared* to gather identification information for all the logical processors sharing the particular cache for which information is being requested.

Identification information about the logical processors sharing the cache is in the return values *proc\_n\_cache\_info1* and *proc\_n\_cache\_info2*. The format of these return values is shown in [Figure 11-9](#) and [Figure 11-10](#).

**Figure 11-9. Layout of *proc\_n\_cache\_info1* Return Value**

- *tid* – Thread id: The thread identifier of the logical processor for which information is being returned. This value will be unique on a per core basis.
- *rv* – Reserved
- *cid* – Core id: The core identifier of the logical processor for which information is being returned. This value will be unique on a per physical processor package basis.
- *rv* – Reserved

There is no guarantee that the core id's and thread id's will be contiguous on a given physical processor package.

**Figure 11-10. Layout of *proc\_n\_cache\_info2* Return Value**

- *la* – Logical address: geographical address of the logical processor for which information is being returned. This is the same value that is returned by the PAL\_FIXED\_ADDR procedure when it is called on the logical processor.
- *rv* – Reserved

This procedure must be supported on all implementations that contain more than one logical processor on a physical processor package and returns an unimplemented procedure error code otherwise.

## PAL\_CACHE\_SUMMARY – Get Cache Hierarchy Summary (4)

**Purpose:** Returns summary information about the hierarchy of caches controlled by the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_CACHE_SUMMARY within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_CACHE_SUMMARY procedure.
	cache_levels	Unsigned 64-bit integer denoting the number of levels of cache implemented by the processor. Strictly, this is the number of levels for which the cache controller is integrated into the processor (the cache SRAMs may be external to the processor).
	unique_caches	Unsigned 64-bit integer denoting the number of unique caches implemented by the processor. This has a maximum of $2 \times \text{cache\_levels}$ , but may be less if any of the levels in the cache hierarchy are unified caches or do not have both instruction and data caches.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** Software is expected to call PAL\_CACHE\_SUMMARY before calling PAL\_CACHE\_INFO to determine the number of times PAL\_CACHE\_INFO should be called and the amount of storage that must be allocated to hold all of the information returned by PAL\_CACHE\_INFO.

## PAL\_CACHE\_WRITE – Write Values into the Processor Cache (260)

**Purpose:** Writes the data and tag of a processor-controlled cache line for diagnostic testing.

**Calling Conv:** Stacked Registers

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_CACHE_WRITE within the list of PAL procedures.
line_id	8-byte formatted value describing where in the cache to write the data.
address	64-bit 8-byte aligned physical address at which the data should be written. The address must be an implemented physical address on the processor model with bit 63 set to 0.
data	unsigned 64-bit integer value to write into the specified <i>part</i> of the cache.

Return Value	Description
status	Return status of the PAL_CACHE_WRITE procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error.
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error.
-7	The operation requested is not supported for this <i>cache_type</i> and <i>level</i> .

**Description:** The value of *data* is written into the specified level, way, and part of the cache. This procedure allows writing cache data, tag, protection, or status bits.

This procedure may also be used to seed errors into a cache line. It calculates the protection bits based on the value of *data*, then inverts a specified bit field before writing *data* to the cache. Bit field inversion is only used for writes to the cache data or tag.

If seeding an error into the instruction cache or seeding an unrecoverable error, then return back to the caller may not be possible.

This procedure call cannot be used where coherency is required.

The *line\_id* argument is an 8-byte quantity in the following format:

**Figure 11-11. Layout of *line\_id* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																																
part								way								level								cache_type																																							
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																																
trigger																length																start																mesi															

- *cache\_type* – Unsigned 8-bit integer denoting whether to write to instruction (1) or data/unified (2) cache. All other values are reserved.
- *level* – Unsigned 8-bit integer specifying which cache within the cache hierarchy to write *data*. This value must be in the range from 0 up to one less than the *cache\_levels* return value from PAL\_CACHE\_SUMMARY.
- *way* – Unsigned 8-bit integer denoting within which cache way to write *data*. If the cache is direct-mapped this argument is ignored.
- *part* – Unsigned 8-bit integer denoting where to write *data* into the cache:

**Table 11-77. *part* Input Values**

Value	Description
0	data
1	tag
2	data protection
3	tag protection
4	combined data and tag protection

All other values of *part* are reserved.

- *mesi* – Unsigned 8-bit integer denoting whether the line should be written as clean or dirty, shared or exclusive. Though there may be multiple calls to PAL\_CACHE\_WRITE to the same cache line, the last call's *mesi* will be in effect. Values are defined as follows:

**Table 11-78. *mesi* Return Values**

Value	Description
0	invalid
1	shared
2	exclusive
3	modified

All other values of *mesi* are reserved.

- *start* – Unsigned 8-bit integer denoting the least-significant bit of the field in *data* to invert. If *length* is 0 or *part* is not 0 or 1, this field is ignored.
- *length* – Unsigned 8-bit integer denoting the number of bits to invert. If *length* is 0, no bits are inverted and *start* is ignored. If *part* is not 0 or 1, this field is ignored.
- *trigger* – Unsigned 8-bit integer denoting whether to trigger the error while in procedure. If *trigger* is 0, the procedure writes *data* and returns. If *trigger* is 1 and *cache\_type* is data/unified, the procedure writes *data* and executes a 64-bit load from *address* before returning. If *trigger* is 1 and *cache\_type* is set to instruction, the procedure writes *data* and branches to the *address*. All other values are reserved.

The *data* argument contains the value to write into the cache. Its contents are interpreted based on the *part* field as follows:

**Table 11-79. Interpretation of *data* Input Field**

Part	Data
0	64-bit data to write to the specified line (with optional bit field inversion).
1	right-justified tag to write into the specified line (with optional bit field inversion).
2	right-justified protection bits corresponding to the 64 bits of data at <i>address</i> . If the cache uses less than 64-bits of data to generate protection, <i>data</i> will contain more than one value. For example if a cache generates parity for every 8-bits of data, this return value would contain 8 parity values. The PAL_CACHE_PROT_INFO call returns information on how a cache generates protection information in order to decode this return value. If a cache uses greater than 64-bits of data to generate protection, <i>data</i> will contain the value to use for the portion of the cache line indicated by <i>address</i> .
3	right-justified protection bits for the cache line tag.
4	right-justified protection bits for the cache line tag and 64 bits of data at <i>address</i> .

To guarantee correct behavior for this procedure, it is required that there shall be no RSE activity that may cause cache side effects.

## PAL\_COPY\_INFO – Return Parameters to Copy PAL Code to Memory (30)

**Purpose:** Returns the parameters needed to copy relocatable PAL code from the firmware address space to memory.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_COPY_INFO within the list of PAL procedures.
copy_type	Unsigned integer denoting type of procedures for which copy information is requested.
Reserved	0
mca_proc_state_info	Unsigned integer denoting the number of bytes that SAL needs for the min-state save area for each processor.

Return Value	Description
status	Return status of the PAL_COPY_INFO procedure.
buffer_size	Unsigned integer denoting the number of bytes of PAL information that must be copied to main memory.
buffer_align	Unsigned integer denoting the starting alignment of the data to be copied.
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** This procedure is called to obtain the information needed to relocate runtime PAL procedures and PAL PMI code from the firmware address space to memory. The information returned in this call is used by SAL to allocate a memory region on the required alignment, and call PAL\_COPY\_PAL to copy the relocatable PAL code.

The *copy\_type* input argument indicates which type of procedure for which copying information is requested. A value of 0 denotes procedures required for SAL, PMI, and Itanium architecture-based operating systems. All other values are reserved. If the *copy\_type* is 0, then SAL shall call PAL\_COPY\_PAL call subsequently to copy the PAL procedures and PAL PMI code to the allocated memory region.

The *buffer\_align* return value must be a power of two between 4 KB and 1 MB.



## PAL\_COPY\_PAL – Copy PAL Code to Memory (256)

**Purpose:** Copy relocatable PAL code from the firmware address space to memory.

**Calling Conv:** Stacked Registers

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_COPY_PAL within the list of PAL procedures.
	target_addr	Physical address of a memory buffer to copy relocatable PAL procedures and PAL PMI code.
	alloc_size	Unsigned integer denoting the size of the buffer passed by SAL for the copy operation.
	copy_option	Unsigned integer indicating whether relocatable PAL code and PAL PMI code should be copied from firmware address space to main memory.

Returns:	Return Value	Description
	status	Return status of the PAL_COPY_PAL procedure.
	proc_offset	Unsigned integer denoting the offset of PAL_PROC in the relocatable segment copied.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This procedure is called to relocate runtime PAL procedures and PAL PMI code from the firmware address space to main memory. A value of 0 for the *copy\_option* indicates that the relocation should be performed; a value of 1 indicates that the relocation should not be performed. This procedure also updates the PALE\_PMI entrypoint in hardware. All other values are reserved.

PAL\_COPY\_INFO should be called first to determine the size and alignment requirements of the memory buffer to which the PAL code will be copied. Bit 63 of *target\_addr* must be set consistently with the cacheability attribute of the memory buffer being copied to. It is PAL's responsibility to ensure that the firmware address space contents that are being copied from, are not in any processor caches. It is the caller's responsibility to ensure that the contents of the memory buffer copied to, are flushed out of the internal processor's data caches if *target\_addr* has a cacheable memory attribute.

If a PAL procedure makes calls to internal PAL functions that execute only out of the firmware address space, that portion of code will continue to execute out of the firmware address space, even though the main procedure has been copied to RAM. This is true only for some PAL procedures that can be called only in physical mode.

PAL\_COPY\_PAL call is mandatory as part of the system boot process. Higher level firmware should guarantee that PAL\_COPY\_PAL is called on all processors before OS launch. This is to guarantee that full processor functionality is available. This procedure can be called more than once.

## PAL\_DEBUG\_INFO – Get Debug Registers Information (11)

**Purpose:** Returns the number of instruction and data debug register pairs.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_DEBUG_INFO within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_DEBUG_INFO procedure.
	i_regs	Unsigned 64-bit integer denoting the number of pairs of instruction debug registers implemented by the processor.
	d_regs	Unsigned 64-bit integer denoting the number of pairs of data debug registers implemented by the processor.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This call returns the number of pairs of registers. Even numbered registers contain breakpoint addresses and odd numbered registers contain breakpoint mask conditions. For example if *i\_regs* is 4, there are 8 instruction debug registers of which 4 are breakpoint address registers (IBR<sub>0,2,4,6</sub>) and 4 are breakpoint mask registers (IBR<sub>1,3,5,7</sub>). The minimum value for both *i\_regs* and *d\_regs* is 4.

On some implementations, a hardware debugger may use two or more debug register pairs for its own use. When a hardware debugger is attached, PAL\_DEBUG\_INFO may return a value for *i\_regs* and/or *d\_regs* less than the implemented number of debug registers. When a hardware debugger is attached, PAL\_DEBUG\_INFO may return a minimum value of 2 for *d\_regs* and a minimum of 2 for *i\_regs*.

## PAL\_FIXED\_ADDR – Get Fixed Geographical Address of Processor (12)

**Purpose:** Returns a unique geographical address of this processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_FIXED_ADDR call within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_FIXED_ADDR procedure.
	address	Fixed geographical address of this processor.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** The *address* return value will contain a unique unsigned integer denoting the position of this processor on its system interconnect. This is an arbitrary number which is expected to have geographical significance and is unique for the system interconnect to which the processor is connected. If the processor is connected to multiple system interconnects, the *address* return value must be unique among all such interconnects. The maximum size of the *address* returned corresponds to the size of the fields (id and eid) in the LID register (CR64).

## PAL\_FREQ\_BASE – Get Processor Base Frequency (13)

**Purpose:** Returns the frequency of the output clock for use by the platform is generated by the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_FREQ_BASE within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_FREQ_BASE procedure.
	base_freq	Base frequency of the platform if generated by the processor chip.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Can not complete call without error

**Description:** If the processor outputs a clock for use by the platform, the *base\_freq* return parameter will be the frequency of this output clock in ticks per second. If the processor does not generate an output clock for use by the platform, this procedure will return with a status of -1.

## PAL\_FREQ\_RATIOS – Get Processor Frequency Ratios (14)

**Purpose:** Returns the ratios of the processor frequency, bus frequency, and interval timer to the input clock of the processor, if the platform clock is generated externally or to the output clock to the platform, if the platform clock is generated by the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

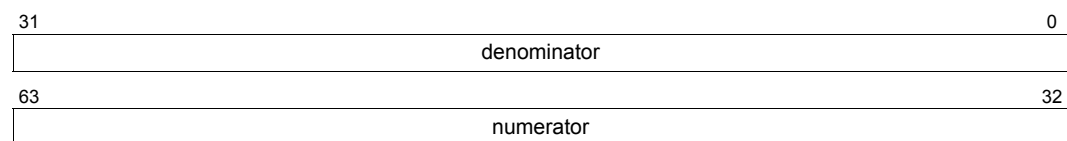
Arguments:	Argument	Description
	index	Index of PAL_FREQ_RATIOS within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_FREQ_RATIOS procedure.
	proc_ratio	Ratio of the processor frequency to the input clock of the processor, if the platform clock is generated externally or to the output clock to the platform, if the platform clock is generated by the processor.
	bus_ratio	Ratio of the bus frequency to the input clock of the processor, if the platform clock is generated externally or to the output clock to the platform, if the platform clock is generated by the processor.
	itc_ratio	Ratio of the interval timer counter rate to input clock of the processor, if the platform clock is generated externally or to the output clock to the platform, if the platform clock is generated by the processor.

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Can not complete call without error

**Description:** Each of the ratios returned is an unsigned 64-bit value, where the upper unsigned 32 bits contain the numerator and the lower unsigned 32 bits contain the denominator of the ratio, as depicted in [Figure 11-12](#). Each ratio is given by dividing the numerator by the denominator.

**Figure 11-12. Return values**



- denominator – Unsigned 32-bit integer
- numerator – Unsigned 32-bit integer

## PAL\_GET\_HW\_POLICY – Retrieve Current Hardware Resource Sharing Policy (48)

**Purpose:** Returns the current hardware resource sharing policy of the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_GET_HW_POLICY within the list of PAL procedures.
proc_num	Unsigned 64-bit integer that specifies for which logical processor information is being requested. This input argument must be zero for the first call to this procedure and can be a maximum value of one less than the number of logical processors impacted by the hardware resource sharing policy, which is returned by the <i>num_impacted</i> return value.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_GET_HW_POLICY procedure.
cur_policy	Unsigned 64-bit integer representing the current hardware resource sharing policy.
num_impacted	Unsigned 64-bit integer that returns the number of logical processors impacted by the <i>policy</i> input argument.
la	Unsigned 64-bit integer containing the logical address of one of the logical processors impacted by policy modification.

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-9	Call requires PAL memory buffer

**Description:** This procedure is used to return information on the current hardware resource sharing policy. This procedure can also be used to identify which logical processors (see [“PAL\\_LOGICAL\\_TO\\_PHYSICAL – Get Information on Logical to Physical Processor Mappings \(42\)”](#) on page 2:404 for a definition of a logical processor) are impacted by the various hardware sharing policies supported on the processor.

The procedure returns information about the current hardware sharing policy, the total number of logical processors impacted by hardware sharing policies and the logical address of one of the processors impacted by the hardware sharing policy.

The definition of the hardware sharing policies that can be returned in the *cur\_policy* value are defined in [Table 11-80](#).

**Table 11-80. Hardware policies returned in *cur\_policy***

Value	Name	Description
0	Performance	The processor has its hardware resources configured to achieve maximum performance across all logical processors that share hardware with the logical processor the procedure was made on.
1	Fairness	The processor has its hardware resources configured to approximately achieve equal sharing of competing hardware resources among all the logical processors that share hardware with the logical processor the procedure was made on.
2	High-priority	The processor has its hardware resources configured such that the logical processor this procedure was called on has a greater share of the competing hardware resources.
3	Exclusive High-priority	The processor has its hardware resources configured such that the logical processor this procedure was called on has a greater share of the competing hardware resources. See <a href="#">“PAL_SET_HW_POLICY – Set Current Hardware Resource Sharing Policy (49)” on page 2:456</a> for differences between high-priority and exclusive high priority.
4	Low-priority	The processor has its hardware resources configured such that the logical processor this procedure was called on has a smaller share of the competing hardware resources. This occurs when a competing logical processor has itself set as high priority or exclusive high priority.
All Other Values		Reserved

The return value *num\_impacted* specifies the number of logical processors impacted by the hardware sharing policy. The return value *la* returns the logical address of one of the logical processors impacted by the hardware sharing policy. The return value *la* is the same value and format of that is returned by the [PAL\\_FIXED\\_ADDR – Get Fixed Geographical Address of Processor \(12\)” on page 2:391](#) for details.

If the caller is interested in identifying all the logical processors impacted by the hardware sharing policy, this procedure will need to be called a number of times equal to the value returned in *num\_impacted* return value. For each subsequent call it needs to increment the 'proc\_num' input argument.

The logical processor this procedure is made on can only return information about how the hardware sharing policy impacts logical processors it is sharing hardware resources with. For example a physical processor package may contain two multi-threaded cores. On this example implementation the hardware sharing policy only impacts the two threads on the core and this procedure would only return the two *la*'s of the threads on that core, but would not return the *la*'s of the threads on the other core. When this procedure was made on the other core, then that procedure call would return the *la*'s of the two threads on that core.

This procedure is only supported on processors that have multiple logical processors sharing hardware resources that can be configured. On all other processor implementations, this procedure will return the Unimplemented procedure return status.

## PAL\_GET\_PSTATE – Return Information on the Performance Index of the Processor (262)

**Purpose:** Returns the performance index of the processor.

**Calling Conv:** Stacked Registers

**Mode:** Physical and Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_GET_PSTATE within the list of PAL procedures.
type	Type of <i>performance_index</i> value to be returned by this procedure.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_GET_PSTATE procedure.
performance_index	Unsigned integer denoting the processor performance for the time duration since the last PAL_GET_PSTATE procedure call was made. The value returned is relative to the performance index of the highest available P-state.
Reserved	0
Reserved	0

Status Value	Description
1	Call completed without error, but accuracy of performance index has been impacted by a thermal throttling event, or a hardware-initiated event.
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-9	Call requires PAL memory buffer

**Description:** This procedure returns a performance index of the processor, and is relative to the highest available P-state, P0. A value of 100 represents the minimum processor performance in the P0 state. For processors that support variable P-state performance, it is possible for a processor to report a number greater than 100, representing that the processor is running at a performance level greater than the minimum P0 performance. The PAL procedure [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)” on page 2:446](#) indicates whether the processor supports variable P-state performance.

The *type* argument allows the caller to select the *performance\_index* value that will be returned. See [Table 11-81](#) below for details.



Table 11-81. PAL\_GET\_PSTATE type Argument

type	Description
0	<p>The <i>performance_index</i> returned will correspond to the target P-state requested by software.</p> <ul style="list-style-type: none"> <li>For SCDD (software-coordinated dependency domain) logical processors, this is the P-state requested by the most recent PAL_SET_PSTATE procedure call made by any logical processor in the domain.</li> <li>For HCDD (hardware-coordinated dependency domain) or HIDD (hardware-independent dependency domain) logical processors, this is simply the P-state requested by the most recent PAL_SET_PSTATE procedure call on this logical processor.</li> </ul> <p>The value returned is not affected by platform power-caps.</p>
1	<p>The <i>performance_index</i> is a weighted-average value of the different P-states that the processor was operating in for the time duration between the current PAL_GET_PSTATE procedure call, and the previous invocation of PAL_GET_PSTATE with <i>type</i>=1. This allows the caller to establish a new starting point for subsequent computation of the weighted-average <i>performance_index</i>. See <a href="#">Section 11.6.1, “Power/Performance States (P-states)” on page 2:315</a> for more details on how the weighted average value is derived.</p>
2	<p>The <i>performance_index</i> is a weighted-average value of the different P-states that the processor was operating in for the time duration between the current PAL_GET_PSTATE procedure call, and the previous invocation of PAL_GET_PSTATE with <i>type</i>=1. This allows the caller to sample the current value of the <i>performance_index</i>, without affecting the starting point used for computing the weighted-average <i>performance_index</i>.</p>
3	<p>The <i>performance_index</i> returned will correspond to the current instantaneous P-state of the dependency domain containing the logical processor, at the time of the procedure call. The value returned is not affected by platform power-caps. When variable P-states performance is supported, the <i>performance_index</i> may be higher than the P-state requested. Please see <a href="#">Section 11.6.1.4, “Variable P-state Performance” on page 2:322</a> for more information about variable P-state performance.</p>
All Other Values	Reserved

For SCDD logical processors, or HIDD logical processors that do not support platform power-caps, note that the *performance\_index* returned for *type*=0 and *type*=3 will have identical values. This is because the most recent PAL\_SET\_PSTATE procedure call that returned a status of 0 will always succeed in transitioning to the requested performance state for these coordination domains (see PAL\_SET\_PSTATE procedure description for additional details).

For SCDD logical processors, the PAL\_GET\_PSTATE procedure should always be called with *type* argument value of 0 or 3. On such processors, calling PAL\_GET\_PSTATE with *type* argument value of 1 or 2 is undefined.

For HIDD logical processors, the *type* argument values of 1 and 2 are supported, since such processors can also support platform power-caps, which affect the weighted-average performance index.

If there was a thermal-throttling or hardware-initiated event (other than a platform power-cap) which affected the processor power/performance for the current time period, and the accuracy of the *performance\_index* value has been impacted by the event, then the procedure will return with *status*=1. The *performance\_index* returned in this case will still have a value that falls within the range of possible *performance\_index* values for this processor implementation (i.e., 0 up to the highest variable p-state *performance\_index* value).

The procedure, when called with *type*=1 or *type*=2, returns a fixed *performance\_index* value of 100 until the procedure has been called with *type*=1 to reset computation of the weighted-average *performance\_index*. For subsequent invocations with *type*=1 or

*type=2*, the procedure will return the *performance\_index* value corresponding to the processor performance in the time duration between the previous call to PAL\_GET\_PSTATE with *type=1* and the current call.

If the processor had transitioned to a HALT state (see [Section 11.6.1, “Power/Performance States \(P-states\)”](#) on page 2:315) in between successive invocations to the PAL\_GET\_PSTATE procedure, the performance index computation returned will not take into account the performance of the processor during the time spent in HALT state (see [Section 11.6.1.5, “Interaction of P-states with HALT State”](#) on page 2:323 for details).

## PAL\_HALT – Halt Processor (28)

**Purpose:** Causes the processor to enter the HALT state, or one of the implementation-dependent low-power states.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_HALT within the list of PAL procedures.
	halt_state	Unsigned 64-bit integer denoting low power state requested.
	io_detail_ptr	8-byte aligned physical address pointer to information on the type of I/O (load/store) requested.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_HALT procedure.
	load_return	Value returned if a load instruction is requested in the <i>io_detail_ptr</i>
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error

**Description:** This call places the processor in a low power state designated by *halt\_state*. This procedure can optionally let the platform know it is about to enter the low power state via an I/O transaction.

*halt\_state* is an unsigned 64-bit integer denoting the low power state requested. The value passed must be a valid halt state in the range from 1 to 7, for which information is returned by PAL\_HALT\_INFO. All other values are reserved.

The processor informs the platform that it has entered the requested low-power state in an implementation-specific manner.

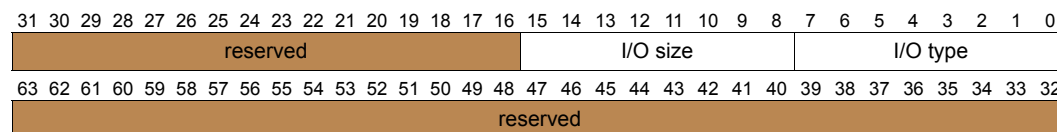
The layout of the information pointed to by the *io\_detail\_ptr* is shown [Table 11-82](#).

**Table 11-82. I/O Detail Pointer Description**

Offset	Description
0x0	I/O size and type information
0x8	Address for I/O
0x10	Data value to store

- I/O size and type information has the format shown in [Figure 11-13](#).

**Figure 11-13. I/O Size and Type Information Layout**



- *I/O type* is an unsigned 8-bit integer denoting the type of I/O transaction to complete.

**Table 11-83. I/O Type Definition**

Value	Description
0	No transaction
1	Perform a load
2	Perform a store

All other values for *I/O type* are reserved.

- *I/O size* is an unsigned 8-bit integer denoting the size of the I/O transaction to complete.

**Table 11-84. I/O Size Definition**

Value	Description
0	No transaction
1	1 byte size
2	2 byte size
4	4 byte size
8	8 byte size

All other values for *I/O size* are reserved.

- Address for the I/O transaction is a physical pointer for the load or store. The address passed should be aligned according to the size of the I/O transaction requested. The most significant bit (63) of the physical address should be set according to the cacheability attribute wanted for the I/O transaction.
- The data value to store is the value that will be stored out if the *io\_type* is 2. If *io\_type* is not equal to a 2, then this value is a don't care.

If an I/O transaction is requested by the caller, the processor will wait until this transaction has been received by the platform before entering the low power state.

On receipt of a PMI, machine check, INIT, reset, or unmasked external interrupt (including NMI), PAL transitions the processor to the normal state. An unmasked external interrupt is defined to be an interrupt that is permitted to interrupt the processor based on the current setting of the TPR.mic and TPR.mmi fields in the TPR control register. PAL sets the value in the *load\_return* return parameter if the *io\_type* is 1, otherwise this value is set to zero.

If the processor transitions to normal state via an unmasked external interrupt, execution resumes to the caller.

If the processor transitions to normal state via a PMI, execution resumes to the caller if PMIs are masked, otherwise execution will resume to the PMI handler.

If the processor transitions to the normal state via a machine check or INIT, execution resumes to the caller if machine checks and INITs are masked, otherwise execution will resume to the corresponding handler.

If the processor transitions to the normal state via a reset event, the processor will reset itself and start execution at the PAL reset address.

For more information on power management, please refer to Section 11.6, "Power Management" on [page 2:313](#).

## PAL\_HALT\_INFO – Get Halt State Information for Power Management (257)

**Purpose:** Returns information about the processor’s power management capabilities.

**Calling Conv:** Stacked Registers

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_HALT_INFO within the list of PAL procedures.
power_buffer	64-bit pointer to a 64-byte buffer aligned on an 8-byte boundary.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_HALT_INFO procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** The power information requested is returned in the data buffer referenced by *power\_buffer*. Power information is returned about the 8 power states. The low power states are LIGHT\_HALT, HALT, plus 6 other low power states. The LIGHT\_HALT state is index 0 in the buffer, and the HALT state is index 1. All 8 low power states need not be implemented

The information returned is in the format of [Figure 11-14](#). The information about the HALT states will be in ascending order of the index values.

**Figure 11-14. Layout of *power\_buffer* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0					
entry_latency																exit_latency																				
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32					
rv	co	im	power_consumption																																	

- *exit\_latency* – 16-bit unsigned integer denoting the minimum number of processor cycles to transition to the NORMAL state.
- *entry\_latency* – 16-bit unsigned integer denoting the minimum number of processor cycles to transition from the NORMAL state.
- *power\_consumption* – 28-bit unsigned integer denoting the typical power consumption of the state, measured in milliwatts.
- *im* – 1-bit field denoting whether this low power state is implemented or not. A value of 1 indicates that the low power state is implemented, a value of 0 indicates that it is not implemented. If this value is 0 then all other fields are invalid.
- *co* – 1-bit field denoting if the low power state maintains cache and TLB coherency. A value of 1 indicates that the low power state keeps the caches and TLBs coherent, a value of 0 indicates that it does not.

The latency numbers given are the minimum number of processor cycles that will be required to transition the states. The maximum or average cannot be determined by PAL due to its dependency on outstanding bus transactions.

For more information on power management, please refer to Section 11.6, "Power Management" on [page 2:313](#).

## PAL\_HALT\_LIGHT – Cause Processor to Enter Coherent Halt State (29)

**Purpose:** Causes the processor to enter the LIGHT HALT state, where prefetching and execution are suspended, but cache and TLB coherency is maintained.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_HALT_LIGHT within the list of PAL procedures.
Reserved	0
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_HALT_LIGHT procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** This call places the processor in the LIGHT HALT state in an implementation-dependent fashion where cache and TLB coherency is maintained, but power consumption is minimized.

The processor acknowledges to the platform that it has entered the LIGHT HALT low-power state in an implementation-specific manner.

On receipt of a PMI, machine check, INIT, reset, or unmasked external interrupt (including NMI), PAL transitions the processor to the normal state. An unmasked external interrupt is defined to be an interrupt that is permitted to interrupt the processor based on the current setting of the TPR.mic and TPR.mmi fields in the TPR control register.

If the processor transitions to normal state via an unmasked external interrupt, execution resumes to the caller.

If the processor transitions to normal state via a PMI, execution resumes to the caller if PMIs are masked, otherwise execution will resume to the PMI handler.

If the processor transitions to the normal state via a machine check or INIT, execution resumes to the caller if machine checks and INITs are masked, otherwise execution will resume to the corresponding handler.

If the processor transitions to the normal state via a reset event, the processor will reset itself and start execution at the PAL reset address.

For more information on power management, please refer to Section 11.6, "Power Management" on [page 2:313](#).

## PAL\_LOGICAL\_TO\_PHYSICAL – Get Information on Logical to Physical Processor Mappings (42)

**Purpose:** Returns information on the logical to physical processor mapping.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_LOGICAL_TO_PHYSICAL within the list of PAL procedures.
proc_number	Signed 64-bit integer that specifies for which logical processor information is being requested. When this input argument is -1, information is returned about the logical processor on which the procedure call is made. This input argument must be in the range of -1 up to one less than the number of logical processors returned by <i>num_log</i> in the <i>log_overview</i> return value.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_LOGICAL_TO_PHYSICAL procedure.
log_overview	The format of <i>log_overview</i> is shown in <a href="#">Figure 11-15</a> .
proc_n_log_info1	The format of <i>proc_n_log_info1</i> is shown in <a href="#">Figure 11-16</a> .
proc_n_log_info2	The format of <i>proc_n_log_info2</i> is shown in <a href="#">Figure 11-17</a> .

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

**Description:** This procedure will return information about the logical processors contained on the physical processor package that the procedure call is made on. A physical processor package can contain one or more logical processors, organized into threads and cores. A logical processor is a compute-capability-centric view of the CPU that allows the physical processor package to execute from more than one instruction stream. A physical processor package that can execute from *n* instruction streams has *n* logical processors. Threads are logical processors that share core pipeline execution resources. Cores are defined as a collection of hardware that implements the main execution pipeline of the processor. Multiple cores on a physical processor package do not share core pipeline resources but may share caches and bus interfaces. A core may support multiple threads of execution.

The *log\_overview* return value provides an overview of the logical processors on the physical processor package this procedure call was made on. The format of the *log\_overview* return argument is shown in [Figure 11-15](#).



**Figure 11-15. Layout of *log\_overview* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
rv								tpc								num_log															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
rv								ppid								rv								cpp							

- *num\_log* – Total number of logical processors on this physical processor package that are enabled.
- *tpc* – Threads per core. Number of threads per core.
- *rv* – Reserved
- *cpp* – Cores per processor. Total number of cores on this physical processor package.
- *rv* – Reserved
- *ppid* – Physical processor package ID. Physical processor package identifier which was assigned at reset by the platform or bus controller. This value may or may not be unique across the entire platform since it depends on the platform vendor's policy.
- *rv* – Reserved

It is not ensured that *num\_log* will always be equal to *cpp* multiplied by *tpc*. This is possible if some logical processors are disabled through implementation specific means.

The caller uses the value returned in *num\_log* to gather additional information about the other logical processors on the same physical processor package. This procedure will need to be called multiple times (equal to the number of logical processors returned in *num\_log*) to gather all additional information about the logical processors on the physical processor package this procedure call was made on. This procedure may be called from any logical processor on the physical processor package to gather information about all the logical processors. It may also be called to get information about the logical processor on which the procedure is running. Information about the logical processors is in the return values *proc\_n\_log\_info1* and *proc\_n\_log\_info2*. The format of these return values is shown in [Figure 11-16](#) and [Figure 11-17](#).

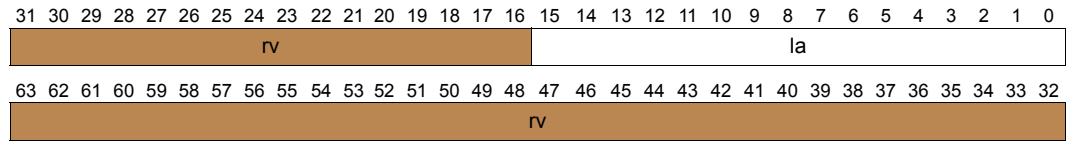
**Figure 11-16. Layout of *proc\_n\_log\_info1* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
rv																tid															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
rv																cid															

- *tid* – Thread id: The thread identifier of the logical processor for which information is being returned. This value will be unique on a per core basis.
- *rv* – Reserved
- *cid* – Core id: The core identifier of the logical processor for which information is being returned. This value will be unique on a per physical processor package basis.
- *rv* – Reserved

There is no guarantee that the core id's and thread id's will be contiguous on a given physical processor package.

**Figure 11-17. Layout of *proc\_n\_log\_info2* Return Value**



- *la* – Logical address: geographical address of the logical processor for which information is being returned. This is the same value that is returned by the PAL\_FIXED\_ADDR procedure when it is called on the logical processor.
- *rv* – Reserved

This procedure must be supported on all implementations that contain more than one logical processor on a physical processor package and returns an unimplemented procedure error code otherwise.

## PAL\_MC\_CLEAR\_LOG – Clear Processor Error Logging Registers (21)

**Purpose:** Clears all processor error logging registers and resets the indicator that allows the error logging registers to be written. This procedure also checks the pending machine check bit and pending INIT bit and reports their states.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_CLEAR_LOG within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

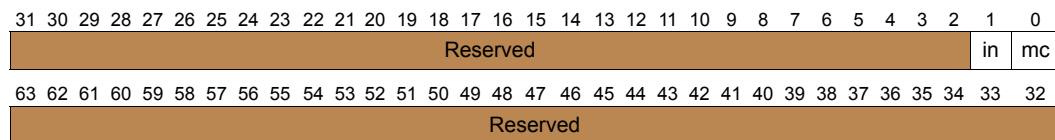
Returns:	Return Value	Description
	status	Return status of the PAL_MC_CLEAR_LOG procedure.
	pending	64-bit vector denoting whether an event is pending.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This procedure is called to clear processor error logging registers after all error information has been obtained. This procedure re-enables the logging registers in the case of a subsequent error. It clears any information that would be returned by either the PAL\_MC\_ERROR\_INFO or PAL\_MC\_DYNAMIC\_STATE procedures.

This procedure does not clear any pending machine checks. The *pending* return parameter returns a value of 0 if no subsequent event is pending, a 1 in bit position 0, if a machine check is pending, and/or a 1 in bit position 1 if an INIT is pending. All other values are reserved.

**Figure 11-18. Pending Return Parameter**



**Table 11-85. Pending Return Parameter Fields**

Field	Description
mc	Pending machine check
in	Pending initialization event

## PAL\_MC\_DRAIN – Complete Outstanding Transactions (22)

**Purpose:** Ensures that all outstanding transactions in a processor are completed or that any MCA due to these outstanding transactions is taken.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_DRAIN within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_MC_DRAIN procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This call causes all outstanding transactions in the processor to be completed. For example:

- Flushes ( $\text{fc}$ ) invalidate the cache, lines that have been modified are written back (issued to the fabric) to memory before invalidation.
- Instruction cache coherence flushes ( $\text{fc.i}$ ) invalidate lines and/or write them back to main memory, if this is required to make the instruction caches coherent with the data caches.
- Loads get their data returned.
- Stores either update the cache or issue transactions to the system fabric.
- Prefetches are either completed or cancelled,

As a result of completing these outstanding transactions Machine Check Aborts (MCAs) may be taken. This call is typically issued by code that needs to guarantee that no MCAs due to outstanding transactions will occur after a given point.

## PAL\_MC\_DYNAMIC\_STATE – Returns Dynamic Processor State (24)

**Purpose:** Returns the Machine Check Dynamic Processor State.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_MC_DYNAMIC_STATE within the list of PAL procedures.
info_type	Unsigned 64-bit value indicating the type of information to return
dy_buffer	64-bit pointer to a buffer aligned on an 8-byte boundary
Reserved	0

Return Value	Description
status	Return status of the PAL_MC_DYNAMIC_STATE procedure.
max_size	Maximum size (in bytes) of the data that can be returned by this procedure for this processor family.
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

**Description:** The *info\_type* input argument designates the type of information the procedure will return. When *info\_type* is 0, the procedure returns the maximum size (in bytes) of processor dynamic state that can be returned for this processor family in the *max\_size* return value.

When *info\_type* is 1, the procedure will copy processor dynamic state into memory pointed to by the input argument *dy\_buffer*. This copy will occur using the addressing attributes used to make the procedure call (physical or virtual) and the caller needs to ensure the *dy\_buffer* input pointer matches this addressing attribute.

The amount of data returned can vary depending on the state of the machine at the time the procedure is called, and may not always return the maximum size for every call. The amount of data returned is provided in the processor state parameter field *dsize*. Please see [Table 11-7](#) for more information on the processor state parameter. The caller of the procedure needs to ensure that the buffer is large enough to handle the *max\_size* that is returned by this procedure.

The contents of the processor dynamic state is implementation dependent. Portions of this information may be cleared by the PAL\_MC\_CLEAR\_LOG procedure. This procedure should be invoked before PAL\_MC\_CLEAR\_LOG to ensure all the data is captured.

## PAL\_MC\_ERROR\_INFO – Get Processor Error Information (25)

**Purpose:** Returns the Processor Machine Check Information

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_ERROR_INFO within the list of PAL procedures.
	info_index	Unsigned 64-bit integer identifying the error information that is being requested. (See <a href="#">Table 11-86</a> ).
	level_index	8-byte formatted value identifying the structure to return error information on. (See <a href="#">Figure 11-19</a> ).
	err_type_index	Unsigned 64-bit integer denoting the type of error information that is being requested for the structure identified in <i>level_index</i> .

Returns:	Return Value	Description
	status	Return status of the PAL_MC_ERROR_INFO procedure.
	error_info	Error information returned. The format of this value is dependant on the input values passed.
	inc_err_type	If this value is zero, all the error information specified by <i>err_type_index</i> has been returned. If this value is one, more structure-specific error information is available and the caller needs to make this procedure call again with <i>level_index</i> unchanged and <i>err_type_index</i> , incremented.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error
	-6	Argument was valid, but no error information was available

**Description:** This procedure returns error information for machine checks as specified by *info\_index*, *level\_index* and *err\_type\_index*. Higher level software is informed that additional machine check information is available when the processor state parameter *mi* bit is set to one. See [Table 11-7, "Processor State Parameter Fields," on page 2:299](#) for more information on the processor state parameter and the *mi* bit description.

The *info\_index* argument specifies which error information is being requested. See [Table 11-86](#) for the definition of the *info\_index* values.

**Table 11-86. *info\_index* Values**

<i>info_index</i>	Error Information Type	Description
0	Processor Error Map	This <i>info_index</i> value will return the processor error map. This return value specifies the processor core identification, the processor thread identification, and a bit-map indicating which structure(s) of the processor generated the machine check. This bit-map has the same layout as the <i>level_index</i> . A one in the structure bit-map indicates that there is error information available for the structure. The layout of the <i>level_index</i> is described in <a href="#">Figure 11-19, “level_index Layout”</a> on page 2:411.
1	Processor State Parameter	This <i>info_index</i> value will return the same processor state parameter that is passed at the PALE_CHECK exit state for a machine check event (provided a valid min-state save area has been registered) or will construct a processor state parameter for a corrected machine check events. This parameter describes the severity of the error and the validity of the processor state when the machine check or CMCI occurred. This procedure will not return a valid PSP for INIT events. The Processor State Parameter is described in <a href="#">Figure 11-11, “Processor State Parameter,”</a> on page 2:299.
2	Structure-specific Error Information	This <i>info_index</i> value will return error information specific to a processor structure. The structure is specified by the caller using the <i>level_index</i> and <i>err_type_index</i> input parameters. The value returned in <i>error_info</i> is specific to the structure and type of information requested.

All other values of *info\_index* are reserved. When *info\_index* is equal to 0 or 1, the *level\_index* and *err\_type\_index* input values are ignored. When *info\_index* is equal to 2, the *level\_index* and *err\_type\_index* define the format of the *error\_info* return value.

The caller is expected to first make this procedure call with *info\_index* equal to zero to obtain the processor error map. This error map informs the caller about the processor core identification, the processor thread identification and indicates which structure(s) caused the machine check. If more than one structure generated a machine check, multiple structure bits will be set. The caller then uses this information to make sub-subsequent calls to this procedure for each structure identified in the processor error map to obtain detailed error information.

The *level\_index* input argument specifies which processor core, processor thread and structure for which information is being requested. See [Table 11-87 on page 2:412](#) for the definition of the *level\_index* fields. This procedure call can only return information about one processor structure at a time. The caller is responsible for ensuring that only one structure bit in the *level\_index* input argument is set at a time when retrieving information, otherwise the call will return that an invalid argument was passed.

**Figure 11-19. *level\_index* Layout**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
erf				ebh				edt				eit				edc				eic				tid				cid			
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32

**Figure 11-19. *level\_index* Layout****Table 11-87. *level\_index* Fields**

Field	Bits	Description
cid	3:0	Processor core ID (default is 0 for processors with a single core)
tid	7:4	Logical thread ID (default is 0 for processors that execute a single thread)
eic	11:8	Error information is available for 1st, 2nd, 3rd, and 4th level instruction caches
edc	15:12	Error information is available for 1st, 2nd, 3rd, and 4th level data/unified caches
eit	19:16	Error information is available for 1st, 2nd, 3rd, and 4th level instruction TLB
edt	23:20	Error information is available for 1st, 2nd, 3rd, and 4th level data/unified TLB
ebh	27:24	Error information is available for the 1st, 2nd, 3rd, and 4th level processor bus hierarchy
erf	31:28	Error information is available on register file structures
ems	47:32	Error information is available on micro-architectural structures
rsvd	63:48	Reserved

The convention for levels and hierarchy in the *level\_index* field is such that the least significant bit in the error information bit-fields represent the lowest level of the structures hierarchy. For example bit 8 if the *eic* field represents the first level instruction cache.

The *erf* field is 4-bits wide to allow reporting of 4 concurrent register related machine checks at one time. One bit would be set for each error. The *ems* field is 16-bits wide to allow reporting of 16-concurrent micro-architectural structures at one time. There is no significance in the order of these bits. If only one register file related error occurred, it could be reported in any one of the 4-bits.

The *err\_type\_index* specifies the type of information will be returned in *error\_info* for a particular structure. See [Table 11-88](#) for the values of *err\_type\_index*

**Table 11-88. *err\_type\_index* Values**

<i>err_type_index</i> value mod 8	Return Value	Description
0	Structure-specific error information specified by <i>level_index</i>	The information returned in <i>error_info</i> is dependant on the structure specified in <i>level_index</i> . See <a href="#">Table 11-89</a> for the <i>error_info</i> return formats.
1	Target address	The target address is a 64-bit integer containing the physical address where the data was to be delivered or obtained. The target address also can return the incoming address for external snoops and TLB shoot-downs that generated a machine check. The structure-specific error information informs the caller if there is a valid target address to be returned for the requested structure.
2	Requester identifier	The requester identifier is a 64-bit integer that specifies the bus agent that generated the transaction responsible for generating the machine check. The structure-specific error information informs the caller if there is a valid requester identifier.



**Table 11-88. *err\_type\_index* Values (Continued)**

<i>err_type_index</i> value mod 8	Return Value	Description
3	Responder identifier	The responder identifier is a 64-bit integer that specifies the bus agent that responded to a transaction that was responsible for generating the machine check. The structure-specific error information informs the caller if there is a valid responder identifier.
4	Precise instruction pointer	The precise instruction pointer is a 64-bit virtual address that points to the bundle that contained the instruction responsible for the machine check. The structure-specific error information informs the caller if there is a valid precise instruction pointer.
5-7	Reserved	Reserved

See [Table 11-89](#) for the format of *error\_info* when structure-specific information is requested.

**Table 11-89. *error\_info* Return Format when *info\_index* = 2 and *err\_type\_index* = 0**

<i>level_index</i> Field Input	<i>error_info</i> Return Format
eic	cache_check return format
edc	cache_check return format
eit	tlb_check return format
edt	tlb_check return format
ebh	bus_check return format
erf	reg_file_check return format
ems	uarch_check return format

The structure specified by the *level\_index* may have the ability to log distinct multiple errors. This can occur if the structure is accessed at the same time by more than one instruction and the processor can log machine check information for each access. To inform the caller of this occurrence, this procedure will return a value of one in the *inc\_err\_type* return value.

It is important to note, that when the caller sees that the *inc\_err\_type* return value is one, it should make a sub-sequent call with the *err\_type\_index* value incremented by 8. If the structure-specific error information returns that there is a valid target address, requester identifier, responder identifier or precise instruction pointer these can be returned as well by incrementing the *err\_type\_index* value in the same manner. Refer to the following example for more information.

For example, to gather information on the first error of a structure that can log multiple errors, *err\_type\_index* would be called with the value of 0 first. The caller examines the information returned in *error\_info* to know if there is a valid target address, requester identifier, responder identifier, or precise instruction pointer available for logging. If there is, it makes sub-sequent calls with *err\_type\_index* equal to 1, 2, 3 and/or 4 depending on which valid bits are set. Additionally if the *inc\_err\_type* return value was set to one, the caller knows that this structure logged multiple errors. To get the second error of the structure it sets the *err\_type\_index* = 8 and the structure-specific information is returned in *error\_info*. The caller examines this *error\_info* to know if there is a valid target address, requester identifier, responder identifier, or precise

instruction pointer available for logging on the second error. If there is, it makes sub-sequent calls with *err\_type\_index* equal to 9, 10, 11, and/or 12 depending on which valid bits are set. The caller continues incrementing the *err\_type\_index* value in this fashion until the *inc\_err\_type* return value is zero.

As shown in Table 11-89, the information returned in *error\_info* varies based on which structure information is being requested on. The next sections describe the *error\_info* return format for the different structures.

**Cache\_Check Return Format:** The cache check return format is returned in *error\_info* when the user requests information on any instruction or data/unified caches in the *level\_index* input argument. The cache\_check return format must be used to report errors in cacheable transactions. These errors may also be reported using the bus\_check return format if the bus structures can detect these errors. The cache\_check return format is a bit-field that is described in Figure 11-20 and Table 11-90.

Figure 11-20. cache\_check Layout

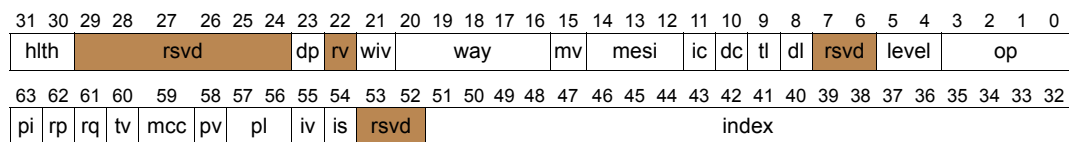


Table 11-90. cache\_check Fields

Field	Bits	Description
op	3:0	Type of cache operation that caused the machine check: 0 – unknown or internal error 1 – load 2 – store 3 – instruction fetch or instruction prefetch 4 – data prefetch (both hardware and software) 5 – snoop (coherency check) 6 – cast out (explicit or implicit write-back of a cache line) 7 – move in (cache line fill) All other values are reserved.
level	5:4	Level of cache where the error occurred. A value of 0 indicates the first level of cache.
rsvd	7:6	Reserved
dl	8	Failure located in the data part of the cache line.
tl	9	Failure located in the tag part of the cache line.
dc	10	Failure located in the data cache
ic	11	Failure located in the instruction cache
mesi	14:12	0 – cache line is invalid. 1 – cache line is held shared. 2 – cache line is held exclusive. 3 – cache line is modified. All other values are reserved.
mv	15	The <i>mesi</i> field in the cache_check parameter is valid.
way	20:16	Failure located in the way of the cache indicated by this value.
wiv	21	The <i>way</i> and <i>index</i> field in the cache_check parameter is valid.
rsvd	22	Reserved
dp	23	An uncorrectable (typically multiple-bit) error was detected and data was poisoned for the corresponding cache line, without any corrupted data being consumed (i.e., no corrupted data has been copied to processor registers).

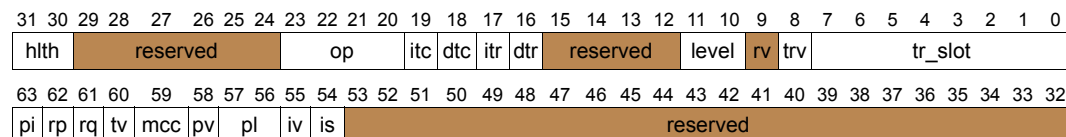
**Table 11-90. cache\_check Fields (Continued)**

Field	Bits	Description
rsvd	29:24	Reserved
hlth	31:30	Health indicator. This field will report if the cache type and level reporting this error supports hardware status tracking and the current status of this cache. 00 – No hardware status tracking is provided for the cache type and level reporting this event. 01 – Status tracking is provided for this cache type and level and the current status is normal status. <sup>a</sup> 10 – Status tracking is provided for the cache type and level and the current status is cautionary. <sup>a</sup> When a cache reports a cautionary status the "hardware damage" bit of the PSP (see Figure 11-11, "Processor State Parameter," on page 2:299) will be set as well. 11 – Reserved
index	51:32	Index of the cache line where the error occurred.
rsvd	53:52	Reserved
is	54	Instruction set. If this value is set to zero, the instruction that generated the machine check was an Intel Itanium instruction. If this bit is set to one, the instruction that generated the machine check was IA-32 instruction.
iv	55	The <i>is</i> field in the cache_check parameter is valid.
pl	57:56	Privilege level. The privilege level of the instruction bundle responsible for generating the machine check.
pv	58	The <i>pl</i> field of the cache_check parameter is valid.
mcc	59	Machine check corrected: This bit is set to one to indicate that the machine check has been corrected.
tv	60	Target address is valid: This bit is set to one to indicate that a valid target address has been logged.
rq	61	Requester identifier: This bit is set to one to indicate that a valid requester identifier has been logged.
rp	62	Responder identifier: This bit is set to one to indicate that a valid responder identifier has been logged.
pi	63	Precise instruction pointer. This bit is set to one to indicate that a valid precise instruction pointer has been logged.

a. Hardware is tracking the operating status of the structure type and level reporting the error. The hardware reports a "normal" status when the number of entries within a structure reporting repeated corrections is at or below a pre-defined threshold. A "cautionary" status is reported when the number of affected entries exceeds a pre-defined threshold.

**TLB\_Check Return Format:** The *tlb\_check* return format is returned in *error\_info* when the user requests information on any instruction or data/unified TLB in the *level\_index* input argument. The *tlb\_check* return format is a bit-field that is described in Figure 11-21 and Table 11-91.

**Figure 11-21. tlb\_check Layout**



**Table 11-91. tlb\_check Fields**

Field	Bits	Description
tr_slot	7:0	Slot number of the translation register where the failure occurred.
trv	8	The <i>tr_slot</i> field in the TLB_check parameter is valid.

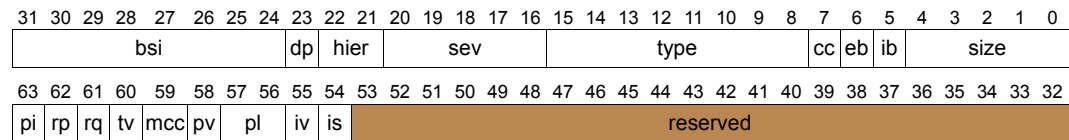
Table 11-91. tlb\_check Fields (Continued)

Field	Bits	Description
rv	9	Reserved
level	11:10	The level of the TLB where the error occurred. A value of 0 indicates the first level of TLB
reserved	15:12	Reserved
dtr	16	Error occurred in the data translation registers
itr	17	Error occurred in the instruction translation registers
dtrc	18	Error occurred in data translation cache
itr	19	Error occurred in the instruction translation cache
op	23:20	Type of cache operation that caused the machine check: 0 – unknown 1 – TLB access due to load instruction 2 – TLB access due to store instruction 3 – TLB access due to instruction fetch or instruction prefetch 4 – TLB access due to data prefetch (both hardware and software) 5 – TLB shoot down access 6 – TLB probe instruction (probe, tpa) 7 – move in (VHPT fill) 8 – purge (insert operation that purges entries or a TLB purge instruction) All other values are reserved.
reserved	29:24	Reserved
hlth	31:30	Health indicator. This field will report if the tlb type and level reporting this error supports hardware status tracking and the current status of this tlb. 00 – No hardware status tracking is provided for the tlb type and level reporting this event. 01 – Status tracking is provided for this tlb type and level and the current status is normal. <sup>a</sup> 10 – Status tracking is provided for the tlb type and level and the current status is cautionary. <sup>a</sup> When a tlb reports a cautionary status the "hardware damage" bit of the PSP (see Figure 11-11, "Processor State Parameter," on page 2:299) will be set as well. 11 – Reserved
reserved	53:32	Reserved
is	54	Instruction set. If this value is set to zero, the instruction that generated the machine check was an Intel Itanium instruction. If this bit is set to one, the instruction that generated the machine check was IA-32 instruction.
iv	55	The <i>is</i> field in the TLB_check parameter is valid.
pl	57:56	Privilege level. The privilege level of the instruction bundle responsible for generating the machine check.
pv	58	The <i>pl</i> field of the TLB_check parameter is valid.
mcc	59	Machine check corrected: This bit is set to one to indicate that the machine check has been corrected.
tv	60	Target address is valid: This bit is set to one to indicate that a valid target address has been logged.
rq	61	Requester identifier: This bit is set to one to indicate that a valid requester identifier has been logged.
rp	62	Responder identifier: This bit is set to one to indicate that a valid responder identifier has been logged.
pi	63	Precise instruction pointer. This bit is set to one to indicate that a valid precise instruction pointer has been logged.

- a. Hardware is tracking the operating status of the structure type and level reporting the error. The hardware reports a "normal" status when the number of entries within a structure reporting repeated corrections is at or below a pre-defined threshold. A "cautionary" status is reported when the number of affected entries exceeds a pre-defined threshold.

**Bus\_Check Return Format:** The bus\_check return format is returned in *error\_info* when the user requests information on any level of hierarchy of the processor bus structures as specified in the *level\_index* input argument. The bus\_check return format must be used to report errors in uncacheable transactions. These errors must not be reported using the cache\_check return format. The bus\_check return format is a bit-field that is described in Figure 11-22 and Table 11-92.

**Figure 11-22. bus\_check Layout**



**Table 11-92. bus\_check Fields**

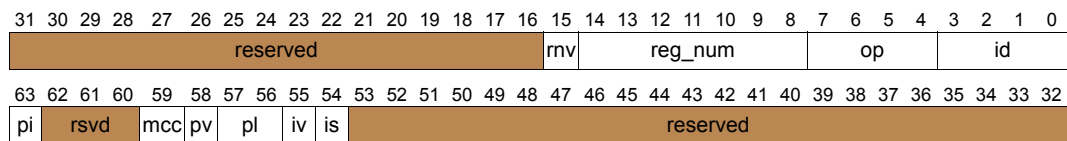
Field	Bits	Description
size	4:0	Size in bytes of the transaction that caused the machine check abort.
ib	5	Internal bus error
eb	6	External bus error
cc	7	Error occurred during a cache to cache transfer.
type	15:8	Type of transaction that caused the machine check abort. 0 – unknown 1 – partial read 2 – partial write 3 – full line read 4 – full line write 5 – implicit or explicit write-back operation 6 – snoop probe 7 – incoming or outgoing ptc.g 8 – write coalescing transactions 9 – I/O space read 10 – I/O space write 11 – inter-processor interrupt message (IPI) 12 – interrupt acknowledge or external task priority cycle All other values are reserved
sev	20:16	Bus error severity. The encodings of error severity are platform specific.
hier	22:21	This value indicates which level or bus hierarchy the error occurred in. A value of 0 indicates the first level of hierarchy.
dp	23	A multiple-bit error was detected, and data was poisoned for the incoming cache line.
bsi	31:24	Bus error status information. It describes the type of bus error. This field is processor bus specific.
reserved	53:32	Reserved
is	54	Instruction set. If this value is set to zero, the instruction that generated the machine check was an Intel Itanium instruction. If this bit is set to one, the instruction that generated the machine check was IA-32 instruction.
iv	55	The <i>is</i> field in the bus_check parameter is valid.
pl	57:56	Privilege level. The privilege level of the instruction bundle responsible for generating the machine check.
pv	58	The <i>pl</i> field of the bus_check parameter is valid.
mcc	59	Machine check corrected: This bit is set to one to indicate that the machine check has been corrected.
tv	60	Target address is valid: This bit is set to one to indicate that a valid target address has been logged.

**Table 11-92. bus\_check Fields (Continued)**

Field	Bits	Description
rq	61	Requester identifier: This bit is set to one to indicate that a valid requester identifier has been logged.
rp	62	Responder identifier: This bit is set to one to indicate that a valid responder identifier has been logged.
pi	63	Precise instruction pointer. This bit is set to one to indicate that a valid precise instruction pointer has been logged.

**Reg\_File\_Check Return Format:** The `reg_file_check` return format is returned in `error_info` when the user requests information on any of the registers as specified in the `level_index` input argument. The `reg_file_check` return format is a bit-field that is described in [Figure 11-23](#) and [Table 11-93](#). When the `reg_file_check` return format is returned, the target address, the requester identifier and the responder identifier will always be invalid.

**Figure 11-23. reg\_file\_check Layout**



**Table 11-93. reg\_file\_check Fields**

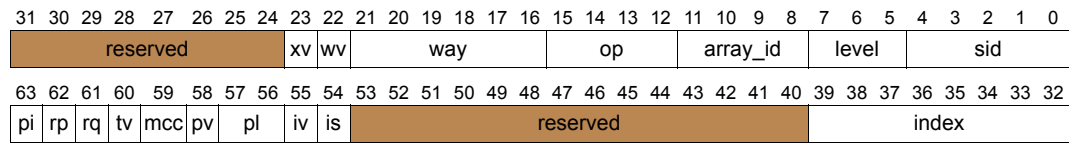
Field	Bits	Description
id	3:0	Register file identifier: 0 – unknown/unclassified 1 – General register (bank1) 2 – General register (bank 0) 3 – Floating-point register 4 – Branch register 5 – Predicate register 6 – Application register 7 – Control register 8 – Region register 9 – Protection key register 10 – Data breakpoint register 11 – Instruction breakpoint register 12 – Performance monitor control register 13 – Performance monitor data register All other values are reserved
op	7:4	Identifies the operation that caused the machine check 0 – unknown 1 – read 2 – write All other values are processor specific
reg_num	14:8	Identifies the register number that was responsible for generating the machine check
rnv	15	Specifies if the <code>reg_num</code> field is valid
reserved	53:16	Reserved
is	54	Instruction set. If this value is set to zero, the instruction that generated the machine check was an Intel Itanium instruction. If this bit is set to one, the instruction that generated the machine check was IA-32 instruction.
iv	55	The <code>is</code> field in the <code>reg_file_check</code> parameter is valid.

**Table 11-93. reg\_file\_check Fields**

Field	Bits	Description
pl	57:56	Privilege level. The privilege level of the instruction bundle responsible for generating the machine check.
pv	58	The <i>pl</i> field of the <i>reg_file_check</i> parameter is valid.
mcc	59	Machine check corrected: This bit is set to one to indicate that the machine check has been corrected.
reserved	62:60	Reserved
pi	63	Precise instruction pointer. This bit is set to one to indicate that a valid precise instruction pointer has been logged.

**Uarch\_Check Return Format:** The *uarch\_check* return format is returned in *error\_info* when the user requests information on any of the micro-architectural structures as specified in the *level\_index* input argument. The *uarch\_check* return format is a bit-field that is described in [Figure 11-24](#) and [Table 11-94](#).

**Figure 11-24. uarch\_check Layout**



**Table 11-94. uarch\_check Fields**

Field	Bits	Description
sid	4:0	Structure identification. These bits identify the micro-architectural structure where the error occurred. The definition of these bits are implementation specific.
level	7:5	Level of the micro-architectural structure where the error was generated. A value of 0 indicates the first level.
array_id	11:8	Identification of the array in the micro architectural structure where the error was generated. 0 – unknown/unclassified All other values are implementation specific
op	15:12	Type of operation that caused the error 0 – unknown 1 – read or load 2 – write or store All other values are implementation specific
way	21:16	Way of the micro-architectural structure where the error was located.
wv	22	The way field in the uarch_check parameter is valid.
xv	23	The index field in the uarch_check parameter is valid.
reserved	31:24	Reserved
index	39:32	Index or set of the micro-architectural structure where the error was located.
reserved	53:40	Reserved
is	54	Instruction set. If this value is set to zero, the instruction that generated the machine check was an Intel Itanium instruction. If this bit is set to one, the instruction that generated the machine check was IA-32 instruction.
iv	55	The is field in the bus_check parameter is valid.
pl	57:56	Privilege level. The privilege level of the instruction bundle responsible for generating the machine check.
pv	58	The pl field of the bus_check parameter is valid.
mcc	59	Machine check corrected: This bit is set to one to indicate that the machine check has been corrected.
tv	60	Target address is valid: This bit is set to one to indicate that a valid target address has been logged.
rq	61	Requester identifier: This bit is set to one to indicate that a valid requester identifier has been logged.
rp	62	Responder identifier: This bit is set to one to indicate that a valid responder identifier has been logged.
pi	63	Precise instruction pointer. This bit is set to one to indicate that a valid precise instruction pointer has been logged.



## PAL\_MC\_ERROR\_INJECT – Inject Processor Error (276)

**Purpose:** Injects the requested processor error or returns information on the supported injection capabilities for this particular processor implementation.

**Calling Conv:** Stacked

**Mode:** Physical and Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_ERROR_INJECT within the list of PAL procedures.
	err_type_info	Unsigned 64-bit integer specifying the first level error information which identifies the error structure and corresponding structure hierarchy, and the error severity.
	err_struct_info	Unsigned 64-bit integer identifying the optional structure specific information that provides the second level details for the requested error.
	err_data_buffer	Unsigned 64-bit integer specifying the address of the buffer providing additional parameters for the requested error. The address of this buffer must be 8-byte aligned.

Returns:	Return Value	Description
	status	Return status of the PAL_MC_ERROR_INJECT procedure.
	capabilities	64-bit vector specifying the supported error injection capabilities for the input argument combination of <i>struct_hier</i> , <i>err_struct</i> and <i>err_sev</i> fields in <i>err_type_info</i> .
	resources	64-bit vector specifying the architectural resources that are used by the procedure.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-4	Call completed with error; the requested error could not be injected due to failure in locating the target location in the specified structure.
	-5	Argument was valid, but requested error injection capability is not supported.
	-9	Call requires PAL memory buffer

**Description:** This procedure enables error injection into processor structures based on information specified by *err\_type\_info*, *err\_struct\_info* and *err\_data\_buffer*. Each invocation of the procedure enables a single error to be injected. The procedure supports error injection for at least one error of each severity type (correctable, recoverable, fatal).

The *err\_type\_info* argument specifies details of the error injection operation that is being requested (see [Figure 11-25](#)). The *err\_struct\_info* and *err\_data\_buffer* specify additional optional information. The format of *err\_struct\_info* is specified for each supported structure type indicated by the *err\_struct* field in *err\_type\_info*. *err\_data\_buffer* is optional, depending on the structure type and whether *trigger* functionality is used. If *err\_data\_buffer* is not required for the error injection, PAL will not attempt to access the memory location specified in this parameter.

**Figure 11-25. *err\_type\_info***

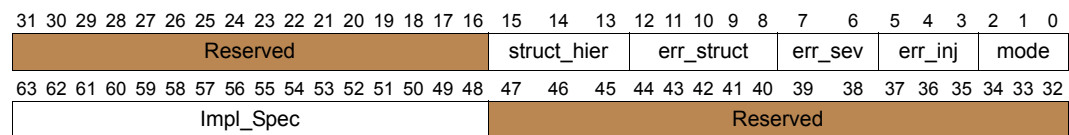


Table 11-95. *err\_type\_info*

Field	Bits	Description
mode	2:0	Indicates the mode of operation for this procedure: 0 – Query mode 1 – Error inject mode ( <i>err_inj</i> should also be specified) 2 – Cancel outstanding trigger. All other fields in <i>err_type_info</i> , <i>err_struct_info</i> and <i>err_data_buffer</i> are ignored. All other values are reserved.
err_inj	5:3	Indicates the mode of error injection: 0 – Error inject only (no error consumption) 1 – Error inject and consume All other values are reserved.
err_sev	7:6	Indicates the severity desired for error injection/query. Definitions of the different error severity types is given in Section 11.8, “PAL Glossary” on page 2:350. 0 – Corrected error 1 – Recoverable error 2 – Fatal error 3 – Reserved
err_struct	12:8	Indicates the structure identification for error injection/query: 0 - Any structure (cannot be used during <i>query mode</i> ). When selected, the structure type used for error injection is determined by PAL. 1 – Cache 2 – TLB 3 – Register file 4 – Bus/System interconnect 5-15 – Reserved 16-31 – Processor specific error injection capabilities. <i>err_data_buffer</i> is used to specify error types. Please refer to the processor specific documentation for additional details.
struct_hier	15:13	Indicates the structure hierarchy for error injection/query: 0 - Any level of hierarchy (cannot be used during <i>query mode</i> ). When selected, the structure hierarchy used for error injection is determined by PAL. 1 – Error structure hierarchy level-1 2 – Error structure hierarchy level-2 3 – Error structure hierarchy level-3 4 – Error structure hierarchy level-4 All other values are reserved.
Reserved	47:16	Reserved
Impl_Spec	63:48	Processor specific error injection capabilities. Please refer to processor specific documentation for additional details.

If *query mode* is selected through the mode bit in the *err\_type\_info* parameter, the return value in the *capabilities* vector indicates which error injection types are *individually* supported on the underlying implementation for the corresponding values of *err\_struct*, *struct\_hier* and *err\_sev* fields in *err\_type\_info*. The caller is expected to iterate through all combinations of *err\_inj*, *err\_sev*, *err\_struct*, and *struct\_hier* to determine the full extent of *individual* error injection types supported by the underlying implementation.

The *capabilities* vector does not indicate which combinations of error injection inputs from *err\_struct\_info* are supported by the implementation. For example, if an implementation supports *tag* error injection only for instruction caches and *data* error injection only for data caches, this cannot be determined by the *capabilities* vector. In this instance, the *capabilities* vector will report *i=1*, *d=1*, *tag=1*, *data=1*, indicating that the error injection is supported *individually* for instruction and data caches, and for *tag* and *data* fields, but not indicating which *combinations* of *i*, *d*, *tag*, and *data* are

supported for error injection. The caller is required to use the *query mode* with appropriate inputs in *err\_struct\_info* to determine which combinations of error injection types are supported. If a given combination is not supported, the procedure returns with status -5.

The procedure supports both an *Error inject* and *Error inject and consume* mode (selectable through the *err\_inj* field in *err\_type\_info*). In the former mode, the procedure performs the requested error injection in the specified structure, but does not perform any additional actions that can lead to consumption of the error and generation of the subsequent machine check. In *Error inject and consume* mode, the procedure will inject the error in the specified structure, and will perform additional operations to ensure that the error condition is encountered resulting in a machine check. Note that in this case, the machine check will be generated within the context of this procedure.

The procedure also provides the ability to set an error injection trigger. In this case, the error injection is delayed until the operation specified by the trigger is encountered and the executing context has the specified privilege level. In the absence of a trigger, the error injection is performed at the time of procedure execution. If an error injection trigger is specified, the mode field in *err\_type\_info* determines whether the error is injected, or injected and consumed when the trigger operation is encountered. There can be only one outstanding trigger programmed at a time. Subsequent procedure calls that use the trigger functionality will overwrite the previous trigger parameters. Once a trigger is programmed it remains active until either the trigger operation is encountered or software cancels the outstanding trigger via this call. Software can cancel outstanding triggers by specifying *Cancel outstanding trigger* via the mode bit in *err\_type\_info*. The *resources* value returned is all zeroes, indicating that the procedure is no longer using any architectural resources (specified in *resources*) for triggering purposes. When using this mode, it is possible that the procedure execution may itself satisfy the trigger conditions while in the process of cancelling the last programmed trigger.

To support triggers, PAL may use existing architectural resources. The *resources* return value defines the list of resources that are being used by PAL (see [Figure 11-26](#)).

In order for triggering to work when PAL is using the IBR or DBR registers, certain PSR bits are required to be set. Software needs to ensure that the PSR.db and the PSR.ic bits are set to one when executing the code that it is targeting with the trigger. If either one of these bits are not set, then triggers will not work as defined.

Procedure operation is undefined if software overwrites or modifies the IBR/DBR resources that PAL indicates it is using for a trigger. The IBR/DBR resources that PAL is not using are available for software to program for their own use.

**Figure 11-26. resources Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
Reserved																								dbr6	dbr4	dbr2	dbr0	ibr6	ibr4	ibr2	ibr0
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
Reserved																															

**Table 11-96. resources Return Value**

Field	Bits	Description
ibr0	0	When 1, indicates that IBR0,1 are being used by the procedure for trigger functionality.
ibr2	1	When 1, indicates that IBR2,3 are being used by the procedure for trigger functionality.
ibr4	2	When 1, indicates that IBR4,5 are being used by the procedure for trigger functionality.
ibr6	3	When 1, indicates that IBR6,7 are being used by the procedure for trigger functionality.
dbr0	4	When 1, indicates that DBR0,1 are being used by the procedure for trigger functionality.
dbr2	5	When 1, indicates that DBR2,3 are being used by the procedure for trigger functionality.
dbr4	6	When 1, indicates that DBR4,5 are being used by the procedure for trigger functionality.
dbr6	7	When 1, indicates that DBR6,7 are being used by the procedure for trigger functionality.

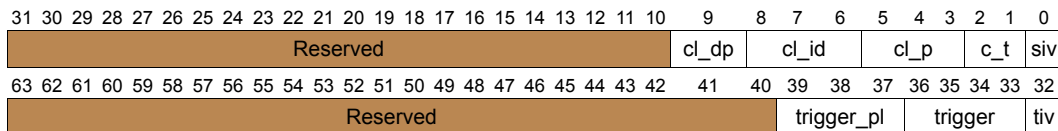
Multiprocessor coherency is not guaranteed when error injection is performed using this procedure. Please refer to the processor-specific documentation for further details regarding possible scenarios which can result in loss of coherency.

In cases where an error cannot be injected due to failure in locating the specified target location (cache line, TC, TR, register number) for the given set of input arguments, the procedure will return with status -4. For example, if the caller requests an error injection in the cache and specifies *cl\_id*=1 (virtual address provided), then PAL will attempt to locate the cache line as indicated by the input virtual address. If the corresponding cache line cannot be found (the cache line could have been evicted from the cache in the time interval between the procedure call and the search process, or the cache line may be in *invalid* state), then the procedure returns with a status value of -4.

The procedure does not check the settings of the error promotion bits (bit 53 and bit 60 in PAL\_PROC\_GET\_FEATURES) before injecting an error in the specified structure. Based on the configuration of these bits, the severity of the error reported may vary.

The detailed descriptions of *err\_struct\_info* and *err\_data\_buffer* are shown below.

**Figure 11-27. err\_struct\_info – Cache**



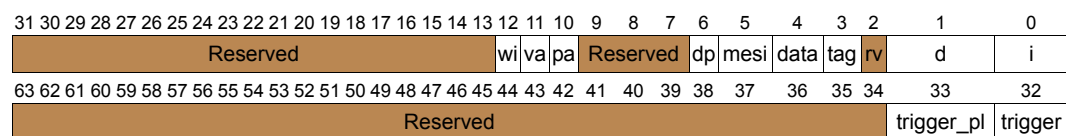
**Table 11-97. err\_struct\_info – Cache**

Field	Bits	Description
siv	0	When 1, indicates that the structure information fields ( <i>c_t</i> , <i>cl_p</i> , <i>cl_id</i> ) are valid and should be used for error injection. When 0, the structure information fields are ignored, and the values of these fields used for error injection are implementation-specific.
c_t	2:1	Indicates which cache should be used for error injection: 0 – Reserved 1 – Instruction cache 2 – Data or unified cache 3 – Reserved
cl_p	5:3	Indicates the portion of the cache line where the error should be injected: 0 – Reserved 1 – Tag 2 – Data 3 – mesi All other values are reserved.

**Table 11-97. *err\_struct\_info* – Cache (Continued)**

Field	Bits	Description
cl_id	8:6	Indicates which mechanism is used to identify the cache line to be used for error injection: 0 – Reserved 1 – Virtual address provided in the <i>inj_addr</i> field of the buffer pointed to by <i>err_data_buffer</i> should be used to identify the cache line for error injection. 2 – Physical address provided in the <i>inj_addr</i> field of the buffer pointed to by <i>err_data_buffers</i> should be used to identify the cache line for error injection. 3 – <i>way</i> and <i>index</i> fields provided in <i>err_data_buffer</i> should be used to identify the cache line for error injection. All other values are reserved.
cl_dp	9	When 1, indicates that a multiple bit, non-correctable error should be injected in the cache line specified by <i>cl_id</i> . If this injected error is not consumed, it may eventually cause a data-poisoning event resulting in a corrected error signal, when the associated cache line is cast out (implicit or explicit write-back of the cache line). The error severity specified by <i>err_sev</i> in <i>err_type_info</i> must be set to 0 ( <i>corrected error</i> ) when this bit is set.
Reserved	31:10	Reserved
tiv	32	When 1, indicates that the trigger information fields ( <i>trigger</i> , <i>trigger_pl</i> ) are valid and should be used for error injection. When 0, the trigger information fields are ignored and error injection is performed immediately.
trigger	36:33	Indicates the operation type to be used as the error trigger condition. The address corresponding to the trigger is specified in the <i>trigger_addr</i> field of the buffer pointed to by <i>err_data_buffer</i> : 0 – Instruction memory access. The trigger match conditions for this operation type are similar to the IBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions”</a> on page 2:154. 1 – Data memory access. The trigger match conditions for this operation type are similar to the DBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions”</a> on page 2:154. All other values are reserved.
trigger_pl	39:37	Indicates the privilege level of the context during which the error should be injected: 0 – privilege level 0 1 – privilege level 1 2 – privilege level 2 3 – privilege level 3 All other values are reserved. If the implementation does not support privilege level qualifier for triggers (i.e. if <i>trigger_pl</i> is 0 in the <i>capabilities</i> vector), this field is ignored and triggers can be taken at any privilege level.
Reserved	63:40	Reserved

**Figure 11-28. *capabilities* vector for cache**



**Table 11-98. *capabilities* vector for cache**

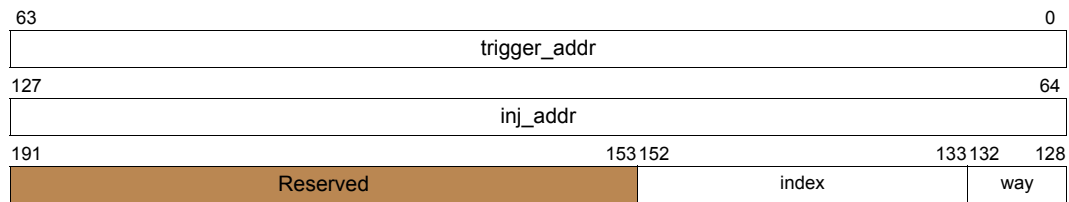
Field	Bits	Description
i	0	Error injection for instruction caches is supported
d	1	Error injection for data caches is supported
rv	2	Reserved

**Table 11-98. capabilities vector for cache (Continued)**

Field	Bits	Description
tag	3	Error injection in <i>tag</i> portion of cache line is supported
data	4	Error injection in <i>data</i> portion of cache line is supported
mesi	5	Error injection in <i>mesi</i> portion of cache line is supported
dp	6	Error injection that results in data poisoning events is supported
Reserved	9:6	Reserved
pa	10	Error injection with physical address input is supported
va	11	Error injection with virtual address input is supported
wi	12	Error injection with <i>way</i> and <i>index</i> input is supported
Reserved	31:13	Reserved
trigger	32	Error injection with trigger is supported
trigger_pl	33	Error injection with privilege level qualifier for trigger is supported
Reserved	63:34	Reserved

*err\_data\_buffer* needs to be specified for *cache* only if *siv* is 1 or *tiv* is 1, in *err\_struct\_info*.

**Figure 11-29. Buffer pointed to by *err\_data\_buffer* – Cache**



**Table 11-99. Buffer pointed to by *err\_data\_buffer* – Cache**

Field	Bits	Description
trigger_addr	63:0	64-bit virtual address to be used by the <i>trigger</i> in the <i>err_struct_info</i> input argument. This field is ignored if <i>tiv</i> in <i>err_struct_info</i> is 0. The field is defined similar to the <i>addr</i> field in the debug breakpoint registers, as specified in <a href="#">Table 7-1, “Debug Breakpoint Register Fields (DBR/IBR)” on page 2:153</a> .
inj_addr	127:64	64-bit virtual or physical address used to identify the cache line to be used for error injection. This field is valid only if <i>cl_id</i> in <i>err_struct_info</i> corresponds to either <i>va</i> or <i>pa</i> (value 1 or 2).
way	132:128	Indicates the <i>way</i> information for error injection. This is used in combination with the <i>index</i> field to identify the cache line for error injection. This field is valid only if <i>cl_id</i> in <i>err_struct_info</i> is 3, else it is ignored.
index	152:133	Indicates the <i>index</i> information for error injection. This is used in combination with the <i>way</i> field to identify the cache line for error injection. This field is valid only if <i>cl_id</i> in <i>err_struct_info</i> is 3, else it is ignored.
Reserved	191:153	Reserved

Figure 11-30. *err\_struct\_info* – TLB

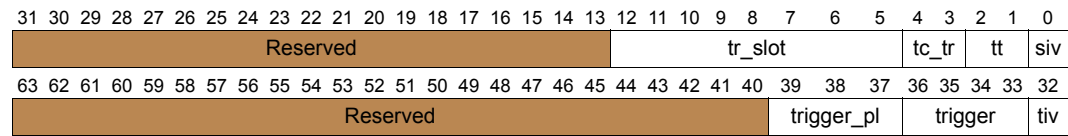
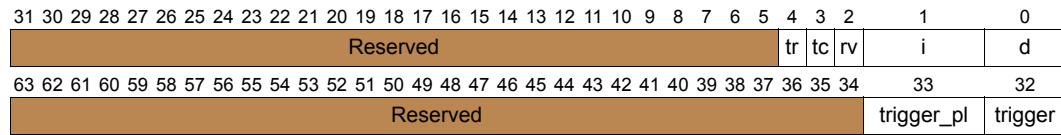


Table 11-100. *err\_struct\_info* – TLB

Field	Bits	Description
siv	0	When 1, indicates that the structure information fields ( <i>tt</i> , <i>tc_tr</i> , <i>tr_slot</i> ) are valid and should be used for error injection. When 0, the structure information fields are ignored, and the values of these fields used for error injection are implementation-specific.
tt	2:1	Indicates which TLB should be used for error injection: 0 – Reserved 1 – Instruction TLB 2 – Data TLB 3 – Reserved
tc_tr	4:3	Indicates which portion of TLB should be used for error injection: 0 – Reserved 1 – tc: error should be injected in a Translation Cache (TC) entry. For TC insertion, the entry is identified by the <i>vpn</i> and <i>rid</i> fields in <i>err_data_buffer</i> 2 – tr: error should be injected in a Translation Register (TR) entry. For TR insertion, the slot number is specified by the <i>tr_slot</i> field. 3 – Reserved
tr_slot	12:5	Indicates the Translation Register (TR) slot number where the error should be injected. This field is valid only when <i>tc_tr</i> is 2, else it is ignored.
Reserved	31:13	Reserved
tiv	32	When 1, indicates that the trigger information fields ( <i>trigger</i> , <i>trigger_pl</i> ) are valid and should be used for error injection. When 0, the trigger information fields are ignored and error injection is performed immediately.
trigger	36:33	Indicates the operation type to be used as the error trigger condition. The virtual address corresponding to the trigger is specified in the <i>trigger_addr</i> field of the buffer pointed to by <i>err_data_buffer</i> . 0 – Instruction memory access. The trigger match conditions for this operation type are similar to the IBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions”</a> on page 2:154. 1 – Data memory access. The trigger match conditions for this operation type are similar to the DBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions”</a> on page 2:154. All other values are reserved.
trigger_pl	39:37	Indicates the privilege level of the context during which the error should be injected 0 – privilege level 0 1 – privilege level 1 2 – privilege level 2 3 – privilege level 3 All other values are reserved. If the implementation does not support privilege level qualifier for triggers (i.e. if <i>trigger_pl</i> is 0 in the <i>capabilities</i> vector), this field is ignored and triggers can be taken at any privilege level.
Reserved	63:40	Reserved

**Figure 11-31. capabilities vector for TLB**

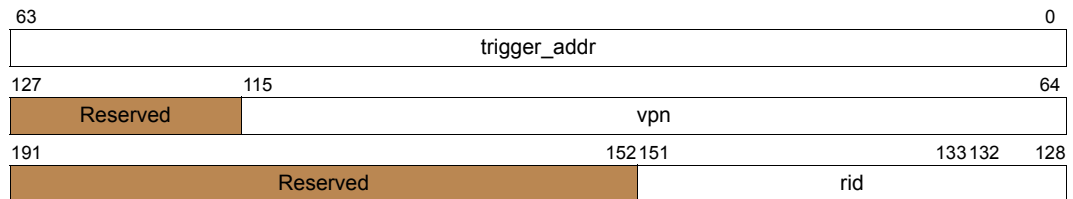


**Table 11-101. capabilities vector for TLB**

Field	Bits	Description
d	0	Error injection for data TLB is supported
i	1	Error injection for instruction TLB is supported
rv	2	Reserved
tc	3	Error injection in TC entries is supported
tr	4	Error injection in TR entries is supported
Reserved	31:5	Reserved
trigger	32	Error injection with trigger is supported
trigger_pl	33	Error injection with privilege level qualifier for trigger is supported
Reserved	63:34	Reserved

*err\_data\_buffer* needs to be specified for TLB only if *tiv* is 1 or if *tc\_tr* value corresponds to *tc*, in *err\_struct\_info*.

**Figure 11-32. Buffer pointed to by *err\_data\_buffer* – TLB**



**Table 11-102. Buffer pointed to by *err\_data\_buffer* – TLB**

Field	Bits	Description
trigger_addr	63:0	64-bit virtual address to be used by the <i>trigger</i> in the <i>err_struct_info</i> input argument. The field is defined similar to the <i>addr</i> field in debug breakpoint registers, as specified in Table 7-1, “Debug Breakpoint Register Fields (DBR/IBR)” on page 2:153.
vpn	115:64	Indicates the Virtual page number. This field is valid only when <i>tc_tr</i> in <i>err_struct_info</i> is 1. <i>vpn</i> used in combination with <i>rid</i> to identify the TC entry for error injection.
Reserved	127:116	Reserved
rid	151:128	Indicates the region identifier. This field is valid only when <i>tc_tr</i> in <i>err_struct_info</i> is 1. <i>rid</i> is used in combination with <i>vpn</i> to identify the TC entry for error injection.
Reserved	191:152	Reserved

**Figure 11-33. *err\_struct\_info* – Register File**

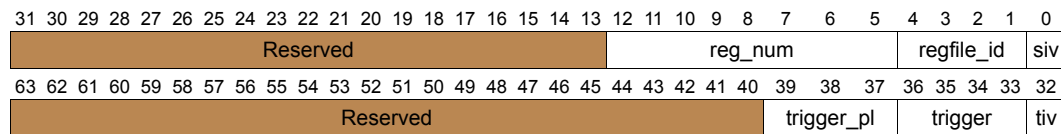
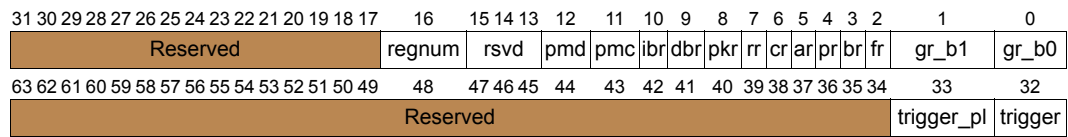




Table 11-103. *err\_struct\_info* – Register File

Field	Bits	Description
siv	0	When 1, indicates that the structure information fields ( <i>regfile_id</i> , <i>reg_num</i> ) are valid and should be used for error injection. When 0, the structure information fields are ignored, and the values of these fields used for error injection are implementation-specific.
regfile_id	4:1	Identifies the register file where the error should be injected: 0 – Any register file type. When selected, the register file used for error injection is determined by PAL. 1 – General register (bank0)(GR16-31) 2 – General register (bank1)(GR0-127) 3 – Floating point register 4 – Branch register 5 – Predicate register 6 – Application register 7 – Control register 8 – Region register 9 – Protection key register 10 – Data breakpoint register 11 – Instruction breakpoint register 12 – Performance monitor control register 13 – Performance monitor data register All other values are reserved.
reg_num	12:5	Indicates the register number where the error should be injected. Procedure operation is undefined if there is a conflict between the register number chosen for error injection, and the registers being used by the procedure for code execution (see PAL calling conventions, Section 11.9.2). 0-127: Specific register number corresponding to <i>regfile_id</i> 128-254: Reserved for future use 255: Any register number. When selected, the actual register number used for error injection is determined by PAL.
Reserved	31:13	Reserved
tiv	32	When 1, indicates that the trigger information fields ( <i>trigger</i> , <i>trigger_pl</i> ) are valid and should be used for error injection. When 0, the trigger information fields are ignored and error injection is performed immediately.
trigger	36:33	Indicates the operation type to be used as the error trigger condition. The address corresponding to the trigger is specified in the <i>trigger_addr</i> field of the buffer pointed to by <i>err_data_buffer</i> . 0 – Instruction memory access. The trigger match conditions for this operation type are similar to the IBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions” on page 2:154</a> 1 – Data memory access. The trigger match conditions for this operation type are similar to the DBR address breakpoint match conditions as outlined in <a href="#">Section 7.1.2, “Debug Address Breakpoint Match Conditions” on page 2:154</a> . All other values are reserved.
trigger_pl	39:37	Indicates the privilege level of the context during which the error should be injected: 0 – privilege level 0 1 – privilege level 1 2 – privilege level 2 3 – privilege level 3 All other values are reserved. If the implementation does not support privilege level qualifier for triggers (i.e. if <i>trigger_pl</i> is 0 in the <i>capabilities</i> vector), this field is ignored and triggers can be taken at any privilege level.
Reserved	63:40	Reserved

**Figure 11-34. capabilities Vector for Register File**

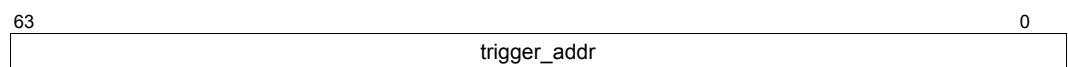


**Table 11-104. capabilities Vector for Register File**

Field	Bits	Description
gr_b0	0	Error injection for General register (bank0) is supported
gr_b1	1	Error injection for General register (bank1) is supported
fr	2	Error injection for Floating point register is supported
br	3	Error injection for Branch register is supported
pr	4	Error injection for Predicate register is supported
ar	5	Error injection for Application register is supported
cr	6	Error injection for Control register is supported
rr	7	Error injection for Region register is supported
pkrr	8	Error injection for Protection key register is supported
dbr	9	Error injection for Data breakpoint register is supported
ibr	10	Error injection for Instruction breakpoint register is supported
pmc	11	Error injection for Performance monitor control register is supported
pmd	12	Error injection for Performance monitor data register is supported
Reserved	15:13	Reserved
regnum	16	Error injection with register number input is supported
Reserved	31:17	Reserved
trigger	32	Error injection with trigger is supported
trigger_pl	33	Error injection with privilege level qualifier for trigger is supported
Reserved	63:34	Reserved

*err\_data\_buffer* needs to be specified for *register file* only if *tiv* in *err\_struct\_info* is 1.

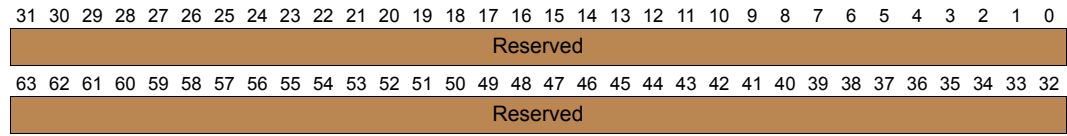
**Figure 11-35. Buffer pointed to by *err\_data\_buffer* – Register File**



**Table 11-105. Buffer pointed to by *err\_data\_buffer* – Register File**

Field	Bits	Description
trigger_addr	63:0	64-bit address to be used by the <i>trigger</i> in the <i>err_struct_info</i> input argument. The field is defined similar to the <i>addr</i> field in the debug breakpoint registers, as specified in Table 7-1, “Debug Breakpoint Register Fields (DBR/IBR)” on page 2:153.

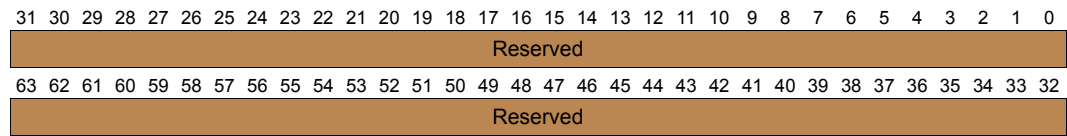
**Figure 11-36. *err\_struct\_info* – Bus/Processor Interconnect**



**Table 11-106. *err\_struct\_info* – Bus/Processor Interconnect**

Field	Bits	Description
Reserved	63:0	Reserved

**Figure 11-37. *capabilities* vector for Bus/Processor Interconnect**



**Table 11-107. *capabilities* vector for Bus/Processor Interconnect**

Field	Bits	Description
Reserved	63:0	Reserved

*err\_data\_buffer* does not need to be specified for *bus/system interconnect*.

## PAL\_MC\_HW\_TRACKING – Query which hardware structures are performing hardware status tracking (51)

**Purpose:** Provide a way to query which hardware structures are performing hardware status tracking for corrected machine check events.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_MC_HW_TRACKING within the list of PAL procedures.
Reserved	0
Reserved	0
Reserved	0

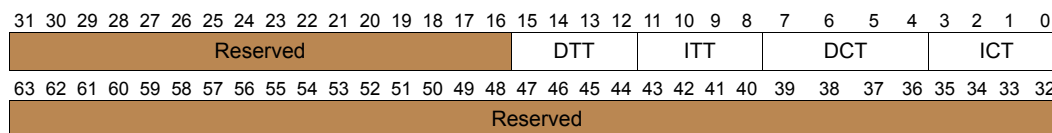
Return Value	Description
status	Return status of the PAL_MC_HW_TRACKING procedure.
hw_track	64-bit vector denoting which hardware structures are providing hardware status tracking. See <a href="#">Figure 11-38</a> .
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

**Description:** This procedure will return information about which hardware structures are providing hardware status tracking for corrected machine check events. This information is also returned in the error logs for corrected machine check events.

The layout of the tracked return value is shown in [Figure 11-38](#).

**Figure 11-38. Layout of hw\_track Return Value**



**Table 11-108. hw\_check Fields**

Field	Bits	Description
ICT	3:0	Instruction cache tracking. This is a 4-bit vector denoting which levels of instruction cache provide hardware tracking.
DCT	7:4	Data cache tracking. This is a 4-bit vector denoting which levels of data/unified caches provide hardware tracking.
ITT	11:8	Instruction TLB tracking. This is a 4-bit vector denoting which levels of the instruction TLB provide hardware tracking.
DTT	15:12	Data TLB tracking. This is a 4-bit vector denoting which levels of data/unified TLB provide hardware tracking.
Reserved	63:16	Reserved

The convention for the levels in the *hw\_track* field is such that the least significant bit in the field represents the lowest level of the structures hierarchy. For example, bit 0 of the ICT field represents the first level instruction cache.

## PAL\_MC\_EXPECTED – Set/Reset Expected Machine Check Indicator (23)

**Purpose:** Informs PALE\_CHECK whether a machine check is expected so that PALE\_CHECK will not attempt to correct any expected machine checks.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_EXPECTED within the list of PAL procedures.
	expected	Unsigned integer with a value of 0 or 1 to set or reset the hardware resource PALE_CHECK examines for expected machine checks.
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_MC_EXPECTED procedure.
	previous	Unsigned integer denoting whether a machine check was previously expected.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** If the argument *expected* contains a value of 1, an implementation-dependent hardware resource is set to inform PALE\_CHECK to expect a machine check. If the argument *expected* is 0, the resource is reset, so that PALE\_CHECK does not expect any following machine checks. All other values of *expected* are reserved.

The implementation-dependent hardware resource should be, by default, in the “not expected” state. Software or firmware should only call PAL\_MC\_EXPECTED immediately prior to issuing an instruction which might generated an expected machine check. It should then immediately reset the bit to the “not expected” state after checking the results of the operation.

The *previous* return parameter indicates the previous state of the hardware resource to inform PALE\_CHECK of an expected machine check. A value of 0 indicates that a machine check was not expected. A value of 1 indicated that a machine check was expected. All other values of *previous* are reserved.

## PAL\_MC\_REGISTER\_MEM – Register Memory with PAL for Machine Check and Init (27)

**Purpose:** Registers a platform dependent location with PAL to which it can save minimal processor state in the event of a machine check or initialization event.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_MC_REGISTER_MEM within the list of PAL procedures.
	address	Physical address of the buffer to be registered with PAL.
	size	Unsigned integer indicating the size in kilobytes (KB) of the buffer passed. This input argument is only required when passing in a size greater than 4KB. The implementation indicates when a size greater than 4KB is required at the reset hand-off. Refer to <a href="#">Section 11.2.2.1, "Definition of SALE_ENTRY State Parameter" on page 2:291</a> for more information.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_MC_REGISTER_MEM procedure.
	req_size	Returns the required size of the min-state save area in kilobytes (KB) if the <i>size</i> input argument did not match the required size for this implementation.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This procedure is used to register with PAL an uncacheable min-state save area memory buffer that is used for machine check and initialization event handling. The size of the min-state save area is either 4KB or a larger size that is indicated in the reset hand-off state described in [Section 11.2.2.1, "Definition of SALE\\_ENTRY State Parameter" on page 2:291](#). The input argument *size* indicates the size of the min-state save buffer in kilobytes (KB) when it is greater than 4KB. If the *size* input argument does not match the required size, the procedure returns an invalid argument return status and a min-state area is not registered. The procedure will also return the required size of the min-state save area in the *req\_size* return value.

The layout of the min-state save area is defined in [Section 11.3.2.4, "Processor Min-state Save Area Layout" on page 2:302](#). The address passed has a minimum alignment requirement of 512-bytes.

## PAL\_MC\_RESUME – Restore Minimal Architected State and Return (26)

**Purpose:** Restores the minimal architectural processor state, sets the CMC interrupt if necessary, and resumes execution.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_MC_RESUME within the list of PAL procedures.
set_cmci	Unsigned 64 bit integer denoting whether to set the CMC interrupt. A value of 0 indicates not to set the interrupt, a value of 1 indicated to set the interrupt, and all other values are reserved.
save_ptr	Physical address of min-state save area used to used to restore processor state.
new_context	Unsigned 64-bit integer denoting whether the caller is returning to a new context. A value of 0 indicates the caller is returning to the interrupted context, a value of 1 indicates that the caller is returning to a new context.

Return Value	Description
status	Return status of the PAL_MC_RESUME procedure <sup>a</sup> .
Reserved	0
Reserved	0
Reserved	0

a. This procedure returns to the caller only in an error situation.

Status Value	Description
-2	Invalid argument
-3	Call completed with error

**Description:** This procedure will restore the processor minimal architected state and optionally set the CMC interrupt.

If the *set\_cmci* argument is set to one, this procedure will set the CMC interrupt and return to the interrupted context. The CMC interrupt handler will be invoked sometime after returning to the interrupted context.

The *save\_ptr* argument specifies the processor min-state save area buffer from which the processor state will be restored. This pointer has the same alignment and size restrictions as the address passed to PAL\_MC\_REGISTER\_MEM procedure on [page 2:435](#).

This procedure is used to resume execution of the interrupted context for both machine check and initialization events. This procedure can resume execution to the same context or a new context. If software attempts to resume execution for these events without using this call, processor behavior is undefined.

If the caller is resuming to the same context, the *new\_context* argument must be set to 0 and the *save\_ptr* argument has to point to a copy of the min-state save area written by PAL when the event occurred.

If the caller is resuming to a new context, the *new\_context* argument must be set to 1 and the *save\_ptr* argument must point to a new min-state save area set up by the caller.

Please see Section 11.3.3, “Returning to the Interrupted Process” on [page 2:305](#) 3for more information on resuming to the interrupted context.



## PAL\_MEM\_ATTRIB – Get Memory Attributes (5)

**Purpose:** Returns the memory attributes implemented by processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

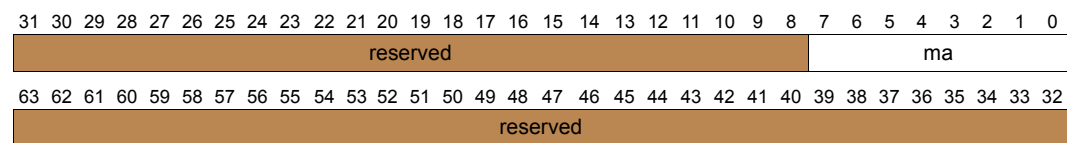
Argument	Description
index	Index of PAL_MEM_ATTRIB within the list of PAL procedures.
Reserved	0
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_MEM_ATTRIB procedure.
attrib	8-bit vector of memory attributes implemented by processor.
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** Returns a 8-bit vector in the low order 8 bits of the return register that specifies the set of memory attributes implemented by the processor. The return register is formatted as follows:

**Figure 11-39. Layout of *attrib* Return Value**



Each bit in the bit field *ma* represents one of the eight possible memory attributes implemented by the processor. The bit field position corresponds to the numeric memory attribute encoding defined in [Section 4.4, “Memory Attributes” on page 2:75](#).

## PAL\_MEMORY\_BUFFER – Allocate a cacheable memory buffer for exclusive PAL usage (277)

**Purpose:** Provides cacheable memory to PAL for exclusive use during runtime.

**Calling Conv:** Stacked

**Mode:** Physical

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_MEMORY_BUFFER within the list of PAL procedures.
base_address	Physical address of the memory buffer allocated for PAL use.
alloc_size	Unsigned integer denoting the size of the memory buffer.
control_word	Formatted bit vector that provides control options for this procedure. See <a href="#">Table 11-109</a> .

Return Value	Description
status	Return status of the PAL_MEMORY_BUFFER procedure.
min_size	Returns the minimum size of the memory buffer required if the <i>alloc_size</i> input argument was not large enough.
Reserved	0
Reserved	0

Status Value	Description
1	Call has not completed a buffer relocation due to a pending interrupt
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

**Description:** This procedure is used to provide PAL firmware a cacheable memory buffer for its exclusive use as well as the ability to relocate this buffer at a later point in time if necessary. PAL provides information at reset hand-off about the minimum buffer size required by this procedure, and also indicates if this procedure is required to be called for correct functionality of the processor. See [Section 11.2.2, “PALE\\_RESET Exit State” on page 2:289](#) for additional information on the reset hand-off state.

The *base\_address* input argument specifies the beginning address for the memory buffer. The *alloc\_size* input argument specifies the size of the memory buffer allocated for PAL use. The minimum alignment requirement for this buffer is 4K. If the *base\_address* is not at least 4K aligned, the procedure will return an invalid argument. If the *alloc\_size* input argument is smaller than the minimum size passed at PAL reset handoff state, the procedure will return an invalid argument and provide the minimum size required in the *min\_size* return argument.

The *control\_word* input argument specifies if this procedure is being used to register the memory buffer or if it is being used to relocate the memory buffer. The format of the *control\_word* is shown in [Table 11-109](#).

**Table 11-109. control\_word Layout**

Field	Bits	Description
reg	0	Value of 0 indicates registration for the first time of the buffer. A value of 1 indicates a relocation of the buffer.
int	1	Value of 1 indicates that the procedure should periodically poll for pending external interrupts. If this value is 0, interrupts will be masked during the execution of the entire procedure.
Reserved	63:2	Reserved

A memory buffer must be allocated for each physical package, and is shared by all logical processors on that package. Software is required to call this procedure on all logical processors on a given package with the same input values. If not, processor operation is undefined.

If the PAL reset hand-off state indicates that the memory buffer is required but no call is made to allocate the memory buffer for a given physical package before calling buffer-dependent PAL procedures on a logical processor on that package, those procedures return an error.

If software would like to relocate this memory buffer at a later point in time, it can do so by setting the value of *reg* field in *control\_word* to one. PAL will copy the contents of the existing buffer to a new buffer. Software is still required to make this call on all logical processors with the same input arguments when relocating the buffer. Once the call has been made on all logical processors in the physical package, the old memory can be reclaimed.

Software can choose if it wants this procedure to periodically poll for interrupts during the execution of the procedure. If an interrupt is seen, the procedure will return a value of 1 and software must re-call this procedure again on the same logical processor, with the same input arguments, until the copy is completed. If this procedure returns with a value of 1, both the old memory buffer and the new memory buffer will be in use by PAL until PAL returns that the procedure has completed execution successfully by setting the return value to 0.

An error will be returned if software calls this procedure with the *reg* value set to one to re-register a buffer and a call has never been made to register the buffer.

It is required that PAL firmware only perform cacheable memory accesses to this buffer.

## PAL\_PERF\_MON\_INFO – Get Processor Performance Monitor Information (15)

**Purpose:** Returns Performance Monitor information about what can be counted and how to configure the monitors to count the desired events.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

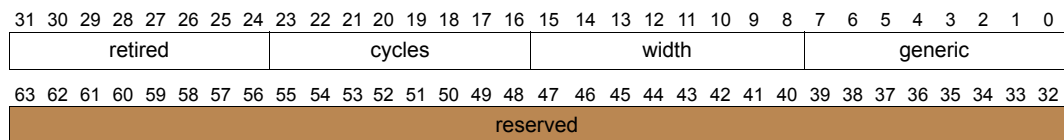
Argument	Description
index	Index of PAL_PERF_MON_INFO within the list of PAL procedures.
pm_buffer	An address to an 8-byte aligned 128-byte memory buffer.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_PERF_MON_INFO procedure.
pm_info	Information about the performance monitors implemented.
Reserved	0
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** PAL\_PERF\_MON\_INFO is called to determine the number of performance monitors and the events which can be counted on the performance monitors. For more information on performance monitoring, see [Section 7.2, “Performance Monitoring” on page 2:155](#). *pm\_info* is a formatted 64-bit return register, as shown in [Figure 11-40](#).

**Figure 11-40. Layout of *pm\_info* Return Value**



**Table 11-110. *pm\_info* Fields**

Field	Description
generic	Unsigned 8-bit number defining the number of generic PMC/PMD pairs.
width	Unsigned 8-bit number in the range 0:60 defining the number of implemented counter bits.
cycles	Unsigned 8-bit number defining the event type for counting processor cycles.
retired	Unsigned 8-bit number defining the event type for retired instruction bundles.

The *pm\_buffer* argument points to a 128-byte memory area where mask information is returned. The layout of *pm\_buffer* is shown in [Table 11-111](#).

**Table 11-111. *pm\_buffer* Layout**

Offset	Description
0x0	256-bit mask defining which PMC registers are implemented.
0x20	256-bit mask defining which PMD registers are implemented.

**Table 11-111. pm\_buffer Layout (Continued)**

Offset	Description
0x40	256-bit mask defining which registers can count cycles.
0x60	256-bit mask defining which registers can count retired bundles.

## PAL\_PLATFORM\_ADDR – Set Processor Interrupt Block Address and I/O Port Space Address (16)

**Purpose:** Specifies the physical address of the processor Interrupt Block and I/O Port Space.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PLATFORM_ADDR within the list of PAL procedures.
	type	Unsigned 64-bit integer specifying the type of block. 0 indicates that the processor interrupt block pointer should be initialized. 1 indicates that the processor I/O block pointer should be initialized.
	address	Unsigned 64-bit integer specifying the address to which the processor I/O block or interrupt block shall be set. The address must specify an implemented physical address on the processor model, bit 63 is ignored.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PLATFORM_ADDR procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error

**Description:** PAL\_PLATFORM\_ADDR specifies the physical address that the processor shall interpret as accesses to the SAPIC memory or the I/O Port space areas.

The default value for the Interrupt block pointer is 0x00000000 FEE00000. If an alternate address is selected by this call, it must be aligned on a 2 MB boundary, else the procedure will return an error status. The address specified must also not overlay any firmware addresses in the 16 MB region immediately below the 4GB physical address boundary.

The default value for the I/O block pointer is to the beginning of the 64 MB block at the highest physical address supported by the processor. Therefore, its physical address is implementation dependent. If an alternate address is selected by this call, it must be aligned on a 64MB boundary, else the procedure will return an error status. The address specified must also not overlay any firmware addresses in the 16 MB region immediately below the 4GB physical address boundary.

The Interrupt and I/O Block pointers should be initialized by firmware before any Inter-Processor Interrupt messages or I/O Port accesses. Otherwise the default block pointer values will be used.

Some processor implementations may not support relocation of the interrupt and I/O block pointers and an unimplemented procedure return status will be returned. In these cases the default address spaces will be used.

## PAL\_PMI\_ENTRYPOINT – Setup SAL PMI Entrypoint in Memory (32)

**Purpose:** Sets the SAL PMI entrypoint in memory.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PMI_ENTRYPOINT within the list of PAL procedures.
	SAL_PMI_entry	256-byte aligned physical address of SAL PMI entrypoint in memory.
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PMI_ENTRYPOINT procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** This procedure is called to set the SAL PMI entrypoint so that the SAL PMI code shall be executed out of main memory instead of the firmware address space. Some processor implementations will allow initialization of the PMI entrypoint only once. Under those situations, this procedure may be called only once after a boot to initialize the PMI entrypoint register. Subsequent calls will return a status of -3. This call must be made before PMI is enabled by SAL.

## PAL\_PREFETCH\_VISIBILITY – Make Processor Prefetches Visible (41)

**Purpose:** Used in the architected sequences for memory attribute transitions described in [Section 4.4.11, “Memory Attribute Transition” on page 2:88](#) to transition a page (or set of pages) from a one memory attribute to another.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PREFETCH_VISIBILITY within the list of PAL procedures.
	trans_type	Unsigned integer specifying the type of memory attribute transition that is being performed
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PREFETCH_VISIBILITY procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	1	Call completed without error; this call is not necessary on remote processors
	0	Call completed without error; this call must also be performed on all remote processors in the coherence domain
	-2	Invalid argument
	-3	Call completed with error

**Description:** This call is intended to be used only in the architected sequences described in [Section 4.4.11, “Memory Attribute Transition” on page 2:88](#).

The *trans\_type* input indicates the type of memory attribute transition the user is making. An input value of 0 is used when transition virtual memory attributes only. A value of 1 is used when transitioning physical memory attributes only, or when transitioning memory that may have a combination of virtual and physical memory attributes. All other values are reserved.

This procedure, when used for transitioning virtual memory attributes, will ensure that all prefetches that were initiated by the processor to the cacheable, speculative memory prior to the call, will either not be cached; have been aborted; or are visible to subsequent  $\epsilon C$  instructions. (from both the local processor and from remote processors).

This procedure when used for transitioning physical memory attributes will ensure that all prefetches that were initiated by the processor to the cacheable, limited speculative memory prior to the call, will either not be cached; have been aborted; or are visible to subsequent  $\epsilon C$  instructions (from both the local processor and from remote processors). It will also terminate the ability for the processor to make speculative references to any limited speculation pages. For the processor to make any speculative reference to a limited speculation page after this call, there must be a verified reference made to that page after this call. See the discussion on limited speculation in [Section 4.4.6.1, “Limited Speculation and the WBL Physical Addressing Attribute” on page 2:81](#).



This procedure, when used to delete a memory range on-line, will ensure that all of the conditions described in both of the preceding paragraphs regarding transition of virtual memory attributes and physical memory attributes are met.

If the processor implementation does not require this procedure call to be made on remote processors in the sequences, this procedure will return a 1 upon successful completion.

A return value of 0 upon successful completion of this procedure is an indication to software that the processor implementation requires that this call be performed on all processors in the coherence domain to make prefetches visible in the sequences.

These return code can be used to tune the architected sequence to the particular system on which is running; see [Section 4.4.11, "Memory Attribute Transition"](#) for details.

## PAL\_PROC\_GET\_FEATURES – Get Processor Dependent Features (17)

**Purpose:** Provides information about configurable processor features.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PROC_GET_FEATURES within the list of PAL procedures.
	Reserved	0
	feature_set	Feature set information is being requested for.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PROC_GET_FEATURES procedure.
	features_avail	64-bit vector of features implemented. See <a href="#">Table 11-112</a> .
	feature_status	64-bit vector of current feature settings. See <a href="#">Table 11-112</a> .
	feature_control	64-bit vector of features controllable by software.

Status:	Status Value	Description
	1	Call completed without error; The <i>feature_set</i> passed is not supported but a <i>feature_set</i> of a larger value is supported
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error
	-8	<i>feature_set</i> passed is beyond the maximum <i>feature_set</i> supported

**Description:** PAL\_PROC\_GET\_FEATURES and PAL\_PROC\_SET\_FEATURES procedure calls are used together to describe current settings of processor features and to allow modification of some of these processor features.

The *feature\_set* input argument for PAL\_PROC\_GET\_FEATURES describes which processor *feature\_set* information is being requested. [Table 11-112](#) describes processor *feature\_set* zero. The *feature\_set* values are split into two categories: architected and implementation-specific. The architected feature sets have values from 0-15. The implementation-specific feature sets are values 16 and above. The architected feature sets are described in this document. The implementation-specific feature sets are described in processor-specific documentation.

This procedure will return an invalid argument if an unsupported architectural *feature\_set* is passed as an input. Implementation-specific feature sets will start at 16 and will expand in an ascending order as new implementation-specific feature sets are added. The return *status* is used by the caller to know which implementation-specific feature sets are currently supported on a particular processor.

For each valid *feature\_set*, this procedure returns which processor features are implemented in the *features\_avail* return argument, the current feature setting is in *feature\_status* return argument, and the feature controllability in the *feature\_control* return argument. Only the processor features which are implemented and controllable can be changed via PAL\_PROC\_SET\_FEATURES. Features for which *features\_avail* are 0 (unimplemented features) also have *features\_status* and *features\_control* of 0.

In [Table 11-112](#), the *class* field indicates whether a feature is required to be available (*Req.*) or is optional (*Opt.*). The *control* field indicates which features are required to be controllable. *Req.* indicates that the feature must be controllable, *Opt.* indicates that

the feature may optionally be controllable, and *No* indicates that the feature cannot be controllable. The *control* field applies only when the feature is available. The sense of the bits is chosen so that for features which are controllable, the default hand-off value at exit from PALE\_RESET should be 0. PALE\_CHECK and PALE\_INIT will not modify these features.

**Table 11-112. Processor Features**

Bit	Class	Control	Scope	Description
63	Opt.	Req.	May <sup>a</sup>	Enable BERR promotion. When 1, the Bus Error (BERR) signal is promoted to the Bus Initialization (BINIT) signal, and the BINIT pin is asserted on the occurrence of each Bus Error. Setting this bit has no effect if BINIT signalling is disabled. (See PAL_BUS_GET/SET_FEATURES)
62	Opt.	Req.	May	Enable MCA promotion. When 1, machine check aborts (MCAs) are promoted to the Bus Error signal, and the BERR pin is assert on each occurrence of an MCA. Setting this bit has no effect if BERR signalling is disabled. (See PAL_BUS_GET/SET_FEATURES)
61	Opt.	Req.	May	Enable MCA to BINIT promotion. When 1, machine check aborts (MCAs) are promoted to the Bus Initialization signal, and the BINIT pin is assert on each occurrence of an MCA. Setting this bit has no effect if BINIT signalling is disabled. (See PAL_BUS_GET/SET_FEATURES)
60	Opt.	Req.	No <sup>b</sup>	Enable CMCI promotion When 1, Corrected Machine Check Interrupts (CMCI) are promoted to MCAs. They are also further promoted to BERR if bit 39, Enable MCA promotion, is also set and they are promoted to BINIT if bit 38, Enable MCA to BINIT promotion, is also set. This bit has no effect if MCA signalling is disabled (see PAL_BUS_GET/SET_FEATURES)
59	Opt.	Req.	May	Disable Cache. When 0, the processor performs cast outs on cacheable pages and issues and responds to coherency requests normally. When 1, the processor performs a memory access for each reference regardless of cache contents and issues no coherence requests and responds as if the line were not present. Cache contents cannot be relied upon when the cache is disabled. WARNING: Semaphore instructions may not be atomic or may cause Unsupported Data Reference faults if caches are disabled.
58	Opt.	Req.	May	Disable Coherency. When 0, the processor uses normal coherency requests and responses. When 1, the processor answers all requests as if the line were not present.
57	Opt.	Req.	May	Disable Dynamic Power Management (DPM). When 0, the hardware may reduce power consumption by removing the clock input from idle functional units. When 1, all functional units will receive clock input, even when idle.
56	Opt.	Req.	May	Disable a BINIT on internal processor time-out. When 0, the processor may generate a BINIT on an internal processor time-out. When 1, the processor will not generate a BINIT on an internal processor time-out. The event is silently ignored.
55	Opt.	Req.	May	Enable external notification when the processor detects hardware errors caused by environmental factors that could cause loss of deterministic behavior of the processor. When 1, this bit will enable external notification, when 0 external notification is not provided. The type of external notification of these errors is processor-dependent. A loss of processor deterministic behavior is considered to have occurred if these environmentally induced errors cause the processor to deviate from its normal execution and eventually causes different behavior which can be observed at the processor bus pins. Processor errors that do not have this effects (i.e., software induced machine checks) may or may not be promoted depending on the processor implementation.

Table 11-112. Processor Features (Continued)

Bit	Class	Control	Scope	Description
54	Opt.	Req.	No	Enable the use of the vmsw instruction. When 0, the vmsw instruction causes a Virtualization fault when executed at the most privileged level. When 1, this bit will enable normal operation of the vmsw instruction. This bit has no effect if virtual machine features are disabled (see bit 40).
53	Opt.	Req.	May	Enable MCA signaling on unconsumed data-poisoning event detection. When 0, a CMCI will be signaled on error detection. When 1, an MCA will be signaled on error detection. Note that the reported error severity depends on which method is chosen for signaling; see <a href="#">Section 11.3.2.3, “Unconsumed Data-Poisoning Event Handling”</a> for details. If this feature is not supported, then the corresponding argument is ignored when calling PAL_PROC_SET_FEATURES. Note that the functionality of this bit is independent of the setting in bit 60 (Enable CMCI promotion), and that the bit 60 setting does not affect CMCI signaling for data-poisoning related events.
52	Opt.	Req.	May	Disable P-states. Provides the ability to disable p-states when they are implemented by the processor. When the feature is available and status is 1 or when the feature is not available, the PAL P-state procedures (PAL_PSTATE_INFO, PAL_SET_PSTATE, PAL_GET_PSTATE) will return with a status of -1 (Unimplemented procedure). When the feature is available and the status is 0, the PAL P-state procedures will operate normally.
51:48	N/A	N/A	N/A	Reserved
47	Opt.	Opt.	May	Disable Dynamic branch prediction. When 0, the processor may predict branch targets and speculatively execute, but may not commit results. When 1, the processor must wait until branch targets are known to execute.
46	Opt.	Opt.	May	Disable Dynamic Instruction Cache Prefetch. When 0, the processor may prefetch into the caches any instruction which has not been executed, but whose execution is likely. When 1, instructions may not be fetched until needed or hinted for execution. (Prefetch for a hinted branch is allowed even when dynamic instruction cache prefetch is disabled.)
45	Opt.	Opt.	May	Disable Dynamic Data Cache Prefetch. When 0, the processor may prefetch into the caches any data which has not been accessed by instruction execution, but which is likely to be accessed. When 1, no data may be fetched until it is needed for instruction execution or is fetched by an lfetch instruction.
44	Opt.	Req.	No	Disable Spontaneous Deferral. When 1, the processor may optionally defer speculative loads that do not encounter any exception conditions, but that trigger other implementation-dependent conditions (e.g., cache miss). This behavior is gated by the programming model described in <a href="#">Section 5.5.5, “Deferral of Speculative Load Faults” on page 2:105</a> . When 0, spontaneous deferral is disabled.
43	Opt.	Opt.	No	Disable Dynamic Predicate Prediction. When 0, the processor may predict predicate results and execute speculatively, but may not commit results until the actual predicates are known. When 1, the processor shall not execute predicated instructions until the actual predicates are known.
42	Opt.	No	RO <sup>c</sup>	XR1 through XR3 implemented. Denotes whether XR1 - XR3 are implemented for machine check recovery. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
41	Opt.	No	RO	XIP, XPSR, and XFS implemented. Denotes whether XIP, XPSR, and XFS are implemented for machine check recovery. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.

**Table 11-112. Processor Features (Continued)**

Bit	Class	Control	Scope	Description
40	Opt.	Opt.	No	Virtual Machine features implemented and enabled. When 1, PSR.vm is implemented and virtual machines features are not disabled. When 0 (features_status) and when the corresponding features_avail bit is 1, virtual machines features are implemented but are disabled. When both the features_avail and features_status bits are 0, virtual machine features are not implemented.  If implemented and controllable, virtual machine features may be disabled by writing this bit to 0 with PAL_PROC_SET_FEATURES. However, virtual machine features cannot be re-enabled except via a power-on; hence, if virtual machine features are disabled, this bit reads as 0 for both features_status and features_control (but still 1 for features_avail).
39	Opt.	Req.	May	Variable P-state performance: A value of 1 indicates that the processor is optimizing performance for the given P-state power budget by dynamically varying the frequency, such that maximum performance is achieved for the power budget. A value of 0 indicates that P-states have no frequency variation or very small frequency variations for their given power budget.
38	Opt.	No	RO	Simple implementation of unimplemented instruction addresses. Denotes how an unimplemented instruction address is recorded in IIP on an Unimplemented Instruction Address trap or fault. When 1, the full unimplemented address is recorded in IIP; when 0, the address is sign extended (virtual addresses) or zero extended (physical addresses). See <a href="#">Section 3.3.5.3, "Interrupt Instruction Bundle Pointer (IIP – CR19)"</a> for details. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
37	Opt.	No	RO	INIT, PMI, and LINT pins present. Denotes the absence of INIT, PMI, LINT0 and LINT1 pins on the processor. When 1, the pins are absent. When 0, the pins are present. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
36	Opt.	No	RO	Unimplemented instruction address reported as fault. Denotes how the processor reports the detection of unimplemented instruction addresses. When 1, the processor reports an Unimplemented Instruction Address fault on the unimplemented address; when 0, it reports an Unimplemented Instruction Address trap on the previous instruction in program order. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
35	Opt.	Req.	May	Disable data speculation and the ALAT. When 1, data speculation checks (chk.a) always fail (i.e., always branch to the target address), thus triggering recovery code; check loads (ld.c) always re-load the target register. When 0, data speculation works as normal.
34	Opt.	No	RO	Interrupt Instruction Bundle interruption registers (IIB0, IIB1) implemented. Denotes whether IIB registers are implemented. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
33	Opt.	No	RO	Interval Timer Offset register (ITO) implemented. Denotes whether ITO register is implemented. This feature may only be interrogated by PAL_PROC_GET_FEATURES. It may not be enabled or disabled by PAL_PROC_SET_FEATURES. The corresponding argument is ignored.
32:0	N/A	N/A	N/A	Reserved

- a. May-span-multiple-logical-processors. Readers should refer to implementation-specific document for details.
- b. Setting this bit affect logical-processor only.
- c. Read-only bit.

## PAL\_PROC\_SET\_FEATURES – Set Processor Dependent Features (18)

**Purpose:** Enables/disables specific processor features.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PROC_SET_FEATURES within the list of PAL procedures.
	feature_select	64-bit vector denoting desired state of each feature (1=select, 0=non-select).
	feature_set	Feature set to apply changes to. See PAL_PROC_GET_FEATURES for more information on feature sets.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PROC_SET_FEATURES procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	1	Call completed without error; The <i>feature_set</i> passed is not supported but a <i>feature_set</i> of a larger value is supported
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error
	-8	<i>feature_set</i> passed is beyond the maximum <i>feature_set</i> supported

**Description:** PAL\_PROC\_GET\_FEATURES should be called to ascertain the implemented processor features and their current setting before calling PAL\_PROC\_SET\_FEATURES. The list of possible processor features is defined in [Table 11-112](#). Any attempt to set processor features which cannot be set will be ignored.

## PAL\_PSTATE\_INFO – Get Information for Power/Performance States (44)

**Purpose:** Returns information about the P-states supported by the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

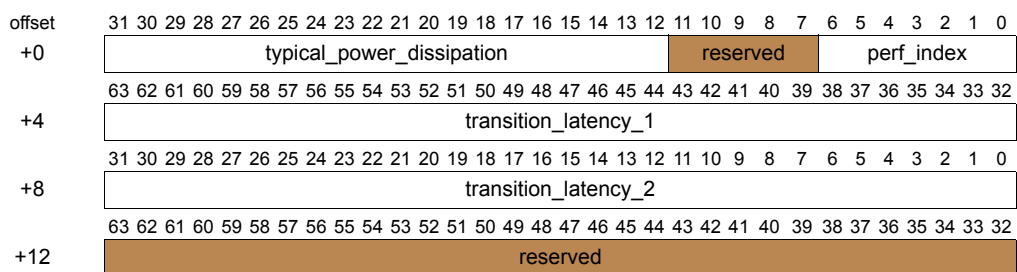
Argument	Description
index	Index of PAL_PSTATE_INFO within the list of PAL procedures.
pstate_buffer	64-bit pointer to a 256-byte buffer aligned on an 8-byte boundary.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_PSTATE_INFO procedure.
pstate_num	Unsigned integer denoting the number of P-states supported. The maximum value of this field is 16.
dd_info	Dependency domain information
Reserved	0

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error

**Description:** Information about available P-states is returned in the data buffer referenced by *pstate\_buffer*. Entries in the buffer are organized in an ascending order. For example, P0 (the highest performance P-state) state information is index 0 in the buffer, P1 state is index 1 in the buffer, and so on. The return argument *pstate\_num* indicates the number of P-states supported on the given implementation. For example, if *pstate\_num* is 4, it indicates that P-states P0-P3 are available for that implementation. Information in *pstate\_buffer* is returned only for entries corresponding to the available P-states. Entries corresponding to unimplemented P-states must be ignored. Figure 11-41 illustrates the format of the *pstate\_buffer*.

**Figure 11-41. Layout of *pstate\_buffer* Entry**



64

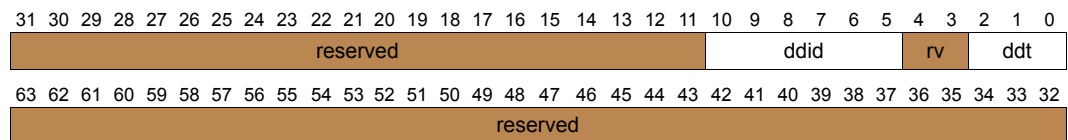
- *typical\_power\_dissipation* is a 20-bit field denoting the typical processor package power dissipation if all logical processors on the package are placed in this P-state, measured in milliwatts.
- *perf\_index* is a 7-bit field denoting the performance index of this P-state, relative to the highest available P-state (P0). This field is enumerated relative to the index of the highest-performing P-state. A value of 100 represents the minimum processor

performance in the P0 state. For example, if the P1-state has a value of 75, and the next P-state (P2) has a value of 50, it implies that P1 performance is 25% lower than P0 performance, and P2 performance is 50% lower than P0 performance.

- *transition\_latency\_1* is a 32-bit field indicating the minimum number of processor cycles required to initiate a transition to this P-state from any other P-state.
- *transition\_latency\_2* is a 32-bit field indicating the minimum recommended number of processor cycles that the caller should wait, before initiating a new P-state transition with a reasonable chance of acceptance. This field is intended to give the caller an estimation of the frequency with which PAL\_SET\_PSTATE procedure calls should be made, without having the transition request be not accepted.

Dependency domain details for the logical processor are returned in *dd\_info*. See [Figure 11-42](#) for *dd\_info* layout.

**Figure 11-42. Layout of *dd\_info* Parameter**



- *ddt* (Dependency Domain Type) is a 3-bit unsigned integer denoting the type of dependency domains that exist on the processor package. The possible values are shown in [Table 11-113](#). See [Section 11.6.1, “Power/Performance States \(P-states\)”](#) on page 2:315 for details of the values in this field.

**Table 11-113. Values for *ddt* Field**

Value	Description
0	Hardware independent (HIDD)
1	Hardware coordinated (HCDD)
2	Software coordinated (SCDD)
3-7	Reserved

- *ddid* (Dependency Domain Identifier) is a 6-bit unsigned integer denoting this logical processor's dependency domain. The *ddid* values are unique only for a given processor package. Software can use the *ddid* field to determine which logical processors belong to the same dependency domain within the package.

For more information on performance states and power management, refer to [Section 11.6.1, “Power/Performance States \(P-states\)”](#) on page 2:315.



## PAL\_PTCE\_INFO – Get PTCE Purge Loop Information (6)

**Purpose:** Returns information required for the architected loop used to purge (initialize) the entire TC.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_PTCE_INFO within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_PTCE_INFO procedure.
	tc_base	Unsigned 64-bit integer denoting the beginning address to be used by the first PTCE instruction in the purge loop.
	tc_counts	Two unsigned 32-bit integers denoting the loop counts of the outer (loop 1) and inner (loop 2) purge loops. count1 (loop 1) is contained in bits 63:32 of the parameter, and count2 (loop 2) is contained in bits 31:0 of the parameter.
	tc_strides	Two unsigned 32-bit integers denoting the loop strides of the outer (loop 1) and inner (loop 2) purge loops. stride1 (loop 1) is contained in bits 63:32 of the parameter, and stride2 (loop 2) is contained in bits 31:0 of the parameter.

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** No explicit hardware support is required by this call. See the purge loop example in the description of the `ptc.e` instruction in [Chapter 2, "Instruction Reference" in Volume 3](#).

## PAL\_REGISTER\_INFO – Return Information about Implemented Processor Registers (39)

**Purpose:** Returns information about implemented Application and Control Registers.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_REGISTER_INFO within the list of PAL procedures.
info_request	Unsigned 64-bit integer denoting what register information is requested.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_REGISTER_INFO procedure.
reg_info_1	64-bit vector denoting information for registers 0-63. Bit 0 is register 0, bit 63 is register 63.
reg_info_2	64-bit vector denoting information for registers 64-127. Bit 0 is register 64, bit 63 is register 127.
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

This procedure is called to obtain information about the implementation of Application Registers and Control Registers. [Table 11-114](#) shows the information that is returned for each request.

**Table 11-114. info\_request Return Value**

info_request	Meaning of Return Bit Vector
0	A 0-bit in the return vector indicates that the corresponding Application Register is not implemented, a 1-bit in the return vector indicates that the corresponding Application Register is implemented.
1	A 0-bit in the return vector indicated that the corresponding Application Register can be read without side effects, a 1-bit in the return vector indicated that the corresponding Application registers may cause side effects when read.
2	A 0-bit in the return vector indicates that the corresponding Control Register is not implemented, a 1-bit in the return vector indicates that the corresponding Control Register is implemented.
3	A 0-bit in the return vector indicated that the corresponding Control Register can be read without side effects, a 1-bit in the return vector indicated that the corresponding Control Register may cause side effects when read.
All others	Reserved.

## PAL\_RSE\_INFO – Get RSE Information (19)

**Purpose:** Returns information about the register stack and RSE for this processor implementation.

**Calling Conv:** Static Registers Only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Argument	Description
index	Index of PAL_RSE_INFO within the list of PAL procedures.
Reserved	0
Reserved	0
Reserved	0

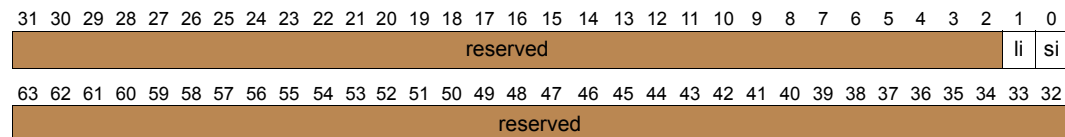
Return Value	Description
status	Return status of the PAL_RSE_INFO procedure.
phys_stacked	Number of physical stacked general registers.
hints	RSE hints supported by processor.
Reserved	0

Status Value	Description
0	Call completed without error
-2	Invalid argument
-3	Call completed with error

**Description:** The return parameter *phys\_stacked* contains a 64-bit unsigned integer that specifies the number of physical registers implemented by the processor for the stacked general registers, r32-r127. *phys\_stacked* will be an integer multiple of 16 greater than or equal to 96.

The return parameter *hints* contains a 2-bit field that specifies which RSE load/store hints are implemented.

**Figure 11-43. Layout of *hints* Return Value**



A bit field value of 1 specifies that the corresponding mode is implemented; a value of 0 specifies that the mode is not implemented. The bit field encodings are:

**Table 11-115. RSE Hints Implemented**

li	si	RSE Hints	Class
0	0	enforced lazy	Required
0	1	eager stores	Optional
1	0	eager loads	Optional
1	1	eager stores and loads	Optional

“Lazy” is the default RSE mode and must be implemented. Hardware is not required to implement any of the other modes.

## PAL\_SET\_HW\_POLICY – Set Current Hardware Resource Sharing Policy (49)

**Purpose:** Sets the current hardware resource sharing policy of the processor.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_SET_HW_POLICY within the list of PAL procedures.
policy	Unsigned 64-bit integer specifying the hardware resource sharing policy the caller is setting.
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_SET_HW_POLICY procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
1	Call completed successfully but could not change the hardware policy since a competing logical processor is set in exclusive high priority
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-9	Call requires PAL memory buffer

**Description:** This procedure is used to set the hardware resource sharing policy on the logical processor it is called on. The setting of this policy will impact other logical processors on the physical processor package. The logical processors impacted is returned by the PAL\_GET\_HW\_POLICY procedure, see “[PAL\\_GET\\_HW\\_POLICY – Retrieve Current Hardware Resource Sharing Policy \(48\)](#)” on page 2:394 for details.

The input argument *policy* selects the hardware policy the caller would like to set. The supported hardware policies are listed in [Table 11-116](#) below. By default the hardware always sets the processor in the performance policy at reset.

**Table 11-116. Processor Hardware Sharing Policies**

Value	Name	Description
0	Performance	The processor has its hardware resources configured to achieve maximum performance across all logical processors.
1	Fairness	The processor configures hardware resources to approximately achieve equal sharing of competing hardware resources among all impacted logical processors.

**Table 11-116. Processor Hardware Sharing Policies (Continued)**

Value	Name	Description
2	High-priority	The processor configures hardware resources to provide the logical processor this procedure was called on a greater share of the competing hardware resources. All competing logical processors will get a smaller share of the competing hardware resources.
3	Exclusive High-priority	The processor configures hardware resources such that the logical processor this procedure was called on has a greater share of the competing hardware resources. All competing logical processors will get a smaller share of the competing hardware resources. This policy also ensures that no other competing logical processor can modify the hardware sharing policy until the logical processor that is in exclusive high priority releases exclusive high-priority by selecting a different policy.
All Other Values		Reserved

The caller must be aware of which logical processors are impacted by hardware policy changes, since making a call on one of the logical processors will impact all logical processors that share the same hardware resources. For example if the caller selects the high-priority policy on one logical processor A and then later in time selects fairness policy on one of the competing logical processors B, the procedure will take away high-priority status from logical processor A and change all impacted logical processors to the fairness policy without an error.

If a caller wants to ensure that high-priority will not be taken away from a logical processor, it can use the exclusive high-priority policy. This policy will return an error if any competing logical processor tries to change the hardware policy. This ensures that the caller can ensure a certain logical processor will retain high-priority status until that status is explicitly released by that logical processor.

This procedure is only supported on processors that have multiple logical processors sharing hardware resources that can be configured. On all other processor implementations, this procedure will return the Unimplemented procedure return status.

## PAL\_SET\_PSTATE – Request Processor to Enter Power/Performance State (263)

**Purpose:** To request a processor transition to a given P-state.

**Calling Conv:** Stacked Registers

**Mode:** Physical and Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_SET_PSTATE within the list of PAL procedures.
p_state	Unsigned integer denoting the processor P-state being requested.
force_pstate	Unsigned integer denoting whether the P-state change should be forced for the logical processor.
Reserved	0

Return Value	Description
status	Return status of the PAL_SET_PSTATE procedure.
Reserved	0
Reserved	0
Reserved	0

Status Value	Description
1	Call completed without error, but transition request was not accepted
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-9	Call requires PAL memory buffer

**Description:** PAL\_SET\_PSTATE is used to request the transition of the processor to the P-state specified by the *p\_state* input parameter. The PAL\_SET\_PSTATE procedure does not wait for the transition to complete before returning back to the caller. The request may either be accepted (*status* = 0) or not accepted (*status* = 1), depending on hardware capabilities and implementation-specific event conditions. The presence of a platform power-cap does not prevent the request from being accepted. (See [Section 11.6.1, “Power/Performance States \(P-states\)” on page 2:315](#) for details.) If the request is not accepted, then no transition is performed, and it is up to the caller to make another PAL\_SET\_PSTATE procedure call to transition to the desired P-state. When the request is accepted, the processor will attempt to initiate a transition to the requested performance state. For SCDD or HIDD logical processors, the procedure will always succeed in transitioning to the requested performance state. For HCDD logical processors, the procedure will make a best-case attempt at fulfilling the transition request, based on the nature of the dependencies that exist between the logical processors in the domain. In such circumstances, the procedure may initiate no transition, partial transition or full transition to the requested P-state.

The *force\_pstate* argument may be used for a HCDD when it is necessary to get a deterministic response for the P-state transition at the expense of compromising the power/performance of other logical processors in the same domain. If the *force\_pstate* argument is non-zero, and if the request is accepted, the procedure will initiate the P-state transition on the logical processor regardless of any dependencies that exist in the dependency domain at the time the procedure is called. Forcing the P-state does not change the P-states requested by other logical processors in the dependency domain, nor the value seen on other logical processors when they do a PAL\_GET\_PSTATE with *type*=0; rather, forcing the P-state effectively suspends hardware

coordination. A subsequent call to PAL\_SET\_PSTATE on any logical processor in the dependency domain (with a *force\_pstate* argument of zero) reinstates hardware coordination. The *force\_pstate* argument is ignored on SCDD and HIDD logical processors.

Calling this procedure on some processor implementations may affect P-states of other processors in the same dependency domain. Please refer to [Section 11.6.1, "Power/Performance States \(P-states\)" on page 2:315](#) and implementation-specific reference manuals for details.

## PAL\_SHUTDOWN – Shutdown the Processor (45)

**Purpose:** Put the logical processor into a low power state which can be exited only by a reset event.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_SHUTDOWN within the list of PAL procedures.
	notify_platform	8-byte aligned physical address pointer providing details on how to optionally notify the platform that the processor is entering a shutdown state.
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_SHUTDOWN procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** This call places the logical processor in a low power state which can be exited only by asserting a reset. This procedure can optionally let the platform know that it is about to shutdown by performing a store operation as specified in the *notify\_platform* input argument.

If the *notify\_platform* input argument is zero, no store operation will be performed. If the *notify\_platform* input argument is non-zero, the layout for this argument is shown in [Table 11-117](#).

**Table 11-117. *notify\_platform* Layout**

Offset	Description
0x0	Size of the store operation to perform (1, 2, 4 or 8 are the only valid values for this field).
0x8	Aligned physical address of the store operation. The most significant bit (63) of the physical address should be set according to the cacheability attribute wanted for the store transaction.
0x10	Data value for the store operation.
All others	Reserved.

If the address value is not naturally aligned to the size selected, this procedure will return an error.

The logical processor will wait until this transaction has been received by the platform before entering the shutdown state.

On receipt of a reset event, the logical processor will reset itself and start execution at the PAL reset address. All other events will be ignored by the logical processor when in shutdown state.



## PAL\_TEST\_INFO – Information for Processor Self-test (37)

**Purpose:** Returns the alignment and size requirements needed for the memory buffer passed to the PAL\_TEST\_PROC procedure as well as information on self-test control words for the processor self-tests.

**Calling Conv:** Static Registers Only

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_TEST_INFO within the list of PAL procedures.
	test_phase	Unsigned integer that specifies which phase of the processor self-test information is being requested on. A value of 0 indicates the phase two of the processor self-test and a value of 1 indicates phase one of the processor self-test. All other values are reserved.
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_TEST_INFO procedure.
	bytes_needed	Unsigned 64-bit integer denoting the number of bytes of main memory needed to perform the second phase of processor self-test.
	alignment	Unsigned 64-bit integer denoting the alignment required for the memory buffer.
	st_control	48-bit wide bit-field indicating if control of the processor self-tests is supported and which bits of the <i>test_control</i> field are defined for use.

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** PAL\_TEST\_INFO returns the size and alignment requirements for the memory buffer that is passed to the PAL\_TEST\_PROC procedure and returns information on the implementation of the self-test control word based on the *test\_phase* input argument. Please see [Section 11.2.3, "PAL Self-test Control Word" on page 2:295](#) for more information on the self-test control word.

When *test\_phase* is equal to zero, information is returned about phase two of the processor self-test. These are the tests that require external memory to execute properly. When *test\_phase* is equal to one, information is returned about phase one of the processor self-test. These are the tests that are normally run during PALE\_RESET and do not require external memory to properly execute. When information is requested about phase one of the processor self-test a memory buffer and alignment argument will be returned as well since these tests may need to save and restore processor state to this memory buffer if executed from the PAL\_TEST\_PROC procedure.

## PAL\_TEST\_PROC – Perform a Processor Self-test (258)

**Purpose:** Performs the second phase of processor self test.

**Calling Conv:** Stacked Registers

PAL\_TEST\_PROC may modify some registers marked unchanged in the Stacked Register calling convention. See additional description below.

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_TEST_PROC within the list of PAL procedures.
	test_address	64-bit physical address of main memory area to be used by processor self-test. The memory region passed must be cacheable, bit 63 must be zero.
	test_info	Input argument specifying the size of the memory buffer passed and the phase of the processor self-test that should be run. See <a href="#">Figure 11-44</a> .
	test_params	Input argument specifying the self-test control word and the allowable memory attributes that can be used with the memory buffer. See <a href="#">Figure 11-45</a> .

Returns:	Return Value	Description
	status	Return status of the PAL_TEST_PROC procedure.
	self-test_state	Formatted 8-byte value denoting the state of the processor after self-test. The format is described in <a href="#">Section 11.2.2.3, "Definition of Self Test State Parameter"</a> on page 2:293.
	Reserved	0
	Reserved	0

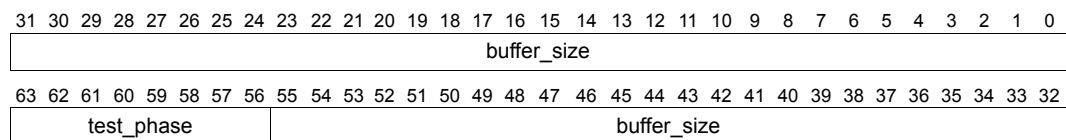
Status:	Status Value	Description
	1	Call completed without error, but hardware failures occurred during self-test
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** The PAL\_TEST\_PROC procedure will perform a phase of the processor self-tests as directed by the *test\_info* and the *test\_control* input parameters.

*test\_address* points to a contiguous memory region to be used by PAL\_TEST\_PROC. This memory region must be aligned as specified by the alignment return value from PAL\_TEST\_INFO, otherwise this procedure will return with an invalid argument return value. The PAL\_TEST\_PROC routine requires that the memory has been initialized and that there are no known uncorrected errors in the allocated memory.

The *test\_info* input parameter specifies the size of the memory buffer passed to the procedure and which phase of the processor self-test is requested to be run (either phase one or phase two).

**Figure 11-44. Layout of *test\_info* Argument**

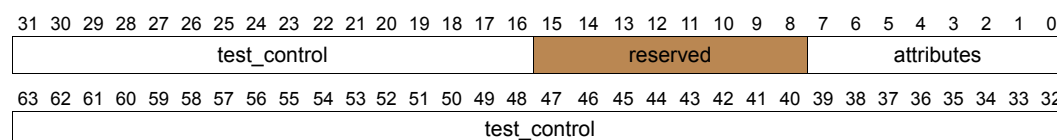


- *buffer\_size* indicates the size in bytes of the memory buffer that is passed to this procedure. *buffer\_size* must be greater than or equal in size to the *bytes\_needed* return value from PAL\_TEST\_INFO, otherwise this procedure will return with an invalid argument return value.

- *test\_phase* defines which phase of the processor self-tests are requested to be run. A value of zero indicates to run phase two of the processor self-tests. Phase two of the processor self-tests are ones that require external memory to execute correctly. A value of one indicates to run phase one of the processor self-tests. Phase one of the processor self-tests are tests run during PALE\_RESET and do not depend on external memory to run correctly. When the caller requests to have phase one of the processor self-test run via this procedure call, a memory buffer may be needed to save and restore state as required by the PAL calling conventions. The procedure PAL\_TEST\_INFO informs the caller about the requirements of the memory buffer.

The *test\_params* input argument specifies which memory attributes are allowed to be used with the memory buffer passed to this procedure as well as the self-test control word. The self-test control word *test\_control* controls the runtime and coverage of the processor self-test phase specified in the *test\_phase* parameter.

**Figure 11-45. Layout of *test\_param* Argument**



- *attributes* specifies the memory attributes that are allowed to be used with the memory buffer passed to this procedure. The *attributes* parameter is a vector where each bit represents one of the virtual memory attributes defined by the architecture. The bit field position corresponds to the numeric memory attribute encoding defined in [Section 4.4, “Memory Attributes” on page 2:75](#). The caller is required to support the cacheable attribute for the memory buffer, otherwise an invalid argument will be returned.
- *test\_control* is the self-test control word corresponding to the *test\_phase* passed. This *test\_control* directs the coverage and runtime of the processor self-tests specified by the *test\_phase* input argument. Information about the self-test control word can be found in [Section 11.2.3, “PAL Self-test Control Word” on page 2:295](#) and information on if this feature is implemented and the number of bits supported can be obtained by the PAL\_TEST\_INFO procedure call. If this feature is implemented by the processor, the caller can selectively skip parts of the processor self-test by setting *test\_control* bits to a one. If a bit has a zero, this test will be run. The values in the unimplemented bits are ignored. If PAL\_TEST\_INFO indicated that the self-test control word is not implemented, this procedure will return with an invalid argument status if the caller sets any of the *test\_control* bits.

PAL\_TEST\_PROC will classify the processor after the self-test in one of four states: CATASTROPHIC FAILURE, FUNCTIONALLY RESTRICTED, PERFORMANCE RESTRICTED, or HEALTHY. These processor self-test states are described in [Figure 11-9 on page 2:293](#). If PAL\_TEST\_PROC returns in the FUNCTIONALLY RESTRICTED or PERFORMANCE RESTRICTED states the *self-test\_status* return value can provide additional information regarding the nature of the failure. In the case of a CATASTROPHIC FAILURE, the procedure does not return.

The procedure will only perform memory accesses to the buffer passed to it using the memory attributes indicated in the *attributes* bit-field. The caller must ensure that the memory region passed to the procedure is in a coherent state.

PAL\_TEST\_PROC may modify PSR bits or system registers as necessary to test the processor. These bits or registers must be restored upon exit from PAL\_TEST\_PROC

## ***PAL\_TEST\_PROC***

with the exception of the translation caches, which are evicted as a result of testing. PAL\_TEST\_PROC is free to invalidate all cache contents. If the caller depends on the contents of the cache, they should be flushed before making this call. PAL\_TEST\_PROC requires that the RSE is set up properly to handle spills and fills to a valid memory location if the contents of the register stack are needed. PAL\_TEST\_PROC requires that the memory buffer passed to it is not shared with other processors running this procedure in the system at the same time. PAL\_TEST\_PROC will use this memory region in a non-coherent manner. PAL\_TEST\_PROC may overwrite floating point registers 32-127 without restoring their values upon exit.

## PAL\_VERSION – Get PAL Version Number Information (20)

**Purpose:** Returns PAL version information.

**Calling Conv:** Static registers only

**Mode:** Physical or Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_VERSION within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VERSION procedure.
	min_pal_ver	8-byte formatted value returning the minimum PAL version needed for proper operation of the processor. See <a href="#">Figure 11-46</a> .
	current_pal_ver	8-byte formatted value returning the current PAL version running on the processor. See <a href="#">Figure 11-46</a> .
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-2	Invalid argument
	-3	Call completed with error

**Description:** PAL\_VERSION provides the caller the minimum PAL version needed for proper operation of the processor as well as the current PAL version running on the processor. The *min\_pal\_ver* and *current\_pal\_ver* return values are 8-byte values in the following format:

**Figure 11-46. Layout of *min\_pal\_ver* and *current\_pal\_ver* Return Values**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
PAL_vendor								Reserved								PAL_B_version															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32
Reserved																PAL_A_version															

- *PAL\_B\_version* is a 16-bit binary coded decimal (BCD) number that provides identification information about the PAL\_B firmware.
- *PAL\_vendor* is an unsigned 8-bit integer indicating the vendor of the PAL code.
- *PAL\_A\_version* is a 16-bit binary coded decimal (BCD) number that provides identification information about the PAL\_A firmware. In the split PAL\_A model, this return value is the version number of the processor-specific PAL\_A. The generic PAL\_A version is not returned by this procedure in the split PAL\_A model.

The version numbers selected for the PAL\_A and PAL\_B firmware is specific to the *PAL\_vendor*. The version numbers selected will always have the property that later versions of firmware will have a higher number than earlier versions of firmware.

## PAL\_VM\_INFO – Get Virtual Memory Information (7)

**Purpose:** Return information about the virtual memory characteristics of the processor implementation.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

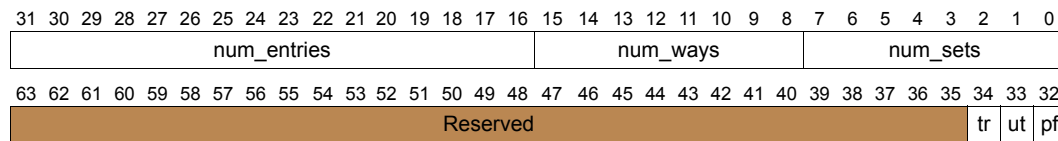
Arguments:	Argument	Description
	index	Index of PAL_VM_INFO within the list of PAL procedures.
	tc_level	Unsigned 64-bit integer specifying the level in the TLB hierarchy for which information is required. This value must be between 0 and one less than the value returned in the <i>vm_info_1.num_tc_levels</i> return value from PAL_VM_SUMMARY.
	tc_type	Unsigned 64-bit integer with a value of 1 for instruction translation cache and 2 for data or unified translation cache. All other values are reserved.
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VM_INFO procedure.
	tc_info	8-byte formatted value returning information about the specified TC.
	tc_pages	64-bit vector containing a bit for each page size supported in the specified TC, where bit position n indicates a page size of 2**n.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error.
	-2	Invalid argument.
	-3	Call completed with error.

**Description:** The *tc\_info return* is an 8-byte quantity in the following format:

**Figure 11-47. Layout of *tc\_info* Return Value**



- *num\_sets* – Unsigned 8-bit integer denoting the number of hash sets for the specified level (1=fully associative)
- *num\_ways* – Unsigned 8-bit integer denoting the associativity of the specified level (1=direct).
- *num\_entries* – Unsigned 16-bit integer denoting the number of entries in the specified TC.
- *pf* – Flag denoting whether the specified level is optimized for the region’s preferred page size (1=optimized). *tc\_pages* indicates which page sizes are usable by this translation cache.
- *ut* – Flag denoting whether the specified TC is unified (1=unified).
- *tr* – Flag denoting whether installed translation registers will reduce the number of entries within the specified TC.

The *num\_entries* will always equal *num\_ways* \* *num\_sets*. For a direct mapped TC, *num\_ways* = 1 and *num\_sets* = *num\_entries*. For a fully associative TC, *num\_sets* = 1 and *num\_ways* = *num\_entries*.

## PAL\_VM\_PAGE\_SIZE – Get Virtual Memory Page Size Information (34)

**Purpose:** Returns page size information about the virtual memory characteristics of the processor implementation.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_VM_PAGE_SIZE within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VM_PAGE_SIZE procedure.
	insertable_pages	64-bit vector containing a bit for each architected page size that is supported for TLB insertions and region registers.
	purge_pages	64-bit vector containing a bit for each architected page size supported for TLB purge operations.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error.
	-2	Invalid argument
	-3	Call completed with error.

**Description:** The values returned from this call are all 64-bit bitmaps. One bit is set for each page size implemented by the processor where bit  $n$  represents a page size of  $2^{*n}$ . Please refer to [Table 4-5 on page 2:58](#) for the minimum page sizes that are supported.

The *insertable\_pages* returns the page sizes that are supported for TLB insertions and region registers.

The *purge\_pages* returns the page sizes that are supported for the TLB purge operations.

## PAL\_VM\_SUMMARY – Get Virtual Memory Summary Information (8)

**Purpose:** Returns summary information about the virtual memory characteristics of the processor implementation.

**Calling Conv:** Static Registers Only

**Mode:** Physical and Virtual

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_VM_SUMMARY within the list of PAL procedures.
	Reserved	0
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VM_SUMMARY procedure.
	vm_info_1	8-byte formatted value returning global virtual memory information.
	vm_info_2	8-byte formatted value returning global virtual memory information.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error.
	-2	Invalid argument
	-3	Call completed with error.

**Description:** The *vm\_info\_1* return is an 8-byte quantity in the following format:

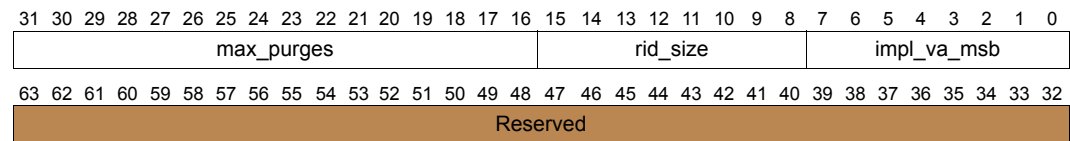
**Figure 11-48. Layout of *vm\_info\_1* Return Value**

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																
hash_tag_id								max_pkr								key_size								phys_add_size								vw															
63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32																
num_tc_levels																num_unique_tcs																max_itr_entry								max_dtr_entry							

- *vw* – 1-bit flag indicating whether a hardware TLB walker is implemented (1 = walker present).
- *phys\_add\_size* – Unsigned 7-bit integer denoting the number of bits of physical address implemented.
- *key\_size* – Unsigned 8-bit integer denoting the number of bits implemented in the PKR.key field.
- *max\_pkr* – Unsigned 8-bit integer denoting the maximum PKR index (number of PKRs-1).
- *hash\_tag\_id* – Unsigned 8-bit integer which uniquely identifies the processor hash and tag algorithm.
- *max\_dtr\_entry* – Unsigned 8 bit integer denoting the maximum data translation register index (number of dtr entries - 1).
- *max\_itr\_entry* – Unsigned 8 bit integer denoting the maximum instruction translation register index (number of itr entries - 1).
- *num\_unique\_tcs* – Unsigned 8-bit integer denoting the number of unique TCs implemented. This is a maximum of  $2 * num\_tc\_levels$ .
- *num\_tc\_levels* – Unsigned 8-bit integer denoting the number of TC levels.

The *vm\_info\_2* return is an 8-byte quantity in the following format:



**Figure 11-49. Layout of *vm\_info\_2* Return Value**

- *impl\_va\_msb* – Unsigned 8-bit integer denoting the bit number of the most significant virtual address bit. This is the total number of virtual address bits - 1.
- *rid\_size* – Unsigned 8-bit integer denoting the number of bits implemented in the RR.rid field.
- *max\_purges* – Unsigned 16 bit integer denoting the maximum number of concurrent outstanding TLB purges allowed by the processor. A value of 0 indicates one outstanding purge allowed. A value of  $2^{16}-1$  indicates no limit on outstanding purges. All other values indicate the actual number of concurrent outstanding purges allowed.

## PAL\_VM\_TR\_READ – Read a Translation Register (261)

**Purpose:** Reads a translation register.

**Calling Conv:** Stacked Registers

**Mode:** Physical

**Buffer:** Not dependent

Arguments:	Argument	Description
	index	Index of PAL_VM_TR_READ within the list of PAL procedures.
	reg_num	Unsigned 64-bit number denoting which TR to read.
	tr_type	Unsigned 64-bit number denoting whether to read an ITR (0) or DTR (1). All other values are reserved.
	tr_buffer	64-bit pointer to the 32-byte memory buffer in which translation data is returned.

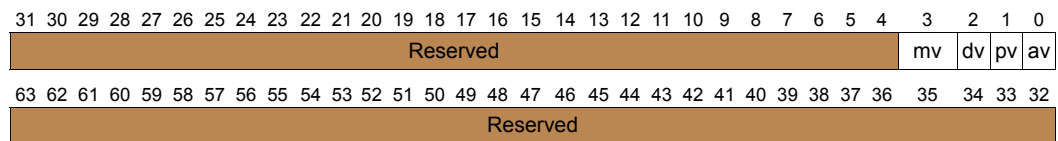
Returns:	Return Value	Description
	status	Return status of the PAL_VM_TR_READ procedure.
	TR_valid	Formatted bit vector denoting which fields are valid. See <a href="#">Figure 11-50</a> .
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error.
	-2	Invalid argument
	-3	Call completed with error.

**Description:** This procedure reads the specified translation register and returns its data in the buffer starting at *tr\_buffer*. The format of the data is returned in Translation Insertion Format, as described in [Figure 4-5, "Translation Insertion Format,"](#) on page 2:54. In addition, bit 0 of the IFA in [Figure 4-5](#) (an ignored field in the figure) will return whether the translation is valid. If bit 0 is 1, the translation is valid.

Some fields of the translation register returned may be invalid. The validity of these fields is indicated by the return argument *TR\_valid*. If these fields are not valid, the caller should ignore the indicated fields when reading the translation register returned in *tr\_buffer*.

**Figure 11-50. Layout of *TR\_valid* Return Value**



- av – denotes that the access rights field is valid
- pv – denotes that the privilege level field is valid
- dv – denotes that the dirty bit is valid
- mv – denotes that the memory attributes are valid.

A value of 1 denotes a valid field. A value of 0 denotes an invalid field. Any value returned in an invalid field must be ignored.

The *tr\_buffer* parameter should be aligned on an 8 byte boundary.

**Note:** This procedure may have the side effect of flushing all the translation cache entries depending on the implementation.

## PAL\_VP\_CREATE – PAL Create New Virtual Processor (265)

**Purpose:** Initializes a new *vpd* for the operation of a new virtual processor in the virtual environment.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_CREATE within the list of PAL procedures
	vpd	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
	host_iva	64-bit host virtual pointer to the host IVT for the virtual processor
	opt_handler	64-bit non-zero host-virtual pointer to an optional handler for virtualization intercepts. See <a href="#">Section 11.7.3, "PAL Intercepts in Virtual Environment" on page 2:332</a> for details.

Returns:	Return Value	Description
	status	Return status of the PAL_VP_CREATE procedure
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** Initializes a new *vpd* for the operation of a new virtual processor within the virtual environment.

The caller must pass a pointer to the new Virtual Processor Descriptor (*vpd*) as argument. The host virtual to host physical translation of the 64K region specified by *vpd* must be mapped with either a DTR or DTC. See [Section 11.10.2.1.3, "Making PAL Procedure Calls in Physical or Virtual Mode" on page 2:359](#) for details on data translation requirements of memory buffer pointers passed as arguments to PAL procedures. The *vac*, *vdc* and *virt\_env\_vaddr* parameters in the VPD must already be initialized before calling this procedure. Invalid argument is returned on unsupported *vac/vdc* combinations. See [Section 11.7.4.4, "Virtualization Optimization Combinations" on page 2:349](#) for details.

The *host\_iva* parameter specifies the host IVT to handle IVA-based interruptions when this virtual processor is running. The VMM can use the same or different *host\_iva* for each virtual processor. The *opt\_handler* specifies an optional virtualization intercept handler. If a non-zero value is specified, all virtualization intercepts are delivered to this handler. If a zero value is specified, all virtualization intercepts are delivered to the Virtualization vector in the host IVT. If the VMM relocates the IVT specified by the *host\_iva* parameter and/or the virtualization intercept handler specified by the *opt\_handler* parameter after this procedure, PAL\_VP\_REGISTER must be called to register the new host IVT and virtualization intercept handler before resuming virtual processor execution or allowing any IVA-based interruptions to occur; otherwise processor operation is undefined.

Upon return, the VMM is responsible for setting up the rest of the VMD state before the new virtual processor is launched (via PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER).

## ***PAL\_VP\_CREATE***

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, "Processor Virtualization" on page 2:44](#) and ["PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)" on page 2:446](#) for details.

## PAL\_VP\_ENV\_INFO – PAL Virtual Environment Information (266)

**Purpose:** Returns the parameters needed to enter a virtual environment.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Argument	Description
index	Index of PAL_VP_ENV_INFO within the list of PAL procedures
Reserved	0
Reserved	0
Reserved	0

Return Value	Description
status	Return status of the PAL_VP_ENV_INFO procedure
buffer_size	Unsigned integer denoting the number of bytes required by the PAL virtual environment buffer during PAL_VP_INIT_ENV
vp_env_info	64-bit vector of virtual environment information. See <a href="#">Table 11-118</a> . for details
Reserved	0

Status Value	Description
0	Call completed without error
-1	Unimplemented procedure
-2	Invalid argument
-3	Call completed with error
-9	Call requires PAL memory buffer

**Description:** This procedure returns the configuration options and the PAL virtual environment buffer size required by PAL\_VP\_INIT\_ENV. This procedure is used by the VMM to setup a virtual environment and determine the amount of memory / resources required. The VMM can then allocate the required amount of physical memory, set up the virtual to physical instruction and data translations that cover the PAL virtual environment buffer in TRs and call PAL\_VP\_INIT\_ENV. The buffer allocated must be at least 4K aligned.

On a multiprocessor system, this procedure need only be invoked once (on any one logical processor) to obtain virtual environment information.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, “Processor Virtualization” on page 2:44](#) and [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)” on page 2:446](#) for details.

**Table 11-118. vp\_env\_info – Virtual Environment Information Parameter**

Field	Bit	Description
Reserved	7:0	Reserved
opcode	8	If 1, hardware support to provide opcode information during PAL intercepts is available. The opcode (and the decoding of cause) passed as parameters to the VMM on intercept will represent the instruction that triggered the intercept. If 0, opcode information during PAL intercepts is provided by PAL. The opcode (and the decoding of cause) passed as parameters to the VMM on intercept will not necessarily represent the instruction that triggered the intercept, but may represent some value that was written to memory between the time the instruction that triggered the intercept was fetched, and when the intercept was triggered.
Reserved	9	Reserved
gitc	10	If 1, guest MOV-from-AR.ITC optimization is supported. <sup>a</sup> If 0, guest MOV-from-AR.ITC optimization is not supported.

**Table 11-118. *vp\_env\_info* – Virtual Environment Information Parameter**

Field	Bit	Description
Reserved	31:11	Reserved
probe	32	If 1, processor supports interception of probe instructions. See <a href="#">Section 11.7.4.2.8, “Probe Instruction Virtualization”</a> on page 2:344 for details on the usage of this control. If 0, intercept of probe instructions is not supported.
tf	33	If 1, guest test feature optimization is supported. If 0, this optimization is not supported. See <a href="#">Section 11.7.4.2.9, “Test Feature Optimization”</a> on page 2:345 for details.
ic_um	34	If 1, guest interruption collection and user mask optimization is supported. If 0, this optimization is not supported. See <a href="#">Section 11.7.4.2.10, “Interruption Collection and User Mask Optimization”</a> on page 2:345 for details.
Reserved	63:35	Reserved

- a. Architecturally, an implementation which supports guest MOV-from-AR.ITC will also support the interval timer offset (ITO) register.

## PAL\_VP\_EXIT\_ENV – PAL Exit Virtual Environment (267)

**Purpose:** Allows a logical processor to exit a virtual environment.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_EXIT_ENV within the list of PAL procedures
	iva	Optional 64-bit host virtual pointer to the IVT when this procedure is done
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VP_EXIT_ENV procedure
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** This procedure allows a logical processor to exit a virtual environment.

Upon successful execution of the PAL\_VP\_EXIT\_ENV procedure and if the *iva* parameter is non-zero, the IVA control register will contain the value from the *iva* parameter.

On a multiprocessor system, the VMM must allow the last logical processor in this environment to complete the procedure before freeing the memory resource allocated to the virtual environment.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, "Processor Virtualization" on page 2:44](#) and ["PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)" on page 2:446](#) for details.

## PAL\_VP\_INFO – PAL Virtual Processor Information (50)

**Purpose:** Returns information about virtual processor features.

**Calling Conv:** Static

**Mode:** Physical

**Buffer:** Not dependent

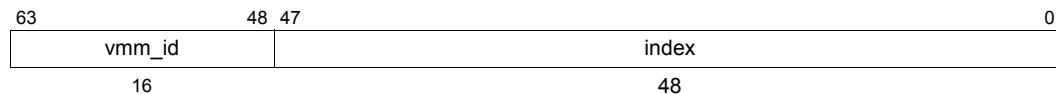
Arguments:	Argument	Description
	index	Index of PAL_VP_INFO within the list of PAL procedures
	feature_set	Feature set information is being requested for.
	vp_buffer	An address to an 8-byte aligned memory buffer (if used).
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VP_INFO procedure
	vp_info	Information about the virtual processor.
	vmm_id	Unique identifier for the VMM.
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-8	Specified feature_set is not implemented

**Description:** The PAL\_VP\_INFO procedure call is used to describe virtual processor features.

The *feature\_set* input argument for PAL\_VP\_INFO describes which virtual-processor *feature\_set* information is being requested, and is composed of two fields as shown:



A *vmm\_id* of 0 indicates architected feature sets, while others are implementation-specific feature sets. Implementation-specific feature sets are described in VMM-specific documentation.

This procedure will return a -8 if an unsupported *feature\_set* argument is passed as an input. The return status is used by the caller to know which feature sets are currently supported on a particular VMM. This procedure always returns unimplemented (-1) when called on physical processors.

For each valid *feature\_set*, this procedure returns information about the virtual processor in *vp\_info*. Additional information may be returned in the memory buffer pointed to by *vp\_buffer*, as needed. Details, for a given implementation-specific *feature\_set*, of whether information is returned in the buffer, the size of the buffer, and the representation of this information in the buffer and in *vp\_info* are described in VMM-specific documentation.

Architected *feature\_set* 0 (*vmm\_id* 0, *index* 0) is defined and required to be implemented (if this procedure is implemented), but there are no architected features defined in it yet, and so all bits in *vp\_info* are reserved for architected *feature\_set* 0. Other architected feature sets (*vmm\_id* 0, *index*>0) are undefined, and return -8 (Specified *feature\_set* is not implemented). Software can call PAL\_VP\_INFO with a *feature\_set* argument of 0 to



get the *vmm\_id*, although *vmm\_id* is also returned for any other implemented feature sets as well. For *feature\_set* 0, the *vp\_buffer* argument is ignored.

## PAL\_VP\_INIT\_ENV – PAL Initialize Virtual Environment (268)

**Purpose:** Allows a logical processor to enter a virtual environment.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_INIT_ENV within the list of PAL procedures
	config_options	64-bit vector of global configuration settings – See <a href="#">Table 11-119</a> . for details
	pbase_addr	Host physical base address of a block of contiguous physical memory for the PAL virtual environment buffer – This memory area must be allocated by the VMM and be 4K aligned. The first logical processor to enter the environment will initialize the physical block for virtualization operations.
	vbase_addr	Host virtual base address of the corresponding physical memory block for the PAL virtual environment buffer – The VMM must maintain the host virtual to host physical data and instruction translations in TRs for addresses within the allocated address space. Logical processors in this virtual environment will use this address when transitioning to virtual mode operations.

Returns:	Return Value	Description
	status	Return status of the PAL_VP_INIT_ENV procedure
	vsa_base	Virtualization Service Address – VSA specifies the virtual base address of the PAL virtualization services in this virtual environment.
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** This procedure allows a logical processor to enter a virtual environment. This call must be made after calling PAL\_VP\_ENV\_INFO and before calling other PAL virtualization procedures and services. All of the logical processors in a virtual environment share the same **PAL virtual environment buffer**. The buffer must be 4K aligned. The first logical processor entering the virtual environment initializes the buffer provided by the VMM. Subsequent processors can enter the virtual environment at any time and will not perform initialization to the buffer.

PAL\_VP\_ENV\_INFO must be called before this procedure to determine the configuration options and size requirements for the virtual environment. The VMM is required to maintain the ITR and DTR translations of the PAL virtual environment buffer throughout this procedure. See "[PAL\\_VP\\_ENV\\_INFO – PAL Virtual Environment Information \(266\)](#)" on page 2:473 for more information on PAL\_VP\_ENV\_INFO.

After this procedure, it is optional for the VMM to maintain the TR mapping for the PAL virtual environment buffer. If the TR translations for the buffer are not installed, the VMM must not make any PAL virtualization service calls; and the VMM must be prepared to handle DTLB faults during any PAL virtualization procedure calls.

[Table 11-119](#) shows the layout of the *config\_options* parameter. The *config\_options* parameter configures the global configuration options and global virtualization optimizations for all the logical processors in the virtual environment. All logical

processors in the virtual environment must specify the same value in the *config\_options* parameter during PAL\_VP\_INIT\_ENV, otherwise processor operation is undefined.

**Table 11-119. *config\_options* – Global Configuration Options**

	Field	Bit	Description
Global Configuration Options	initialize	0	If 1, this procedure will initialize the PAL virtual environment buffer for this virtual environment. If 0, this procedure will not initialize the PAL virtual environment buffer. On a multiprocessor system, the VMM must wait until this procedure completes on the first logical processor before calling this procedure on additional logical processors; otherwise processor operation is undefined.
	fr_pmc	1	If 1, for virtualization intercepts the performance counters are disabled by setting PSR.up and pp to 0, see <a href="#">Section 11.7.3.1, “PAL Virtualization Intercept Handoff State”</a> on page 2:333 for details on PSR settings at virtualization intercepts; for all other IVA-based interruptions PSR.pp and up are set according to Interruption State column described in Processor Status Field table described in <a href="#">Table 3-2, “Processor Status Register Fields”</a> on page 2:24. The VMM must have DCR.pp equal to 0 when the <i>fr_pmc</i> option is 1, whenever the IVA control register on the logical processor is set to point to the per-virtual-processor host IVT. See <a href="#">Section 11.7.2, “Interruption Handling in a Virtual Environment”</a> on page 2:331 and <a href="#">Table 11-21, “IVA Settings after PAL Virtualization-related Procedures and Services”</a> on page 2:332 for details on per-virtual-processor host IVT. If 0, PSR.pp and up are set according to Interruption State column described in Processor Status Field table described in <a href="#">Table 3-2, “Processor Status Register Fields”</a> on page 2:24
	be	2	Big-endian – Indicates the endian setting of the VMM. If 1, the values in the VPD are stored in big-endian format and the PAL services calls are made with PSR.be bit equal to 1. If 0, the values in the VPD are stored in little-endian format and the PAL services calls are made with PSR.be bit equal to 0. The VMM must match DCR.be with the value set in this field when the IVA control register on the logical processor is set to point to the per-virtual-processor host IVT. See <a href="#">Section 11.7.2, “Interruption Handling in a Virtual Environment”</a> on page 2:331 and <a href="#">Table 11-21, “IVA Settings after PAL Virtualization-related Procedures and Services”</a> on page 2:332 for details on per-virtual-processor host IVT.
	Reserved	7:3	Reserved.

**Table 11-119. *config\_options* – Global Configuration Options (Continued)**

	Field	Bit	Description
Global Virtualization Optimizations	opcode	8	This bit must be set to 1 – opcode information will be provided to the VMM during PAL intercepts within the virtual environment. This opcode may or may not be guaranteed to be the opcode that triggered the intercept. See <a href="#">Table 11-118, “vp_env_info – Virtual Environment Information Parameter”</a> on page 2:473 for details. This procedure returns an error if this bit is not set to 1.
	cause	9	If 1, the causes of virtualization intercepts will be provided to the VMM during PAL intercept handoffs within the virtual environment. No information will be provided if 0. See <a href="#">Section 11.7.3.1, “PAL Virtualization Intercept Handoff State”</a> on page 2:333 for details of virtualization intercept handoffs.
	gitc	10	If 1, enables guest MOV-from-AR.ITC optimization. For details see <a href="#">Section 11.7.4.1.3, “Guest MOV-from-AR.ITC Optimization”</a> on page 2:337 and <a href="#">Section 3.3.4.4, “Interval Timer Offset (ITO – CR4)”</a> on page 2:34. This bit is reserved if guest MOV-from-AR.ITC optimization is not supported.
	Reserved	62:11	Reserved.
	impl	63	Implementation-specific configuration option. This field is ignored if not implemented. Please refer to processor-specific documentation for details.

The *fr\_pmc* bit in the global *config\_options* parameter specifies whether the performance counters will be frozen when the Virtualization optimizations specified in the Virtualization Acceleration Control (*vac*) and Virtualization Disable Control (*vdc*) are running. When a virtual processor is running, the *vac* field in the corresponding VPD specifies whether a certain virtualization accelerations are enabled. If the *fr\_pmc* in the virtual environment was also enabled, the performance counters will be frozen when the enabled virtualization optimizations are running. See [Section 11.7.4, “Virtualization Optimizations”](#) on page 2:335 for details on Virtualization Acceleration Control (*vac*) and Virtualization Disable Control (*vdc*).

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, “Processor Virtualization”](#) on page 2:44 and [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)”](#) on page 2:446 for details.

## PAL\_VP\_REGISTER – PAL Register Virtual Processor (269)

**Purpose:** Register a different host IVT and/or a different optional virtualization intercept handler for the virtual processor specified by *vpd*.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_REGISTER within the list of PAL procedures
	vpd	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
	host_iva	64-bit host virtual pointer to the host IVT for the virtual processor
	opt_handler	64-bit non-zero host-virtual pointer to an optional handler for virtualization intercepts. See <a href="#">Section 11.7.3, “PAL Intercepts in Virtual Environment” on page 2:332</a> for details.

Returns:	Return Value	Description
	status	Return status of the PAL_VP_REGISTER procedure
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** PAL\_VP\_REGISTER registers a different host IVT and/or a different optional virtualization intercept handler specific to the virtual processor specified by *vpd*. On creation of a virtual processor by PAL\_VP\_CREATE, the VMM specifies a host IVT specific to the virtual processor. This procedure allows the VMM to specify a host IVT different from the one specified during PAL\_VP\_CREATE.

The host virtual to host physical translation of the 64K region specified by *vpd* must be mapped with either a DTR or DTC. See [Section 11.10.2.1.3, “Making PAL Procedure Calls in Physical or Virtual Mode” on page 2:359](#) for details on data translation requirements of memory buffer pointers passed as arguments to PAL procedures. The *virt\_env\_vaddr* parameter in the VPD must be setup with the host virtual address of the PAL virtual environment buffer before calling this procedure.

The *host\_iva* parameter specifies the host IVT to handle IVA-based interruptions when this virtual processor is running. The VMM can use the same or different *host\_iva* for each virtual processor. The *opt\_handler* specifies an optional virtualization intercept handler. If a non-zero value is specified, all virtualization intercepts are delivered to this handler. If a zero value is specified, all virtualization intercepts are delivered to the Virtualization vector in the host IVT. Upon completion of this procedure, the VMM must not relocate the IVT specified by the *host\_iva* parameter and/or the virtualization intercept handler specified by the *opt\_handler* parameter. The VMM can call this procedure again in case it wishes to associate a different host IVT and/or virtualization intercept handler with the virtual processor.

PAL\_VP\_REGISTER returns invalid argument on unsupported virtualization optimization combinations in *vpd*. See [Section 11.7.4.4, “Virtualization Optimization Combinations” on page 2:349](#) for details.

This procedure can be used by the VMM to:

## **PAL\_VP\_REGISTER**

- Relocate the host IVT associated with the virtual processor.
- Specify a different optional virtualization intercept handler for the virtual processor.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, "Processor Virtualization" on page 2:44](#) and ["PAL\\_PROC\\_GET\\_FEATURES - Get Processor Dependent Features \(17\)" on page 2:446](#) for details.

## PAL\_VP\_RESTORE – PAL Restore Virtual Processor (270)

**Purpose:** Restores virtual processor state for the specified *vpd* on the logical processor.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_RESTORE within the list of PAL procedures.
	vpd	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD.)
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VP_RESTORE procedure.
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** PAL\_VP\_RESTORE performs an implementation-specific restore operation of the virtual processor specified by the *vpd* parameter on the logical processor. The host virtual to host physical translation of the 64K region specified by *vpd* and the PAL virtual environment buffer must be mapped by instruction and data translation registers (TR). The instruction and data translation must be maintained until after the next invocation of PAL\_VP\_SAVE or PAL\_VPS\_SAVE and a different host IVT is set up by the VMM by writing to the IVA control register. PAL\_VP\_RESTORE configures the logical processor to run the specified virtual processor by loading implementation-specific virtual processor context from the VPD, and returns control back to the VMM.

This procedure performs an implicit PAL\_VPS\_SYNC\_WRITE; there is no need for the VMM to invoke PAL\_VPS\_SYNC\_WRITE unless the VPD values are modified before resuming the virtual processor. After the procedure, the caller is responsible for restoring all of the architectural state before resuming to the new virtual processor through PAL\_VPS\_RESUME\_NORMAL or PAL\_VPS\_RESUME\_HANDLER.

Upon completion of this procedure, the IVA-based interruptions will be delivered to the host IVT associated with this virtual processor.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, “Processor Virtualization” on page 2:44](#) and [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)” on page 2:446](#) for details.

## PAL\_VP\_SAVE – PAL Save Virtual Processor (271)

**Purpose:** Saves virtual processor state for the specified *vpd* on the logical processor.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_SAVE within the list of PAL procedures
	vpd	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
	Reserved	0
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VP_SAVE procedure
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** PAL\_VP\_SAVE performs an implementation-specific save operation of the virtual processor specified by the *vpd* parameter on the logical processor. The host virtual to host physical translation of the 64K region specified by *vpd* must be mapped by instruction and data translation registers (TR).

This procedure performs an implicit PAL\_VPS\_SYNC\_READ; there is no need for the VMM to invoke PAL\_VPS\_SYNC\_READ to synchronize the implementation-specific control resources before this procedure.

Upon completion of this procedure, the IVA-based interruptions will continue to be delivered to the host IVT associated with this virtual processor. After this procedure, the VMM can setup the IVA control register to use a different host IVT.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, "Processor Virtualization" on page 2:44](#) and ["PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)" on page 2:446](#) for details.



## PAL\_VP\_TERMINATE – PAL Terminate Virtual Processor (272)

**Purpose:** Terminates operation for the specified virtual processor.

**Calling Conv:** Stacked Registers

**Mode:** Virtual

**Buffer:** Dependent

Arguments:	Argument	Description
	index	Index of PAL_VP_TERMINATE within the list of PAL procedures
	vpd	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
	iva	Optional 64-bit host virtual pointer to the IVT when this procedure is done
	Reserved	0

Returns:	Return Value	Description
	status	Return status of the PAL_VP_TERMINATE procedure
	Reserved	0
	Reserved	0
	Reserved	0

Status:	Status Value	Description
	0	Call completed without error
	-1	Unimplemented procedure
	-2	Invalid argument
	-3	Call completed with error
	-9	Call requires PAL memory buffer

**Description:** Terminates operation of the virtual processor specified by *vpd* on the logical processor. The host virtual to host physical translation of the 64K region specified by *vpd* must be mapped by instruction and data translation registers (TR). See [Section 11.10.2.1.3, “Making PAL Procedure Calls in Physical or Virtual Mode” on page 2:359](#) for details on data translation requirements of memory buffer pointers passed as arguments to PAL procedures. All resources allocated for the execution of the virtual machine are freed.

Upon successful execution of PAL\_VP\_TERMINATE procedure and if the *iva* parameter is non-zero, the IVA control register will contain the value from the *iva* parameter.

This procedure returns unimplemented procedure when virtual machine features are disabled. See [Section 3.4, “Processor Virtualization” on page 2:44](#) and [“PAL\\_PROC\\_GET\\_FEATURES – Get Processor Dependent Features \(17\)” on page 2:446](#) for details.

## 11.11 PAL Virtualization Services

In order to support efficient handling of interruptions when PSR.vm was 1, a set of PAL virtualization services is defined to allow certain high-frequency PAL functions to be performed in a low-latency and low-overhead manner.

Upon successful completion of PAL\_VP\_INIT\_ENV, the virtual base address of the PAL virtualization services (VSA) is returned to the VMM. VMM can invoke PAL services by branching to the defined offsets from the virtual base address. See [Table 11-120](#) for the defined services. See [Section 11.11, “PAL Virtualization Services” on page 2:486](#) for details on PAL virtualization services.

These PAL virtualization services will only make references to the PAL virtual environment buffer. The VMM is required to maintain the ITR and DTR translations of the PAL virtual environment buffer during any PAL virtualization service calls.

**Table 11-120. PAL Virtualization Services**

Offset	PAL Service
0x0000	PAL_VPS_RESUME_NORMAL
0x0400	PAL_VPS_RESUME_HANDLER
0x0800	PAL_VPS_SYNC_READ
0x0c00	PAL_VPS_SYNC_WRITE
0x1000	PAL_VPS_SET_PENDING_INTERRUPT
0x1400	PAL_VPS_THASH
0x1800	PAL_VPS_TTAG
0x1c00	PAL_VPS_RESTORE
0x2000	PAL_VPS_SAVE
All other offsets	Reserved

### 11.11.1 PAL Virtualization Service Invocation Convention

This section describes the required parameters applicable to all PAL Virtualization Services. Additional parameters are listed in the description section of specific PAL Virtualization Services. Architectural state not listed in this section is managed by the VMM and can contain both VMM and/or virtual processor state. The architectural state not listed is unchanged by PAL virtualization services.

The state of the processor on handing off to any PAL Virtualization Service is:

- GR24-31: Parameters for PAL virtualization services.
- BRs:
  - BR0: Scratch, the VMM will use BR0 to specify the 64-bit host virtual address of the PAL Virtualization Service being invoked.
- Predicates: The predicates are preserved by the PAL virtualization services.
- PSR State (see [Table 11-121](#) for details):
  - PSR.be, i, cpl, is, ss, db, tb, vm must be 0.
  - PSR.dt, rt and it must be 1.
  - All other values are don't cares.

**Table 11-121. State Requirements for PSR for PAL Virtualization Services**

PSR Bit	Description	Value
be	big-endian memory access enable	_a
up	user performance monitor enable	-
ac	alignment check	-
mfl	floating-point registers f2-f31 written	-
mfh	floating-point registers f32-f127 written	-
ic	interruption state collection enable	0 <sup>b</sup> _c
i	interrupt enable	0
pk	protection key validation enable	-
dt	data address translation enable	1
dfi	disabled FP register f2 to f31	-
dfh	disabled FP register f32 to f127	-
sp	secure performance monitors	-
pp	privileged performance monitor enable	-
di	disable ISA transition	-
si	secure interval timer	-
db	debug breakpoint fault enable	0
lp	lower-privilege transfer trap enable	-
tb	taken branch trap enable	0
rt	register stack translation enable	1
cpl	current privilege level	0
is	instruction set	0
mc	machine check abort mask	-
it	instruction address translation enable	1
id	instruction debug fault disable	-
da	data access and dirty-bit fault disable	-
dd	data debug fault disable	-
ss	single step trap enable	0
ri	restart instruction	-
ed	exception deferral	-
bn	register bank	_d 0 <sup>e</sup>
ia	instruction access-bit fault disable	-
vm	processor virtualization	0

a. PAL services can be called with PSR.be bit equal to 0 or 1. The behavior is undefined if PSR.be setting does not match the *be* parameter during PAL\_VP\_INIT\_ENV. See “PAL\_VP\_INIT\_ENV – PAL Initialize Virtual Environment (268)” on page 2:478 for details.

b. Most PAL services are invoked with PSR.ic equal to 0.

- c. Specific PAL services can be invoked with PSR.ic equal to 1 or 0. See the description of specific PAL services for details.
- d. Most PAL services can be invoked with PSR.bn equal to 1 or 0.
- e. Specific PAL services must be invoked with PSR.bn equal to 0. See the description of specific PAL services for details.

### **11.11.2 PAL Virtualization Service Specifications**

The following pages provide detailed interface specifications for each of the PAL Virtualization Services.

## PAL\_VPS\_RESUME\_NORMAL – Resume Virtual Processor Normal (0x0000)

**Purpose:** Resumes the current virtual processor. This service is used when vpsr.ic is 1. This service can also be used independent of the state of vpsr.ic if all virtualization accelerations and disables are disabled.

Argument	Description
GR24	VBR0
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
GR26	Reserved
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

**Returns:** PAL\_VPS\_RESUME\_NORMAL does not return to the VMM.

**Description:** On interruptions or intercepts, PAL\_VPS\_RESUME\_NORMAL allows the VMM to resume the same virtual processor where the vpsr.ic is 1. PAL\_VP\_RESTORE can be used to restore the state of a different virtual processor.

The VMM specifies the VBR0 of the virtual processor in GR24 and the 64-bit virtual pointer to the VPD in GR25.

The VMM is responsible for setting up all the required virtual processor state in the architectural registers as well as in the VPD prior to invoking this service. See [Table 11-122, “Virtual Processor Settings in Architectural Resources for PAL\\_VPS\\_RESUME\\_NORMAL and PAL\\_VPS\\_RESUME\\_HANDLER” on page 2:489](#) for details.

PAL\_VPS\_RESUME\_NORMAL must be called with PSR.bn equal to 0.

If all virtualization accelerations and disables are disabled, PAL\_VPS\_RESUME\_NORMAL can also be used to resume to the guest independent on the state of vpsr.ic.

**Table 11-122. Virtual Processor Settings in Architectural Resources for PAL\_VPS\_RESUME\_NORMAL and PAL\_VPS\_RESUME\_HANDLER**

Resource	Description
Bank 1 GRs	Contains state of bank 0/1 GRs of the virtual processor (depends on vpsr.bn.)
FRs	Contains floating-point register state of the virtual processor.
Predicate Register	Contains the predicates of the virtual processor.
Branch Registers	BR1-BR7 contains the state of the virtual processor. BR0 of the virtual processor resides in bank 0 GR24.
Application Registers	Contains application register state of the virtual processor.
Interval Timer Offset Register <sup>a</sup>	If guest MOV-from-AR.ITC optimization is enabled, this register contains an offset, programmed by the VMM, to ensure that guest reads of ITC get the proper value.
Interruption Control Registers	IIP, IPSR and IFS contains the IP, PSR and CFM of the virtual processor. See <a href="#">Table 11-123</a> for the PSR settings for the execution of the virtual processor. The rest of the interruption control registers are don't cares. For PAL_VPS_RESUME_HANDLER, the virtual interruption control registers are specified in the VPD. See <a href="#">Section 11.7.4, “Virtualization Optimizations” on page 2:335</a> for synchronization of VPD resources before resuming the virtual processor.

**Table 11-122. Virtual Processor Settings in Architectural Resources for PAL\_VPS\_RESUME\_NORMAL and PAL\_VPS\_RESUME\_HANDLER**

Resource	Description
External Interrupt Control Registers	The external interrupt control registers contain the state of the virtual processor if d_extint in Virtualization Disable Control (vdc) is 1. Otherwise the external interrupt control registers are virtualized by the VMM and contain VMM state.
Data/Instruction Breakpoint Registers	The data/instruction breakpoint registers contain the state of the virtual processor if d_jbr_dbr in Virtualization Disable Control (vdc) is 1. Otherwise the data/instruction breakpoint registers are virtualized by the VMM and contain VMM state.
Performance Monitor Configuration Registers	The performance monitor configuration registers contain the state of the virtual processor if d_pmc in Virtualization Disable Control (vdc) is 1. Otherwise the performance monitor configuration registers are virtualized by the VMM and contain VMM state.
Performance Monitor Data Registers	Contain the state of the virtual processor.

- a. Interval Timer Offset register is not supported on all processor implementations. See [Section 3.3.4.4, “Interval Timer Offset \(ITO – CR4\)”](#) on page 2:34 for details.

**Table 11-123. Processor Status Register Settings for Virtual Processor Execution**

Field	Bits	Description
User Mask = PSR{5:0}		
rv	0	Reserved
be	1	Contain user mask of the virtual processor.
up	2	
ac	3	
mfl	4	
mfh	5	
System Mask = PSR{23:0}		
ic	13	Must be 1.
i	14	VMM-specific.
pk	15	
rv	12:6, 16	Reserved
dt	17	Must be 1.
dfl	18	VMM-specific.
dfh	19	
sp	20	
pp	21	
di	22	
si	23	
PSR.I = PSR{31:0}		
db	24	VMM-specific.
lp	25	Contains the lp bit of the virtual processor.
tb	26	Contains the tb bit of the virtual processor.
rt	27	Must be 1.
rv	31:28	Reserved
PSR{63:0}		

**Table 11-123. Processor Status Register Settings for Virtual Processor Execution (Continued)**

Field	Bits	Description
cpl	33:32	Contains the cpl field of the virtual processor.
is	34	VMM-specific.
mc	35	VMM-specific.
it	36	Must be 1.
id	37	VMM-specific.
da	38	VMM-specific.
dd	39	VMM-specific.
ss	40	VMM-specific.
ri	42:41	Contains the ri field of the virtual processor.
ed	43	Contains the ed bit of the virtual processor.
bn	44	Must be 1.
ia	45	VMM-specific.
vm	46	Must be 1.
rv	63:47	Reserved

PAL\_VPS\_RESUME\_NORMAL performs the following actions:

- Perform any implementation-specific setup to run a virtual processor.
- Re-enable performance counters if the value of the *fr\_pmc* field in the *config\_options* parameter passed to PAL\_VP\_INIT\_ENV was 1.
- Resume the virtual processor.

## PAL\_VPS\_RESUME\_HANDLER – Resume Virtual Processor Handler (0x0400)

**Purpose:** Resumes the current virtual processor. This service is used when `vpsr.ic` is 0.

Argument	Description
GR24	VBR0
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
GR26	Virtualization Acceleration Control ( <i>vac</i> ) field from the VPD specified in GR25 and CFLE setting at the target instruction.
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

**Returns:** PAL\_VPS\_RESUME\_HANDLER does not return to the VMM.

**Description:** On interruptions or intercepts, PAL\_VPS\_RESUME\_HANDLER allows the VMM to resume to the same virtual processor where the `vpsr.ic` is 0<sup>1</sup>.

GR24 specifies the BR0 of the virtual processor; GR25 specifies the 64-bit virtual pointer to the VPD; GR26 specifies the *vac* field of the VPD argument specified in GR25; bit 63 of GR26 specifies the value of CFLE setting at the target instruction. Behavior is undefined if the *vac* in GR26 does not match the *vac* field in the VPD argument specified in GR25.

The VMM is responsible for setting up all the required virtual processor state in the architectural registers as well as in the VPD prior to invoking this service. See [Table 11-122, “Virtual Processor Settings in Architectural Resources for PAL\\_VPS\\_RESUME\\_NORMAL and PAL\\_VPS\\_RESUME\\_HANDLER” on page 2:489](#) for details.

PAL\_VPS\_RESUME\_HANDLER must be called with `PSR.bn` equal to 0.

PAL\_VPS\_RESUME\_HANDLER performs the following actions:

- Perform any implementation-specific setup to run a virtual processor.
- Re-enable performance counters if the value of the *fr\_pmc* field in the *config\_options* parameter passed to PAL\_VP\_INIT\_ENV was 1.
- Resume the virtual processor.

---

1. PAL\_VP\_RESTORE can be used to restore the state of a different virtual processor.



## PAL\_VPS\_SYNC\_READ – Synchronize VPD State for Reads (0x0800)

**Purpose:** Synchronize VPD with the latest implementation-specific virtual architectural state.

Arguments:	Argument	Description
	GR24	64-bit host virtual return address
	GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
	GR26	Reserved
	GR27	Reserved
	GR28	Reserved
	GR29	Reserved
	GR30	Reserved
	GR31	Reserved

Returns:	Return Value	Description
	GR24	Scratch
	GR25	Scratch
	GR26	Scratch
	GR27	Scratch
	GR28	Scratch
	GR29	Scratch
	GR30	Scratch
	GR31	Scratch

**Description:** On processor implementations that support virtualization accelerations, implementation-specific control resources can be provided to enhance performance of virtual processors. When a specific acceleration is enabled, after interruptions and intercepts which occur when `PSR.vm` was 1, the VMM must invoke this service to synchronize the related resources before reading the value from the VPD. For the accelerations that are disabled, the corresponding resources in the VPD are unchanged.

The synchronization requirements of the related resources for each acceleration are described in the corresponding sections for each acceleration in [Section 11.7.4.2, “Virtualization Accelerations”](#) on page 2:337.

PAL\_VPS\_SYNC\_READ performs the following actions:

- Copy implementation-specific control resources of the enabled accelerations into VPD.
- Return to VMM by an indirect branch specified in the GR24 parameter.

## PAL\_VPS\_SYNC\_WRITE – Synchronize VPD State for Writes (0x0c00)

**Purpose:** Synchronize the implementation-specific virtual architectural state with VPD.

Argument	Description
GR24	64-bit host virtual return address.
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD.)
GR26	Reserved
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

Return Value	Description
GR24	Scratch
GR25	Scratch
GR26	Scratch
GR27	Scratch
GR28	Scratch
GR29	Scratch
GR30	Scratch
GR31	Scratch

**Description:** On processor implementations that support virtualization accelerations, implementation-specific control resources can be provided to enhance performance of virtual processors. When a specific acceleration is enabled, the VMM must invoke this service to synchronize the related resources after modifying the value in the VPD and before resuming the virtual processor. For the accelerations that are disabled, the corresponding resources in the VPD are ignored.

The synchronization requirements of the related resources for each acceleration are described in the corresponding sections for each acceleration in [Section 11.7.4.2, “Virtualization Accelerations” on page 2:337](#).

PAL\_VPS\_SYNC\_WRITE performs the following actions:

- Copy values of the enabled accelerations in the VPD into implementation-specific control resources.
- Return to VMM by an indirect branch specified in the GR24 parameter.

## PAL\_VPS\_SET\_PENDING\_INTERRUPT – Register Highest Priority Pending Interrupt (0x1000)

**Purpose:** Register highest priority pending interrupt of the running virtual processor.

Argument	Description
GR24	64-bit host virtual return address
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
GR26	Reserved
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

Return Value	Description
GR24	Scratch
GR25	Scratch
GR26	Scratch
GR27	Scratch
GR28	Scratch
GR29	Scratch
GR30	Scratch
GR31	Scratch

**Description:** PAL\_VPS\_SET\_PENDING\_INTERRUPT allows the VMM to register the highest priority pending interrupt for the virtual processor. The virtual highest priority pending interrupt is specified in the vhpi field in the VPD. See [Table 11-124, “vhpi – Virtual Highest Priority Pending Interrupt”](#) on page 2:495 for details.

PAL\_VPS\_SET\_PENDING\_INTERRUPT can be called with PSR.ic equal to 1 or 0.

**Table 11-124. vhpi – Virtual Highest Priority Pending Interrupt**

Value	Description
0	Nothing pending.
1	Class 1 interrupt pending.
2	Class 2 interrupt pending.
3	Class 3 interrupt pending.
4	Class 4 interrupt pending.
5	Class 5 interrupt pending.
6	Class 6 interrupt pending.
7	Class 7 interrupt pending.
8	Class 8 interrupt pending.
9	Class 9 interrupt pending.
10	Class 10 interrupt pending.
11	Class 11 interrupt pending.
12	Class 12 interrupt pending.
13	Class 13 interrupt pending.
14	Class 14 interrupt pending.
15	Class 15 interrupt pending.
16	ExtINT pending.
17-31	Reserved.
32	NMI pending.
33+	Reserved.

## ***PAL\_VPS\_SET\_PENDING\_INTERRUPT***

PAL\_VPS\_SET\_PENDING\_INTERRUPT performs the following actions:

- Copy the virtual highest priority pending interrupt from the VPD into implementation-specific resources.
- Return to VMM by an indirect branch specified in the GR24 parameter.

## PAL\_VPS\_THASH – Compute Long Format VHPT Entry Address (0x1400)

**Purpose:** Compute a long format VHPT entry address.

Argument	Description
GR24	64-bit host virtual return address
GR25	64-bit virtual address used to compute the hash entry address
GR26	Region register value used to compute the hash entry address
GR27	Virtual PTA
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

Return Value	Description
GR24	Scratch
GR25	Scratch
GR26	Scratch
GR27	Scratch
GR28	Scratch
GR29	Scratch
GR30	Scratch
GR31	64-bit VHPT entry address

**Description:** PAL\_VPS\_THASH computes a long format Virtual Hashed Page Table (VHPT) entry address based on the input arguments and the result is returned in GR31. The format of the region register parameter (GR26) is defined in [Section 4.1.2, “Region Registers \(RR\)” on page 2:58](#), the *ve* field is ignored by the service. The format of the Virtual PTA parameter (GR27) is defined in [Section 3.3.4.6, “Page Table Address \(PTA – CR8\)” on page 2:35](#), the *vf* field is ignored by the service.

PAL\_VPS\_THASH returns the same long format VHPT entry address given the same input arguments across different implementations. The long format VHPT entry address returned may not be the same as the long format VHPT entry address generated by the `thash` instruction of the processor.

PAL\_VPS\_THASH can be called with PSR.ic equal to 1 or 0.

## PAL\_VPS\_TTAG – Compute Translated Hashed Entry Tag (0x1800)

**Purpose:** Compute the long format translated hashed entry tag.

Arguments:	Argument	Description
	GR24	64-bit host virtual return address
	GR25	64-bit virtual address used to compute the hash entry tag
	GR26	Region register value used to compute the hash entry tag
	GR27	Reserved
	GR28	Reserved
	GR29	Reserved
	GR30	Reserved
	GR31	Reserved

Returns:	Return Value	Description
	GR24	Scratch
	GR25	Scratch
	GR26	Scratch
	GR27	Scratch
	GR28	Scratch
	GR29	Scratch
	GR30	Scratch
	GR31	64-bit VHPT entry tag

**Description:** PAL\_VPS\_TTAG computes the tag value of the long format Virtual Hashed Page Table (VHPT) based on the input arguments and the result is returned in GR31. The format of the region register parameter (GR26) is defined in [Section 4.1.2, "Region Registers \(RR\)"](#) on page 2:58, the ve field is ignored by the service.

PAL\_VPS\_TTAG returns the same tag value given the same input arguments across different implementations. The tag value returned may not be the same as the tag value generated by the `ttag` instruction of the processor.

PAL\_VPS\_TTAG can be called with PSR.ic equal to 1 or 0.

## PAL\_VPS\_RESTORE – Fast Restore Virtual Processor State (0x1c00)

**Purpose:** Performs an implementation-specific light-weight restore operation for the specified VPD on the logical processor.

Argument	Description
GR24	64-bit host virtual return address
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
GR26	Skip implicit synchronization
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

Return Value	Description
GR24	Scratch
GR25	Scratch
GR26	Scratch
GR27	Scratch
GR28	Scratch
GR29	Scratch
GR30	Scratch
GR31	Scratch

**Description:** PAL\_VPS\_RESTORE performs an implementation-specific light-weight restore operation of the virtual processor specified by the VPD parameter (GR25) on the logical processor. The host virtual to host physical translation of the 64K region specified by the VPD parameter (GR25) and the PAL virtual environment buffer must be mapped by instruction and data translation registers (TR). The instruction and data translation must be maintained until after the next invocation of PAL\_VP\_SAVE or PAL\_VPS\_SAVE and a different host IVT is set up by the VMM by writing to the IVA control register. PAL\_VPS\_RESTORE configures the logical processor to run the specified virtual processor by loading the minimal implementation-specific virtual processor context from the VPD, and returns control back to the VMM.

If GR26 is zero, this service performs an implicit PAL\_VPS\_SYNC\_WRITE; there is no need for the VMM to invoke PAL\_VPS\_SYNC\_WRITE to synchronize the implementation-specific control resources before this service. If GR26 is one (0x1), no implicit synchronization will be performed by this service.

Upon completion of this service, the IVA-based interruptions will be delivered to the host IVT associated with this virtual processor.

This service does not restore any PAL procedure implementation-specific state<sup>1</sup>. The caller of this service is responsible to manage the difference in settings for the PAL procedures between the VMM and virtual processors.

---

1. PAL\_VP\_RESTORE can be used to restore PAL procedure implementation-specific state. See [“PAL\\_VP\\_RESTORE – PAL Restore Virtual Processor \(270\)” on page 2:483](#) for details.

## PAL\_VPS\_SAVE – Fast Save Virtual Processor State (0x2000)

**Purpose:** Performs an implementation-specific light-weight save operation for the specified VPD on the logical processor.

Argument	Description
GR24	64-bit host virtual return address
GR25	64-bit host virtual pointer to the Virtual Processor Descriptor (VPD)
GR26	Skip implicit synchronization
GR27	Reserved
GR28	Reserved
GR29	Reserved
GR30	Reserved
GR31	Reserved

Return Value	Description
GR24	Scratch
GR25	Scratch
GR26	Scratch
GR27	Scratch
GR28	Scratch
GR29	Scratch
GR30	Scratch
GR31	Scratch

**Description:** PAL\_VPS\_SAVE performs an implementation-specific light-weight save operation of the virtual processor specified by the VPD parameter (GR25) on the logical processor. The host virtual to host physical translation of the 64K region specified by the VPD parameter (GR25) must be mapped by instruction and data translation registers (TR).

If GR26 is zero, this service performs an implicit PAL\_VPS\_SYNC\_READ; there is no need for the VMM to invoke PAL\_VPS\_SYNC\_READ to synchronize the implementation-specific control resources before this service. If GR26 is one (0x1), no implicit synchronization will be performed by this service.

Upon completion of this service, the IVA-based interruptions will continue to be delivered to the host IVT associated with this virtual processor. After this service, the VMM can setup the IVA control register to use a different host IVT.

This service does not save any PAL procedure implementation-specific state<sup>1</sup>. The caller of this service is responsible to manage the difference in settings for the PAL procedures between the VMM and virtual processors.

§

1. PAL\_VP\_SAVE can be used to save PAL procedure implementation-specific state. See “PAL\_VP\_SAVE – PAL Save Virtual Processor (271)” on page 2:484 for details.



## ***Part II: System Programmer's Guide***



*Part II: System Programmer's Guide* is intended as a companion section to the information presented in [Part I: "System Architecture Guide"](#). While *Part I* provides a crisp and concise architectural definition of the Itanium instruction set, *Part II* provides insight into programming and usage models of the Itanium system architecture. This section emphasizes how the various architecture features fit together and explains how they contribute to high performance system software.

The intended audience for this section is system programmers who would like to better understand the Itanium system architecture. The goal of this document is to:

- Familiarize system programmers with Itanium system architecture principles and usage models.
- Provide recommendations, code examples, and performance guidelines.

This section does not re-define the Itanium instruction set. Please refer to [Part I: "System Architecture Guide"](#) as the authoritative definition of the system architecture.

The reader is expected to be familiar with the contents of *Part I* and is expected to be familiar with modern virtual memory and multiprocessing concepts. Furthermore, this document is platform architecture neutral (i.e. no assumptions are made about platform architecture capabilities, such as busses, chipsets, or I/O devices).

## 1.1 Overview of the System Programmer's Guide

The Itanium architecture provides numerous performance enhancing features of interest to the system programmer. Many of these instruction set features focus on reducing overhead in common situations. The chapters outlined below discuss different aspects of the Itanium system architecture.

[Chapter 2, "MP Coherence and Synchronization"](#) describes Itanium architecture-based multiprocessing synchronization primitives and the Itanium memory ordering model. This chapter also discusses programming rules for self- and cross-modifying code. This chapter is useful for application and system programmers who write multi-threaded code.

[Chapter 3, "Interruptions and Serialization"](#) discusses how the Itanium architecture, despite its explicitly parallel instruction execution semantics, provides the system programmer with a precise interruption model. This chapter describes how the processor serializes execution around interruptions and what state is preserved and made available to low-level system code when interruptions are taken. This chapter introduces the interrupt vector table and describes how low-level kernel code is expected to transfer control to higher level operating system code written in a high-level programming language. This chapter is useful for operating system and firmware programmers.

[Chapter 4, “Context Management”](#) describes how operating systems need to preserve Itanium register contents. In addition to spilling and filling a register’s data value, the Itanium architecture also requires software to preserve control and data speculative state associated with that register, i.e. its NaT bit and ALAT state. This chapter also discusses system architecture mechanisms that allow an operating system to significantly reduce the number of registers that need to be spilled/filled on interruptions, system calls, and context switches. These optimizations improve the performance of an Itanium architecture-based operating system by reducing the amount of required memory traffic. This chapter is useful for operating system programmers.

[Chapter 5, “Memory Management”](#) introduces various memory management strategies in the Itanium architecture: region register model, protection keys, and the virtual hash page table usage models are described. This chapter is of interest to virtual memory management software developers.

[Chapter 6, “Runtime Support for Control and Data Speculation”](#) describes the operating system support that is required for control and data speculation. This chapter describes various speculation software models and their associated operating system implications. This chapter is of interest to operating system developers and compiler writers.

[Chapter 7, “Instruction Emulation and Other Fault Handlers”](#) describes a variety of instruction emulation handlers that Itanium architecture-based operating systems are expected to support. This chapter is useful for operating system developers.

[Chapter 8, “Floating-point System Software”](#) discusses how processors based on the Itanium architecture handle floating-point numeric exceptions and how the Itanium architecture-based software stack provides complete IEEE-754 compliance. This includes a discussion of the floating-point software assist firmware, the FP SWA EFI driver. This chapter also describes how Itanium architecture-based operating systems are expected to support IEEE floating-point exception filters. This chapter is useful for operating system developers and floating-point numerics experts.

[Chapter 9, “IA-32 Application Support”](#) outlines how software needs to perform instruction set transitions, and what low-level kernel handlers are required in an Itanium architecture-based operating system to support IA-32 applications. This chapter is useful for operating system developers.

[Chapter 10, “External Interrupt Architecture”](#) describes the external interrupt architecture with a focus on how external asynchronous interrupt handling can be controlled by software. Basic interrupt prioritization, masking, and harvesting capabilities are discussed in this chapter. This chapter is of interest to operating system developers and to device driver writers.

[Chapter 11, “I/O Architecture”](#) describes the I/O architecture with a focus on platform considerations and support for the existing IA-32 I/O port space platform infrastructure. This chapter is of interest to operating system developers and to device driver writers.

[Chapter 12, “Performance Monitoring Support”](#) describes the performance monitor architecture with a focus on what kind of operating system support is needed from Itanium architecture-based operating systems. This chapter is of interest to operating system and performance tool developers.

Chapter 13, “Firmware Overview” introduces the firmware model and how various firmware layers (PAL, SAL, UEFI, ACPI) work together to enable processor and system initialization and operating system boot. This chapter also discusses how firmware layers and the operating system work together to provide error detection, error logging, as well as fault containment capabilities. This chapter is of interest to platform firmware and operating system developers.

## 1.2 Related Documents

The following documents are referred to fairly often in this document. For more details on software conventions and platform firmware, please consult these manuals (available at <http://developer.intel.com>).

[SWC] **Intel® Itanium® Software Conventions and Runtime Architecture Guide**

[UEFI] *Unified Extensible Firmware Interface Specification*

[SAL] **Intel® Itanium® Processor Family System Abstraction Layer Specification**

§



This chapter describes how to enforce an ordering of memory operations, how to update code images, and presents examples of several simple multiprocessor synchronization primitives on a processor based on the Itanium architecture. These topics are relevant to anyone who writes either user- or system-level software for multiprocessor systems based on the Itanium architecture.

The chapter begins with a brief overview of Itanium memory access instructions intended to summarize the behaviors that are relevant to later discussions in the chapter. Next, this chapter presents the Itanium memory ordering model and compares it to a sequentially-consistent ordering model. It then explores versions of several common synchronization primitives. This chapter closes by describing how to correctly update code images to implement self-modifying code, cross-modifying code, and paging of code using programmed I/O.

## 2.1 An Overview of Intel® Itanium® Memory Access Instructions

The Itanium architecture provides load, store, and semaphore instructions to access memory. In addition, it also provides a memory fence instruction to enforce further ordering relationships between memory accesses. As [Section 4.4.7, “Memory Access Ordering” on page 1:73](#) describes, memory operations in the Itanium architecture come with one of four semantics: unordered, acquire, release, or fence. [Section 2.2 on page 2:510](#) describes how the memory ordering model uses these semantics to indicate how memory operations can be ordered with respect to each other.

[Section 2.1.1](#) defines the four memory operation semantics. [Section 2.2](#), [Section 2.3](#), and [Section 2.4](#) present brief outlines of load and store, semaphore, and memory fence instructions in the Itanium architecture. Refer to [Chapter 2, “Instruction Reference”](#) for more information on the behavior and capabilities of these instructions.

### 2.1.1 Memory Ordering of Cacheable Memory References

The Itanium architecture has a relaxed memory ordering model which provides unordered memory opcodes, explicitly ordered memory opcodes, and a fencing operation that software can use to implement stronger ordering. Each memory operation establishes an ordering relationship with other operations through one of four semantics:

- *Unordered* semantics imply that the instruction is made visible in any order with respect to other orderable instructions.
- *Acquire* semantics imply that the instruction is made visible prior to all subsequent orderable instructions.
- *Release* semantics imply that the instruction is made visible after all prior orderable instructions.

- *Fence* semantics combine acquire and release semantics (i.e. the instruction is made visible after all prior orderable instructions and before all subsequent orderable instructions).

In the above definitions “prior” and “subsequent” refer to the program-specified order. An “orderable instruction” is an instruction that the memory ordering model can use to establish ordering relationships<sup>1</sup>. The term “visible” refers to all architecturally-visible (from the standpoint of multiprocessor coherency) effects of performing an instruction. Specifically,

- Accesses to uncacheable or write-coalescing memory regions are visible when they reach the processor bus.
- Loads from cacheable memory regions are visible when they hit a non-programmer-visible structure such as a cache or store buffer.
- Stores to cacheable memory regions are visible when they enter a snooped (in a multiprocessor coherency sense) structure.

Memory access instructions typically have an ordered and an unordered form (i.e. a form with unordered semantics and a form with either acquire, release, or fence semantics). The Itanium architecture does not provide all possible combinations of instructions and ordering semantics. For example, the Itanium instruction set does not contain a store with fence semantics.

[Section 4.4.7, “Memory Access Ordering” on page 1:73](#) and [Section 4.4.7, “Sequentiality Attribute and Ordering” on page 2:82](#) discuss ordering, orderable instructions, and visibility in greater depth.

[Section 2.2 on page 2:510](#) describes how the ordering semantics affect the Itanium memory ordering model.

## 2.1.2 Loads and Stores

In the Itanium architecture, a load instruction has either unordered or acquire semantics while a store instruction has either unordered or release semantics. By using acquire loads (`ld.acq`) and release stores (`st.rel`), the memory reference stream of an Itanium architecture-based program can be made to operate according to the IA-32 ordering model. The Itanium architecture uses this behavior to provide IA-32 compatibility. That is, an Itanium acquire load is equivalent to an IA-32 load and an Itanium release store is equivalent to an IA-32 store, from a memory ordering perspective.

Loads can be either speculative or non-speculative. The speculative forms (`ld.s`, `ld.sa`, and `ld.a`) support control and data speculation.

## 2.1.3 Semaphores

The Itanium architecture provides a set of three semaphore instructions: exchange (`xchg`), compare and exchange (`cmpxchg`), and fetch and add (`fetchadd`). Both `cmpxchg` and `fetchadd` may have either acquire or release semantics depending on the

1. The ordering semantics of an instruction *do not* imply the orderability of the instruction. Specifically, unordered ordering semantics alone *do not* make an instruction unordered; there are orderable instructions with each of the four ordering semantics.



specific opcode chosen. The `xchg` instruction always has acquire semantics. These instructions read a value from memory, modify this value using an instruction-specific operation, and then write the modified value back to memory. The read-modify-write sequence is atomic by definition.

### 2.1.3.1 Considerations for using Semaphores

The memory location on which a semaphore instruction operates on must obey two constraints. First, the location must be cacheable (the `fetchadd` instruction is an exception to this rule; it may also operate on exported uncacheable locations, UCE). Thus, with the exception of `fetchadd` to UCE locations, the Itanium architecture does not support semaphores in uncacheable memory. Second, the location must be naturally-aligned to the size of the semaphore access. If either of these two constraints are not met, the processor generates a fault.

The exported uncacheable memory attribute, UCE, allows a processor based on the Itanium architecture to export fetch and add operations to the platform. A processor that does not support exported `fetchadd` will fault when executing a `fetchadd` to a UCE memory location. If the processor supports exported `fetchadd` but the platform does not, the behavior is undefined when executing a `fetchadd` to a UCE memory location.

Sharing locks between IA-32 and Itanium architecture-based code does work with the following restrictions:

- Itanium architecture-based code can only manipulate an IA-32 semaphore if the IA-32 semaphore is aligned.
- Itanium architecture-based code can only manipulate an IA-32 semaphore if the IA-32 semaphore is allocated in write-back cacheable memory.

An Itanium architecture-based operating system can emulate IA-32 uncacheable or misaligned semaphores by using the technique described in the next section.

### 2.1.3.2 Behavior of Uncacheable and Misaligned Semaphores

A processor based on the Itanium architecture raises an Unsupported Data Reference fault if it executes a semaphore that accesses a location with a memory attribute that the semaphore does not support.

If the alignment requirement for Itanium architecture-based semaphores is not met, a processor based on the Itanium architecture raises an Unaligned Data Reference fault. This fault is taken regardless of the setting of the user mask alignment checking bit, `UM.ac`.

The `DCR.lc` bit controls how the processor behaves when executing an atomic IA-32 memory reference under an external bus lock. When the `DCR.lc` bit (see [Section 3.3.4.1, "Default Control Register \(DCR - CR0\)"](#)) is 1 and an IA-32 atomic memory reference requires a non-cacheable or misaligned read-modify-write operation, an `IA_32_Intercept(Lock)` fault is raised. Such memory references require an external bus lock to execute correctly. To preserve `LOCK` pin functionality, an Itanium architecture-based operating system can virtualize the bus lock by implementing a shared cacheable global `LOCK` variable.

To support existing IA-32 atomic read-modify-write operations that require the `LOCK` pin, an Itanium architecture-based operating system can use the `DCR.lc` bit to intercept all external IA-32 read-modify-write operations. Then, the `IA_32_Intercept(Lock)` handler can emulate these operations by first acquiring a cacheable virtualized `LOCK` variable, then performing the required memory operations non-atomically, and then releasing the virtualized `LOCK` variable. This emulation allows the read-modify-write sequence to appear atomic to other processors that use the semaphore.

### 2.1.4 Memory Fences

The memory fence instruction (`mf`) is the only instruction in the Itanium instruction set with fence semantics. This instruction serializes the set of memory accesses before the memory fence in program order with respect to the set of memory accesses that follow the fence in program order.

## 2.2 Memory Ordering in the Intel® Itanium® Architecture

Understanding a system's memory ordering model is key to writing either user- or system-level multiprocessor software that uses shared memory to communicate between processes and also that executes correctly on a shared-memory multiprocessor system. For a general introduction to memory ordering models, see Adve and Gharachorloo [AG95].

Four factors determine how a processor or system based on the Itanium architecture orders a group of memory operations with respect to each other:

- *Data dependencies* define the relationship between operations from the same processor that have register or memory dependencies on the same address<sup>1</sup>. This relationship need only be honored by the local processor (i.e. the processor that executes the operations).
- The *memory ordering semantics* define the relationship between memory operations from a particular processor that reference different addresses. For cacheable references, this relationship is honored by *all* observers in the coherence domain.
- Aligned *release stores* and *semaphore operations* (both require and release forms) become visible to all observers in the coherence domain in a single total order except each processor may observe its own release stores (via loads or acquire loads) prior to their being observed globally<sup>2</sup>.
- Non-programmer-visible state, such as *store buffers*, *processor caches*, or any logically-equivalent structure, may satisfy read requests from loads or acquire loads on the local processor before the data in the structure is made globally visible to other observers.

- 
1. That is, A precedes B in program order and A produces a value that B consumes. This relationship is transitive.
  2. Consequently, each such operation appears to become visible to each observer in the coherence domain at the same time, with the exception that a release store can become visible to the storing processor before others.

In the Itanium architecture, dependencies between operations by a processor have implications for the ordering of those operations at that processor. The discussion in [Section 2.2.1.6](#) on [page 2:515](#) and [Section 2.2.1.7](#) on [page 2:516](#) explores this issue in greater depth.

The following sections examine the Itanium ordering model in detail. [Section 2.2.1](#) presents several memory ordering executions to illustrate important behaviors of the model. [Section 2.2.2](#) discusses how memory attributes and the ordering model interact. Finally, [Section 2.2.3](#) describes how the Itanium memory ordering model compares with other memory ordering models.

## 2.2.1 Memory Ordering Executions

Multiprocessor software that uses shared memory to communicate between processes often makes assumptions about the order in which other agents in the system will observe memory accesses. As [Section 2.1.1](#) on [page 2:507](#) describes, the Itanium architecture provides a rich set of ordering semantics that allows software to express different ordering constraints on a memory operation, such as a load. Writing correct multiprocessor software requires that the programmer (or compiler) select the ordering semantic appropriate to enforce the expected behavior.

For example, an algorithm that requires two store operations A and B become visible to other processors in the order {A, B} will use stores with different ordering semantics than an algorithm that does not require any particular ordering of A and B. Although it is always safe to enforce stricter ordering constraints than an algorithm requires, doing so may lead to lower performance. If the ordering of memory operations is not important, software should use unordered ordering semantics whenever possible for best possible performance.

This section presents multiprocessor executions to demonstrate the ordering behaviors that the Itanium architecture allows and to contrast the Itanium ordering model with other ordering models. The executions consist of sequences of memory accesses that execute on two or more processors and highlight outcomes that the Itanium memory ordering model either allows or disallows once all accesses on all processors complete. A programmer can use these executions as a guide to determine which Itanium memory ordering semantics are appropriate to ensure a particular visibility order of memory accesses.

[Section 2.2.1.1](#) presents the assumptions and notational conventions that the upcoming discussions use to examine the executions. The remaining eleven sections each explore a different facet of the Itanium ordering model:

- Relaxed ordering of unordered memory operations ([Section 2.2.1.2](#)).
- Using acquire and release semantics to order operations ([Section 2.2.1.3](#)).
- Loads may pass stores ([Section 2.2.1.4](#)) and how to prevent this behavior ([Section 2.2.1.5](#)).
- When dependencies do or do not establish memory ordering ([Section 2.2.1.6](#) and [Section 2.2.1.7](#)).
- Satisfying loads from store buffers ([Section 2.2.1.8](#)) and how to prevent this behavior ([Section 2.2.1.9](#)).
- Semaphore operations and local bypass ([Section 2.2.1.10](#)).

- Global visibility order of memory operations ([Section 2.2.1.11](#) and [Section 2.2.1.12](#)).

This presentation is organized to begin with simple behaviors and move to increasingly complex behaviors.

### 2.2.1.1 Assumptions and Notation

The discussions of the multiprocessor executions in the upcoming sections adopt two main notational conventions.

First, the memory accesses in the executions in this document are written using a pseudo-Itanium architecture-based assembly language that allows a store to write an immediate operand to memory. All memory locations are cacheable and aligned. Unless stated otherwise, memory locations do not overlap. Initially, all registers and memory locations contain zero.

Second, given two different memory operations  $X$  and  $Y$ ,  $X \gg Y$  specifies that  $X$  precedes  $Y$  in program order and  $X \rightarrow Y$  indicates that  $X$  is visible if  $Y$  is visible (i.e.  $X$  becomes visible before  $Y$ ).

Using this notation, [Figure 2-1](#) expresses the Itanium ordering semantics from [Section 2.1.1, “Memory Ordering of Cacheable Memory References” on page 2:507](#) and also [Section 4.4.7, “Memory Access Ordering” on page 1:73](#). There are no implications regarding the ordering of the visibility for the following pairs of operations: a release followed by an unordered operation; a release followed by an acquire; an unordered operation followed by another; or an unordered operation followed by an acquire.

**Figure 2-1. Intel® Itanium® Ordering Semantics**

Acquire $\gg$ X $\Rightarrow$ Acquire $\rightarrow$ X
X $\gg$ Release $\Rightarrow$ X $\rightarrow$ Release
X $\gg$ Fence $\Rightarrow$ X $\rightarrow$ Fence
Fence $\gg$ Y $\Rightarrow$ Fence $\rightarrow$ Y

In [Figure 2-1](#), “Acquire”, “Release”, and “Fence” represent an orderable instruction with the corresponding memory ordering semantics whereas “X” and “Y” indicate any orderable instruction.

### 2.2.1.2 The Intel® Itanium® Architecture Provides a Relaxed Ordering Model

The Itanium memory ordering model is a relaxed model. As a result, the Itanium architecture permits any outcome when executing the code shown in [Table 2-1](#).

**Table 2-1. Intel® Itanium® Architecture Provides a Relaxed Ordering Model**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld	r1 = [y]	// M3
st	[y] = 1	// M2	ld	r2 = [x]	// M4

*Outcomes: all are allowed*

Because all of the operations in [Table 2-1](#) are unordered, the Itanium memory ordering model does not place any constraints on the order in which a processor based on the Itanium architecture makes the operations visible.

Observing a particular value in `r2`, for example, does not allow any inferences to be made about the value of `r1` because the pair of stores on Processor #0 may become visible in any order. Therefore, all outcomes are possible as the system may interleave M1, M2, M3, and M4 in any order without violating the memory ordering constraints.

### 2.2.1.3 Enforcing Basic Ordering

Using acquire and release ordering semantics enforces an ordering between both the Processor #0 operations M1 and M2 and the Processor #1 operations M3 and M4 from the [Table 2-1](#) execution as shown in [Table 2-1](#).

**Table 2-2. Acquire and Release Semantics Order Intel® Itanium® Memory Operations**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld.acq	r1 = [y]	// M3
st.rel	[y] = 1	// M2	ld	r2 = [x]	// M4

*Outcome: only r1 = 1 and r2 = 0 is not allowed*

The Itanium ordering model only disallows the outcome `r1 = 1` and `r2 = 0` in this execution. The release semantics on M2 and acquire semantics on M3 affect the following ordering constraints:

M1 → M2  
M3 → M4

Given the code in [Table 2-2](#), these two ordering constraints along with the assumption that the outcome is `r1 = 1` and `r2 = 0` together imply that:

$$r1 = 1 \Rightarrow M2 \rightarrow M3 \Rightarrow M1 \rightarrow M4 \text{ (because } M1 \rightarrow M2 \text{ and } M3 \rightarrow M4) \Rightarrow r2 = 1$$

This contradicts the postulated outcome `r1 = 1` and `r2 = 0` and thus the Itanium ordering model disallows the `r1 = 1` and `r2 = 0` outcome.

In operational terms, if Processor #1 observes M2, the release store to `y` (i.e. `r1` is 1), it must have also observed M1, the unordered store to `x` (i.e. `r2` is 1 as well), given the ordering constraints. Therefore, the Itanium ordering model must disallow the outcome `r1 = 1` and `r2 = 0` in this execution as this outcome violates these constraints.

Stronger ordering models that do not relax load-to-load and store-to-store ordering, such as sequential consistency, impose these same ordering constraints on M1, M2, M3, and M4 and therefore also do not allow the outcome `r1 = 1` and `r2 = 0`.

### 2.2.1.4 Allow Loads to Pass Stores to Different Locations

The Itanium memory ordering model allows loads to pass stores as shown in the execution sequence in [Table 2-3](#). Permitting this behavior can improve performance by allowing the processor to complete loads that follow a store that misses the cache.

The Itanium ordering semantics always allow a processor to make operations that follow a release visible before the release and to make operations that precede an acquire visible after the acquire.

**Table 2-3. Loads May Pass Stores to Different Locations**

Processor #0			Processor #1		
st.rel	[x] = 1	// M1	st.rel	[y] = 1	// M3
ld.acq	r1 = [y]	// M2	ld.acq	r2 = [x]	// M4

Outcomes: all are allowed

Like the execution shown in Table 2-1, the Itanium memory ordering model does not place any constraints on the ordering of the operations on each processor in this execution either.

Therefore, for reasons similar to those given in Section 2.2.1.2 for the execution shown in Table 2-1, the Itanium memory ordering model allows any outcome in this execution as well. Further, the Itanium memory ordering model also allows all outcomes in similar executions that differ only in the ordering semantics of the load and store operations (e.g. those that replace M1 with an unordered store, etc.). There is no combination of legal ordering semantics on these operations (recall that the Itanium instruction set does not provide stores with acquire or fence semantics) that enforce either  $M1 \rightarrow M2$  or  $M3 \rightarrow M4$ .

### 2.2.1.5 Preventing Loads from Passing Stores to Different Locations

The only way to prevent the loads from moving ahead of the stores in the Table 2-3 execution is to separate them with a memory fence as the execution in Table 2-4 illustrates.

**Table 2-4. Loads May Not Pass Stores in the Presence of a Memory Fence**

Processor #0			Processor #1		
st	[x] = 1	// M1	st	[y] = 1	// M4
mf		// M2	mf		// M5
ld	r1 = [y]	// M3	ld	r2 = [x]	// M6

Outcome: only  $r1 = 0$  and  $r2 = 0$  is not allowed

The Itanium memory ordering model only disallows the outcome  $r1 = 0$  and  $r2 = 0$  in this execution. The memory fences on Processor #0 and Processor #1 (operations M2 and M5) force the load and store memory accesses to be made visible in program order; no re-ordering is permitted across the fence. Thus, the following ordering constraints must be met:

$$M1 \rightarrow M2 \rightarrow M3$$

$$M4 \rightarrow M5 \rightarrow M6$$

Given the code in Table 2-4, these two constraints along with the assumption that the outcome is  $r1 = 0$  and  $r2 = 0$  together imply that

$$r1 = 0 \Rightarrow M3 \rightarrow M4 \Rightarrow M3 \rightarrow M6 \text{ because } M4 \rightarrow M5 \rightarrow M6$$

$$r1 = 0 \Rightarrow M1 \rightarrow M3 \text{ because } M1 \rightarrow M2 \rightarrow M3$$

$$M1 \rightarrow M3 \text{ and } M3 \rightarrow M6 \Rightarrow M1 \rightarrow M6 \Rightarrow r2 = 1$$

This contradicts the postulated outcome  $r1 = 0$  and  $r2 = 0$  and thus the Itanium memory ordering model disallows the  $r1 = 1$  and  $r2 = 0$  outcome. Specifically, if M3 reads 0, then M4, M5, and M6 may not yet be visible but M1 and M2 must be visible. Thus, when M6 becomes visible it must observe  $x = 1$  because M1 is already visible.

### 2.2.1.6 Data Dependency Does Not Establish MP Ordering

The dependency rules define the relationship between memory operations that access the same address. Specifically, the Itanium architecture resolves read-after-write (RAW), write-after-read (WAR), and write-after-write (WAW) dependencies through memory in program order on the local processor. As [Section 2.2](#) discusses, dependencies are fundamentally different from the ordering semantics even though both affect ordering relationships between groups of memory accesses.

The execution shown in [Table 2-5](#) illustrates this difference.

**Table 2-5. Dependencies Do Not Establish MP Ordering (1)**

Processor #0			Processor #1		
st	[x] = 1 ;;	// M1	ld.acq	r2 = [y]	// M4
ld	r1 = [x] ;;	// M2	ld	r3 = [x]	// M5
st	[y] = r1 ;;	// M3			

Outcomes:  $r1 = 1$ ,  $r2 = 1$ , and  $r3 = 0$  is allowed

The following discussion focuses on the outcome  $r1 = 1$ ,  $r2 = 1$ , and  $r3 = 0$ . This outcome is allowed only because the Itanium architecture treats data dependencies and the ordering semantics differently.

The ordering semantics require  $M4 \rightarrow M5$ , but do not place any constraints on the relative order of operations M1, M2, or M3. Due to the register and memory dependencies between the instructions on Processor #0, these operations complete *in program order* on Processor #0 and also become *locally* visible in this order. However, the operations need *not* be made visible to remote processors in program order. In this outcome it appears to Processor #0 as if  $M1 \rightarrow M3$  while to Processor #1 it appears that  $M3 \rightarrow M1$ . There are two things to note here. First, the behavior is another example of the local bypass behavior that [Section 2.2.1.8](#) presents on [page 2:518](#). Second, there are no dependencies *directly* between M1 and M3 that requires them to become globally visible in program order.

**Note:** All processors will observe the order established by a particular processor in case of a WAW memory dependency to the same location. For example, all processors in the coherence domain eventually see a value of 1 in location x in the following code:

```

st      [x] = 0      // M1: set [x] to 0
st      [x] = 1      // M2: set [x] to 1,
                        // cannot move above M1 due to WAW

```

because there is a WAW memory dependency between from M2 to M1 and the Itanium architecture requires that the local processor resolves RAW, WAR, and WAW dependencies between its memory accesses in program order. Thus,  $M1 \rightarrow M2$  even though the ordering semantics do not place any constraints on the relative ordering of M1 and M2.

### 2.2.1.7 Data Dependency Establishes Local Ordering

In the Itanium architecture, a dependency (e.g., a later operation reading the value written by an earlier operation) can imply a local ordering relationship between the two operations. This section focuses on dependencies through registers only.

Section 2.2.1.6 discusses dependencies and MP ordering.

The execution shown in Table 2-6 illustrates how data dependency and memory ordering interact in a simple “pointer chase.”

**Table 2-6. Memory Ordering and Data Dependency**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld	r1 = [y] ;;	// M3
st.rel	[y] = x	// M2	ld	r2 = [r1]	// M4

Outcome: r1 = x and r2 = 0 is not allowed

In this example, Processor #0 could be executing code that updates a shared object with M1 and then publishes a pointer to the object with M2. Processor #1 then loads the pointer and dereferences it to read the contents of the shared object. The outcome r1 = x and r2 = 0 implies that Processor #1 observes the new value of the object pointer, y, but the old value of the data field, x.

The ordering semantics require M1 → M2 but place no requirements on the relative ordering of M3 and M4.

Thus, the memory semantics alone would allow the outcome r1 = x and r2 = 0 in the absence of other constraints. Using an acquire load for M3 can avoid this outcome as doing so forces M3 → M4 and thus prevents the outcome. However, this use of acquire is non-intuitive given the RAW dependency through register r1 between M3 and M4. That is, M3 produces a value that M4 requires in order to execute so how should it be possible for them to go out of order? Further, using an acquire in this case prevents any memory operation following M3 from moving above M3, even if they are completely independent of M3.

To avoid this potential confusion and performance issue, the Itanium architecture treats data dependency and memory ordering in the same fashion on the local processor. That is, if A » B and A produces a value that B consumes, then A → B on the local processor. This relationship is also transitive as the execution in Table 2-7 illustrates.

**Table 2-7. Memory Ordering and Data Dependency Through a Predicate Register**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld	r1 = [y]	// M3
st.rel	[y] = x	// M2	cmp.eq	p1, p2 = r1, x ;;	// C1
			(p1)ld	r2 = [x]	// M4

Outcome: r1 = x and r2 = 0 is not allowed

The Processor #0 code is the same as in Table 2-6. The Processor #1 now performs the following operation: if the pointer value y is equal to x, load a value from x.



The Itanium architecture does not allow the outcome  $r1 = x$  and  $r2 = 0$  in this execution either. Unlike the execution in [Table 2-6](#), there is no *direct* dependency between the values that M3 produces and the values that M4 consumes. However, there is a RAW through register r1 from M3 to C1 and a RAW through register p1 from C1 to M4. Thus, by transitivity,  $M3 \rightarrow M4$ .

The execution in [Table 2-8](#) illustrates a similar construct but introduces a control dependency.

**Table 2-8. Memory Ordering and Data and Control Dependencies**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld	r1 = [y];;	// M3
st.rel	[y] = x	// M2	cmp.eq	p1, p2 = r1, x	// C1
			(p2)br	t	// B1
			ld	r2 = [x]	// M4
			t:		

Outcome:  $r1 = x$  and  $r2 = 0$  is not allowed

This execution is semantically the same as the execution in [Table 2-7](#); however, this execution uses a control dependency rather than predication to conditionally execute M4. As a result, the outcome  $r1 = x$  and  $r2 = 0$  is not allowed in the [Table 2-8](#) execution.

The execution of the load M4 is data-dependent on the value of p2 that the branch B1 uses to resolve. Further, p2 is dependent on the value of r1 that the load M3 produces through the compare C1. Thus,  $M3 \rightarrow M4$ .

The execution in [Table 2-9](#) is a variation on the execution from [Table 2-8](#) where the loads are truly independent.

**Table 2-9. Memory Ordering and Control Dependency**

Processor #0			Processor #1		
st	[x] = 1	// M1	ld	r1 = [y]	// M3
st.rel	[y] = x	// M2	cmp	p1, p2 = r3, x	// C1
			(p2) br	t	// B1
			ld	r2 = [x]	// M4
			t:		

Outcome: all are allowed

In this execution, there is no dependency between M3 and M4, and thus, there are no constraints on the relative ordering of M3 and M4. Like the execution in [Table 2-8](#), M4 is data-dependent on the value of p2 that the branch B1 uses to resolve. However, p2 is *independent* of the value that the load M3 produces (specifically, because the compare does not use the value of register r1 that the load produces). Thus, there is no chain of dependencies between M3 and M4 and therefore there are no constraints on the relative ordering of M3 and M4. As a result, all outcomes are allowed in this execution.

### 2.2.1.8 Store Buffers May Satisfy Local Loads

In the Itanium memory ordering model, store buffers (or other logically-equivalent structures) may satisfy local read requests from loads or acquire loads even if the stored data is not yet visible to other agents in the coherence domain. Such bypassing must honor any ordering semantics in the memory reference stream. Table 2-10 and Table 2-11 that Section 2.2.1.9 presents illustrate this behavior.

**Table 2-10. Store Buffers May Satisfy Loads if the Stored Data is Not Yet Globally Visible**

Processor #0			Processor #1		
st.rel	[x] = 1	// M1	st.rel	[y] = 1	// M4
ld.acq	r1 = [x]	// M2	ld.acq	r3 = [y]	// M5
ld	r2 = [y]	// M3	ld	r4 = [x]	// M6

Outcome: r1 = 1, r3 = 1, r2 = 0, and r4 = 0 is allowed

In this sequence, each processor bypasses its locally-written value from a store buffer before the value becomes visible to the other processor. This behavior may make accesses of different sizes that have overlapping memory addresses appear to complete non-atomically.

The following discussion focuses on the outcome r1 = 1, r3 = 1, r2 = 0, and r4 = 0 because this outcome is allowed if and only if store buffers can satisfy local loads (other outcomes are allowed but do not depend on being able to satisfy local loads from a store buffer).

The Itanium memory ordering semantics only require that M2 → M3 and M5 → M6. There are no constraints on the relative ordering of M1 and M2 or M3 nor on the relative ordering of M4 and M5 or M6.

Remember that both dependencies and the memory ordering model place requirements on the manner in which a processor based on the Itanium architecture may re-order accesses. Even though the Itanium memory ordering model allows loads to pass stores, a processor based on the Itanium architecture cannot re-order the following sequence:

```
st.rel    [x] = r0    // M1: store 0 to [x]
ld.acq   r1 = [x]   // M2: cannot move above st.rel due to RAW
```

This is because there is a RAW dependency through memory between M1 and M2 and the Itanium memory ordering model requires that the local processor resolve RAW, WAR, and WAW dependencies between its memory accesses in program order. Thus, M1 → M2 even though the ordering semantics place no constraints on the relative ordering of M1 and M2.

Because there is a RAW dependency through memory between M1 and M2 and between M4 and M5, the ordering constraints *effectively* become:<sup>1</sup>

```
M1 → M2 → M3
M4 → M5 → M6
```

---

1. That is, the store operations must become visible to the local processors before their loads that read the stored value.

to account for both the memory ordering semantics and dependencies. It is important to keep in mind that the observance of a dependency between two operations does not imply an ordering relationship (from the standpoint of the memory ordering model) between the operations as [Section 2.2.1.6](#) describes.

Assuming that a processor can bypass locally-written values before they are made globally-visible implies that there is a local and a global visibility points for a memory operation where a value always becomes locally visible before it becomes globally visible. Since M1 and M4 can have local visibility with respect to M2 and M5 as well as global visibility,

$$\begin{aligned} m1 \rightarrow M2 \rightarrow M3; m1 \rightarrow M1 \\ m4 \rightarrow M5 \rightarrow M6; m4 \rightarrow M4 \end{aligned}$$

where m1 and M1 represent local and global visibility of memory operation 1, respectively. There are two things to note. First, the ordering of the local visibilities of operations M1 and M4 (m1 and m4, respectively) allow each processor to honor its data dependencies. That is, Processor #2 honors the RAW dependency through memory between M1 and M2 by requiring m1 to become visible before M2. Second, that these requirements do not place any constraints on the relative ordering perceived by a *remote* observer of operation M1 with M2 and M3 or of operation M4 with M5 and M6 (as the local visibilities meet the *local* ordering constraints that the dependencies impose).

The code in [Table 2-10](#) and these constraints together imply that

$$\begin{aligned} r1 = 1 \Rightarrow m1 \rightarrow M2 \\ r3 = 1 \Rightarrow m4 \rightarrow M5 \\ r2 = 0 \Rightarrow M3 \rightarrow M4 \Rightarrow m1 \rightarrow M6 \text{ because } m1 \rightarrow M3 \text{ and } M3 \rightarrow M4 \text{ and } M4 \rightarrow M6 \\ r4 = 0 \Rightarrow M6 \rightarrow M1 \\ m1 \rightarrow M6 \text{ and } M6 \rightarrow M1 \Rightarrow m1 \rightarrow M1 \end{aligned}$$

Thus, the outcome  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  is allowed because these statements are consistent with our definition of local and global visibility. Specifically, a value becomes locally visible before it becomes globally visible. Similar reasoning can show that the constraints also imply that  $m4 \rightarrow M4$ .

### 2.2.1.9 Preventing Store Buffers from Satisfying Local Loads

In the code shown in [Table 2-10](#) from [Section 2.2.1.8](#), there are no ordering constraints between the store and acquire load from the standpoint of memory ordering semantics (however, there is a RAW dependency through memory that forces the acquire load to follow the store). Bypassing may not occur if doing so violates the memory ordering constraints of memory operations between the store and the bypassing read.

[Table 2-11](#) presents a variation on the execution in [Table 2-10](#) from [Section 2.2.1.8](#) that illustrates this behavior.

**Table 2-11. Preventing Store Buffers from Satisfying Local Loads**

Processor #0			Processor #1		
st	[x] = 1	// M1	st	[y] = 1	// M5
mf		// M2	mf		// M6
ld.acq	r1 = [x]	// M3	ld.acq	r3 = [y]	// M7
ld	r2 = [y]	// M4	ld	r4 = [x]	// M8

*Outcome:*  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  is not allowed

Like [Section 2.2.1.8](#), the discussion in this section focuses on the outcome  $r1 = 1$ ,  $r3 = 1$ ,  $r2 = 0$ , and  $r4 = 0$  because it is allowed if and only if store buffers can satisfy local loads. The line of reasoning to show that the outcome  $r1 = 1$ ,  $r3 = 1$ ,  $r2 = 0$ , and  $r4 = 0$  is not allowed in [Table 2-11](#) is similar to the reasoning used to show that this outcome is allowed in the [Table 2-10](#) execution from [Section 2.2.1.8](#) on [page 2:518](#).

By the definition of the Itanium memory ordering semantics,

$$\begin{aligned} M1 \rightarrow M2 \rightarrow M3 \rightarrow M4 \\ M5 \rightarrow M6 \rightarrow M7 \rightarrow M8 \end{aligned}$$

By allowing local and global visibility of operations M1 and M5 (similar to the discussion in [Section 2.2.1.8](#)), this assumption, along with the above constraints, together imply that,

$$\begin{aligned} m1 \rightarrow M1 \Rightarrow m1 \rightarrow M2 \rightarrow M3 \rightarrow M4 \\ m5 \rightarrow M5 \Rightarrow m5 \rightarrow M6 \rightarrow M7 \rightarrow M8 \end{aligned}$$

Consider these constraints on the Processor #0 operations m1, M1, M2, M3, and M4. Making m1 visible before M2, M3, and M4 correctly honors the data dependency through memory on Processor #0. However, unless it constrains the global visibility of M1 to occur before M2, M3, and M4, Processor #0 violates the Itanium ordering semantics. Specifically, the memory fence M2 must always be made visible after the store M1. Allowing global and local visibilities of M1 in this case violates this constraint, and thus, is not allowed. Essentially, by allowing M1 to become locally visible early, M3 would see M1 before the fence semantics for M2 were met (namely, that M1 be visible before M2 and thus M3). Without local and global visibility of M1 and M5, the ordering constraints are as this example originally postulated.

The code in [Table 2-11](#) and these constraints together imply that

$$r2 = 0 \Rightarrow M4 \rightarrow M5 \Rightarrow M1 \rightarrow M8 \text{ because } M1 \rightarrow M4 \text{ and } M4 \rightarrow M5 \text{ and } M5 \rightarrow M8 \Rightarrow r4 = 1$$

This contradicts the  $r1 = 1$ ,  $r3 = 1$ ,  $r2 = 0$ , and  $r4 = 0$  outcome. The visibility of the memory fence, M2, implies that all prior operations including the store to x, M1, are globally visible. Thus, the load from x on Processor #1, M8, must observe the new value of x and  $M1 \rightarrow M8$  but the outcome requires  $M8 \rightarrow M1$ .

### 2.2.1.10 Semaphores Do Not Locally Bypass

As [Section 2.2.1.8](#) and [Section 2.2.1.9](#) discuss, loads and acquire loads may be satisfied with values placed in local store buffers (or other logically-equivalent structures) by stores or release stores before the stored data becomes visible to other agents in the coherence domain. The Itanium architecture explicitly prohibits such local bypass either to or from semaphore operations. That is, semaphore operations cannot be satisfied in this way nor can the data they store be used to satisfy loads or acquire loads in this way.

The execution in [Table 2-12](#) illustrates a variation on the execution in [Table 2-10](#) where the acquire loads have been replaced with exchange semaphore operations (which also have acquire semantics).

**Table 2-12. Bypassing to a Semaphore Operation**

Processor #0			Processor #1		
mov	r5 = 2		mov	r6 = 2	
st.rel	[x] = 1	// M1	st.rel	[y] = 1	// M4
xchg	r1 = [x], r5	// M2	xchg	r3 = [y], r6	// M5
ld	r2 = [y]	// M3	ld	r4 = [x]	// M6

Outcome: r1 = 1, r3 = 1, r2 = 0, and r4 = 0 is not allowed

Although each semaphore operation can be decomposed into a read access followed by a write access, the Itanium architecture does *not* allow a read request by a semaphore to be satisfied from a store buffer (or other logically-equivalent structure). As a result, the outcome r1 = 1, r3 = 1, r2 = 0, and r4 = 0 is not allowed. The reasoning is similar to that presented in [Section 2.2.1.9](#).

Specifically, by the definition of the Itanium memory ordering semantics, M2 → M3 and M5 → M6. The relative ordering between operation M1 and operations M2 or M3 is not constrained. Likewise, the relative ordering between operation M4 and operations M5 and M6.

Now, assume the outcome r1 = 1, r3 = 1, r2 = 0, and r4 = 0. Given that r1 = 1, r3 = 1, and r2 = 0, we observe the following:

$$\begin{aligned}
 r1 = 1 &\Rightarrow M1 \rightarrow M2 \\
 r3 = 1 &\Rightarrow M4 \rightarrow M5 \\
 r2 = 0 &\Rightarrow M3 \rightarrow M4 \\
 M3 \rightarrow M4 &\Rightarrow M1 \rightarrow M6 \text{ because } M1 \rightarrow M3 \rightarrow M4 \rightarrow M6 \\
 M1 \rightarrow M6 &\Rightarrow r4 = 2
 \end{aligned}$$

This conclusion contradicts the assumed outcome where r4 = 0 and thus the outcome r1 = 1, r3 = 1, r2 = 0, and r4 = 0 is not allowed. Because M1 and M4 cannot become locally-visible to M2 and M5 before they become globally-visible to M6 and M3 (as read accesses from semaphores may not bypass from store buffers or other logically-equivalent structures), it is not possible to avoid this contradiction.

The Itanium architecture also prohibits local bypass from a semaphore operation to a local read access from a load or acquire load as shown in the execution in [Table 2-13](#).

**Table 2-13. Bypassing from a Semaphore Operation**

Processor #0			Processor #1		
fetchadd.rel	r5 = [x], 1	// M1	fetchadd.rel	r6 = [y], 1	// M4
ld.acq	r1 = [x]	// M2	ld.acq	r3 = [y]	// M5
ld	r2 = [y]	// M3	ld	r4 = [x]	// M6

Outcome: r1 = 1, r3 = 1, r2 = 0, r4 = 0, r5 = 0, and r6 = 0 is not allowed

A store buffer may not provide a local read operation early access to a value written by a semaphore operation. Therefore, the outcome  $r1 = 1, r3 = 1, r2 = 0, r4 = 0, r5 = 0,$  and  $r6 = 0$  in the [Table 2-13](#) execution is not allowed. The reasoning is similar to that used in the previous execution.

### 2.2.1.11 Ordered Cacheable Operations are Seen in the Same Order by All Observers

The Itanium memory ordering model requires that release stores and semaphore operations (both acquire and release forms) become visible to all observers in the coherence domain in a single total order with the exception that each processor may observe (via loads or acquire loads) its own update early. Thus, each observer in the coherence domain sees the same interleaving of release stores and semaphores (both acquire and release forms) from the other processors in the coherence domain except that each processor may observe its own release stores (via loads or acquire loads) prior to their being observed globally. [Table 2-14](#) illustrates this behavior.

**Table 2-14. Enforcing the Same Visibility Order to All Observers in a Coherence Domain**

Processor #0	Processor #1	Processor #2	Processor #3
st.rel [x] = 1// M1	ld.acq r1 = [x]//M2 ld r2 = [y]//M3	st.rel [y] = 1// M4	ld.acq r3 = [y]//M5 ld r4 = [x]//M6

Outcome: only  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  is not allowed

The Itanium memory ordering model only disallows the outcome  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  in this execution. By the definition of the Itanium memory ordering semantics,

$$\begin{aligned} M2 &\rightarrow M3 \\ M5 &\rightarrow M6 \end{aligned}$$

The Itanium memory ordering model does not permit the  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  outcome as this would require that Processors #1 and #3 observe the release stores to x and y in different orders. Specifically, assuming that the outcome is  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$ :

$$\begin{aligned} r1 = 1 &\Rightarrow M1 \rightarrow M2 \\ r3 = 1 &\Rightarrow M4 \rightarrow M5 \\ r2 = 0 &\Rightarrow M3 \rightarrow M4 \Rightarrow M1 \rightarrow M4 \text{ because } M1 \rightarrow M2, M2 \rightarrow M3, \text{ and } M3 \rightarrow M4 \\ r4 = 0 &\Rightarrow M6 \rightarrow M1 \Rightarrow M4 \rightarrow M1 \text{ because } M4 \rightarrow M5, M5 \rightarrow M6, \text{ and } M6 \rightarrow M1 \end{aligned}$$

The final two statements are inconsistent since both  $M1 \rightarrow M4$  and  $M4 \rightarrow M1$  cannot be true unless Processors #1 and #3 are allowed to see the release stores to x and y in different orders.

The Itanium memory ordering model allows the  $r1 = 1, r3 = 1, r2 = 0,$  and  $r4 = 0$  outcome if either one or both of the release stores M1 and M4 are unordered since unordered operations need not be seen in the same total order by all observers in the coherence domain. Thus, in a version of the execution shown in [Table 2-14](#) with unordered stores, Processor #2 observes  $M1 \rightarrow M4$  while Processor #4 observes  $M4 \rightarrow M1$ .

The Itanium memory ordering model also allows this outcome if the release stores M1 and M4 are replaced with a memory fence followed by an unordered store. From the standpoint of a single processor, a release store has equivalent ordering semantics on the local processor to a memory fence followed by an unordered store. However, because the store in the memory fence/unordered store pair is unordered, it does not have any ordering requirements with respect to a remote processor. Even when processors are allowed to construct different interleavings, the ordering of an individual processor's memory references within the interleaving must always respect the ordering constraints placed on those references.

### 2.2.1.12 Obeying Causality

As noted in [Section 2.2.1.11](#), the Itanium memory ordering model requires that release stores and semaphore operations (both acquire and release forms) become visible to all observers in the coherence domain in a single total order with the exception that each processor may observe (via loads or acquire loads) its own update early. Thus, each observer in the coherence domain sees the same interleaving of release stores, and semaphores operations from the other processors in the coherence domain.

A consequence of this is the fact that the Itanium memory ordering model respects causality in a certain way. Specifically, if a release store or semaphore operation causally precedes any store or semaphore operation, then the two operations will become visible to all processors in the causality order. [Table 2-1](#) illustrates this behavior. Suppose that M2 reads the value written by M1. In this case, there is a causal relationship from M1 to M3 (a control dependency could also establish such a relationship). The fact that the store to x is a release store implies that, since there is a causal relationship from M1 to M3, M1 must become visible to processor #2 before M3.

**Table 2-15. Intel® Itanium® Architecture Obeys Causality**

Processor #0	Processor #1	Processor #2
st.rel [x] = 1 // M1	ld.acq r1 = [x] // M2 st [y] = 1 // M3	ld.acq r2 = [y] // M4 ld r3 = [x] // M5

*Outcome: only r1 = 1, r2 = 1, and r3 = 0 is not allowed*

The Itanium memory ordering model disallows the outcome  $r1 = 1, r2 = 1, \text{ and } r3 = 0$  in this execution (all other outcomes are allowed). To see this, we note the following. If  $r1 = 1$ , then  $M1 \rightarrow M2$  at Processor #1. Because M2 is an acquire load and  $M2 \gg M3$ ,  $M2 \rightarrow m3$ , where  $m3$  represents the local visibility of memory operation 1 (see [Section 2.2.1.8](#)). Thus,  $M1 \rightarrow m3$ . Since M1 is a release store, it appears to become visible to all processors at the same time. This fact and  $m3 \rightarrow M3$  together imply  $M1 \rightarrow M3$ .

If  $r2 = 1$ ,  $M3 \rightarrow M4$ . Because M4 is an acquire load,  $M4 \rightarrow M5$ . If  $r3 = 0$ , then  $M5 \rightarrow M1$ . Together, these imply  $M3 \rightarrow M1$ , which contradicts the observation from the previous paragraph. Thus, the outcome  $r1 = 1, r2 = 1, \text{ and } r3 = 0$  is disallowed.

The indicated outcome would also be disallowed if M1 were a semaphore operation because, like release stores, each semaphore must appear to become visible at all processors at the same time. The indicated outcome would be allowed if M1 were a weak store, as a weak store may appear to become visible at different times to different processors.

## 2.2.2 Memory Attributes

In addition to the ordering semantics and data dependencies, the memory attributes of the page that is being referenced also influence access ordering and visibility. Using memory attributes allows the Itanium architecture to match the performance and the usage model to the type of device (e.g. main memory, memory-mapped I/O device, frame buffer, locations with side-effects, etc.) that backs a page of memory. Typically, memory with side-effects is mapped uncacheable while memory without side-effects is mapped as write-back cacheable.

[Section 4.4, “Memory Attributes”](#) describes memory attributes in the Itanium architecture in greater depth.

Memory with the uncacheable UC or UCE attributes is sequential by definition. A processor based on the Itanium architecture ensures that accesses to sequential memory locations reach a peripheral domain (a platform-specific collection of uncacheable locations, colloquially known as “a device”) in program order with respect to all other accesses to sequential locations in the same peripheral domain. The sequential behavior of UC or UCE memory is independent of the ordering semantics (i.e. acquire, release, fence, or unordered) attached to the accesses.

Other observers (e.g. processors or other peripheral domains) need not see references to UC or UCE memory in sequential order if at all. When multiple agents are writing to the same device, it is up to software to synchronize the accesses to the device to ensure the proper interleaving.

The ordering semantics of an access to sequential memory determines how the access becomes visible to the peripheral domain with respect to other operations. For example, consider the code sequence shown in [Figure 2-2](#).

**Figure 2-2. Interaction of Ordering and Accesses to Sequential Locations**

```
sequential_example:
    st    [data_0] = 0      // M1: put data in cacheable mem
    st    [data_1] = 0      // M2: put data in cacheable mem
    st.rel [ready] = 1     // M3: tell device to get ready
    st    [start] = 1      // M4: tell device to start
```

In this code, assume that `data_0` and `data_1` are cacheable locations and `start` and `ready` are an uncacheable UC or UCE locations.

Sequentiality ensures that M3 and M4 reach the peripheral domain in program order (i.e. M3 before M4). Further, the release semantics on M3 ensures that it is not made visible to the peripheral domain until after M1 and M2 are made visible to the coherence domain. The M1 and M2 accesses may become visible to the coherence domains in any order as they both have unordered semantics. Even though the memory ordering semantics allow M4 to become visible before M3, the processor must make M3 visible before M4 because both `ready` and `start` are sequential locations.



### 2.2.3 Understanding Other Ordering Models: Sequential Consistency and IA-32

To provide a point of reference, it is helpful to understand other memory ordering models. These ordering models affect not only the programmer's view of the system, but also the overall system performance and design. Processors with relaxed memory ordering models may achieve higher performance than those with strict ordering models.

The most intuitive memory ordering model is "sequential consistency" (SC) which Lamport formally defines in [L79]. In sequential consistency, all processors see the memory references from a given processor in program order, and, in addition, all processors see the same system-wide interleaving of memory references from each processor.

The SC model precludes many common optimizations made in modern microprocessors to enhance performance. For example, in an SC system, a load may not pass a prior store until that store becomes globally visible (because all memory operations must become visible in program order). This requirement prevents the SC system from using a store buffer to hide the latency of store traffic by allowing loads that hit the cache to be serviced under a prior store that miss the cache.

To address such performance issues, many memory ordering models have been developed that relax the constraints of sequential consistency. Adve categorizes these memory models by noting how they relax the ordering requirements between reads and writes and if they allow writes to be read early [AG95]. The Itanium architecture allows for relaxed ordering between reads and writes and also allows writes to be read early under certain circumstances.

Aside from disallowing any relaxation of memory references, sequential consistency has two other subtle differences from the Itanium memory ordering model. First, it requires a total order of operations whereas the Itanium memory ordering model only requires a total order for release stores and semaphores. Second, remote processors must always honor data dependencies since the local processor does not have the option of re-ordering such accesses as can occur.

The IA-32 memory ordering relaxes write to read ordering and allows a processor to read its own writes before they are globally visible. Further, IA-32 allows each processor in the coherence domain to interleave the reference streams from other processors in the coherence domain in a different order. The per-processor orders must meet some additional constraints to ensure they are consistent with each other (enumerating and explaining these constraints is beyond the scope of this document). For more information on the IA-32 ordering model see [Section 6.2.3.2, "IA-32 Segmentation" on page 1:131](#).

## 2.3 Where the Intel® Itanium® Architecture Requires Explicit Synchronization

The Itanium architecture requires a memory synchronization (`sync.i`) and a memory fence (`mf`) during a context switch to ensure that all memory operations prior to the context switch are made visible before the context changes. Without this requirement, the ordering constraints may be violated if the process migrates to a different processor. For example, consider the example shown in [Figure 2-3](#).

**Figure 2-3. Why a Fence During Context Switches is Required in the Intel® Itanium® Architecture**

```
// Process A begins executing on Processor #0...

    ld.acq    r1 = [x]           // load executes on processor #0

// 1) Context switch occurs
// 2) O/S migrates Process A from Processor #0 to Processor #1
// 3) Process A resumes at the instruction following the ld.acq

    st        [y] = r2          // store executes on processor #1
```

In this example, Processor #1 may make the unordered store visible to the coherence domain before Processor #0 makes the acquire load visible. This violates the ordering constraints. Executing a memory fence during the context switch handler ensures that this violation can not occur.

See [Section 4.5, “Context Switching” on page 2:557](#) on context management in a processor based on the Itanium architecture.

Interruptions do not affect memory ordering. On entry to an interrupt handler, memory operations from the interrupted program may still be in-flight and not yet visible to other processors in the coherence domain. A handler that expects that all memory operations that precede the interruption to be visible must enforce this requirement by executing a memory fence at the beginning of the handler.

## 2.4 Synchronization Code Examples

There are many synchronization primitives that software uses in multiprocessor or multi-threaded environments to coordinate the activities of different code streams. In this section, we present several typical examples to illustrate how some common constructs translate to the Itanium instruction set. In addition, the discussions identify special considerations with various implementations.

The examples use the syntax “[`foo`]” to indicate the memory location that holds the variable `foo`. Actual Itanium architecture-based assembly language would first move the address of `foo` into a register and then use this register as an operand to a memory access instruction. The alternate syntax is chosen to simplify and clarify the examples.

## 2.4.1 Spin Lock

Software commonly uses spin locks to guard access to a critical region of code. In these locks, the software “spins” while waiting for a shared lock variable to indicate that the critical region can be safely accessed. Typically, the lock code uses atomic operations such as compare and exchange or fetch and add to update the shared lock variable. Figure 2-4 shows a spin lock based on the `cmpxchg` instruction.

**Figure 2-4. Spin Lock Code**

```
// available. If it is 1, another process is in the critical section.
//
spin_lock:
    mov     ar.ccv = 0           // cmpxchg looks for avail (0)
    mov     r2 = 1             // cmpxchg sets to held (1)

spin:
    ld8     r1 = [lock] ;;      // get lock in shared state
    cmp.ne  p1, p0 = r1, r2     // is lock held (ie, lock == 1)?
    (p1)   br.cond.spnt  spin ;; // yes, continue spinning

    cmpxchg8.acq  r1 = [lock], r2, ar.ccv ;; // attempt to grab lock
    cmp.ne  p1, p0 = r1, r2     // was lock empty?
    (p1)   br.cond.spnt  spin ;; // bummer, continue spinning

cs_begin:
    // critical section code goes here...
cs_end:
    st8.rel     [lock] = r0 ;; // release the lock
```

The spin lock code first initializes `ar.ccv` and a register with the values that indicate that the lock is available and held, respectively. A compare and exchange obtains the lock by exchanging `lock` with 1 if it currently holds 0. Next, the first loop ensures that the code spins in cache while the lock is held by someone else. Once this loop finds that the lock is available, a compare and exchange instruction attempts to obtain the lock. If this instruction fails (e.g. because someone else obtained the lock in the meantime), the code resumes spinning in the first loop.

Spinning using only the `cmpxchg/cmp/br` loop may generate excessive coherency traffic. For example, if the `cmpxchg` always stores to memory (even if the comparison fails) and the lock is highly-contested, the platform may have to generate a number of read for ownership transactions causing `lock` to move around the system. Using the first `ld8/cmp/br` loop avoids this problem by obtaining `lock` in a shared state. In the worst case, when `lock` is not contested, this loop adds only the overhead of the additional compare and branch.

The initial `ld8` need not be an acquire load because of the control-flow in the spin loop: this load must become visible before the `cmpxchg8` because the load must return data in order for the compare and branch to resolve. Further, the store that relinquishes the lock after the critical section uses release semantics to prevent memory references from the critical from moving after the reference that releases the lock. Finally, the branches use “static predict not taken” hints to optimize for the case where the lock is not highly contested.

## 2.4.2 Simple Barrier Synchronization

A barrier is a common synchronization primitive used to hold a set of processes at a particular point in the program (the barrier) until all processors reach the location. Once all processes arrive at the barrier, they may all continue to execute. Figure 2-5 shows a sense-reversing barrier synchronization based on the `fetchadd` instruction from Hennessy and Patterson [HP96].

This type of barrier prevents a process that races ahead to the next instance of the barrier from trapping other (slow) processors that are in the process of leaving the barrier.

**Figure 2-5. Sense-reversing Barrier Synchronization Code**

```
// The total shared variable is one less than the number of processors
// that wait at the barrier.
// The release shared variable indicates if the processor must wait at
// the barrier (initially, this variable is 0).
// local_sense is a per-processor local variable that indicates the
// "sense" of the barrier (initially, this variable is 0).

sr_barrier:
    fetchadd8.acq r1 = [count], 1           // update counter
    ld8          r2 = [total]              // get number of procs - 1
    ld8          r3 = [local_sense] ;;     // get local "sense" variable
    xor          r3 = 1, r3                // local_sense != local_sense
    cmp.eq       p1, p2 = r1, r2 ;;        // p1 => last proc to arrive
    st8          [local_sense] = r3        // save new value of local_sense
(p1) st8        [count] = r0               // last resets count to 0
(p1) st8.rel    [release] = r3 ;;         // last allows other to leave

wait_on_others:
(p2) ld8        r1 = [release] ;;          // p2 => more procs to come
(p2) cmp.ne.and p0, p2 = r1, r3           // have all arrived yet?
(p2) br.cond.sptk wait_on_others ;;       // nope, continue waiting

    // This mf prevents memory operations that follow the barrier code
    // from moving ahead of memory operations that precede the barrier
    // code
    mf ;;
```

The barrier code begins by atomically updating the number of processors that are waiting at the barrier, `count`, using a `fetchadd` instruction. For the last processor that reaches the barrier, the `fetchadd` instruction returns the same value as the `total` shared variable, which is one less than the number of processors that wait at the barrier. Other processors each get a unique value on the interval  $[0, total)$  based on the order in which they arrive at the barrier.

All processors except the last processor wait in the `wait_on_others` loop for the signal that all have arrived at the barrier. The last processor to arrive at the barrier provides this signal.

The signal to leave the barrier is deduced from the value of the `release` shared variable and the `local_sense` local variable. Upon entering the barrier, each processor complements the value in its private `local_sense` variable. Once in the barrier, all processors always have the same value in their `local_sense` variables. This variable

indicates the value that `release` must have before the processor can leave the barrier. The last processor to arrive at the barrier releases the other processors by setting `release` to the new `local_sense` value.

The `mf` instruction in [Figure 2-5](#) is necessary only if the programmer wishes to ensure that memory operations performed before the barrier code are visible to memory operations performed by any processor after the barrier code.

### 2.4.3 Dekker's Algorithm

Dekker's algorithm [D65] is a common synchronization construct that arbitrates for a resource through the use of several shared variables that indicate which processor is using the resource. Each processor has its own flag variable that it shares with all other processors in the system. When a processor attempts to enter the critical section, it sets its flag to one and checks to make sure the flags for the other processors are all zero.

The code in [Figure 2-6](#) illustrates the core of this algorithm for a two-way multiprocessor system. In this example, a processor makes a single attempt to acquire the resource; typically, this code would appear in a loop. Although there is an array of per-processor flag variables, the code uses `flag_me` and `flag_you` to indicate to the flag variables for the processor attempting to obtain the resource and the other remote processor, respectively.

Dekker's algorithm assumes a sequential consistency ordering model. Specifically, it assumes that loading zero from `flag_you` implies that a processor's load and stores to the flag variables occur before the other processor's load and store to the flag variables. If this is not the case, both processors can enter the critical section at the same time.

Using unordered loads or stores to access the `flag_me` and `flag_you` variables does not guarantee correct behavior as the processor may re-order the accesses as it sees fit. Using an acquire load and release store is also not sufficient to ensure correct behavior because the ordering semantics always allow acquire loads to move earlier and release stores to move later. In the absence of the `mf`, it is possible for the load from `flag_you` to occur before the store to `flag_me`; even with acquire and release operations.

The first `ld8` need not be an acquire load because of the control-flow that skips the critical section: this load must become visible before any memory operations in the critical section because the load must return data in order for the compare and branch to resolve.

**Figure 2-6. Dekker’s Algorithm in a 2-way System**

```
// The flag_me variable is zero if we are not in the
// synchronization and critical section code and non-zero
// otherwise; flag_you is similarly set for the other processor.
// This algorithm does not retry access to the
// resource if there is contention.
//
dekker:
    mov    r1 = 1 ;;           // my flag = 1 (i want access!)
    st8    [flag_me] = r1
    mf ;;                     // make st visible first
    ld8    r2 = [flag_you] ;; // is other's flag 0?
    cmp.ne p1, p0 = 0, r2
    (p1) br.cond.spnt cs_skip ;; // if not, resource in use

cs_begin:
    // critical section code goes here...
cs_end:

cs_skip:
    st8.rel [flag_me] = r0 ;; // release lock
```

### 2.4.4 Lamport’s Algorithm

Like Dekker’s algorithm, Lamport’s algorithm [L85] also provides mutual exclusion for critical sections of code. Lamport’s algorithm is very simple and, in the case of non-contested locks, only requires two read and two write memory accesses to enter the critical section. The algorithm uses two shared variables,  $x$  and  $y$ , and a shared array,  $b$ , that identify the process entering and using the critical section. Figure 2-7 presents Lamport’s algorithm 2 [L85].

Lamport’s algorithm expects that a processor that enters the critical section performs the set of operations:  $S = \{\text{store } x, \text{load } y, \text{store } y, \text{load } x\}$ <sup>1</sup>. To enforce this ordering, the Itanium architecture requires a memory fence in the middle of the  $\{\text{store } x, \text{load } y\}$  sequence and the  $\{\text{store } y, \text{load } x\}$  sequence. No combination of ordered semantics on the operations in each of these sequences will guarantee the correct ordering.

It is not possible for the store  $y$  in the second sequence to pass the load  $y$  in the first sequence because of the data dependency from the load  $y$  to the compare and branch. If the processor reaches the store  $y$  in the second sequence, the load of  $y$  from the first sequence must be visible. Likewise, it is not possible for memory operations in the critical section to move ahead of the final load  $x$  because of the data dependency between this load and the compare and branch that guards the critical section.

The accesses to the  $b$  array allow the algorithm to correctly handle contention for the lock. In such cases, the algorithm backs off and re-tries.

---

1. There are some additional operations on the  $b$  array that are interposed in this sequence when contention for the resource occurs.

**Figure 2-7. Lamport's Algorithm**

```
// The proc_id variable holds a unique, non-zero id for the process that
// attempts access to the critical section. x and y are the synchronization
// variables that indicate who is in the critical section and who is
// attempting entry. ptr_b_1 and ptr_b_id point at the 1'st and id'th
// element of b[].
//
lamport:
    ld8            r1 = [proc_id] ;;           // r1 = unique process id
start:
    st8            [ptr_b_id] = r1           // b[id] = "true"
    st8            [x] = r1                  // x = process id
    mf             // MUST fence here!
    ld8            r2 = [y] ;;
    cmp.ne         p1, p0 = 0, r2;;         // if (y != 0) then...
(p1) st8          [ptr_b_id] = r0           // ... b[id] = "false"
(p1) br.cond.sptk wait_y                   // ... wait until y == 0

    st8            [y] = r1                  // y = process id
    mf             // MUST fence here!
    ld8            r3 = [x] ;;
    cmp.eq         p1, p0 = r1, r3 ;;       // if (x == id) then...
(p1) br.cond.sptk cs_begin                 // ... enter critical section

    st8            [ptr_b_id] = r0           // b[id] = "false"
    ld8            r3 = [ptr_b_1]           // r3 = &b[1]
    mov            ar.lc = N-1 ;;           // lc = number of processors - 1
wait_b:
    ld8            r2 = [r3] ;;
    cmp.ne         p1, p0 = r1, r2         // if (b[j] != 0) then...
(p1) br.cond.spnt wait_b ;;                // ... wait until b[j] == 0
    add            r3 = #8, r3              // r3 = &b[j+1]
    br.cloop.sptk wait_b ;;                // loop over b[j] for each j

    ld8            r2 = [y] ;;
    cmp.ne         p1, p0 = r2, r1 ;;       // if (y != id) then...
(p1) br.cond.sptk cs_begin                 // ... enter critical section
wait_y:
    ld8            r2 = [y] ;;              // wait until y == 0
    cmp.ne         p1, p2 = 0, r2
(p1) br.cond.spnt wait_y
    br             start                     // back to start to try again

cs_begin:
    // critical section code goes here...
cs_end:

    st8            [y] = r0                 // release the lock
    st8.rel        [ptr_b_id] = r0;;        // b[id] = "false"
```

## 2.5 Updating Code Images

There are four general techniques for updating code images in order to modify the code stream of a local or remote processor.

- Self-modifying code or code that modifies its own image.
- Cross-modifying code or code that modifies the image of code running concurrently on another processor.

- Programmed I/O for paging of code pages.
- DMA for paging of code pages.

The next four sections discuss these techniques in greater depth.

To illustrate the code sequences for self- and cross-modifying code, the examples in this section use the syntax `st [foo] = new` to represent a group of aligned stores that change the instruction at address `foo` to the instruction `new`. The Itanium architecture requires that the instruction stream see aligned stores atomically. In addition, the syntax `fc.i foo` represents a group of flush cache instructions that ensures the cache line addressed by `foo` is coherent with all the instruction caches. Updating more than one instruction simply requires the appropriate store/flush “pair” for each updated instruction<sup>1</sup>.

## 2.5.1 Self-modifying Code

Figure 2-8 presents the Itanium instruction sequence necessary to update a code image location on the local processor only.

**Figure 2-8. Updating a Code Image on the Local Processor**

```
patch_local:
    st      [code] = new_inst      // write new instruction
    fc.i   code ;;                // flush new instruction
    sync.i ;;                      // sync i stream with store
    srlz.i ;;                      // serialize

    // Local caches and pipeline are now coherent with new_inst...
```

This code fragment changes the instruction at the address `code` to the new instruction `new_inst`. After executing this code, the change is visible to both the local processor’s caches and its pipeline.

The `st` instruction updates the code image and the `fc.i` instruction ensures the value stored is coherent with the instruction cache. The `fc.i` is necessary because the Itanium architecture does not require instruction caches to be coherent with data stores for Itanium architecture-based code. Next, the `sync.i` ensures that the code update is visible to the instruction stream of the local processor and orders the cache flush with respect to subsequent operations by waiting for the prior `fc.i` instructions to be made visible. Finally, the `srlz.i` instruction forces the pipeline to re-initiate any instruction group fetches it performed after the `srlz.i` and also waits for the `sync.i` to complete; effectively making the pipeline coherent with the updated code image.

The serialization instruction is not necessary if software can *guarantee* that the processor encounters an event that re-initiates code fetches performed after the `sync.i`, such as an interruption or an `rfi`, before executing the new code. Events such as an interrupt or `rfi` both perform an instruction serialization which in this example waits for the `sync.i` to complete and then re-initiates code fetches.

- 
1. This description hides some of the complexity involved. Specifically, the flush and store operations have different sizes. Whereas multiple store instructions are necessary to update a 16 byte instruction, a single cache line flush invalidates at least two 16 byte instructions.



## 2.5.2 Cross-modifying Code

Consider a multi-threaded program for a multiprocessor system that dynamically updates some procedure that any processor in the system may execute. The program maintains several disjoint buffers to hold the new code and requires a processor to execute an IP-relative branch instruction at some address  $x$  to reach the code. In this scenario, the program updates the procedure by emitting the new code into a different buffer and then patching the branch at address  $x$  to target this new buffer. By carefully writing the update code, software can ensure that any processor in the system sees either:

- The original branch at address  $x$  that targets the original code in the old buffer along with the original code, or
- The new branch at address  $x$  that targets the new code in the new buffer along with the new code.

The code in [Figure 2-9](#) illustrates an optimized Itanium architecture-based code sequence that implements the cross-modifying code for this example.

**Figure 2-9. Supporting Cross-modifying Code without Explicit Serialization**

```
patch:
    st    [new_code] = new_inst // write new instruction
    fc.i  new_code ;;           // flush new instruction
    sync.i ;;                   // sync i stream with store

// Update the target of the branch that jumps to the updated code.
// This branch MUST be ip-relative. Before executing the following
// store, the branch jumps to somewhere other than "new_code".
//
    st.rel [x] = "branch <new_code>"

// If it is desired to propagate "branch <new_code>" to both
// the local processor and remote processor now, the following
// code is also necessary:
//
    fc.i  x ;;                 // flush branch
    sync.i ;;                 // sync i stream with store
    mf ;;                     // fence
```

To reach the new code at `new_code`, the processor executes the branch instruction at  $x$ . Initially, this branch jumps to an address other than `new_code`.

**Note:** The programmer needs to ensure that the branch to `new_code` is updated atomically. If an 8-byte store is used to update the branch, then the programmer needs to ensure that the branch to `new_code` is either in the first or last slot of the bundle.

The release store ensures a processor cannot see the new branch at address  $x$  and the original code at address `new_code`. That is, if a processor encounters "branch <new\_code>" at address  $x$ , then the processor's instruction cache must be coherent with the code image updates applied before the release store that updates the branch.

If remote processors may see either the old or new code sequence, the final three instructions in [Figure 2-9](#) are not necessary. In this case, the remote processors see the code image updates at some point in the future. In the meantime, they continue to execute the old code.

The release store ensures that the code image updates are made visible to the remote processors in the proper order (i.e. `new_code` is updated before the branch at address `x` is updated). Using the final three instructions ensures that the remote processors will see the new code the next time they execute the branch at address `x`.

On the local processor, the branch at address `x` also serves to force the pipeline to be coherent with the code image update the machine without requiring an interrupt, `rfi` instruction, or `srlz.i` instruction. Table 2-16 enumerates the potential pipeline behaviors to illustrate this point.

**Table 2-16. Potential Pipeline Behaviors of the Branch at `x` from Figure 2-9**

Pipeline Operation	Scenario #1	Scenario #2	Scenario #3	Scenario #4
Fetch branch at <code>x</code>	Old branch	Old branch	New branch	New branch
Predict branch at <code>x</code>	Old target	New target	Old target	New target
Code at target	Old instruction	“New” instruction (but could be stale)	Old instruction	New instruction
Retire branch at <code>x</code>	Old retires	Must flush due to misprediction	Must flush due to misprediction	New retires

In the first and fourth scenarios, the pipeline fetches and executes either the old branch and old target instruction or the new branch and new target instruction. Note that if the pipeline sees the new branch, it must also see the new target instruction by virtue of the way the code in Figure 2-9 is written. Either of these behaviors is consistent.

In the second and third scenarios, the pipeline obtains a mix of the old or new branch and the old or new target instruction. In these cases, the pipeline must flush because the predicted target will not agree with the branch instruction.

This behavior is not guaranteed unless the branch at address `x` is IP-relative and taken. The branch must be IP-relative to ensure that both the instruction and target address can be atomically updated (this is only possible with an IP-relative branch because in this type of branch, the target address is part of the instruction).

### 2.5.3 Programmed I/O

Programmed I/O requires that the CPU copy data from the device controller to main memory using load instructions to read from the device and store instructions to write data into cacheable memory (page-in).

To ensure correct operation, Itanium architecture-based software must exercise care in the presence of Programmed I/O due to two features of the architecture. First, the Itanium architecture does not require an implementation to maintain coherency between local instruction and data caches for Itanium architecture-based code. Second, the Itanium architecture allows aggressive instruction prefetching. Specifically, an implementation can move any location from a cacheable page into its instruction cache(s) any time a translation for the location indicates that the page is present (i.e. the `p` bit of the translation is set).

A system that performs Programmed I/O can use a sequence similar to that shown in Figure 2-8 to perform the data movement. Figure 2-10 presents a code sequence that updates a code image on both the local and remote processors.

**Figure 2-10. Updating a Code Image on a Remote Processor**

```
patch_l_and_r:
    st    [code] = new_inst    // write new instruction
    fc.i  code ;;              // flush new instruction
    sync.i ;;                  // sync i stream with store

// If the local processor must ensure that remote processors see
// the preceding memory updates before any subsequent memory
// operations, the following code is also necessary.
//
    mf ;;                      // make store visible to others

// If the local processor is going to execute the code and cannot
// cannot ensure instruction stream serialization, the following
// code is also necessary,
//
    srlz.i ;;                  // serialize my pipeline

// Local caches and pipeline are now coherent with new_inst, remote
// caches are now coherent with new_inst...
```

This code fragment changes the instruction at the address `code` to the new instruction `new_inst`. After executing this code, the change is visible to the local and remote processor's caches and to the local processor's pipeline, but may not be visible to remote processor's pipelines.

The sequence in [Figure 2-10](#) is similar to the code from [Figure 2-8](#) except an `mf` instruction occurs between the `sync.i` and `srlz.i` instructions. The fence is necessary if software must ensure that the code image update is made visible to all remote processors before any subsequent memory operations from the local processor. Although the `sync.i`, which orders the `st/fc.i` pair, has unordered semantics, it is an orderable operation and thus obeys the release or fence semantics of subsequent instructions (unlike an `fc.i` instruction; see [Section 4.4.7, "Sequentiality Attribute and Ordering"](#) for more information).

Because the pipeline is not snooped, the code in [Figure 2-10](#) cannot ensure that a remote processor's pipeline is coherent with the code image update. In the local case shown in [Figure 2-8](#), the `srlz.i` instruction enforces this coherency. As a result, the remote processor must serialize its instruction stream before it executes the updated code in order to ensure that a stale copy of some of the updated code is not present in the pipeline. This can be accomplished by explicitly executing a `srlz.i` before executing the updated code or by forcing an event that re-initiates any code fetches performed after the `fc.i` is observed to occur, such as an interruption or `rfi`.

Several optimizations to this code are possible depending on how software uses the updated code. Specifically, the `mf` and `srlz.i` can be eliminated under certain circumstances.

The `srlz.i` is not necessary if the local processor that updates the code image does not ever execute the new code. In this case, the local processor does not require its pipeline to be coherent with the changes to the code image. The fence is not necessary if the code image update can be made visible to remote processors in any relationship with subsequent memory operations from the local processor.

Finally, software may also eliminate the `mf` or `srlz.i` instructions if it *guarantees* that these operations will take place elsewhere (e.g. in the operating system) before the processor attempts to execute the updated code. For example, context switch routines must contain a memory fence (see [Section 2.3](#) on page [page 2:526](#)). Thus, the fence is not required if a context switch *always* occurs before any program can use the updated code.

## 2.5.4 DMA

Unlike Programmed I/O, which requires intervention from the CPU to move data from the device to main memory, data movement in DMA occurs without help from the CPU. A processor based on the Itanium architecture expects the platform to maintain coherency for DMA traffic. That is, the platform issues snoop cycles on the bus to invalidate cacheable pages that a DMA access modifies. These snoop cycles invalidate the appropriate lines in both instruction and data caches and thus maintain coherency. This behavior allows an operating system to page code pages without taking explicit actions to ensure coherency.

Software must maintain coherency for DMA traffic through explicit action if the platform does not maintain coherency for this traffic. Software can provide coherency by using the flush cache instruction, `fc`, to invalidate the instruction and data cache lines that a DMA transfer modifies. Code such as that shown in [Figure 2-8](#) on page [page 2:532](#) and [Figure 2-10](#) on page [page 2:535](#) accomplish this task.

## 2.6 References

- [AG95] S. V. Adve and K. Gharachorloo. "Shared memory consistency models: A Tutorial," Rice University ECE Technical Report 9512, September 1995.
- [L79] L. Lamport. "How to make a multiprocessor computer that correctly executes multiprocess programs," *IEEE Transactions on Computers*, C-28(9):690-691, September 1979.
- [HP96] J. L. Hennessy and D. A. Patterson. *Computer Architecture: A Quantitative Approach*, second edition, Morgan-Kaufmann, 1996.
- [D65] E. W. Dijkstra. "Cooperating sequential processes," Eindhoven, the Netherlands, Technological University Technical Report EWD-123, 1965.
- [L85] L. Lamport. "A Fast Mutual Exclusion Algorithm," Compaq Systems Research Center Technical Report 7, November 1985.

This chapter discusses the interruption and serialization model. Although the Itanium architecture is an explicitly parallel architecture, faults and traps are delivered in program order based on IP, and from left-to-right in each instruction group. In other words, faults and traps are reported precisely on the instruction that caused them.

## 3.1 Terminology

In the Itanium architecture, an **interruption** is an event which causes the hardware automatically to stop execution of the current instruction stream, and start execution at the instruction address corresponding to the **interruption handler** for that interruption. When this happens, we say that an interruption has been **delivered** to the processor core.

There are two classes of interruptions in the Itanium architecture. **IVA-based interruptions** are handled by the operating system (OS), at an address determined by the location of the interrupt vector table (IVT) and the particular interruption that has occurred. **PAL-based interruptions** are handled by the processor firmware. PAL-based interruptions are not visible to the OS, though PAL may notify the OS that a PAL-based interruption has occurred; see [Section 13.3, “Event Handling in Firmware” on page 2:632](#).

The architecture supports several different types of interruptions. These are defined below:

- A **fault** occurs when OS intervention is required before the current instruction can be executed. For example, if the current instruction misses the TLBs on a data reference, a Data TLB Miss fault may be delivered by the processor. Faults are delivered precisely on the instruction that caused the fault. The faulting instruction and all subsequent instructions do not update any architectural state (with the possible exception of subsequent instructions which violate a resource dependency<sup>1</sup>). All instructions executed prior to the faulting instruction update all their architectural state before the fault handler begins execution.
- A **trap** occurs when OS intervention is required after the current instruction has completed. For example, if the last instruction executed was a branch and PSR.tb is 1, a Taken Branch trap will be delivered after the instruction completes. Traps are delivered precisely on the instruction following the trapping instruction. The trapping instruction and all prior instructions update all their architectural state before the trap handler begins execution. All instructions subsequent to the trapping instruction do not update any architectural state.<sup>1</sup>

---

1. When an interruption is delivered on an instruction whose instruction group contains one or more illegal dependency violations, instructions which follow the interrupted instruction in program order and which violate the resource dependency may appear to complete before the interruption handler begins execution. Software cannot rely upon the value(s) written to the resource(s) whose dependencies have been violated; the value(s) are undefined. For details refer to [Section 3.4, “Instruction Sequencing Considerations” on page 1:39](#).

- When an external or independent agent (I/O device, timer, another processor) requires attention from the processor, an **interrupt** occurs. There are several types of interrupts. An initialization interrupt occurs when the processor has received an initialization request. A **Platform Management Interrupt** (PMI) can be generated by the platform to request features such as power management. Initialization interrupts and PMIs are PAL-based interruptions. An **external interrupt** occurs when an agent in the system requires the OS to perform some service on its behalf. External interrupts are IVA-based interruptions. Interrupts are delivered asynchronously with respect to program execution. The instruction upon which an interrupt is delivered may or may not be related to the interrupt itself.
- An **abort** is generated by the processor when a malfunction (Machine Check) is detected, or when a processor reset occurs. Aborts are asynchronous with respect to program execution. If caused by a particular instruction, an abort may be delivered sometime after that instruction completes. Aborts are PAL-based interruptions.

An interruption handler returns from interruption when it executes an `rfi` instruction. The `rfi` instruction copies state from specific control registers known as **interruption registers** into their corresponding architectural state (e.g. IIP is copied into IP and execution begins at that instruction address). Whether or not the state that is restored by the `rfi` is the same state that was captured when the interruption occurred is up to the operating system.

## 3.2 Interruption Vector Table

The Interruption Vector Address (IVA) control register defines the base address of the interruption vector table (IVT). Each IVA-based interruption has its own architected offset into this table as defined in [Section 5.7, “IVA-based Interruption Vectors” on page 2:113](#). For the remainder of this section, “interruption” refers to an IVA-based interruption, unless otherwise noted.

When an interruption occurs, the processor stops execution at the current IP, sets the current privilege level to 0, and begins fetching instructions from the address of the entry point to the interruption handler for the particular interruption that occurred. The address of this entry point is defined by the base address of the IVT contained in the IVA register and the architected offset into the table according to the interruption that occurred.

The IVT is 32Kbytes long and contains the code for the interruption handlers. Execution of the interruption handler begins at the entry point. The interruption handler may be contained entirely in the IVT, or the handler may branch to code outside the IVT if more space is needed.

When an interruption occurs, if the processor is operating with instruction address translation enabled (PSR.it is 1), then the address in IVA is treated as a virtual address; otherwise, it is treated as a physical address. Whenever an interruption may occur (i.e. whenever external interrupts are not masked or disabled, or whenever an instruction may raise a fault or trap), the software must ensure that the processor can safely reference the IVT. As a result, the IVT must be permanently resident in physical memory. If instruction address translation is enabled, the IVT must be mapped by an instruction translation register and must point at a valid physical page frame. When

instruction address translation is disabled, the IVA register should contain the physical address of the base of the IVT. Software must further ensure that instruction and memory references from low-level interruption handlers do not generate additional interruptions until enough state has been saved and interruption collection can be re-enabled.

There are many more interruptions than there are interruption vectors in the IVT. As specified in [Section 5.6, “Interruption Priorities”](#) there is a many-to-one relationship between interruptions and interruption vectors. The interruptions that share a common interruption vector (and hence, the code for an interruption handler) can determine which interruption occurred by reading the Interruption Status Register (ISR) control register. See [Chapter 8, “Interruption Vector Descriptions”](#) and [Chapter 9, “IA-32 Interruption Vector Descriptions”](#) for details of the specific ISR settings for each unique interruption.

## 3.3 Interruption Handlers

### 3.3.1 Execution Environment

As defined in [Section 5.5, “IVA-based Interruption Handling” on page 2:101](#), the processor automatically clears the PSR.i and PSR.ic bits when an interruption is delivered. This disables external interrupts and interrupt state collection, respectively. PMI delivery is also disabled while PSR.ic is 0; other PAL-based interruptions can be delivered at any point during the execution of the interruption handler, regardless of the state of PSR.i and PSR.ic.

In addition to clearing the PSR.i and PSR.ic bits, the processor also automatically clears the PSR.bn bit when an interruption is delivered, switching to bank 0 of general registers GR16 - GR31. This provides the interruption handler with its own set of registers which can be used without spilling any of the interrupted context’s register state, effectively saving GR16 - GR31 of the interrupted context. (This assumes PSR.bn is 1 at the time of interruption; see [Section 3.4.3, “Nested Interruptions” on page 2:546](#) for how to deal with the case where PSR.bn is 0 at the time of interruption.)

As specified in [Section 3.3.7, “Banked General Registers” on page 2:42](#), GR24 - GR31 **of bank 0** should not be used while PSR.ic is 1. By firmware convention, PAL-based interruption handlers may use these registers without preserving their values when PSR.ic is 1. When PSR.ic is 0, software may safely use GR24 - GR31 of bank 0 as scratch register.

Several other PSR bits and the RSE.CFLE are modified by the hardware when an interruption is delivered. [Table 3-1](#) summarizes the execution environment that interruption handlers operate in, and what each PSR bit and the RSE.CFLE values mean for the interruption handler.

**Table 3-1. Interruption Handler Execution Environment (PSR and RSE.CFLE Settings)**

PSR Bit	New Value	Effect on Low-level Interruption Handler
be	DCR.be	Byte order used by handler is determined by be-bit in DCR register.
ic & i	0	Disables interruption collection and external interrupts. Bank 0 is made active bank. This is discussed above
bn	0	
dt, rt, it, pk	unchanged	Instruction/Data/RSE address translation and protection key setting remain unchanged.
dfi & dfh	0	Floating-point registers are made accessible. This allows handlers to spill FP registers without having to toggle FP disable bits first. Modified bits indicate which registers were touched. See <a href="#">Section 4.2.2, "Preservation of Floating-point State in the OS"</a> on page 2:553 for details.
mfi, mfh	unchanged	
pp	DCR.pp	Privileged Monitoring is determined by pp-bit in DCR register. By default, user counters are enabled and performance monitors are unsecured in handlers. See <a href="#">Chapter 12, "Performance Monitoring Support"</a> for details.
up	unchanged	
sp	0	
di	0	Instruction set transitions are not intercepted.
si	0	Interval timer is unsecured.
ac	0	No alignment checks are performed.
db, lp, tb, ss	0	Debug breakpoints, lower-privilege interception, taken branch and single step trapping are disabled.
cpl	0	Current privilege level becomes most privileged.
is	0	Intel Itanium Instruction set. Handlers execute Intel Itanium instructions.
id, da, ia, dd, ed	0	Instruction/data debug, access bit and speculation deferral bits are disabled. For details, refer to <a href="#">Section 5.5.4, "Single Instruction Fault Suppression"</a> on page 2:104 and <a href="#">Section 5.5.5, "Deferral of Speculative Load Faults"</a> on page 2:105.
ri	0	Interrupt handler starts at first instruction is bundle.
mc	unchanged	Software can mask delivery of some machine check conditions by setting PSR.mc to 1, but the processor hardware does not set this bit upon delivery of an IVA-based interruption. Delivery of resets and BINITs cannot be masked.
RSE.CFLE (not a PSR bit)	0	Allows interruption handler to service faults in presence of an incomplete current register stack frame. This can happen when a mandatory RSE load takes an exception during when RSE is servicing a register stack underflow. For details refer to <a href="#">Section 6.6, "RSE Interruptions"</a> on page 2:144.

### 3.3.2 Interruption Register State

The Itanium architecture provides a set of hardware registers which, if interruption collection is enabled, capture relevant interruption state when an interruption occurs. The state of the PSR.ic bit at the time of an interruption controls whether collection is enabled. In this section, it is assumed that interruption collection is enabled (PSR.ic is 1); see [Section 3.4.3, "Nested Interruptions"](#) on page 2:546 for details on handling interruptions when collection is disabled (PSR.ic is 0). For details on collection of interruption resources for each interruption vector refer to [Chapter 8, "Interruption Vector Descriptions"](#) and [Chapter 9, "IA-32 Interruption Vector Descriptions."](#)



A processor based on the Itanium architecture provides the following interruption registers for collecting information about the latest interruption or the state of the machine at the time of the interruption:

- **IPSR** – A copy of the processor status register (PSR) at the moment the interruption occurred. The OS can use the IPSR to determine the value of any PSR bit when the interruption occurred. The contents of IPSR are restored into the PSR when the OS executes an `rfi` instruction. If the OS wishes to change the PSR state of the interrupted process (e.g. to step over an instruction debug fault), it can do so by modifying the IPSR contents before executing the `rfi`. When an interruption occurs, the processor sets IPSR.ri to the slot number (0, 1, or 2) of the instruction that was interrupted.
- **IIP** – A copy of the instruction pointer (IP) where the interruption occurred. The instruction bundle address contained in IIP, along with the IPSR.ri field, defines the instruction whose execution was interrupted. This instruction has not completed (i.e. it has not retired), so when the OS returns to the interrupted context, typically this is the instruction at which execution of the interrupted context resumes<sup>1</sup>. When the OS executes an `rfi` instruction, the contents of IIP are copied into the IP register and the processor begins fetching instructions from this address.
- **ISR** – Contains extra information about the specific interruption that occurred. This register is useful for determining exactly which interruption occurred for interruptions which share the same IVT vector.
- **IFA** – Faults related to addressing (e.g. Data TLB fault) materialize the faulting address in this register.
- **ITIR** – Faults related to addressing materialize the default page size and permission key for the region to which the faulting address belongs in this register.
- **IIPA** – Contains the instruction bundle address of the last instruction to retire successfully while PSR.ic was 1. In conjunction with ISR.ei, IIPA can be used by software to locate the instruction that caused a trap or that was executed successfully prior to a fault or interrupt.
- **IIM** – Instructions that take a Speculation fault (e.g. `chk`) or a Break Instruction fault (e.g. `break.i`) write this register with their immediate field when taking these faults. For these cases, the IIM register can be used to emulate the instruction, or to pass information to the fault handler; for example, software can use a particular immediate field value in a break instruction to indicate to the operating system that a system call is being performed.
- **IHA** – Faults related to the VHPT place the VHPT hash address in this register. See [Section 5.3, “Virtual Hash Page Table” on page 2:571](#) for details.
- **IFS** – This register can be used by software to save a copy of the interrupted context’s PFS register, but an interruption handler must do this explicitly; hardware only clears the valid bit (IFS.v) upon interruption. See below for details.
- **IIB0, IIB1** – Contain the 16-byte instruction bundle related to the interruption. Note that the IIB registers do not provide bundle information for all interruptions and are not supported on all processor implementations; please refer to [Chapter 8](#),

---

1. When an instruction faults because it requires emulation by the OS, the OS will normally skip the emulated instruction by returning to the instruction bundle address and slot number that follows IIP in program order. It does so by writing the next in-order bundle address and slot number into IIP and IPSR.ri, respectively, before executing an `rfi` instruction. Details on emulation handlers is in [Chapter 7, “Instruction Emulation and Other Fault Handlers.”](#)

[“Interruption Vector Descriptions”](#) for details. Software can use the instruction bundle information for debug and emulation purposes.

No other architectural state is modified when an interruption occurs. Note that only IIP, IPSR, ISR, and IFS are written by all interruptions (assuming PSR.ic is 1 at the time of interruption); the other interruption control registers are only written by certain interruptions, and their values are undefined otherwise. For details on which faults update which interruption resources refer to [Chapter 8, “Interruption Vector Descriptions”](#) and [Chapter 9, “IA-32 Interruption Vector Descriptions.”](#)

### 3.3.3 Resource Serialization of Interrupted State

As defined in [Section 3.2, “Serialization” on page 2:17](#), Itanium control register updates do not take effect until software explicitly serializes the processor’s data or instruction stream with a `srlz.d` or a `srlz.i` instruction, respectively. Control register updates that change a control register’s value and that have not yet been serialized are termed “in-flight.” Refer to [Section 3.2.3, “Definition of In-flight Resources” on page 2:19](#) for a precise definition.

When an interruption is delivered and before execution begins in the interruption handler, the processor hardware automatically performs an instruction and data serialization on all “in-flight” resources. As described in [Section 3.3.1](#) and [Section 3.3.2](#) above, the following resources determine the execution environment of the interruption handler:

- CR[IVA] – determines new IP
- CR[DCR].be – determines new value of PSR.be
- CR[DCR].pp – determines new value of PSR.pp
- PSR.ic – determines whether interruption collection is enabled
- RR[7:0] – determines new value of CR[ITIR] and CR[IHA]
- CR[PTA] – determines new value of CR[IHA]

Although these resources are guaranteed to be serialized prior to interruption handler execution, there is no guarantee that they will be serialized prior to the determination of the handler’s execution environment. If there is a value in-flight for any of these resources at the time of interruption delivery, either the old or new value may be used to generate the values of IP, PSR, CR[ITIR] and CR[IHA] seen by the handler.

As a result, if the handler requires the latest value of the listed resources to determine its execution environment, software must ensure that external interrupts are disabled and that no instruction or data references will take an exception until the resource updates have been appropriately serialized. Typically, the code toggling these resources is mapped by an instruction translation register to avoid TLB related faults.

Note that CR[IPSR] is guaranteed to get the latest value of the PSR on an interruption, even if there are PSR updates in-flight that have not been previously serialized by software.

For example, assume that GR2 contains the new value for IVA and that PSR.i is 1. To modify the IVA register, software would perform the following code sequence, where the code page is mapped by an instruction translation register or instruction translation is disabled:

```
rsm psr.i          // external interrupts disabled upon next instruction
mov cr[iva] = r2
;;
srlz.i            // writing IVA requires instruction serialization
;;
ssm psr.i        // external interrupts will be re-enabled after next srlz
```

### 3.3.4 Resource Serialization upon rfi

An `rfi` instruction also performs an instruction and a data serialization operation when it is executed. Any values that were written to processor register resources by instructions in an earlier instruction group than the `rfi` will be observed by the returned-to instruction, except for those register resources which are also written by the `rfi` itself, in which case the value written by the `rfi` will be observed. This makes the interruption handler more efficient by avoiding additional data and instruction serialization operations before returning to the interrupted context.

## 3.4 Interruption Handling

The Itanium architecture-based operating systems need to distinguish the following interruption handler types:

- **Lightweight interruptions:** Lightweight interruption handlers are allocated 1024 bytes (192 instructions) per handler in the IVT. These are discussed in [Section 3.4.1](#).
- **Heavyweight interruptions:** Heavyweight interruption handlers are allocated only 256 bytes (48 instructions) per handler in the IVT. These are discussed in [Section 3.4.2](#).
- **Nested interruptions:** If an interruption is taken when PSR.ic was 0 or was in-flight, a nested interruption occurs. Nested interruptions are discussed in [Section 3.4.3](#).

### 3.4.1 Lightweight Interruptions

Lightweight interruption handlers are allocated 1024 bytes (192 instructions) per handler in the IVT. Typically, lightweight handlers are written in Itanium architecture-based assembly code, and run in their entirety with interruption collection turned off (PSR.ic = 0) and external interrupts disabled (PSR.i = 0). Because these lightweight handlers are usually very short and performance-critical, they are intended to fit entirely in the space allocated to them in the IVT. An example of a lightweight interruption handler is the Data TLB vector (offset 0x0800). The first 20 vectors in the IVT, offsets 0x0000 (VHPT Translation vector) through 0x4c00 (reserved), are lightweight vectors. Typical lightweight handlers deal with instruction, data or VHPT TLB Misses, protection key miss handling, and page table dirty or access bit updates.

A typical lightweight interruption handler can operate completely out of register bank 0. If the bank 0 registers provide sufficient storage for the handler, none of the interrupted context's register state need be saved to memory, and the handler does not need to use stacked registers. Assuming no stacked registers are needed, the lightweight interruption handler can operate with an incomplete current register stack frame, obviating the need for `cover` and `alloc` instructions in the handler. This also allows the TLB related handlers to service TLB misses that result from mandatory RSE loads to the current frame.

### 3.4.2 Heavyweight Interruptions

Heavyweight interruption handlers are allocated only 256 bytes (48 instructions) per handler in the IVT. This stub provides enough space to save minimal processor state, re-enable interruption collection and external interrupts, and branch to another routine to handle the interruption. Unlike a lightweight interruption handlers described above, heavyweight interruption handlers use general register bank 0 only until they can establish a safe memory context for spilling the interrupted context's state. This allows heavyweight handlers to be interruptible and to take exceptions.

A heavyweight handler stub (i.e. the portion of the handler that is located in the IVT) should determine exactly which type of interruption has occurred based on its offset in the IVT and the contents of the ISR control register. It can then branch out of the IVT to the actual interruption handler. For some heavyweight interruptions (e.g. Data Debug fault), these handlers are typically written in a high-level programming language; for others (e.g. emulation handlers) the interruption can be handled efficiently in Itanium architecture-based assembly code.

The sequence given below illustrates the steps that an Itanium architecture-based heavyweight handler needs to perform to save the interrupted context's state to memory and to create an interruptible execution environment. These steps assume that the low-level kernel code, the kernel backing store, and the kernel memory stack are pinned in the TLB (using a translation register), so that no TLB misses arise from referencing those memory pages. The ordering of the steps below is approximate and other operating system strategies are possible.

1. Copy the interruption resources (IIP, IPSR, IIPA, ISR, IFA, IIB0-1) into bank 0 of the banked registers. To avoid conflicts with processor firmware, use registers GR24-31 for this purpose. Both register bank 0 and the interruption control registers are accessible, since, as described in [Section 3.3.1](#), the processor hardware, upon an interruption always switches to register bank 0, and clears PSR.ic and PSR.i.
2. Preserve the interrupted the predicate registers into bank 0 of the banked registers.
3. Determine whether interruption occurred in the operating system kernel or in user space by inspecting both IPSR.cpl and the memory stack pointer (GR12).
  - a. If IPSR.cpl is zero and the interrupted context was already executing on a kernel stack, then no memory stack switch is required.
  - b. Otherwise, software needs to switch to a kernel memory stack by preserving the interrupted memory stack pointer to a banked register in bank 0, and setting up a new kernel memory stack pointer in GR12.

4. Allocate a “trap frame” to store the interrupted context’s state on the kernel memory stack, and move the interruption state (IIP, IPSR, IIPA, ISR, IFA, IFS, IIB0-1), the interrupted memory stack pointer and the interrupted predicate registers from the banked registers to the trap frame.
5. Save register stack and RSE state by following the steps outlined in [Section 6.11.1, “Switch from Interrupted Context” on page 2:148](#).
  - a. If IPSR.cpl is zero and the interrupted context was not executing on a kernel backing store (determined by inspecting BSPSTORE), then the new kernel BSPSTORE needs to be allocated such that enough space is provided for the RSE to spill all stacked registers. The architectural required maximum RSE spill area is 16KBytes. As a result, BSPSTORE should be offset from the base of the kernel backing store base by at least 16KBytes. This offset can be reduced if the kernel queries PAL for the actual implementation-specific number of stacked physical registers (RSE.N\_STACK\_PHYS). Based on RSE.N\_STACK\_PHYS, the required minimum offset in bytes is:

$$8 * (RSE.N\_STACK\_PHYS + 1 + \text{truncate}((RSE.N\_STACK\_PHYS + 62)/63))$$

Otherwise, the interrupted context was already executing on the kernel backing store. In this case, no new BSPSTORE pointer needs to be setup. The sequence in [Section 6.11.1, “Switch from Interrupted Context” on page 2:148](#), is still required, however, step 6 in that sequence can be omitted.

In either case, the interrupted register stack and RSE state (RSC, PFS, IFS, BSPSTORE, RNAT, and BSP) needs to be preserved, and should be saved either to the trap frame on the kernel memory stack, or to a newly allocated register stack frame.

6. Switch banked register to bank one and re-enable interruption collection as follows:

```

ssm 0x2000 // Set PSR.ic
bsw.1;;    // Switch to register bank 1
srlz.d     // Serialize PSR.ic update

```

With interruptions collection re-enabled, the kernel may now branch to paged code and may reference paged data structures.

7. Preserve branch register and application register state according to operating system conventions.
8. Preserve general and floating-point register state. If this is an involuntary interruption, e.g. an external interrupt or an exception, then software must save the interrupted context’s volatile general register state (scratch registers) to the “trap frame” on the kernel memory stack, or to the newly allocated register stack frame. If this is a voluntary system call then there is no volatile register state. Preserved registers may or may not be spilled depending on operating system conventions. Additionally, the Itanium architecture provides mechanisms to reduce the amount of floating-point register spills and fills. More details on preservation of register context are given in [Section 4.2, “Preserving Register State in the OS” on page 2:551](#).
9. At this point enough context has been saved to allow complete restoration of the interrupted context. Re-enable taking of external interrupts using the ssm instruction as follows:

```
ssm 0x4000 ;; // Set PSR.i
```

There is no need to explicitly serialize the PSR.i update, unless there is a requirement to force sampling of external interrupts right away. Without the serialization, the PSR.i update will occur at the very latest when the next exception causes an implicit instruction serialization to occur.

10. Dispatch interruption service routine (can be high-level programming language routine).
11. Return from interruption service routine.
12. Disable external interrupts as follows:

```
rsm 0x4000 ;; // Clear PSR.i
```

There is no need to explicitly serialize the PSR.i update, since clearing of the PSR.i bit with the `rsm` instruction takes effect at the next instruction group. For details refer to the `rsm` instruction page in [Chapter 2, "Instruction Reference" in Volume 3](#).

13. Restore general and floating-point register state saved in step 8 above.
14. Restore branch register and application register state saved in step 7 above.
15. Disable collection of interruption resources and switch banked register to bank zero as follows:

```
rsm 0x2000 // Clear PSR.ic  
bsw.0;; // Switch to register bank 0  
srlz.d // Serialize PSR update
```

16. Restore register stack and RSE state by following the steps outlined in [Section 6.11.2, "Return to Interrupted Context" on page 2:148](#).
17. Restore interrupted context's interruption state (e.g., IIP, IPSR, IFS) from the "trap frame" on the kernel memory stack.
18. Restore interrupted context's memory stack pointer and predicate registers from the trap frame on the kernel memory stack. This step essentially deallocates the trap frame from the kernel memory stack.
19. Return from interruption using the `rfi` instruction.

Many of the steps shown above are identical for different heavyweight interruptions, so unless there is a specific need to create a different handler for a particular interruption, a common handler can be used. Because external interrupt handlers use the Itanium external interrupt control registers to determine the specific external interrupt vector that needs servicing and to mask off other external interrupt vectors, an external interrupt handler looks somewhat different. Refer to [Section 10.4, "External Interrupt Delivery" on page 2:606](#) for details on writing external interrupt handlers.

### 3.4.3 Nested Interruptions

The Itanium architecture provides a single set of interruption registers whose updates are controlled by PSR.ic. When an IVA-based interruption is delivered and PSR.ic is 0 or in-flight (e.g. during a lightweight interruption handler, or at the beginning of a

heavyweight interruption handler), we say that a nested interruption has occurred. On a nested interruption (other than a Data Nested TLB fault) only ISR is updated by the hardware. All other interruption registers preserve their pre-interruption contents.

With the exception of the Data Nested TLB fault, the Itanium architecture does not support nested interruptions. Data Nested TLB faults are special and are discussed in [Section 5.4.4, “Data Nested TLB Vector” on page 2:576](#). The remainder of this section does not apply to Data Nested TLB faults.

When a nested interruption occurs, the processor will update ISR as defined in [Chapter 8, “Interruption Vector Descriptions”](#) and it will set the ISR.ni bit to 1. A value of 1 in ISR.ni is the only indication to an interruption handler that a nested interruption has occurred. Since all other interruption registers are not updated, there is generally no way for the OS to recover from nested interruptions; the handler for the nested interruption has no context other than ISR for handling the nested interruption. If a nested interruption is detected, it is often useful for the handler to call some function in the OS that logs the state of ISR, IIP, and any other relevant register state to aid in debugging the problem.

## §





This chapter discusses specific context management considerations in the Itanium architecture. With 128 general registers and 128 floating-point registers, the architecture provides a comparatively large amount of state. This chapter discusses various context management and state preservation rules. This chapter introduces some architectural features that help an operating system limit the amount of register spill/fill and gives recommendations to system programmers as to how to use some of the instruction set features.

## 4.1 Preserving Register State across Procedure Calls

The Itanium Software and Runtime Architecture Conventions [SWC] define a contract on register preservation between procedures as follows:

- Scratch Registers (Caller Saves): GR2-3, GR8-11, GR14-GR15, and GR16-31 in register bank 1, FR6-15, and FR32-127. Code that expects scratch registers to hold their value across procedure calls is required to save and restore them.
- Preserved Registers (Callee Saves): GR4-7, FR2-5, and FR16-31. Procedures using these registers are required to preserve them for their callers.
- Stacked Registers: GR32-127, when allocated, are preserved by the RSE.
- Constant Register: GR0 is always 0. FR0 is always +0.0. FR1 is always +1.0.
- Special Use Registers: GR1, GR12, and GR13 have special uses.

Additional architectural register usage conventions apply to GR16-31 in register bank 0 which are used by low-level interrupt handlers and by processor firmware. For details refer to [Section 3.3.1](#).

Itanium general registers and floating-point registers contain three state components: their register value, their control speculative (NaT/NaTVal) state, and their data speculative (ALAT) state. When software saves and restores these registers, all three state components need to be preserved. As described in [Table 4-1](#), software is required to use different state preservation methods depending on the type of register. More details on register preservation are provided in the next two sections.

**Table 4-1. Preserving Intel® Itanium® General and Floating-point Registers**

State Components	General Registers		Floating-point Registers
	GR1-31 (static)	GR32-127 (stacked)	FR2-127
Register Value	<code>st8.spill &amp; ld8.fill</code> preserve register value.	RSE automatically preserves register value.	<code>stf.spill &amp; ldf.fill</code> preserve register value.
Control Speculative State (NaT/NaTVal)	<code>st8.spill &amp; ld8.fill</code> preserve register NaT.	RSE automatically preserves register NaT.	<code>stf.spill &amp; ldf.fill</code> preserve NaTVal.
Data Speculative State (ALAT)	Software must <code>invalidate</code> a register's ALAT state when restoring the register.	RSE and ALAT manage stacked register's ALAT state automatically.	Software must <code>invalidate</code> a register's ALAT state when restoring the register.

### 4.1.1 Preserving General Registers

The Itanium general register file is partitioned into two register sets: GR0-31 are termed the **static general registers** and GR32-127 are termed the **stacked general registers**. Typically, `st8.spill` and `ld8.fill` instructions are used to preserve the static GRs, and the processor's register stack engine (RSE) automatically preserves the stacked GRs.

Using the `st8.spill` and `ld8.fill` instructions, the general register value and its NaT bit are always preserved and restored in unison. However, these instructions do not save and restore a register's data speculative state in the Advanced Load Address Table (ALAT). To maintain the correct ALAT state, software is therefore required to explicitly invalidate a register's ALAT entry using the `invala.e` instruction when restoring a general register. The Itanium calling conventions avoid such explicit ALAT invalidations by disallowing data speculation to preserved registers (GR4-7) across procedure calls.

Spills and fills of general registers using `st8.spill` and `ld8.fill` cause implicit collection and restoration of the accompanying NaT bits to/from the User NaT collection application register (UNAT). The UNAT register needs to be preserved by software explicitly. The spill and fill instructions derive the UNAT bit index of a spilled/filled NaT bit from the spill/fill memory address and not from the spilled/filled register index. As a result, software needs to ensure that the 512-byte alignment offset<sup>1</sup> of the spill/fill memory address is preserved when a general register is restored. This can be an issue particularly for user context data structures that may be moved around in memory (e.g. a `setjmp()` jump buffer).

Unlike the `st8.spill` and `ld8.fill` instructions, the register stack engine (RSE) preserves not only register values and register NaT bits, but it also manages the stacked register's ALAT state by invalidating ALAT that could be reused by software when the physical register stack wraps. This automatic management of ALAT state across procedure calls permits compilers to use speculative advanced loads (`ld.sa`) to perform cross-procedure call control and data speculation in stacked general registers (GR32-127). Whenever software changes the virtual to physical register mapping of the stacked registers, the ALAT needs to be invalidated explicitly using the `invala` instruction. Typically this happens during process/thread context switches or in `longjmp()` when the register stack is reloaded with a new BSPSTORE. Refer to [Section 4.5.1.1, "Non-local Control Transfers \(setjmp/longjmp\)" on page 2:557](#).

The RSE collects the NaT bits of the stacked general registers within the RNAT application register and automatically saves and restores accumulated RNAT collections to/from fixed locations within the register stack backing store. RNAT collections are placed on the backing store whenever BSPSTORE bits{8:3} are all one, which results in one RNAT collection for every 63 registers. When software copies a backing store to a new location, it is required to maintain the backing store's 512-byte alignment offset<sup>2</sup> to ensure that the RNAT collections get placed at the proper offset.

- 
1. The specific requirement is that  $(\text{fill\_address} \bmod 512)$  must be equal to  $(\text{spill\_address} \bmod 512)$ .
  2. The specific requirement is that  $(\text{old\_bspstore} \bmod 512)$  must be equal to  $(\text{new\_bspstore} \bmod 512)$ .

### 4.1.2 Preserving Floating-point Registers

The Itanium architecture encodes a floating-point register's control speculative state as a special unnormalized floating-point number called NaTVal. As a result, Itanium floating-point registers do not have a NaT bit. The architecture provides the `stf.spill` and `ldf.fill` instructions to save and restore floating-point register values and control speculative state. These instructions always generate a 16-byte memory image regardless of the precision of the floating-point number contained in the register.

Preservation of data speculative state associated with floating-point registers needs to be managed by software. As with the general registers, software is required to explicitly invalidate a register's ALAT entry using the `invalidate` instruction when restoring a floating-point register. The Itanium calling conventions avoid such explicit ALAT invalidations by disallowing data speculation to preserved floating-point registers (FR2-5, FR16-31) across procedure calls.

## 4.2 Preserving Register State in the OS

The software calling conventions described in the previous section apply to state preservation across procedure call boundaries. When entering the operating system kernel either voluntarily (for a system call) or involuntarily (for handling an exception or an external interrupt) additional concerns arise because the interrupted user's context needs to be preserved in its entirety.

The Itanium architecture defines a large register set: 128 general registers and 128 floating-point registers account for approximately 1 KByte and 2 KBytes of state, respectively. The architecture provides a variety of mechanisms to reduce the amount of state preservation that is needed on commonly executed code paths such as system calls and high frequency exceptions such as TLB miss handlers.

Additionally, Itanium architecture-based operating systems have opportunities to reduce the amount of context they need to save by distinguishing various kernel entry and exit points. For instance, when entering the kernel on behalf of a voluntary system call, the kernel need only preserve registers as outlined by the calling conventions. Furthermore, the operating system can be sensitive to whether the preserved context is coming from the IA-32 or Itanium instruction set, especially since the IA-32 register context is substantially smaller than the full Itanium register set. Ideally, an Itanium architecture-based operating system should use a single state storage structure which contains a field that indicates the amount of populated state.

Table 4-2 summarizes several key operating system points at which state preservation is needed.

Scratch GRs and FRs, the bulk of all state, only need to be preserved at involuntary interruptions resulting from unexpected external interrupts or from exceptions that need to call code written in a high-level programming language. The demarcation of floating-point registers FR32-127 as "scratch" along with architectural support for lazy state save/restore of the floating-point register file allows software to substantially reduce the overhead of preserving the scratch FRs. See [Section 4.2.2](#) for details.

In principal, preserved GRs and FRs need not be spilled/filled when entering the kernel. Whatever function is called from the low-level interruption handler or the system call entry point will itself observe the calling conventions and preserve the registers. The only occasion when preserved registers need to be spilled/filled is on a process or thread context switch. However, many operating systems provide `get_context()` functions that provide user context upon demand. Although such functions are called infrequently, many operating systems prefer to pay the penalty of spilling preserved registers at system call and at interruption entry points to avoid the complexity of piecing together user state from various potentially unknown kernel stack locations on demand. Fortunately, the amount of preserved Itanium general register state is relatively small, and the Itanium architecture provides additional mechanisms for lazy floating-point state management. See [Section 4.2.2](#) for details.

**Table 4-2. Register State Preservation at Different Points in the OS**

Register Type	Number of Registers	System Call (Voluntary)	Lightweight Interruptions <sup>a</sup> (Involuntary)	Heavyweight Interruptions <sup>b</sup> (Involuntary)	Process/Thread Context Switch (Voluntary)
Scratch GRs	23	no spill/fill required	Untouched (use banked registers)	spill/fill required	no spill/fill required (done at interruption)
Preserved GRs	4	no spill/fill required	Untouched (use banked registers)	no spill/fill required	spill/fill required
Stacked GRs	96	Backing Store Switch	Untouched	Backing Store Switch	Synchronous Backing Store Switch using <code>flushrs</code> <sup>c</sup>
Scratch FRs	106	no spill/fill required	Untouched	spill/fill required	no spill/fill required (done at interruption)
Preserved FRs	20	no spill/fill required	Untouched	no spill/fill required	spill/fill required

- a. For details on lightweight interruption handlers refer to [Section 3.4.1, “Lightweight Interruptions” on page 2:543](#).
- b. For details on heavyweight interruption handlers refer to [Section 3.4.2, “Heavyweight Interruptions” on page 2:544](#).
- c. Refer to [Section 6.11.3, “Synchronous Backing Store Switch”](#) for details.

Stacked GRs are managed by the register stack engine (RSE). On process/thread context switches the operating system is required to completely flush the register stack to its backing store in memory (using the `flushrs` instruction). In cases where the operating system knows that it will return to the user process along the same path, e.g. in system calls and exception handling code, the Itanium architecture allows operating systems to switch the register stack backing store without having to flush all stacked registers to memory. This allows such kernel entry points to switch from the user’s to the kernel’s backing store without causing any memory traffic, as described in the next section.

### 4.2.1 Preservation of Stacked Registers in the OS

A switch from a thread of execution into the operating system kernel, whether on behalf of an involuntary interruption or a voluntary system call, requires preservation of the stacked registers. Instead of flushing all dirty stacked register’s to memory, the RSE can be used to automatically preserve the stacked registers of the interrupted context.

Automatic preservation offers performance benefits: the register stack may contain only a handful of dirty registers, system call parameters can be passed on the register stack, and, upon return to the interrupted context the `loadrs` instruction only needs to restore registers that were actually spilled to memory. Since system call rates scale with processor performance, the RSE offers a key method for reducing the kernel's execution time of a system call.

To ensure operating system integrity the RSE requires a valid backing store (i.e. one with a valid page mapping). The validity of the current backing store depends on the interrupted context. If the interrupted context is itself a kernel thread, then its backing store is in a known state, and no backing store switch is required (assuming that kernel interruptions are nested). If the interrupted context is a user process, then the backing store could be pointing at an invalid region of memory, and software is required to redirect the RSE at a kernel backing store. [Section 6.11.1, "Switch from Interrupted Context" on page 2:148](#) describes the code sequence to switch the RSE backing store without causing memory traffic.

If the kernel redirects the backing store to a kernel memory region, then the kernel must restore the backing store of the interrupted context prior to resumption of the interrupted context. The kernel must also restore the register stack to its interrupted state by manually pulling the spilled registers from the backing store. The kernel uses the `loadrs` instruction to restore stacked registers from the backing store. The `loadrs` instruction requires the backing store pointer to align with any registers spilled from the interrupted context. Thus the kernel should have paired all function calls (`br.call` instructions) with function returns (`br.ret` instructions), or manually manipulated the kernel backing store pointer, so that all kernel contents have been removed from the kernel backing store prior to the `loadrs`. After loading the stacked registers, the kernel can switch to the backing store of the interrupted frame. This code sequence is described in [Section 6.11.1, "Switch from Interrupted Context" on page 2:148](#).

The kernel may occasionally gather the complete interrupted user context, such as to satisfy a debugger request or to provide extended information to a user signal handler. To provide the preserved register stack contents, including NaT values, the kernel must extract the user context values from its backing store.

## 4.2.2 Preservation of Floating-point State in the OS

A full preservation of Itanium floating-point register file requires approximately 2 KBytes of memory. To reduce the frequency of such large register spills and fills, the Itanium architecture offers additional mechanisms for lazy floating-point state management. These features allow the system programmer to eliminate many unnecessary floating-point state spills and fills especially around voluntary and involuntary entries into the kernel, e.g. around system calls, external interrupts and exceptions. Lazy state preservation can provide a significant reduction of memory traffic and hence faster interrupt handlers and system calls, especially since most interrupt handlers and much system code rarely perform floating-point computations.

The 126 non-constant floating-point registers are architecturally divided into the lower set (FR2-31) and the higher set (FR32-127). The Itanium architecture provides two floating-point register set "modified" bits, `PSR.mfl` and `PSR.mfh`, which are set by hardware upon a write to any register in the lower and higher sets, respectively. The "modified" bits are accessible to a user process through the user mask. Additionally,

two “disabled” bits, PSR.dfl and PSR.dfh, are accessible to the privileged software alone. Setting a “disabled” bit causes a fault into the disabled-fp vector upon first use (read or write) of the corresponding register set.

As mentioned earlier, an involuntary kernel entry (e.g. interruption) needs to preserve all scratch floating-point registers. Instead of blindly always spilling all registers, state spills can be conditionalized upon the “modified” bits in the PSR. Additionally, the “disabled” bits allow a deferred, or lazy, approach to both spills and fills. This is particularly useful for “on demand” state motion in an involuntary interruption handler that does not use many floating-point registers. To perform deferred spills on the high set, the handler sets PSR.dfh immediately upon entry. Any reference to a floating-point register in the high set will then fault into the disabled-fp vector which spills the corresponding state to a prearranged store before allowing use within the handler. Lazy state restoration is performed in a similar manner: the handler sets the “disabled” bit just before exit, causing the first reference by the interrupted context to the disabled set to fault into the kernel’s disabled floating-point vector which can then restore the appropriate state. Note the importance of agreeing upon prearranged stores for deferred spill/fill policies and the need for a mechanism to communicate a past fill or spill.

At process or thread context switches all preserved floating-point registers need to be context switched. The higher (scratch) set is also managed here if the context-switch was occasioned by an involuntary interruption (e.g. timer interrupt) which did not already spill the higher set. Use of the “modified” bits by the OS to determine if the appropriate register set is “dirty” with previously unsaved data can help avoid needless spills and fills.

The “modified” bits are intentionally accessible through the user mask so that a user process can provide hints to the OS code about its register liveness requirements. Clearing PSR.mfh, for instance, suggests that the user process does not see the higher register set as containing useful data anymore.

## 4.3 Preserving ALAT Coherency

As described in [Section 4.4.5.3, “Detailed Functionality of the ALAT and Related Instructions” on page 1:65](#), software is required to explicitly invalidate the entire ALAT using the `invala` instruction whenever the virtual to physical register mapping is changed. Typically this occurs when the `clrrb` instruction is used, when a synchronous backing store switch is performed (e.g. in a user-level or kernel thread context switch), or when software “discontinuously” remaps the register to backing store mapping by resetting BSPSTORE (e.g. by calling `longjmp()`).

When returning to a user-process after servicing an involuntary interruptions, an Itanium architecture-based operating system is required to invalidate the entire ALAT using the `invala` instruction. This is required because the operating system may have targeted advanced loads at scratch registers, and thereby altered the user-visible ALAT state.

When returning from a system call, however, full ALAT invalidations can be avoided by using `invala.e` instructions to selectively invalidate ALAT entries of all preserved registers (GR4-7, FR2-5, and FR16-31), or by ensuring that these registers where

never accessible to software during the system call (see [Section 4.2.2](#) for details). This works, because at the system call entry user-code may not have any dependencies on the state of the scratch registers.

## 4.4 System Calls

Reducing the overhead associated with system calls becomes more important as processor efficiency increases. As processor frequencies and pipeline lengths increase, the typical overhead associated with flushing the processor pipeline to effect privilege domain crossings is increased. To reduce system call overhead, the Itanium architecture provides an efficient “enter privileged code” (`epc`) instruction ([page 3:53](#)) that can be paired with the demoting branch return. Additionally, the Itanium architecture provides the traditional `break` instruction ([page 3:29](#)) to enter privileged mode, that is typically paired with the `rfi` instruction ([page 3:236](#)) to return to user mode.

The `epc` instruction offers higher efficiency than the `break` instruction for invoking a kernel system call. Whereas a `break` instruction will always cause a pipeline flush to change privilege level, the `epc` is designed not to. The `break` instruction also passes the system call number as a parameter, and requires a table lookup with an indirect branch to the system call. With the `epc` instruction, the user application can directly branch to the system call code.

More information about `epc`-based system calls is provided in [Section 4.4.1](#). More information about `break`-based system calls is provided in [Section 4.4.2](#). Regardless of whether the `epc` or `break` instruction are used, an Itanium architecture-based operating system needs to check the integrity of system call parameters. In addition to traditional integrity checking of the passed parameter values, the system call handler should inspect system call parameters for set NaT bits as described in [Section 4.4.3](#).

### 4.4.1 `epc`/Demoting Branch Return

To execute a system call with `epc`, a user system call stub branches to an execute-only kernel page containing the system call, using the `br.call` instruction. The kernel page executes an `epc` to raise the privilege level. The privilege level is raised to the privilege level of the page mapping corresponding to the instruction address of the `epc` instruction. The page mapping must be execute-only (see [Section 4.1.1.6](#), “Page Access Rights” for details).

After the kernel completes its system call, it returns to the user system call stub with a `br.ret` instruction. The `br.ret` demotes the privilege level, by restoring the privilege level contained within the PFS application register (PFS.ppl). To ensure operating system integrity `epc` checks that the PFS.ppl field is no greater than the PSR.cpl at the time the `epc` is executed.

As described in [Section 4.2.1](#), interruptions and system calls in a typical Itanium architecture-based operating system need to switch to the kernel register stack backing store upon kernel entry. The `epc` instruction does not disable interrupts nor does it switch the processor to the kernel backing store. As a result, code directly following the `epc` instruction that runs at increased privilege level is still running on the caller’s backing store. It is recommended that software disable external interrupts right after



the `epc` until the switch to the kernel backing store has been completed. Additionally, low-level operating system handlers should not only use `IPSR.cpl`, but should also check `BSPSTORE`, to determine whether they are running on the kernel backing store (imagine an external interrupt being delivered on the first instruction after the `epc`).

## 4.4.2 `break/rfi`

The `break` instruction, when issued in the *i*, *f*, and *m* syllables, specifies an arbitrary 21-bit immediate value. The kernel can choose a specific `break` immediate value to differentiate system calls from other usage of the `break` instruction (such as debug). The `break` instruction jumps to the `break` fault handler, which should be a valid address mapping for each user application, and raises the privilege mode to the most privileged level.

The system call number is an additional parameter passed to the kernel when invoking a system call via the `break` instruction. The system call number must reside in a fixed location. If stored within GR32, then the system call stub must rearrange its input parameters to map to the register stack starting at GR33. This register jostling can be avoided by passing the system call number through a scratch static general register or by using the `break` immediate itself. Additionally, the system call can utilize all eight input registers of the register stack for system call parameters.

## 4.4.3 NaT Checking for NaTs in System Calls

In addition to regular range/value checking on system call arguments, Itanium architecture-based operating systems need to additionally ensure that system call arguments passed in by a user application do not have any NaT bits set. The following code fragment can be used:

```
        mov mask = 0xff
        clrrrb
        ;;
// create register stack frame with only output registers for system call args
        alloc tmp = ar.pfs, 0, 0, 8, 0
        shl mask = mask, syscall_arg_count
        ;;
        mov pr = mask, 0xff00          // define p8 .. p15
        ;;
        cmp.eq p7 = r0, r0            // set p7 to true
        ;;
// test for NaT bits in the input arguments
(p8)   cmp.eq.and p7 = r32, r32      // and type compare clears p7 if r32 is NaT
(p9)   cmp.eq.and p7 = r33, r33
(p10)  cmp.eq.and p7 = r34, r34
(p11)  cmp.eq.and p7 = r35, r35
(p12)  cmp.eq.and p7 = r36, r36
(p13)  cmp.eq.and p7 = r37, r37
(p14)  cmp.eq.and p7 = r38, r38
(p15)  cmp.eq.and p7 = r39, r39
(p7)   br.cond.sptk ok_arguments    // No NaTs found
        ;;
// p7 was cleared by at least one NaT argument
```



## 4.5 Context Switching

This section discusses context switching at the user and kernel levels.

### 4.5.1 User-level Context Switching

#### 4.5.1.1 Non-local Control Transfers (`setjmp/longjmp`)

A non-local control transfer such as the C language `setjmp()/longjmp()` pair requires software to correctly handle the register stack and the RSE. The register stack provides the BSP application register which always contains the backing store address of the current GR32. This permits execution of a `setjmp()` without having to manipulate any register stack or RSE state. All register stack and RSE manipulation is postponed to the much less frequent `longjmp()`.

In `setjmp()` only the RSC, PFS and BSP application registers have to be preserved. This can be accomplished by reading these registers, and without having to disable the RSE. The preserved values will be referred to as `setjmp_rsc`, `setjmp_pfs`, and `setjmp_bsp` further on.

In `longjmp()` restoration of the appropriate register stack and RSE state is more involved, and software needs to take the following steps:

1. Stop RSE by setting RSC.mode bits to zero.
2. Read current BSPSTORE (referred to as `current_bspstore` further down).
3. Find `setjmp()`'s RNAT collection (`rnat_value`).
  - a. Compute the backing store location of `setjmp()`'s RNAT collection as follows:

```
rnat_collection_address{63:0} = setjmp_bsp{63:0} | 0x1F8
```

The RNAT location is computed by setting bits{8:3} of `setjmp()`'s BSP to all ones. This is where `setjmp()`'s RNAT collection will have been spilled to memory.
  - b. If (`current_bspstore > rnat_collection_address`), then the required RNAT collection has already been spilled to the backing store.
  - c. Otherwise if (`current_bspstore <= rnat_collection_address`), the required RNAT collection is incomplete and is still contained in the register stack. To materialize the complete RNAT collection, flush the register stack to the backing store using a `flushrs` instruction.
  - d. Finally, load `rnat_value` from `rnat_collection_address` in memory.
4. Invalidate the contents of the register stack as follows:
  - a. Allocate a zero size register stack frame using the `alloc` instruction.
  - b. Write RSC.loadrs field with all zeros and execute a `loadrs` instruction.
  - c. Invalidate the ALAT using the `invala` instruction.
5. Restore `setjmp()`'s register stack and RSE state as follows:
  - a. Write BSPSTORE with `setjmp_bsp`.
  - b. Write RNAT with `rnat_value`.

- c. Write RSC with `setjmp_rsc`.
  - d. Write PFS with `setjmp_bsp`.
6. Restore `setjmp()`'s return IP into BR7.
  7. Return from `longjmp()` into `setjmp()`'s caller using `br.ret` instruction.

#### 4.5.1.2 User-level Co-routines

The following steps need to be taken to execute a voluntary user-level thread switch.

1. Save all preserved register state of outgoing thread to memory stack. Refer to [Section 4.1](#) for details on preservation of general and floating-point registers.
2. Preserve predicate, branch, and application registers.
3. Flush outgoing register stack to backing store, and switch to incoming thread's backing store as described in [Section 6.11.3, "Synchronous Backing Store Switch" on page 2:148](#). This code sequence includes ALAT invalidation.
4. Switch thread memory stack pointers.
5. Restore incoming thread's predicate, branch, and application registers.
6. Restore incoming thread's preserved register state.

### 4.5.2 Context Switching in an Operating System Kernel

#### 4.5.2.1 Thread Switch within the Same Address Space

To switch between different threads in the same address space the following steps are required:

1. Application architecture state associated with each thread (GRs, FRs, PRs, BRs, ARs) are saved and restores as if this were a user-level coroutine. This is described in [Section 4.5.1.2](#).
2. Memory Ordering: to preserve correct memory ordering semantics the context switch routine needs to fence all memory references and flush cache (`fc`, `fc.i`) operations by executing a `sync.i` and `mf` instruction. More details on memory ordering are given in [Section 2.3](#).

#### 4.5.2.2 Address Space Switching

When an operating system switches address spaces it needs to perform the same steps as a same address space thread switch (described in the previous section). Additionally, however between the saves of the outgoing and the restores of the incoming process, the operating system context switch handler is required to:

1. Save the contents of the protection key registers associated with the outbound context, and then invalidate the protection key registers.
2. Save the default control register (DCR) of the outbound context (if the DCR is maintained on a per-process basis).
3. Save the region registers of the outbound address space.
4. Restore the region registers of the inbound address space.

5. Restore the default control register (DCR) of the inbound context (if the DCR is maintained on a per-process basis).
6. Restore the contents of the protection key registers associated with the inbound context.

§



This chapter introduces various memory management mechanisms of the Itanium architecture: region register model, protection keys, and the virtual hash page table usage models are described. This chapter also discusses usage of the architecture translation registers and translation caches. Outlines are provided for common TLB and VHPT miss handlers.

## 5.1 Address Space Model

The Itanium architecture provides a byte-addressable 64-bit virtual address space. The address space is divided into 8 equally-sized sections called regions. Each region is  $2^{61}$  bytes in size and is tagged with a unique region identifier (RID). As a result, the processor TLBs can hold translations from many different address spaces concurrently, and need not be flushed on address switches. The regions provide the basic virtual memory architecture to support multiple address space (MAS) operating systems.

Additionally, each translation in the TLB contains a protection key that is matched against a set of software maintained protection key registers. The protection keys are orthogonal to the region model and allow efficient object sharing between different address spaces. The protection key registers provide the basic virtual memory architecture to support single address space (SAS) operating systems.

### 5.1.1 Regions

For each of the eight regions, there is a corresponding region register (RR), which contains a RID for that region. The operating system is responsible for managing the contents of the region registers. RIDs are between 18 and 24 bits wide, depending on the processor implementation. This allows an Itanium architecture-based operating system to uniquely address up to  $2^{24}$  address spaces each of which can be up to  $2^{61}$  bytes in virtual size. An address space is made accessible to software by loading its RID into one of the eight region registers.

**Address Translation:** The upper 3 bits of a 64-bit virtual address (bits 63:61) identify the region to which the address belongs; these are called the virtual region number (VRN) bits. When a virtual address is translated to a physical address, the VRN bits select a region register which provides the RID used for this translation. Each TLB entry contains the RID tag bits for the translation it maps; these are matched against the RID bits from the selected region register when the TLB is looked up during address translation. Address translation only succeeds if the RID and VPN bits from the virtual address match the RID and VPN bits from the TLB entry. Note that the VRN bits are used only to select the region register, are not matched against the TLB entries.

**Inserting/Purging of Translations:** When a translation is inserted into the processor TLBs (either by software, or by the processor's hardware page walker), the VRN bits of the virtual address translation being inserted are used only to index the corresponding

region register; they are not inserted into the TLB. Likewise, when software purges a translation from the processor's TLBs, the VRN bits of the address used for the purge are used only to index the corresponding region register and are not used to find a matching translation. Only the RID and VPN bits are used to find overlapping translations in the TLBs.

The fact that the VRN bits are not contained in the processor TLB allows the same address space (identified by a RID) to be referenced through any of the eight region registers. In other words, the combination of RID and VPN establishes a unique 85-bit virtual address, regardless of which VRN (and region register) was used to form the pair. Independence of VRN allows easy creation of temporary virtual mappings of an address space and can accelerate cross-address space copying as described in [Section 5.1.1.3](#).

### 5.1.1.1 RID Management

Before a RID that has been used for one address space can be reused for another address space, all TLB entries relating to the first address space have to be purged. In general, this will require a complete flush of the TLBs of all processors in the system. This can be accomplished by performing an IPI to all processors and executing the `ptc.e` loop described in [Section 5.2.2.2](#) on each processor in the TLB coherence domain.

A more efficient alternative, depending on the size of the defunct address space, might be to perform a series of `ptc.ga` operations on one processor to tear down just the translations used by the recycled RID. Some processor implementations support an efficient region-wide purge page size such that this can be accomplished with a single `ptc.ga` operation.

The frequency of these global TLB flushes can be reduced by using a RID allocation strategy that maximizes the time between use and reuse of a RID. For example, RIDs could be assigned by using a counter that is as wide as the number of implemented RID bits and that is incremented after every assignment. Only when the RID counter wraps around it is necessary to do a global TLB flush. After the flush the operating system can either remember the in-use RIDs or it can re-assign new RIDs to all currently active address spaces.

### 5.1.1.2 Multiple Address Space Operating Systems

Multiple address space (MAS) operating systems provide a separate address space for each process. Typically, only when a process is running is its address space visible to software.

The application view of the virtual address space in the MAS OS model is a contiguous 64-bit address space, though normally not all of this virtual address space is accessible by the application. At least one of the 8 regions must be used to map the OS itself so that the OS can handle interruptions and system services invoked by the application.

The OS chooses a region ID and a region (e.g. region 7) into which to map itself during the boot process and usually does not change this mapping after enabling address translation. The other seven regions may be used to map process-private code and data; code and data that are shared amongst multiple processes; to map large files; temporary mappings to allow efficient cross-address space copies (see [Section 5.1.1.3](#)); and, for operating systems which use it, the long format VHPT.

In a MAS OS, the RID bits act as an address space identifier or tag. For each process-private region, a unique RID is assigned to that process by the OS. If a process needs multiple process-private regions (e.g. the process requires a private 64-bit address space), the OS assigns multiple unique RIDs for each such region. Because each translation in the processor's TLBs is tagged with its RID, the TLBs may contain translations from many different address spaces (RIDs) concurrently. This obviates the need for the OS to purge the processor's TLBs upon an address space switch. When the OS performs a context switch from process A to process B, the OS need only remove process A's private RIDs from the CPU's region registers and replace them with process B's private RIDs.

### 5.1.1.3 Cross-address Space Copies in a MAS OS

The use of regions, region registers, and RIDs provides a mechanism for efficient address space-to-address space copies. Because translations are tied to RIDs and not to a particular static region, a MAS OS can easily copy a memory range from one address space to another by temporarily remapping the target memory location to another region. This remapping is accomplished simply by placing the RID to which the target location belongs into a different region register and then performing the copy from source to target directly.

For example, assume a MAS OS wishes to copy an 8-byte buffer from virtual address 0x000000000A00000 of the currently executing process (process A) to virtual address 0x000000000A00000 of another process (process B):

```

    movl r2 = (2 << 61)
    mov  r3 = process_b_rid
    movl r4 = 0x000000000A00000
    movl r5 = 0x4000000000A00000;; // reference process B through RR[2]
    mov  rr[r2] = r3 ;; // put process B RID into RR[2]
    srlz.d // serialize RR write
copyloop:
    ld8  r6 = [r4] ;; // read buffer from process A addr space
    st8  [r5] = r6 // store buffer into process B addr
space
    (p4)br copyloop // loop until done
    mov  r3 = original_rr2_rid ;;
    mov  rr[r2] = r3 ;; // restore RR[2] RID
    srlz.d // serialize RR write

```

When the OS switches to process B and places process B's RID into RR[0] and resumes execution of process B, the process can reference the message via virtual address 0x000000000A00000. Note that no new translations need to be created to make the sequence shown above work; because translations are tagged by RID and not by region, all existing translations for process B's address space are visible regardless of which region the reference is made to, as long as the region register for that region contains the correct process B RID. Note that the sequence shown above is intended for illustrative purposes only; the OS may need to perform other steps as well to perform a cross-address space copy.

## 5.1.2 Protection Keys

The Itanium architecture provides two mechanisms for applying protection to pages. The first mechanism is the access rights bits associated with each translation. These bits provide privilege level-granular access to a page. The second mechanism is the protection keys. Protection keys permit domain-granular access to a page. These are especially useful for mapping shared code and data segments in a globally shared region, and for implementing domains in a single address space (SAS) operating system.

Protection key checking is enabled via the PSR.pk bit. When PSR.pk is 1, instruction, data, and RSE references go through protection key access checks during the virtual-to-physical address translation process.

All processors based on the Itanium architecture implement at least 16 protection key registers (PKRs) in a protection key register cache. The OS is responsible for maintaining this cache and keeping track of which protection keys are present in the cache at any given time.

Each protection key register contains the following fields:

- v – valid bit. When 1, this register contains a valid key, and is checked during address translation whenever protection keys are enabled (PSR.pk is 1).
- wd – write disable. When 1, write permission is denied to translations which match this protection key, even if the data TLB access rights permit the write.
- rd – read disable. When 1, read permission is denied to translations which match this protection key, even if the data TLB access rights permit the read.
- xd – execute disable. When 1, execute permission is denied to translations which match this protection key, even if the instruction TLB access rights give execute permission.
- key – protection key. An 18- to 24-bit (depending on the processor implementation) unique key which tags a translation to a particular protection domain.

When protection key checking is enabled, the protection key tagged to a referenced translation is checked against all protection keys found in the protection key register cache. If a match is found, the protection rights specified by that key are applied to the translation. If the access being performed is allowed by the matching key, the access succeeds. If the access being performed is not allowed by the matching key (e.g. instruction fetch to a translation tagged with a key marked 'xd'), a Protection Key Permission fault is raised by the processor. The OS may then decide whether to terminate the offending program or grant it the requested access.

If no match is found, a Protection Key Miss fault is raised by the processor, and the OS must insert the correct protection key into the PKRs and retry the access.

Protection keys can be used to provide different access rights to shared translations to each process. For example, assume a shared data page is tagged with a protection key number of 0xA. Two processes share this data page: one is the producer of the data on this page, and the other is only a consumer. When the producer process is running, the OS will insert a valid PKR with the protection key 0xA and the 'wd' and 'rd' bits cleared, to allow this process to both read and write this page. When the consumer process is



running, the OS will insert a valid PKR with the protection key 0xA and the 'rd' bit cleared, to allow this process to read from the page. However, the 'wd' bit for this PKR will be set when the consumer process is running to prevent it from writing the page.

The processor hardware has no notion of which protection keys belong to which process. The only check the hardware performs is to compare the protection key from the translation to any valid protection keys in the PKR cache. On a context switch, the OS must purge any valid protection keys from the PKRs which would provide access rights to the switched-to context that are not allowed. The OS may purge an existing PKR by performing a move to PKR instruction with the same key as the existing PKR, but with the PKR valid bit set to 0.

Protection keys can be read from the processor's data TLBs via the `tak` instruction. However, instruction TLB key values cannot be read directly. Software must keep track of these values in its own data structures.

### 5.1.2.1 Single Address Space Operating Systems

Processes in a single address space (SAS) OS all cohabit a global address space. SAS operating systems running on a processor based on the Itanium architecture can view the RID bits as effectively extending the single virtual address space to between 79 and 85 bits (depending on the number of RID bits implemented by the processor). This address space is then divided into between  $2^{18}$  and  $2^{24}$  61-bit regions, up to eight of which may be accessed concurrently.

Note that there is no "SAS OS" or "MAS OS" mode in the Itanium architecture. The processor behavior is the same, regardless of the address space model used by the OS. The difference between a SAS OS and a MAS OS is one of OS policy: specifically how the RIDs and protection keys are managed by the OS, and whether different processes are permitted to share RIDs for their private code and data. Multiple, unrelated processes in a SAS OS may share the same RID for their private pages; it is the responsibility of the OS to use protection keys and the protection key registers (PKRs) to enforce protection. In a MAS OS, the unique per-process RIDs enforce this protection.

Hybrid SAS/MAS models that combine unique RIDs for process-private regions and shared RIDs with protection keys for per-page memory protection in shared regions are also possible.

## 5.2 Translation Lookaside Buffers (TLBs)

All processors based on the Itanium architecture implement one or more translation lookaside buffers (TLBs) for fast virtual-to-physical address translation. The architecture provides instructions for managing instruction and data TLBs as separate structures.

Both the instruction and data TLBs are further divided into a set of translation registers (TRs), which are managed exclusively by software and are "locked down" to pin critical address translations (e.g. kernel memory); and a set of translation cache entries (TCs), which can be managed by both software and the processor hardware. The TRs are divided into slots, each of which are individually addressable on insertion by software.

The TCs are treated as a set associative cache and are not addressable by software. The TC replacement policy is determined by software. All processor models implement at least 8 instruction and 8 data TRs, and at least 1 instruction and 1 data TC entry.

Software inserts translations into the TLBs via insertion instructions. There are four variants of insertion instructions. `itr.i` and `itr.d` insert a translation into the specified instruction or data TR slot, respectively. `itc.i` and `itc.d` insert a translation into a hardware-selected instruction or data TC entry, respectively.

Software TR purge instructions also distinguish between the instruction and data TRs (`ptr.i`, `ptr.d`). TC purge instructions do not.

## 5.2.1 Translation Registers (TRs)

Once a translation is inserted by software into a TR, it remains in that TR until either the translation is overwritten by software, or the translation is purged. TRs are used by the OS to pin critical address translations; all memory references made to a TR translation will always hit the TLB and will never cause the processor's hardware page walker to walk the VHPT or raise a fault. Examples of memory areas that the OS might cover with one or more TRs are the Interruption Vector Table, critical interruption handlers not contained completely in the Interruption Vector Table, the root-level page table entries, the long format VHPT, and any other non-pageable kernel memory areas.

Two address translations are said to overlap when one or more virtual addresses are mapped by both translations. Software must ensure that translations in an instruction TR never overlap other instruction TR or TC translations; likewise, software must ensure that translations in a data TR never overlap other data TR or TC translations. If an overlap is created, the processor will raise a Machine Check Abort.

The processor hardware will never overwrite or purge a valid TR. TRs that are currently unused may be used by the processor hardware as extra TC entries, but if software subsequently inserts a translation into an unused a TR, the TC translation will be purged when the insertion is executed.

### 5.2.1.1 TR Insertion

To insert a translation into a TR, software performs the following steps:

1. If `PSR.ic` is 1, clear it and execute a `srlz.d` instruction to ensure the new value of `PSR.ic` is observed.
2. Place the base virtual address of the translation into the IFA control register.<sup>1</sup>
3. Place the page size of the translation into the `ps` field of the ITIR control register. If protection key checking is enabled, also place the appropriate translation key into the key field of the ITIR control register. See below for an explanation of protection keys.
4. Place the slot number of the instruction or data TR into which the translation is be inserted into a general register.
5. Place the base physical address of the translation into another general register.

---

1. The upper 3 bits (VRN) of this address specify a region register whose contents are inserted along with the rest of the translation. See [Section 5.1.1](#) for details.

6. Using the general registers from steps 4 and 5, execute the `itr.i` or `itr.d` instruction.

A data or instruction serialization operation must be performed after the insert (for `itr.d` or `itr.i`, respectively) before the inserted translation can be referenced.

Software may insert a new translation into a TR slot already occupied by another valid translation. However, software must perform a TR purge to ensure that the overwritten translation is no longer present in any of the processor's TLB structures.

Instruction TR inserts will purge any instruction TC entries which overlap the inserted translation, and may purge any data TC entries which overlap it. Data TR inserts will purge any data TC entries which overlap the inserted translation and may purge any instruction TC entries which overlap it.

Software may insert the same (or overlapping) translation into both the instruction TRs and the data TRs. This may be desirable for locked pages which contain both code and data, for example.

### 5.2.1.2 TR Purge

To purge a TR from the TLBs, software performs the following steps:

1. Place the base virtual address of the translation to be purged into a general register.<sup>1</sup>
2. Place the address range in bytes of the purge into bits {7:2} of a second general register.
3. Using these two GRs, execute the `ptr.d` or `ptr.i` instruction.

A data or instruction serialization operation must be performed after the purge (for `ptr.d` or `ptr.i`, respectively) before the translation is guaranteed to be purged from the processor's TLBs.

**Note:** The TR purge instruction operates independently of the slot into which the translation was originally inserted.

A `ptr.d` instruction will never purge an overlapping translation in an instruction TR, but may purge an overlapping translation in an instruction TC; likewise, a `ptr.i` instruction will never purge an overlapping translation in a data TR, but may purge an overlapping translation in a data TC.

A TR purge does not modify the page tables nor any other memory location, nor does it affect the TLB state of any processor other than the one on which it is executed.

## 5.2.2 Translation Caches (TCs)

The TC array acts as a cache of the dynamic working set for data and instruction translations. It is managed by software (via `itc` and `ptc` instructions) and, optionally by hardware, if the processor provides a hardware page walker (HPW) and the walker is enabled. See [Section 5.3](#) below.

---

1. The upper 3 bits (VRN) of this address specify a region register whose contents are used as part of the translation to be purged. See [Section 5.1.1](#) for details.

The size, associativity, and replacement policy of the TC array are implementation-dependent. With the exception of the forward progress rules defined in [Section 4.1.1.2, "Translation Cache \(TC\)" on page 2:49](#), software cannot depend on the existence or life-span of a TC translation, as a TC entry may be replaced or invalidated by the hardware at any time.

### 5.2.2.1 TC Insertion

To insert a TC entry, software performs the following steps:

1. If PSR.ic is 1, clear it and execute a `srlz.d` instruction to ensure the new value of PSR.ic is observed.
2. Place the base virtual address of the translation into the IFA control register.<sup>1</sup>
3. Place the page size of the translation into the ps field of the ITIR control register. If protection key checking is enabled, also place the appropriate translation key into the key field of the ITIR control register. See below for an explanation of protection keys.
4. Place the base physical address of the translation into a general register.
5. Using the general register from step 4, execute the `itc.i` or `itc.d` instruction.

A data or instruction serialization operation must be performed after the insert (for `itc.d` or `itc.i`, respectively) before the inserted translation can be referenced.

Instruction TC inserts always purge overlapping instruction TCs and may purge overlapping data TCs. Likewise, data TC inserts always purge overlapping data TCs and may purge overlapping instruction TCs.

### 5.2.2.2 TC Purge

There are several types of TC purge instructions. Unlike the other TLB management instructions, the TC purge instructions do not distinguish between instruction and data translations; they will purge any matching translations in either the data or instruction TC arrays.

#### 5.2.2.2.1 ptc.l

The most basic TC purge is the local TC purge instruction (`ptc.l`). To purge a TC from the local processor TLBs, software performs the following steps:

1. Place the base virtual address of the translation to be purged into a general register.<sup>2</sup>
2. Place the address range in bytes of the purge into bits {7:2} of a second general register.
3. Using these two GRs, execute the `ptc.l` instruction.

---

1. The upper 3 bits (VRN) of this address specify a region register whose contents are inserted along with the rest of the translation. See [Section 5.1.1](#) for details.

2. The upper 3 bits (VRN) of this address specify a region register whose contents are used as part of the translation to be purged. See [Section 5.1.1](#) for details.

A data or instruction serialization operation must be performed after the `ptc.l` before the translation is guaranteed to be no longer visible to the local data or instruction stream, respectively.

The `ptc.l` instruction does not modify the page tables nor any other memory location, nor does it affect the TLB state of any processor other than the one on which it is executed.

The `ptc.l` instruction ensures that all prior stores are made locally visible before the actual purge operation is performed. Consider the following code sequence:

```
st8 [VHPT] = <new_translation>
ptc.l <old_translation>
srlz.i
```

The `ptc.l` instruction will purge the translation only after the local store update is seen. If there was a hardware-initiated VHPT walk for the same translation, it would either insert the *old\_translation* in the TLB before the `ptc.l` executes and then get purged by the `ptc.l`, or insert the *new\_translation* after both the local store update and `ptc.l` purge are complete.

#### 5.2.2.2.2 `ptc.e`

To purge all TC entries from the local processor's TLBs, software uses a series of `ptc.e` instructions. Software must call the `PAL_PTCE_INFO` PAL routine at boot time to determine the parameters needed to use the `ptc.e` instruction. Specifically, `PAL_PTCE_INFO` returns:

- `tc_base` – an unsigned 64-bit integer denoting the beginning address to be used by the first `ptc.e` instruction in the purge loop.
- `tc_counts` – two unsigned 32-bit integers packed into a 64-bit parameter denoting the loop counts of the outer and inner purge loops. `count1` (outer loop) is contained in bits {63:32} of the parameter, and `count2` (inner loop) is contained in bits {31:0} of the parameter.
- `tc_strides` – two unsigned 32-bit integers packed into a 64-bit parameter denoting the loop stride of the outer and inner purge loops. `stride1` (outer loop) is contained in bits {63:32} of the parameter, and `stride2` (inner loop) is contained in bits {31:0} of the parameter.

Software then executes the following sequence:

```
disable_interrupts();
addr = tc_base;
for (i = 0; i < count1; i++) {
    for (j = 0; j < count2; j++) {
        ptc.e addr;
        addr += stride2;
    }
    addr += stride1;
}
enable_interrupts();
```

A data or instruction serialization operation must be performed after the sequence shown above before the translations are guaranteed to be no longer visible to the local data or instruction stream, respectively.

The `ptc.e` instruction does not modify the page tables nor any other memory location, nor does it affect the TLB state of any processor other than the one on which it is executed.

### 5.2.2.2.3 ptc.g, ptc.ga

The Itanium architecture supports efficient global TLB shutdowns via the `ptc.g` and `ptc.ga` instructions. These instructions obviate the need for performing inter-processor interrupts to maintain TLB coherence in a multiprocessor system. A TLB coherence domain is defined as a group of processors in a multiprocessor system which maintain TLB coherence via hardware.

For the remainder of this section, `ptc.g` refers to both the `ptc.g` and `ptc.ga` instructions, except where otherwise noted.

The number of `ptc.g` operations that can be in progress at any time is implementation dependent, and can be determined from the `max_purges` return parameter of `PAL_VM_SUMMARY`. Attempting to execute more than the maximum allowed number of simultaneous `ptc.g` purges will have undefined effects, including possibly raising a Machine Check Abort on one or more processors. Software should implement some semaphoring mechanism to ensure that not more than the maximum `ptc.g` purges allowed are in flight at any one time.

A `ptc.g` instruction is a release operation; all memory references that precede a `ptc.g` in program order are made visible to all other processors before the `ptc.g` is made visible. To guarantee visibility of the `ptc.g` prior to a particular point in program execution, software must use another release operation or a memory fence.

To purge a translation from all TLBs in the coherence domain, software performs the following steps:

1. Acquire the semaphore.
2. Place the base virtual address of the translation to be purged into a general register.
3. Place the address range in bytes of the purge into bits {7:2} of a second general register.
4. Using these two GRs, execute the `ptc.g` instruction. Note that the `ptc.g` instruction must be followed by a `stop`.
5. Release the semaphore.

Global purges can be batched together by performing multiple `ptc.g` instructions prior to releasing the lock.

A data or instruction serialization operation must be performed after the sequence shown above before the translations are guaranteed to be no longer visible to the local data or instruction stream, respectively. To guarantee the translations are no longer visible on remote processors, a release operation or memory fence instruction is required after the `ptc.g` instruction.

The `ptc.g` instruction does not modify the page tables nor any other memory location. It affects both the local and all remote TC entries in the TLB coherence domain. It does not remove translations from either local or remote TR entries. If a `ptc.g` overlaps a translation contained in a TR on the local processor, the local processor will raise a Machine Check Abort; if the `ptc.g` overlaps a translation contained in a TR on any remote processor in the coherence domain, no Machine Check Abort is raised.

The `ptc.ga` variant of the global purge instruction behaves just like the `ptc.g` variant, but it also removes any ALAT entries which fall into the address range specified by the global shutdown from all remote processors' ALATs. The `ptc.ga` variant is intended to be used whenever a translation is remapped to a different physical address to ensure that any stale ALAT entries are invalidated. Note that the `ptc.ga` is not guaranteed to affect the issuing processor's ALAT; processor implementations may optionally remove matching entries from the local ALAT, therefore software must perform a local ALAT invalidation via the `invala` instruction on the processor issuing the `ptc.ga` to ensure the local ALAT is coherent.

Note that processors based on the Itanium architecture may support one or more implementation-dependent purge sizes; some implementations may include a region-wide purge. The `PAL_VM_PAGE_SIZE` firmware call returns the supported page sizes for purges for a particular processor implementation. Refer to [Section 11.10.1, "PAL Procedure Summary"](#) for details. When software wishes to purge an address range that is much larger than the largest supported purge size from all TCs in the coherence domain, performance may be enhanced by issuing inter-processor interrupts to all processors and using the `ptc.e` loop described in [Section 5.2.2.2](#) on each processor, instead of issuing many `ptc.g` instructions from one processor.

`ptc.g` instructions do not apply to processors outside the coherence domain of the processor issuing the `ptc.g` instruction. Systems with multiple coherence domains must use a platform-specific method for maintaining TLB coherence across coherence domains.

## 5.3 Virtual Hash Page Table

The Itanium architecture defines a data structure that allows for the insertion of TLB entries by a hardware mechanism. The data structure is called the "virtual hash page table" (VHPT) and the hardware mechanism is called the VHPT walker.

Unlike the IA-32 page tables, the Itanium VHPT itself is virtually mapped, i.e. VHPT walker references can take TLB faults themselves. Virtual mapping of the page tables is needed because the page tables for  $2^{64}$  address space are quite large and typically do not fit into physical memory.

The Itanium architecture prescribes the format of a leaf-node page table entry (PTE) seen by the VHPT walker, but does not impose an OS page table data structure itself. As summarized in [Table 5-1](#), the architecture support two different VHPT formats:

- **Short** format uses 8-byte PTEs, and is a linear page table. The short format VHPT does not contain protection key information (there are not enough PTE bits for that). Short format is a per-region linear page table, i.e. the PTEs and hash function are independent of the RID. The short format prefers use of a self-mapped page table. The short format VHPT is an efficient representation for address spaces that contain only a few large clusters of pages, like the text, data, and stack segments of applications running on a MAS operating system.
- **Long** format uses 32-byte PTEs, and is a hashed page table. The hash function embedded in hardware. The long format supports protection keys and the use of multiple page sizes in a region. The long format hash and tag functions incorporate the RID, and allows multiple address space translations to be present in the same VHPT. The long format is expected to be used either as a cache of the real OS page

tables, or as a primary page table with collision chains. The long format VHPT is a much better representation for address spaces that are sparsely populated, since the short format VHPT has a linear layout and would consume a large amount of memory. Single address space operating systems may prefer the long format VHPT for this reason.

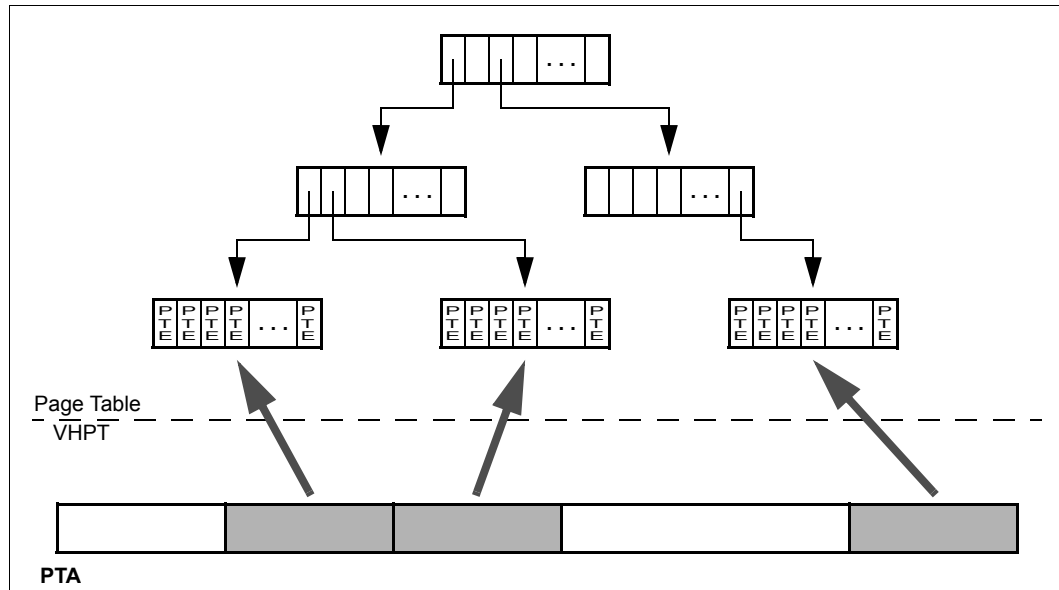
**Table 5-1. Comparison of VHPT Formats**

Attribute	Short Format	Long Format
Entry Size	8 Byte	32 Byte
Lookup	Linear	Hashed
Protection Keys	No	Yes
Page Size	per region	per entry

### 5.3.1 Short Format

The short format VHPT is a per-region linear table that contains translation entries for every page in the region’s virtual address space. This makes the VHPT very large, but since the VHPT itself lives in virtual address space only those parts of the VHPT that actually contain valid translation entries have to be present in physical memory. If the operating system’s page table is a hierarchical data structure and the last level of the hierarchy is a linear list of translations, the VHPT can be mapped directly onto the page table as shown in Figure 5-1.

**Figure 5-1. Self-mapped Page Table**



If the VHPT walker tries to access a location in the VHPT for which no translation is present in the TLB, a VHPT Translation fault is raised. The original address for which the VHPT walker was trying to find an entry in the VHPT is supplied to the fault handler in the IFA register. The fault handler can use this address to traverse the page table and insert a translation into the TLB that maps the address the VHPT walker tried to access (in IHA) to the page that contains the corresponding leaf page table.



### 5.3.2 Long Format

The long format VHPT is organized as a hash table which contains a subset of all translation entries. The long format VHPT entries contain a 8-byte field that is ignored by the VHPT walker and can be used by the operating system to link VHPT entries to software-walkable hash collision chains if it uses the VHPT as its primary page table. The size of the long format VHPT is usually kept small enough to keep a mapping for it in one of the translation registers (TRs), so it is not necessary to handle VHPT translation faults.

The long format hash algorithm is based on the per-region preferred page size, but a translation for a larger page can still be entered into the VHPT by subdividing the large page into multiple smaller pages with the preferred page size and placing an entry for the large page at all VHPT locations that correspond to the smaller pages.

### 5.3.3 VHPT Updates

Visibility of VHPT updates to a VHPT walker on another processor follows the rules outlined in [Section 4.1.7, “VHPT Environment” on page 2:67](#). Since a global TLB purge has release semantics, prior modifications to the VHPT will be visible to operations that occur after the TLB purge operation.

Atomic updates to short format VHPT entries can easily be done through 8-byte stores. For atomic updates of long format VHPT entries, the “ti” flag in bit 63 of the tag field can be utilized as follows:

- Set the “ti” bit to 1.
- Issue a memory fence.
- Update the entry.
- Clear the “ti” bit through a store with release semantics.

## 5.4 TLB Miss Handlers

The Itanium architecture enables lightweight TLB fault handlers by providing individual entry points for different excepting conditions and by pre-setting the translation insertion registers for the various types of TLB faults. The following subsections list the typical steps for resolving each kind of fault.

### 5.4.1 Data/Instruction TLB Miss Vectors

These faults occur when the data or instruction TLB required for a data access or instruction fetch is not found in the processor TLBs, the VHPT walker is enabled, and:

- Either the VHPT walker aborted the walk (for any reason and at any time), or
- The VHPT walker found the translation but the insert failed (due to tag mismatch in the long format or badly formed PTE), or
- The walker is not implemented on this processor.

There is a separate vector for each fault type (data and instruction).

Since the VHPT walker may abort a walk at any time and raise these faults, software must always be able to handle all TLB faults, even when the VHPT walker is enabled. Upon entry to these fault handlers, the IHA, ITIR, and IFA control registers are initialized by the hardware as follows:

- IHA – contains the virtual address of the hashed page table address corresponding to the reference which raised the fault.
- ITIR – contains the default translation information for the reference which raised the fault (i.e. for the virtual address contained in IFA). The access key field is set to the region ID from the RR corresponding to the faulting address. The page size field is set to the preferred page size (RR.ps) from the RR corresponding to the faulting address.
- IFA – the virtual address of the bundle (for instruction faults) or data reference (for data faults) which missed the TLB.

The fault handler for a short format VHPT performs the following steps, at a minimum, to handle the fault:

1. Move IHA into a general register, chosen by convention to match the register expected by the nested TLB fault handler.
2. Perform an 8-byte load into another general register from the address contained in this general register to grab the VHPT entry. Note that the format of these first 8 bytes is identical to the format required for TLB insertion. If the VHPT is not mapped by a TR, software must be prepared to handle a nested TLB fault when performing this load.
3. Using the general register from step 2 that holds the contents of the VHPT entry, perform a TC insert (`itc.i` for instruction faults, `itc.d` for data faults).
4. In an MP environment, reload the VHPT entry from step 2 into a third general register and compare the value to the one loaded in step 2. If the values are not the same, then the VHPT has been modified by another processor between steps 2 and 3, and the entry will have to be re-inserted. In this case, purge the entry just inserted using a `ptc.l` instruction. The fault will re-occur after the `rfi` in step 5 (unless the VHPT walker succeeds on the next TLB miss) and the fault handler will re-attempt the insertion. (Uniprocessor environments may skip this step.)
5. `rfi`.

For a long format VHPT, additional steps are required to load bytes 16-23 of the VHPT entry and check for the correct tag (the correct tag for the reference can be generated using the `ttag` instruction). If the tags do not match, this indicates a VHPT collision, and the handler must proceed to walk the operating system's collision chain manually to find the correct entry. The handler may then choose to swap places between the correct entry and the VHPT entry. Note that the pointers for a collision chain can be stored in bytes 24-31 of the VHPT entry format since these bytes are ignored by the VHPT walker.

If the default page size and key are not sufficient, the handler must also perform additional steps to load the correct page size and key into the ITIR register before performing the TC insert in step 3 of the sequence shown above.

## 5.4.2 VHPT Translation Vector

Processors based on the Itanium architecture does not perform recursive TLB hardware page walks. Since the VHPT is itself a virtually addressed structure, each reference performed by the walker itself goes through the TLBs and may miss. These faults are raised when the VHPT walker is enabled, but the walker misses the TLBs when attempting to service a TLB miss caused by the program.

There is a separate vector for each fault type (data and instruction).

Upon entry to this fault handler, the IHA, IFA, and ITIR control registers are initialized by the hardware as follows:

- IHA – contains the virtual address of the hashed page table address corresponding to the reference which raised the fault.
- ITIR – contains the default translation information for the VHPT address which missed the TLBs (i.e. for the virtual address contained in IHA). The access key field is set to the region ID from the RR corresponding to the VHPT address. The page size field is set to the preferred page size (RR.ps) from the RR corresponding to the VHPT address.
- IFA – contains the original faulting address that the VHPT walker was attempting to resolve.

The fault handler for a short format VHPT performs the following steps, at a minimum, to handle the fault:

1. Move the IHA register into a general register.
2. Perform a thash instruction using the general register from step 1 This will produce, in the target register, the VHPT address of the VHPT entry that maps the VHPT entry corresponding to the original faulting address (i.e. the address in IFA).
3. Using the target general register of the thash from step 2 as the load address, perform an 8-byte load from the VHPT. Note that the format of these first 8 bytes is identical to the format required for TLB insertion. Software must be prepared to take a nested TLB fault if this load misses the TLBs.
4. Move the IHA value from the general register written in step 1 into the IFA register.
5. Using the general register from step 3 that holds the contents of the VHPT entry, perform a data TC insert using the `itc.d` instruction. (VHPT references always go through the data TLBs.)
6. In an MP environment, reload the VHPT entry from step 3 into a different general register and compare the value to the one loaded in step 3. If the values are not the same, then the VHPT has been modified by another processor between steps 3 and 4, and the entry will have to be re-inserted. In this case, purge the entry just inserted using a `ptc.l` instruction. The fault will re-occur after the `rfi` in step 7 (unless the VHPT walker succeeds on the next TLB miss) and the fault handler will re-attempt the insertion. (Uniprocessor environments may skip this step.)
7. `rfi`.

For a long format VHPT, additional steps are required to load bytes 16-23 of the VHPT entry and check for the correct tag; see [Section 5.4.1](#) for more details.

A separate structure other than the VHPT may be used to back VHPT translations, in which case the handler would not use the thash instruction to generate the address of the translation mapping the VHPT entry corresponding to the original faulting address. Instead, the handler would use the operating system's own mechanism for finding VHPT back-mappings. Other schemes for handling VHPT misses are also possible, but are beyond the scope of this document.

### 5.4.3 Alternate Data/Instruction TLB Miss Vectors

These faults are raised when an instruction or data reference misses the processor's TLBs and the VHPT walker is not enabled for the faulting address, i.e. TLB misses are handled entirely in software. Operating systems which do not wish to use the VHPT walker can disable the walker and use these fault vectors for software TLB fill handlers. The OS may also choose to enable the walker on a per-region basis and use these vectors to handle misses in regions where the walker is disabled.

Upon entry to these fault handlers, the IFA and ITIR registers are initialized by the hardware as follows:

- ITIR – contains the default translation information for the reference which raised the fault (i.e. for the virtual address contained in IFA). The access key field is set to the region ID from the RR corresponding to the faulting address. The page size field is set to the preferred page size (RR.ps) from the RR corresponding to the faulting address.
- IFA – the virtual address of the bundle (for instruction faults) or data reference (for data faults) which missed the TLB.

The OS needs to lookup the PTE for the faulting address in the OS page table, convert it to the architected insertion format (see [Section 4.1.1.5, "Translation Insertion Format"](#)), and insert it into the TLB. The mechanism used to handle these faults is OS specific and is beyond the scope of this document.

### 5.4.4 Data Nested TLB Vector

To enable efficient handling of software TLB fills, the Itanium architecture provides a dedicated Data Nested TLB fault vector. The Data Nested TLB fault handler is intended to be used by the Data TLB fault handler, which allows the OS to page the page tables themselves. When PSR.ic is 0, any data reference that misses the TLB and would normally raise a Data TLB Miss fault (e.g. a load performed by the Data TLB fault handler to the page tables) will vector to the Data Nested TLB fault handler instead. Because IFA is not updated when PSR.ic is 0, the Data Nested TLB fault handler must get the faulting address from the general register used as the load address in the Data TLB fault handler<sup>1</sup>. Unlike other nested interruptions, the hardware does *not* update ISR when a Data Nested TLB fault is delivered.

---

1. This requires a register usage convention between all TLB miss handlers and the Data Nested TLB miss handler.

The processor will not deliver a Data Nested TLB fault when PSR.ic is in-flight; Data Nested TLB faults are only delivered when PSR.ic is 0. If PSR.ic is in-flight, any data references which miss the TLB and trigger a fault will raise a Data TLB fault, and the processor will set ISR.ni to 1.

### 5.4.5 Dirty Bit Vector

The operating system is expected to lookup the PTE for the faulting address in the OS page table and load the PTE into a general register  $r_x$ . It can then set the “dirty” bit in  $r_x$  and write the updated PTE back to the page table. To continue execution, the OS must insert the updated PTE into the data TLB or update the PTE memory image and let the VHPT walker perform the insertion.

### 5.4.6 Data/Instruction Access Bit Vector

The operating system is expected to lookup the PTE for the faulting address in the OS page table and load the PTE into a general register  $r_x$ . It can then set the “access” bit in  $r_x$  and to continue execution, the OS must either:

- Write the updated PTE back to the page table, and have the VHPT walker pick it up, or
- Insert the updated PTE into the TLB using `itc.i`  $r_x$  for instruction pages, and `itc.d`  $r_x$  for data pages, or
- Step over the instruction/data access bit fault by setting the IPSR.ia or IPSR.da bits prior to performing an `rfi`.

### 5.4.7 Page Not Present Vector

Forward the fault to the operating system’s virtual memory subsystem.

### 5.4.8 Data/Instruction Access Rights Vector

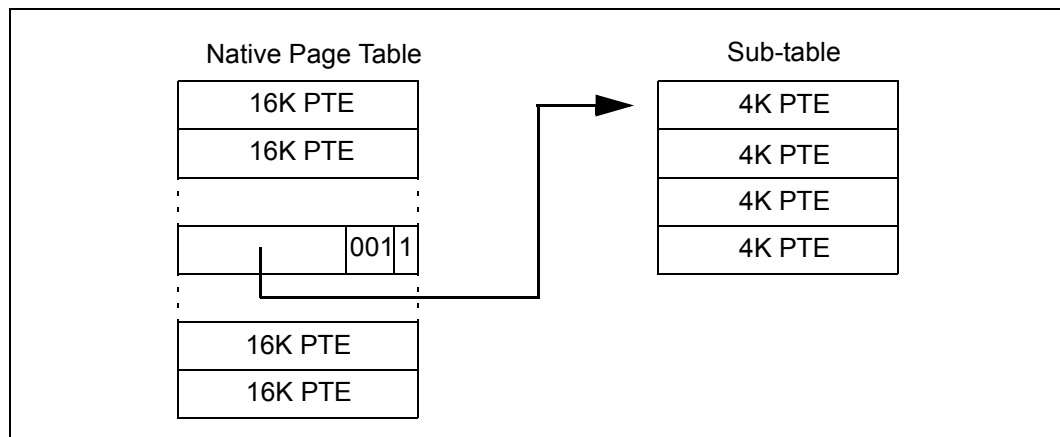
Forward the fault to the operating system’s virtual memory subsystem.

## 5.5 Subpaging

The native page size an Itanium architecture-based operating system will choose for its page tables is likely be larger than the architectural minimum page size of 4 KB. Some legacy IA-32 applications, however, expect a page protection granularity of 4 KB. The following technique allows support for these applications with minimal impact on the native, larger page size paging mechanism.

A special type of entry is used in the native page table to mark pages that are subdivided into smaller 4 KByte units. The entry must have its memory attribute field set to the architecturally “software reserved” encoding (binary 001), and it carries a pointer to an array of 4 KB subentries in its most significant 59 bits. An example using a native page size of 16 KB is shown in [Figure 5-2](#). The use of the “software reserved” memory attribute prevents the VHPT walker from attempting to insert the entry into the TLB.

**Figure 5-2. Subpaging**



When one of the subdivided pages is referenced and does not have a translation in the TLB, a TLB miss will occur. The handler for this fault can then use the faulting address to calculate the appropriate offset into the sub-table and insert the corresponding 4KByte PTE into the TLB.

Some care is required to ensure forward progress for IA-32 instructions. Each IA-32 instruction can reference up to 8 distinct memory pages during its execution (see also [Section 10.6.3, "IA-32 TLB Forward Progress Requirements"](#)). This means that the fault handler not only has to insert the PTE for the current fault into the TLB, but also the PTEs for up to seven faults that occurred before, if these faults originate from the same IA-32 instruction. This can be accomplished by maintaining a buffer for the most recent faulting IIP and for the parameters of up to 7 TLB insertions. If a TLB fault occurs while executing in IA-32 mode and the IIP matches the most recent IIP, all TLB insertions in the buffer have to be repeated and the parameters for the new TLB fault must be added to the buffer. Otherwise, the buffer can be cleared out and the most recent IIP can be updated. The buffer also has to be cleared out when a TLB purge occurs.

§

An Itanium architecture-based operating system needs to handle exceptions generated by control speculative loads (`ld.s` or `ld.sa`), data speculative loads (`ld.a`) and architectural loads (`ld`) in different ways.

Software does not have to worry about control or data speculative loads potentially hitting uncacheable memory with side-effects, since `ld.s`, `ld.sa`, and `ld.a` instructions to non-speculative memory are always deferred by the processor for details refer to [Section 4.4.6, “Speculation Attributes” on page 2:79](#). As a result, compilers can freely use control and data speculation to all program variables.

Control speculative loads require special exception handling and the Itanium architecture provides a variety of deferral mechanisms for handling of control speculative exception handling. This is discussed in [Section 6.1](#).

The Itanium architecture supports different control speculation recovery models. These are discussed in [Section 6.2](#).

Handling of exceptions caused by architectural and data speculative loads is the same, except for emulation of unaligned data speculative references, which require special unaligned emulation handling. This is discussed in [Section 6.3.1](#).

## 6.1 Exception Deferral of Control Speculative Loads

Exceptions that occur on control speculative loads (`ld.s` or `ld.sa`) can be handled by the operating system in different ways. The operating system can configure a processor based on the Itanium architecture in three ways:

- **Hardware-Only Deferral:** automatic hardware deferral of all control speculative exceptions. In this case, the processor hardware will always defer excepting control speculative loads without invoking the operating system.
- **Combined Hardware/Software Deferral:** automatic deferral of some control speculative exceptions, but deliver others to software. In this case, some exceptions will result in hardware deferral as described above, other exceptions will be reported to the operating system. The operating system fault handlers can identify that an exception has been caused by a control speculative load (`ISR.sp` will be 1). Furthermore, OS handlers can software-defer an exception on a control speculative load by setting `IPSR.ed` to 1 prior to `rfi`-ing back to the `ld.s` or `ld.sa`. This allows an operating system to service “cheap” non-fatal exceptions (e.g. simple TLB misses), while software-deferring both “expensive” non-fatal (e.g. page faults) as well as fatal exceptions (e.g. non-recovery protection violation).
- **Software-Only Deferral:** processor is configured to deliver all control speculative exceptions to software. In this case, operating system software handles all non-fatal control speculative exceptions, and software-defers all fatal control speculative exceptions.

Details on these three models are discussed in the next three sections as well as in [Section 5.5.5, “Deferral of Speculative Load Faults” on page 2:105.](#)

### 6.1.1 Hardware-only Deferral

Hardware only deferral is configured by setting all speculation deferral bits in the DCR register (dd, da, dr, dx, dk, dp and dm) to 1. All excepting control speculative loads are automatically deferred by the processor. As a result, all excepting control speculative loads that hit non-fatal exceptions, e.g. a TLB miss or a page fault, will be deferred by the processor hardware, and will cause speculation recovery code to be invoked. This can cause speculation recovery code to be invoked more often than strictly necessary.

### 6.1.2 Combined Hardware/Software Deferral

Setting of a DCR deferral bit to 1 results in hardware deferral by the processor, whereas clearing of a deferral bit causes exceptions to be delivered to software. The operating system may want to configure the processor to deliver control speculative exceptions to its handlers for certain non-fatal faults such as TLB misses or protection key misses. Early handling of these exceptions avoids unnecessary invocation of speculation recovery code, and the associated performance penalty. This is especially useful for exceptions handlers whose overhead is small. Note that handlers will also be invoked for excepting control speculative loads that have been hoisted from not taken paths, and therefore are not needed. As a result, software handling of control speculative exceptions is recommended only for statistically infrequent light weight fault handlers such as TLB miss or protection key miss handlers. If, while handling the exception, the operating system determines that this instance of the exception may require too much effort, e.g. a TLB miss turns out to be a page fault, the handler still has the choice of software-deferring the exception.

### 6.1.3 Software-only Deferral

Software only deferral is configured by clearing all speculation deferral bits in the DCR register (dd, da, dr, dx, dk, dp and dm) to 0. Control speculative loads that hit any Debug, Access Bit, Access Rights, Key Permissions, Key Miss, or Not Present fault, or that suffer a TLB miss or a VHPT Translation fault will be delivered to software.

## 6.2 Speculation Recovery Code Requirements

As described by [Table 6-1](#), code generators for the Itanium architecture are not always required to generate speculation recovery code for all forms of speculation. Compilers and operating systems can collaborate to provide two models for handling of recovery from failed control speculation:

- `ITLB.ed=1` (application with recovery code – the default): The compiler generates appropriate recovery code for all `ld.s` instructions, as well as for `ld.sa` and `ld.a` instructions that have speculatively executed uses. Speculation failure of `ld.sa` and `ld.a` instructions that have no speculatively executed uses can be recovered by a `ld.c` instruction, and hence do not require recovery code. The operating system may defer non-fatal exceptions.



- `ITLB.ed=0` (no control speculative recovery code): The compiler generates recovery code only for `ld.sa` and `ld.a` instructions that have speculatively executed uses. Speculation failure of `ld.sa` and `ld.a` instructions that have no speculatively executed uses can be recovered by a `ld.c` instruction, and hence do not require recovery code. Speculation failure of `ld.s` instructions does not require recovery code, because, in this model, the operating system must guarantee that only fatal exceptions will be deferred. This requires software-only deferral of all potential non-fatal exceptions. The motivation for this model is that the absence of `chk.s` instructions and their associated recovery code may make for shorter and more compact in-line code, especially in loops with tight instruction schedules.

**Table 6-1. Speculation Recovery Code Requirements**

Usage Model	OS May Defer Non-fatal Exceptions on Control Speculative Loads (ITLB.ed=1)	OS Must Not Defer Non-fatal Exceptions on Control Speculative Loads (ITLB.ed=0)
<b>No Speculative Load Uses</b>		
<code>ld.s</code>	Recovery code required; Invoked by <code>chk.s</code> or non-speculative use of speculative value recovers from failed control speculation.	No recovery code required; OS handles all non-fatal exceptions speculatively.
<code>ld.sa, ld.a</code>	No recovery code required; <code>ld.c</code> recovers from failed data speculation.	
<b>With Speculative Load Uses</b>		
<code>ld.s</code>	Recovery code required; invoked by <code>chk.s</code> or non-speculative use of speculative value recovers from failed control speculation.	No recovery code required; OS handles all non-fatal exceptions speculatively.
<code>ld.sa, ld.a</code>	Recovery code required; <code>chk.a</code> recovers from failed data speculation.	

Presence or lack of control speculation recovery code is communicated from the compiler and the runtime system to the operating system by marking the code page's page table entry `ed`-bit appropriately (this bit is referred to as `ITLB.ed`). When `ITLB.ed` is 1, the operating system will expect recovery code to be present; when `ITLB.ed` is 0 no recovery code is expected. When a control speculative load takes an exception, the code page's `ITLB.ed` bit is copied into `ISR.ed` and is made available to the operating system exception handler. Furthermore, a set `ISR.sp` bit indicates that an exception was caused by a control speculative load.

## 6.3 Speculation Related Exception Handlers

### 6.3.1 Unaligned Handler

Misaligned control and data speculative loads, as well as architectural loads, are not required to be handled by the processor. As a result, the operating system's unaligned reference handler has to be prepared to emulate such misaligned memory references, especially in cases where the application has not provided any recovery code (see [Section 6.2](#) for details). Furthermore, misaligned data speculative loads (`ld.sa` or `ld.a`) must be forced failed by the unaligned emulation handler, because the ALAT cannot track all sizes of misalignment for store conflict detection.

The following pseudo code outlines the basic steps for an unaligned reference handler:

1. Ensure that only `ISR.r` is 1, and that `ISR.w`, `ISR.x`, and `ISR.na` are 0.
2. Inspect the `ISR.sp` and `ISR.ed`. If both are 1, then defer this control speculative load by setting `IPSR.ed` and `rfi-ing`.
3. Crack the instruction opcode to determine:
  - a. Size of the load: 1, 2, 4, 8, 10 bytes
  - b. Type of the load: `ld.sa`, `ld.s`, `ld.a`, `ld.c.clr`, `ld.c.nc` or `ld`
  - c. Target, source and post-increment registers of the load
4. If this is a data speculative load (`ld.sa`, or `ld.a`), invalidate the target register's ALAT entry using an `invala.e` instruction, and `rfi`.
5. If this is a `ld.c.clr` instruction invalidate the target register's ALAT entry using an `invala.e` instruction.
6. Emulate the memory read of the load instruction by updating the target register as follows:
  - a. Validate that emulated code has the access rights to the target memory location at the privilege level that it was running prior to taking the alignment fault. The `regular_form probe` instruction can be used on the first and the last byte of the unaligned memory reference. If both probes succeed the memory reference may proceed.
  - b. Using architectural `ld` instructions if the emulated operation is a `ld` or a `ld.c` (either clear or no clear flavor).
  - c. Using `ld.s` instructions if the emulated operation is a `ld.s`. The result in the target register may end up with its NaT bit or NaTVal set, if one of the parts of emulation causes an exception. If `ITLB.ed` is 0 (no control speculation recovery code), then the misaligned `ld.s` may only be deferred if a fatal exception occurred on either half or the `ld.s` emulation.
7. If this is a post-increment load, compute the new value for the source register.

## §

This chapter introduces several common emulation handlers that an Itanium architecture-based operating system must support. A general overview is given for:

- Unaligned Reference Handler – emulation of misaligned memory references that the processor hardware cannot handle, or has been configured to fault on.
- Unsupported Data Reference Handler – emulation of memory operations that the processor hardware does not support. Examples are `semaphore`, `ldfe` or `stfe` operations to uncacheable memory.
- Illegal Dependency Fault Handler – this is a fatal condition that operating system needs to provide error logging functionality for.
- Long Branch Handler – the Itanium processor does not implement the long branch instruction. When encountered on the Itanium processor, long branches must be emulated by the operating system.

Floating-point software assist emulation handlers are not discussed here, but are presented in [Chapter 8, “Floating-point System Software.”](#) Additionally, [Section 5.5.1, “Efficient Interruption Handling”](#) on page 2:102 discusses more details about emulation code in the Itanium architecture.

## 7.1 Unaligned Reference Handler

Misaligned memory references that are not supported by the processor cause Unaligned Reference Faults. This behavior is implementation specific but typically occurs in cases where the access crosses a cache line or page boundary. In cases where the operating system chooses to emulate misaligned operations, some special cases need to be considered:

- Emulation of control and data speculative loads as well as advanced check and “regular” loads requires special attention. For details consult [Section 6.3.1, “Unaligned Handler”](#) on page 2:581.
- Emulation of unaligned semaphores, especially when interacting with IA-32 code require special attention. For details consult [Section 2.1.3.2, “Behavior of Uncacheable and Misaligned Semaphores”](#) on page 2:509.

IA-32 programs do not use the Itanium architecture-based handler to support unaligned references. The hardware that supports IA-32 execution provides the appropriate behavior if alignment checking is disabled through `EFLAGS.ac`. If an unaligned reference occurs in IA-32 code when `EFLAGS.ac` is set to enable alignment checking, alignment faults are delivered to a different vector from the unaligned reference handler. Specifically they are delivered to the `IA_32_Exception(AlignmentCheck)` vector; see [Chapter 9, “IA-32 Interruption Vector Descriptions”](#) for details.

## 7.2 Unsupported Data Reference Handler

Processors based on the Itanium architecture do not support all types of memory references to all memory attributes. In particular:

- Semaphore operations to uncacheable memory are not supported. For details consult [Section 2.1.3.2, “Behavior of Uncacheable and Misaligned Semaphores”](#) on page 2:509.
- A 10-byte memory access, e.g. `ldfe` or `stfe`, to uncacheable memory are not supported by all implementations.

The handler for 10-byte memory accesses must go through the following steps to emulate the `ldfe` or `stfe` instructions:

- Determine that the opcode at the faulting address is an `ldfe` or `stfe`. On control-speculative flavors of these instructions (`ldfe.s` or `ldfe.sa`) processor hardware always defers the unsupported data reference fault. In other words, software does not have to emulate control-speculative fault deferral.
- If the instruction is an advanced load `ldfe.a` then the emulation handler should invalidate the ALAT entry of the appropriate floating-point target register using the `invala.e` instruction. Furthermore, a zero should be returned in the floating-point target register.
- If the instruction is a regular `ldfe` or `stfe`, then software must emulate the load or store behavior of the instruction taking the appropriate faults if necessary.
- If the instruction is the base register update form, update the appropriate base register.

A number of these steps may require the use of self-modifying code to patch instructions with the appropriate operands (for example, the target register of the `inval.e` must be patched to the destination register of the `ldfe` or `stfe`). See [Section 2.5, “Updating Code Images”](#) on page 2:531 for more information.

## 7.3 Illegal Dependency Fault

The Itanium instruction sequencing rules specify that, generally speaking, instructions within an instruction group are free of dependencies as described in [Section 3.4, “Instruction Sequencing Considerations”](#) on page 1:39. A dependency violation occurs anytime a program violates read-after-write (RAW), write-after-write (WAW) or write-after-read (WAR) resource dependency rules within an instruction group.

As [Section 3.4.4, “Processor Behavior on Dependency Violations”](#) on page 1:44 describes, an implementation may provide hardware to detect and report dependency violations. It is important to note that the presence and capabilities of such hardware is implementation specific. A processor based on the Itanium architecture reports dependency violations through the General Exception Vector with an `ISR.code` of 8.

It is recommended that operating systems log the dependency violation and then terminate the offending application, as hardware behavior is undefined when a dependency violation occurs.

## 7.4 Long Branch

The Itanium architecture supports “long” branches with a 64-bit offset. This provides IP-relative conditional- and call-type branches that can reach any address in a 64-bit address space. These instructions use the MLX template, and similar to the move long instruction (`movl`), they encode their immediate in the L and the X slot of the bundle.

The Intel Itanium processor does not support the long branch instruction, `brl`, and requires the operating system to emulate its behavior. When an Itanium processor encounters a `brl` instruction, it vectors to the Illegal Operation Fault handler, regardless of the branches’ qualifying predicate. This handler is expected to emulate the long branch instruction in software. A general outline of the long branch emulation handler is as follows:

- The emulation handler reads the IIP, IPSR, and predicates at the time of the fault.
- If the fault occurred in IA-32 code or if the fault did not occur in slot 2 of a bundle (IPSR.ri is not 2), the handler passes the fault to regular illegal operation fault handler.
- Two floating-point registers are spilled into the integer register file to get ready to load the bundle.
- The emulation handler speculatively loads the 128-bit bundle at the faulting IP using the integer form of the floating-point load pair instruction. This instruction is chosen because it operates atomically (see [Section 4.5, “Memory Datum Alignment and Atomicity”](#)). Using two 64-bit integer loads would require the handler to ensure that another agent does not update the bundle between the two reads.
- If the speculation fails, the recovery code re-issues the load. Before re-issuing an architectural load, the processor must first re-enable PSR.ic to be able to handle potential TLB misses when reading the opcode from memory. In other words, this becomes a heavyweight handler. For details see [Section 3.4.2, “Heavyweight Interruptions” on page 2:544](#). Once the opcode has been read from memory successfully flow of the emulation continues at the next step.
- The 128-bit bundle is moved from the FP register file into two integer registers and the FP registers are restored to their contents at the time of the fault.
- The handler extracts the fields necessary to decode the instruction (specifically, the qp, template, major opcode, and btype or b<sub>1</sub> fields of slot 2). It also determines the value of the qualifying predicate of the instruction in slot 2 from the contents of the predicate register at the time of the fault. Itanium instructions are always stored in memory in little-endian memory format. When extracting bit fields from the loaded opcode current processor endianness (PSR.be) must be taken into account.
- The emulation handler passes the fault off to the regular illegal operation fault handler if the bundle is not an MLX or if the faulting instruction is not a `brl.cond` or `brl.call`.
- If the faulting instruction is a not-taken `brl.cond` or `brl.call`, the code prepares to change the IIP to the address of the sequential successor of the faulting branch (i.e. IIP + 16) and jumps ahead to the trap detection code mentioned below.
- If the faulting instruction is a taken `brl.call`, the handler emulates the appropriate behavior of the call. The code uses a `br.call` to move the appropriate values into CFM and AR[PFS]. There are several details, however. First, the branch register update from the call must be backed out (as it is not the correct update for the `brl.call`). Second, AR[PFS].ppl must be set based on the cpl at the time of the fault (which is given by IPSR.cpl). Finally, the code must update the branch register

specified in the `brl.call` instruction with the IP of the successor of the `brl.call` (predication helps here as the Itanium instruction set does not provide an indirect move to branch register instruction).

- The handler forms the 60-bit immediate IP-offset for the `brl` target from the `i` and `imm20` fields from the X syllable of the bundle (the `brl` instruction) and the `imm39` field from the L syllable of the bundle.
- The handler checks to see if there are any traps to be taken. Specifically, it verifies that the next IP is at an implemented address (the specific test depends on whether the processor was in virtual or physical mode at the time of the fault as `IPSR.it` indicates), that taken branch traps are not enabled if the branch is taken, and that single stepping is not enabled.
- If a trap condition is detected, the `ISR.code` and `ISR.vector` fields are set up as appropriate and the handler jumps to the appropriate operating system entry point after restoring the predicates at the time of the fault and setting the IIP to the appropriate address.
- If no trap occurs, the handler restores the predicates and returns to the faulting code at the appropriate IP.

A processor based on the Itanium architecture typically does not fault on instructions with false qualifying predicates. However, an implementation may take an Illegal Operation Fault on an MLX instruction with a false predicate; the Itanium processor is such an implementation. This implies that the `brl` emulation handler must also provide the means to skip the faulting instruction when its qualifying predicate is false.

## §

This chapter details the way floating-point exceptions are handled in the Itanium architecture and how the architecture can be used to implement the ANSI/IEEE Std. 754-1985 for Binary Floating-point Arithmetic (IEEE-754). It is useful in creating and maintaining floating-point exception handling software by operating system writers.

## 8.1 Floating-point Exceptions in the Intel® Itanium® Architecture

Floating-point exception handling in the Itanium architecture has two major responsibilities. The first responsibility is to assist a hardware implementation to conform to the Itanium floating-point architecture specification. The Floating-point Software Assistance (FP SWA) Exception handler supports this conformance and is included as a driver in the Unified Extensible Firmware Interface (UEFI). The second responsibility is to provide conformance to the IEEE-754 standard. The IEEE Floating-point Exception Filter (IEEE Filter) supports providing this conformance.

When a floating-point exception occurs, a minimal amount of processor state information is saved in interruption control registers. Additional information is contained in the Floating-point Status Register (FPSR), i.e. application register (AR40). This register contains the IEEE exception enable controls, the IEEE rounding controls, the IEEE status flags, and information to determine the dynamic precision and range of the result to be produced.

When a floating-point exception occurs, execution is transferred to the appropriate interruption vector, either the Floating-point Fault Vector (at vector address 0x5c00) or the Floating-point Trap Vector (at vector address 0x5d00.) There the operating system may handle the exception or save additional processor information and arrange for handling of the exception elsewhere in the operating system. Floating-point exception faults must be handled differently than other faults. Correcting the condition that caused the fault (e.g. a page not present is brought into memory) and re-executing the instruction is how most other faults are handled. For floating-point faults, software is required to emulate the operation and continue execution at the next instruction as is normally done for traps. Part of this emulation needs to include a check for any lower priority traps that would have been raised if the instruction hadn't faulted, e.g. a single-step trap.

### 8.1.1 Software Assistance Exceptions (Faults and Traps)

There are three categories of Software Assistance (SWA) exceptions that must be handled by the operating system. The first two categories, SWA Faults and SWA Traps, are implementation dependent and could be generated by any Itanium floating-point arithmetic instruction that contains a status field specifier in the instruction's encoding. An implementation may choose to raise a SWA Fault as needed. The SWA Trap can only be raised under special circumstances. The third category, architecturally mandated

SWA Faults, is limited to the scalar reciprocal and scalar reciprocal square-root approximation instructions and is not implementation dependent. It is required for the correctness of the divide and square root algorithms.

### 8.1.1.1 SWA Faults

The Itanium architecture allows an implementation to raise SWA faults as required. Therefore an implementation-independent operating system must be able to emulate the architectural behavior of all FP instructions that can raise a floating-point exception. However, hardware implementations will limit the cases that raise SWA Faults for performance reasons. The most likely cases would be for the consumption of denormalized or unnormalized operands and production of denormalized results.

The general flow of the SWA Fault handler is as follows:

1. From the interruption instruction bundle pointer (IIP) and faulting instruction index (IPSR.ri), determine the FP instruction that faulted.
2. From the instruction, decode the opcode, static precision, status field and input/output register specifiers.
3. Read the data from the input registers.
4. From the opcode and the FPSR's status field, decode the result range and precision.
5. From the ISR.code, determine that a SWA Fault has occurred, if not go to the last step.
6. From the FPSR, determine if the trap disabled or trap enabled result is wanted.
7. Emulate the Itanium instruction to produce the Itanium architecture specified result.
8. Place the result(s) in the correct FR and/or PR registers, if required.
9. Update the flags in the appropriate status field of the FPSR, if required.
10. Update the ISR.code if required. (This is required if the SWA fault has been translated into an IEEE fault or trap.)
11. Check to see if an IEEE fault or trap needs to be raised. If so, then queue it to the IEEE Filter, otherwise continue checking for lower priority traps that may need to be raised and if required invoke their handler. When finished, continue execution at the next instruction.

### 8.1.1.2 SWA Traps

SWA traps are allowed in the Itanium architecture as an optimization for cases when the hardware implementation has produced the result of the first (exponent unbounded) IEEE rounding<sup>1</sup> and can't continue with the second (exponent bounded) IEEE rounding to produce the final result. One option for the implementation would be to throw away the first IEEE rounding result and raise the SWA Fault. The SWA Fault handler would then have to redo the computation of the first IEEE rounding. A potentially more efficient option would be for the implementation to return the first IEEE rounding result and raise a SWA trap. Returning the first IEEE rounded result is

---

1. ANSI/IEEE Std 754-1985 sections 7.3 Overflow and 7.4 Underflow.

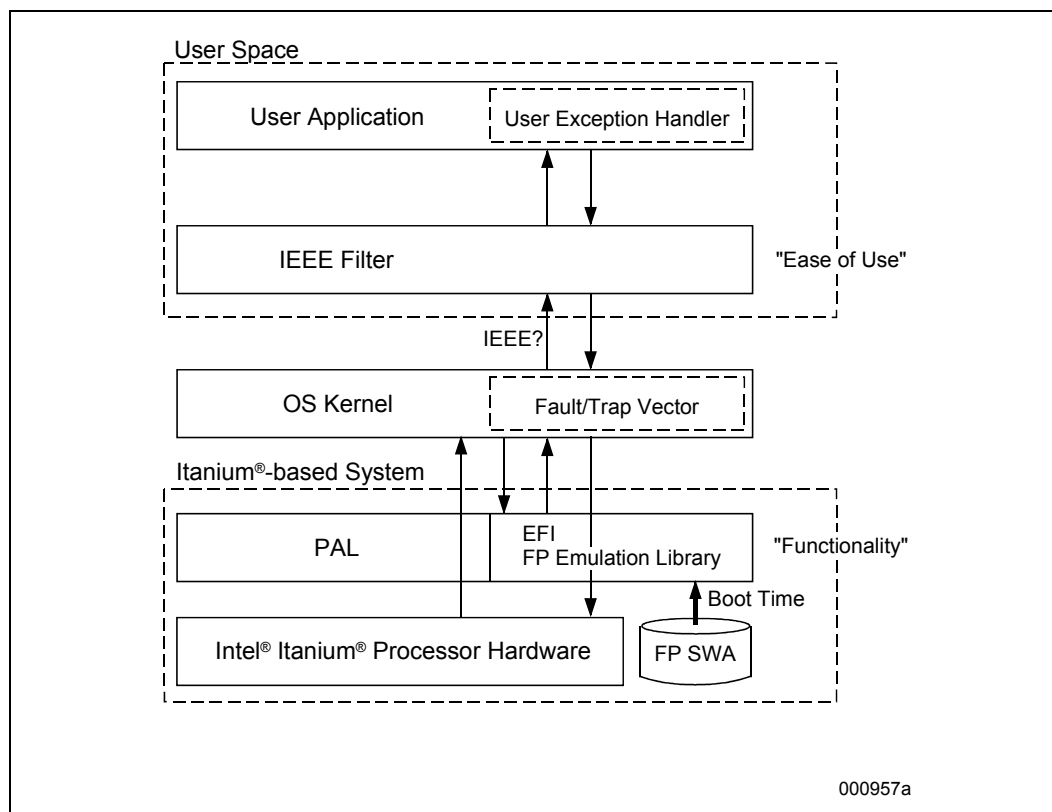


the same as what is done when the IEEE Overflow or Underflow exceptions are enabled. However, hardware implementations will limit the cases that raise SWA Traps for performance reasons. The most likely case would be for the production of denormalized results.

For tiny<sup>1</sup> results, the SWA Trap handler has the simpler task of taking the intermediate result of the first IEEE rounding, the ISR.fpa and ISR.i status bits and producing the correctly rounded and signed minimum normal, denormal or zero. For huge<sup>2</sup> results, the SWA Trap handler has the even simpler task of taking the intermediate result of the first rounding and producing the correctly signed maximum representable normal or infinity, based on the sign of the result, the rounding direction, and the result precision and range.

**Note:** The Itanium architecture also allows for SWA Traps to be raised when the result is just Inexact. This is a trivial case for the SWA Trap handler, since result of the second IEEE rounding is identical to the first IEEE rounding.

**Figure 8-1. Overview of Floating-point Exception Handling in the Intel® Itanium® Architecture**



The general flow of the SWA Trap handler is as follows:

1. From the interruption instruction previous address (IIPA) and exception instruction index (ISR.ei), determine the FP instruction that trapped.
2. From the instruction, decode the opcode, static precision, status field and

- 
1. Tiny numbers are non-zero values with a magnitude smaller than the smallest normal floating-point number.
  2. Huge numbers have values larger in magnitude than the largest normal floating-point number.

input/output register specifiers.

3. From the `ISR.code` and `FPSR` trap enable controls, determine if a SWA Trap has occurred, if not go to the last step.
4. Read the first IEEE rounded result from the FR output register.
5. From the opcode and the status field, decode the result range and precision.
6. From the `ISR.code`'s `FPA`, `O`, `U`, and `I` status bits and the intermediate result, produce the Itanium architecture specified result.
7. Place the result in the output FR register.
8. Update the flags in the appropriate status field of the `FPSR`, if required.
9. Update the `ISR.code` if required. (This is required if the SWA trap has been translated into an IEEE trap.)
10. Check to see if an IEEE trap needs to be raised. If so, then queue it to the IEEE Filter, otherwise continue checking for lower priority traps that may need to be raised and if required invoke their handler. When finished, continue execution at the next instruction.

### 8.1.1.3 Approximation Instructions and Architecturally Mandated SWA Faults

The scalar approximation instructions, `frcpa` and `frsqrta`, can raise architecturally mandated SWA Faults. This occurs when their input operands are such that they are potentially prevented from generating the correct result by the usual software algorithms that are employed for divide and square root. The reasons for this are that these algorithms may suffer from underflow, overflow, or loss of precision, because the inputs or result are at the extremes of their range. For these special cases, the SWA Fault handler must use alternate algorithms to provide the correct quotient or square root and place that result in the floating-point destination register. The predicate destination register is also cleared to indicate the result is not an approximation that needs to be improved via the iterative algorithm.

The parallel approximation instructions `fprcpa` and `fprsqrta` have situations similar to the scalar approximation instruction's architecturally mandated SWA Faults. This occurs when their input operands are such that they are potentially prevented from generating the correct result by the usual software algorithms that are employed for divide and square root. For these special cases, instead of generating a SWA Fault, the parallel approximation instructions indicate that software must use alternate algorithms to provide the correct reciprocal or square-root reciprocal by clearing the destination predicate register. The cleared predicate is the indication to the parallel IEEE-754 divide and square root software algorithms that alternative algorithms are required to produce the correct IEEE-754 quotient or square root.

## 8.1.2 The IEEE Floating-point Exception Filter

The Itanium architecture supports the reporting of the five IEEE-754 standard floating-point exceptions and the IA-32 Denormal Operand exception. In the Itanium architecture the Denormal Operand exception is expanded to the Denormal/Unnormal Operand exception. When referring to the IEEE-754 exceptions in the Itanium architecture the Denormal/Unnormal Operand exception is included.

At the application level, a user floating-point exception handler could handle the Itanium floating-point exception directly. This is the traditional operating system approach of providing a signal handler with a pointer to a machine-dependent data structure. It would be more convenient for the application developer if the operating system were to first transform the results to make them IEEE-754 conforming and then present the exception to the user in an abstracted manner. It is recommended that the operating system include such a software layer to enable application developers that want to handle floating-point exceptions in their application. The IEEE Floating-point Exception Filter provides this convenience to the developer through three functions.

- The first function of the IEEE Filter is to map the Itanium architecture's result to the IEEE-754 conforming result. This includes the wrapping of the exponent for Overflow and Underflow exceptions. The Itanium architecture keeps the exponent in the 17-bit format, which is not wrapped (i.e. scaled) with the appropriate value for the destination precision.
- The second function of an IEEE Filter is to transform the interruption information to a format that is easier to interpret and to invoke a user handler for the exception. The user's handler may then provide a value to be substituted for the IEEE default result, based on the operation, exception and inputs.
- The third function of the filter is to hide the complexities of the parallel instructions from the user. If a floating-point fault occurs in the high half of a parallel floating-point instruction and there is a user handler provided, the parallel instruction is split into two scalar instructions. The result for the high half comes from the user handler, while the low half is emulated by the IEEE Filter. The two results are combined back into a parallel result and execution is continued. More complicated cases can also occur with multiple faults and/or traps occurring in the same instruction.

**Note:** Usage of the IEEE Filter should not be compulsory – the user should be able to choose to handle enabled floating-point exceptions directly. The IEEE filter just hides the details of the instruction set and frees the user handler from having to emulate instructions directly and potentially incorrectly.

### 8.1.2.1 Invalid Operation Exception (Fault)

The exception-enabled response of an Itanium floating-point arithmetic instruction to an Invalid Operation exception is to leave the operands unchanged and to set the V bit in the ISR.code field of the ISR register. The operating system kernel, reached via the floating-point fault vector, will then invoke the user floating-point exception handler, if one has been registered.

### 8.1.2.2 Divide by Zero Exception (Fault)

The exception-enabled response of an Itanium floating-point arithmetic instruction to a Divide-by-Zero exception is to leave the operands unchanged and to set the Z bit in the ISR.code field of the ISR register. The operating system kernel, reached via the floating-point fault vector, will then invoke the user floating-point exception handler, if one has been registered.

### 8.1.2.3 Denormal/Unnormal Operand Exception (Fault)

The exception-enabled response of the Itanium arithmetic instruction to a Denormal/Unnormal Operand exception is to leave the operands unchanged and to set the D bit in the `ISR.code` field of the `ISR` register. The operating system kernel, reached via the floating-point fault vector, will then invoke the user floating-point exception handler, if one has been registered.

### 8.1.2.4 Overflow Exception (Trap)

The exception-enabled response of an Itanium floating-point arithmetic instruction to an Overflow exception is to deliver the first (exponent unbounded) IEEE rounded result, and to set the O bit (and possibly the I and FPA bits) in the `ISR.code` field of the `ISR` register and the Overflow flags (and possibly the Inexact flag) in the appropriate status field of the `FPSR` register.

The IEEE-754 standard requires that, when raising an overflow exception, the user handler should be provided with the result rounded to the destination precision with the exponent range unbounded. For the huge result to fit in the destination's range, it must be scaled down by a factor equal to  $2 \cdot 0^a$  (with  $a$  equal to  $3 \cdot 2^{n-2}$ , where  $n$  is the number of bits in the exponent of the floating-point format used to represent the result.) This scaling down will bring the result close to the middle of the range covered by the particular format. The exponent adjustment factors to do the scaling for the various formats are determined as follows:

- 8-bit (single) exponents are adjusted by  $3 \cdot 2^6 = 0xc0 = 192$ .
- 11-bit (double) exponents are adjusted by  $3 \cdot 2^9 = 0x600 = 1536$ .
- 15-bit (double-extended) exponents are adjusted by  $3 \cdot 2^{13} = 0x6000 = 24576$ .
- 17-bit (register) exponents are adjusted by  $3 \cdot 2^{15} = 0x18000 = 98304$ .

The actual scaling of the result is not performed by the Itanium architecture. The IEEE filter that is invoked before calling the user floating-point exception handler typically performs the scaling.

### 8.1.2.5 Underflow Exception (Trap)

The exception-enabled response of an Itanium floating-point arithmetic instruction to an Underflow exception is to deliver the first (exponent unbounded) IEEE rounded result, and to set the U bit (and possibly the I and FPA bits) in the `ISR.code` field of the `ISR` register and the Underflow flag (and possibly the Inexact flag) in the appropriate status field of the `FPSR` register.

The IEEE-754 standard requires that, when raising an underflow exception, the user handler should be provided with the result rounded to the destination precision with the exponent range unbounded. For the tiny result to fit in the destination's range, it must be scaled up by a factor equal to  $2 \cdot 0^a$  (with  $a$  equal to  $3 \cdot 2^{n-2}$ , where  $n$  is the number of bits in the exponent of the floating-point format used to represent the result). The scaling up will bring result close to the middle of the range covered by the particular format. The exponent adjustment factors to do this scaling for the various formats are the same as those for enabled overflow exceptions, listed above.

Just as for overflow, the actual scaling of the result is not performed by the Itanium architecture. It is typically performed by the IEEE Filter, which is invoked before calling the user floating-point exception handler.

#### **8.1.2.6 Inexact Exception (Trap)**

The exception-enabled response of an Itanium arithmetic instruction to an Inexact exception is to set the I bit (and possibly the FPA bit) in the ISR.code field of the ISR register and the Inexact flag in the appropriate status field of the FPSR register. The operating system kernel, reached via the floating-point fault vector, will then invoke the user floating-point exception handler, if one has been registered.

## **8.2 IA-32 Floating-point Exceptions**

IA-32 floating-point exceptions may occur when executing code in IA-32 mode. When this happens, execution is transferred to the Itanium interruption vector for IA-32 Exceptions (at vector address 0x6900.) For classic IA-32 floating-point instructions, they are raised via the "IA\_32\_Exception(FPError) – Pending Floating-point Error." For SSE instructions, they are raised via the "IA\_32\_Exception(StreamingSIMD) – SSE Numeric Error Fault." The operating system may schedule Itanium architecture-based and/or IA-32 exception handlers for these exceptions.

§



The Itanium architecture enables Itanium architecture-based operating systems to host IA-32 applications, Itanium architecture-based applications, as well as mixed IA-32/Itanium architecture-based applications. Unless the operating system explicitly intercepts ISA transfers (using the PSR.di), user-level code can transition between the two instruction sets without operating system intervention. This allows IA-32 programs to call Itanium architecture-based subroutines or vice-versa. Itanium architecture-based and IA-32 code can share data through registers and/or memory. Multi-threaded IA-32 and Itanium architecture-based applications can easily communicate with each other or the Itanium architecture-based operating system using shared memory. The Itanium architecture does not support execution of Itanium architecture-based programs on an IA-32 operating system. While the architecture does not prevent IA-32 code from executing as part of an Itanium architecture-based operating system, it is strongly recommended that Itanium architecture-based operating systems do **not** contain IA-32 code.

One of the most compelling motivations for executing IA-32 code on an Itanium architecture-based operating system is the ability to run existing unmodified IA-32 application binaries. Because IA-32 performs 32-bit instruction/memory references that are zero-extended into 64-bit virtual addresses, Itanium architecture-based operating systems must ensure that all IA-32 code and data is located in the lower 4GBytes of the virtual address space. Compute intensive IA-32 applications can improve their performance substantially by migrating compute kernels from IA-32 to Itanium architecture-based code while preserving the bulk of the application's IA-32 binary code. If mixed IA-32/Itanium architecture-based applications are supported, care has to be taken that the data accessible to IA-32 portions of the application is located in the lower 4GBytes of the virtual address space.

While processors based on the Itanium architecture are capable of supporting a wide range of Itanium architecture-based/IA-32 code mixing, Itanium architecture-based operating systems need to provide a software support infrastructure to enable full interoperability between the IA-32 and Itanium instruction set. Most Itanium architecture-based operating systems are expected to support user-level IA-32 applications, and, as a result, must be able to provide the full range of operating system services through a 32-bit system call interface. However, different operating systems and runtime conventions may reduce the set of interoperability modes as desired by the operating system vendor.

While it is an interesting topic, this chapter does not discuss 32-bit application binary interfaces provided by specific operating systems. Instead, this chapter focusses on what services are required from an Itanium architecture-based operating system by a processor based on the Itanium architecture that is executing IA-32 code. In other words, the focus of this chapter is the low-level processor / operating system interface rather than the IA-32 software / operating system (application binary) interface.

## 9.1 Transitioning between Intel® Itanium® and IA-32 Instruction Sets

As mentioned earlier, user-level code can transition from Itanium to IA-32 (or back) instruction sets without operating system intervention. As described in [Chapter 6, “IA-32 Application Execution Model in an Intel® Itanium® System Environment” in Volume 1](#), two instructions are provided for this purpose: `br.ia` (an Itanium unconditional branch), and `JMPE` (an IA-32 register indirect and absolute jump). Prior to executing any IA-32 instructions, however, the Itanium architecture-based operating system needs to setup an execution environment for executing IA-32 code.

### 9.1.1 IA-32 Code Execution Environments

Processors based on the Itanium architecture are capable of executing IA-32 code in real mode, VM86 mode or protected mode. When segmentation is enabled both 16 and 32-bit code are supported. Prior to transferring control to IA-32 code, an Itanium architecture-based application and/or operating system is expected to setup the complete IA-32 execution environment in Itanium registers.

In particular, Itanium architecture-based software must setup IA-32 segment descriptor and selector registers in Itanium application registers, and must ensure that code and stack segment descriptors (CSD, SSD) are pointing at valid and correctly aligned memory areas. It is also worth noting that the IA-32 GDT and LDT descriptors are maintained in GR30 and GR31, and are unprotected from Itanium architecture-based user-level code. For more details on the IA-32 execution environment please refer to [Section 6.2.2, “IA-32 Application Register State Model” on page 1:113](#).

Some IA-32 execution environments may need support from an Itanium architecture-based operating system. Which IA-32 software environments are supported by an Itanium architecture-based operating system is determined by the operating system vendor. Itanium architecture-based platform firmware (SAL) provides a runtime environment that allows execution of real-mode IA-32 code found in PCI configuration option ROMs.

### 9.1.2 `br.ia`

`br.ia` is an unconditional indirect branch that transitions from Itanium to IA-32 instruction set. Prior to entering IA-32 code with `br.ia`, software is also required to flush the register stack. `br.ia` sets the size of the current register stack frame to zero. The register stack is disabled during IA-32 code execution. Because IA-32 code execution uses Itanium registers, much of the Itanium register state is overwritten and left in an undefined state when IA-32 code is run. As a result, software can not rely on the value of such registers across an instruction set transition. Execution of IA-32 code also invalidates the ALAT. For more details refer to [Table 6-2, “IA-32 Segment Register Fields” on page 1:118](#).



For best performance, the following code sequence is recommended for transitioning from Itanium to IA-32 instruction set:

```
    {.mii
      flushrs          // flush register stack
      mov b7 = rTarget // Setup IA-32 target address
      nop.i            // nop.i or other instruction
      ;;
    {.mib
      nop.m            // nop.m or other instruction
      nop.i            // nop.i or other instruction
      br.ia.sptk b7    // branch to IA-32 target defined by
                       // lower 32-bits of branch register b7
      ;;
```

Key to performance is that the register stack flush (`flushrs`) and the `br.ia` instruction are separated by a single cycle, and that the `br.ia` instruction is the first B-slot in the bundle directly following the `flushrs`. The `nop` instruction slots in the code example may be used for other instructions.

### 9.1.3 JMPE

JMPE is an IA-32 instruction that comes in a register indirect and absolute branch flavors. The code segment descriptor base is held in the CSD application register (`ar.csd`).

- JMPE `reg16/32` computes the target of the Itanium instruction set as  
 $IP = ([reg16/32] + CSD.base) \& 0xffffffff0$
- JMPE `disp16/32` computes the target of the Itanium instruction set as  
 $IP = (disp16/32 + CSD.base) \& 0xffffffff0$

Targets of the IA-32 JMPE instruction are forced to be 16-byte aligned, and are constrained to the lower 4Gbytes of the 64-bit virtual address space. The JMPE instruction leaves the IA-32 return address (address of the IA-32 instruction following the JMPE itself) in IA\_64 register GR1.

### 9.1.4 Procedure Calls between Intel® Itanium® and IA-32 Instruction Sets

If procedure call linkage is required between Itanium architecture-based and IA-32 subroutines, software needs to perform additional work as described in the next two sections.

#### 9.1.4.1 Itanium® Architecture-based Caller to IA-32 Callee

This section outlines what steps an Itanium architecture-based caller of an IA-32 procedure needs to perform. The ordering of the steps is approximate and need not be executed exactly in the order presented.

1. Setup IA-32 execution environment, if not already done (see [Section 9.1.2](#) for details). Ensure that no NaTed registers are used to setup IA-32 environment nor that they are passed as procedure call arguments to IA-32 code.
2. Marshall arguments from the register stack to memory stack according to IA-32 software conventions.
3. Set up exception handle unwind data structures according to OS convention.

4. Make sure JMPE knows where to return to, e.g. deposit return address for the JMPE on memory stack or pass it in an IA-32 visible register.
5. Setup IA-32 branch target in branch register.
6. Flush register stack, but no other RSE updates.
7. `br.ia` is an indirect branch to IA-32 code. There is no need to preserve Itanium only application registers, since IA-32 code execution leaves them unmodified.
8. Run in the IA-32 callee until it executes a JMPE instruction.
9. JMPE instruction is an unconditional jump to Itanium architecture-based code. JMPE should use the return address specified in step 4.
10. Move return values from memory stack to static Itanium register used for procedure return value according to Itanium calling conventions.
11. Ensure that IA-32 code correctly unwound memory stack, and that memory stack pointer is correctly aligned.
12. Update exception handle unwind data structures according to OS convention.
13. `br.ret` returns to Itanium architecture-based caller.

#### 9.1.4.2 IA-32 Caller to Itanium® Architecture-based Callee

This section outlines what steps an IA-32 caller of an Itanium architecture-based procedure needs to perform. The ordering of the steps is approximate and need not be executed exactly in the order presented.

1. Caller deposits arguments on memory stack, and calls Itanium architecture-based transition stub using the JMPE instruction.
2. Execute JMPE instruction as an unconditional branch to Itanium architecture-based code. The JMPE instruction will leave the address of the IA-32 instruction following the JMPE itself in Itanium register GR1. This address may be used as a return address later.
3. Allocate a register stack frame with the `alloc` instruction.
4. Load procedure arguments from memory stack into Itanium stacked registers. Preserve IA-32 return address in memory or register stack.
5. Set up exception handle unwind data structures according to OS convention.
6. `br.call` to target Itanium architecture-based callee.
7. Execute Itanium architecture-based code until it returns using `br.ret`.
8. Move return value from static Itanium register to memory stack.
9. Load IA-32 return address from step 4 into branch register.
10. Instead of flushing the register stack to memory, the contents of the register stack can be discarded at this point since IA-32 code execution will overwrite it anyway. Invalidate register stack by:
  - a. Allocating a zero-size stack frame using the `alloc` instruction.
  - b. Writing zero into RSC application register, and executing a `loadrs` instruction.
  - c. Restore RSC application register to its original value in preparation for the next call from IA-32 to Itanium instruction set.

11. Ensure memory stack pointer is correctly aligned prior to returning to IA-32 code.
12. `br.ia` returns to IA-32 caller.

## 9.2 IA-32 Architecture Handlers

An Itanium architecture-based operating system needs to be prepared to handle exceptions from Itanium architecture-based and IA-32 code. Depending on the exception cause, exception vectors can be:

- Shared Itanium/IA-32 Exception Vectors: all virtual memory related instruction and data reference faults share a common exception vector, regardless of whether they were caused by Itanium architecture-based or IA-32 code.
- Unique Itanium Exception vectors: these are conditions that only Itanium architecture-based code can cause. Examples are: Instruction Breakpoint fault, Illegal Operation fault, Illegal Dependency fault, Unimplemented Data Address fault, etc.
- Unique IA-32 Exception Vectors: these conditions can occur only from IA-32 instructions.

A detailed break-down of which exceptions occur on which interruption vector and from which instruction set is given in [Table 5-6](#). [Table 9-1](#) shown below summarizes all IA-32 related exceptions that an Itanium architecture-based operating system needs to be ready to handle. These IA-32 specific interrupts are grouped into three vectors: the IA-32 Exception vector, the IA-32 Intercept, and the IA-32 Interrupt vector. Within each of these vectors the interrupt status register (ISR) provides detailed codes as to the origin of this exception. Details on the IA-32 vectors is provided in [Chapter 9, "IA-32 Interruption Vector Descriptions."](#) More details on debug related IA-32 exceptions is given in the following section of this document.

**Table 9-1. IA-32 Vectors that need Itanium® Architecture-based OS Support**

Vector (IVA offset)	Exception Name	Exception Related To	Expected OS Behavior
IA-32 Exception vector (0x6900)	IA-32 Instruction Debug fault	Debug	Relay to debugger.
	IA-32 Code Fetch fault	Segmentation	Signal application.
	IA-32 Instruction Length > 15 bytes fault	Bad Opcode	Signal application.
	IA-32 Device Not Available fault	Numeric	Signal application.
	IA-32 FP Error fault	Numeric	Signal application.
	IA-32 Segment Not Present fault	Segmentation	Signal application.
	IA-32 Stack Exception fault	Segmentation	Signal application.
	IA-32 General Protection fault	Segmentation	Signal application.
	IA-32 Divide by Zero fault	Numeric	Signal application.
	IA-32 Alignment Check fault	Misaligned IA-32 Memory Reference with alignment checking enabled.	Depends on convention.
	IA-32 Bound fault	Segmentation	Signal application.
	IA-32 SSE Numeric Error Fault	Numeric	Signal application.
	IA-32 INTO Overflow trap	Numeric	Signal application.
	IA-32 Breakpoint (INT 3) trap	Software Breakpoint	Depends on convention.
	IA-32 Data Breakpoint trap	Debug	Relay to debugger.

**Table 9-1. IA-32 Vectors that need Itanium® Architecture-based OS Support (Continued)**

Vector (IVA offset)	Exception Name	Exception Related To	Expected OS Behavior
	IA-32 Taken Branch trap	Debug	Relay to debugger.
	IA-32 Single Step trap	Debug	Relay to debugger.
	IA-32 Invalid Opcode fault	Bad Opcode	Signal application.
IA-32 Intercept vector (0x6a00)	IA-32 Instruction Intercept fault	Attempted to access IA-32 paging, MTRRs, IDT, IA-32 control registers, IA-32 debug registers or attempted to execute IA-32 privileged instructions.	This is not supported on an Itanium architecture-based OS. Signal application.
	IA-32 Locked Data Reference fault	Attempt to reference misaligned or uncacheable semaphore.	Emulation handler if needed. Refer to <a href="#">Section 2.1.3.2, "Behavior of Uncacheable and Misaligned Semaphores"</a> on page 2:509.
	IA-32 System Flag Intercept trap	System Flag intercept	Depends on convention.
	IA-32 Gate Intercept trap	Gate/Task transfer intercept	Depends on convention.
IA-32 Interrupt vector (0x6b00)	IA-32 Software Interrupt (INT) trap	Software Interrupt	Depends on convention.
Cannot happen in Itanium architecture-based operating system	IA-32 Double Fault IA-32 Invalid TSS Fault, IA-32 Page Fault, IA-32 Machine Check	N/A	Don't worry,

## 9.3 Debugging IA-32 and Itanium® Architecture-based Code

Itanium architecture-based operating systems that want to provide debug support for both IA-32 and Itanium architecture-based applications, need to be aware of the differences between taking instruction and data breakpoint exceptions as well as single step or taken branch traps on Itanium and IA-32 instructions.

### 9.3.1 Instruction Breakpoints

If an Itanium instruction matches an instruction breakpoint register (IBR) then an Instruction Debug Fault is delivered on the Itanium Debug vector. To step across a single Itanium instruction, `IPSR.id` must be set to one. An IA-32 instruction, however, that matches an IBR causes an IA-32 Instruction Breakpoint fault which is delivered to the IA-32 Exception vector (Debug). To step across a single IA-32 instruction, either `IPSR.id` or `EFLAGS.rf` must be set to one.

### 9.3.2 Data Breakpoints

If an Itanium memory reference matches a data breakpoint register (DBR) then a Data Debug Fault is delivered on the Itanium Debug vector. To step across a single data breakpoint, `IPSR.dd` must be set to one. An IA-32 instruction, however, that matches a DBR causes an IA-32 Data Breakpoint *trap* which is delivered to the IA-32 Exception vector (Debug). In other words, the debugger only gets control after the instruction

making the reference has completed. Since IA-32 instruction can make multiple memory references, a single IA-32 instruction may cause multiple data break points to trigger. Details on how this is communicated to software in the interrupt status register (ISR) is given in [Section 9.1, “IA-32 Trap Code” on page 2:213](#). Since IA-32 data breakpoints are traps, there is no need to step over them.

### 9.3.3 Single Step Traps

When PSR.ss enables single stepping of Itanium architecture-based applications, each instruction that is stepped will stop at the Single Step trap handler. When PSR.ss or EFLAG.tf enable single stepping of IA-32 applications, an IA\_32\_Exception(Debug) trap is taken after each IA-32 instruction. For more details refer to [Section 9.1, “IA-32 Trap Code” on page 2:213](#).

### 9.3.4 Taken Branch Traps

When PSR.tb enables taken branch trapping on Itanium architecture-based applications, each taken branch will transfer control to the Taken Branch Trap handler. When PSR.tb is set, taken IA-32 branches transfer control to the IA\_32\_Exception(Debug) trap handler taken after each IA-32 instruction. For more details refer to [Section 9.1, “IA-32 Trap Code” on page 2:213](#).

§



The Itanium architecture provides a high performance external interrupt architecture. While IA-32 processors commonly use a three wire shared APIC bus, processors based on the Itanium architecture utilize a high performance, message-based, point-to-point protocol between processors and multiple I/O interrupt controllers. To ensure that processors based on the Itanium architecture can fully leverage the large set of existing platform infrastructure and I/O devices, compatibility with existing platform infrastructure is provided in the form of direct support for Intel 8259A compatible interrupt controllers and limited support for level sensitive interrupts.

This chapter introduces the basic external interrupt mechanism provided by the architecture, while [Section 5.8, “Interrupts”](#) provides the complete architectural definition for the Itanium external interrupt architecture.

## 10.1 External Interrupt Basics

Interrupts are identified by their vector number. The vector number implies interrupt priority, and also determines whether the interrupt is delivered to processor firmware as a “PAL-based” interrupt, or whether it is delivered to the operating system as an “IVA-based” external interrupt.

This chapter discusses asynchronous external interrupts only. PAL-based platform management interrupts (PMI) are not discussed here. External interrupts are IVA-based and are delivered to the operating system by transferring control to code located at address CR[IVA]+0x3000. This code location is also known as the external interrupt vector and is described on [page 2:186](#).

Software can distinguish interrupts based on their vector number. Vector numbers range from 0 to 255. Vector numbers also establish interrupt priorities as follows:

- Vector numbers below 16 are special, and are architecturally defined in [Section 5.8.1, “Interrupt Vectors and Priorities” on page 2:118](#). The non-maskable interrupt (NMI) is always vector 2 and is higher priority than all in-service external interrupts. ExtINT, Intel 8259A compatible external interrupt controller interrupt, is always vector 0. Vector numbers below 16 have higher priority than vectors above 16. Vector 15 is used to indicate that the highest priority pending interrupt in the processor is at a priority level that is currently masked or there are no pending external interrupts.
- For vector numbers between 16 and 255, higher vector numbers imply higher priority. In this range, vectors are freely assignable by software. This is achieved by programming of interrupt controllers and the processor internal interrupt configuration registers.

## 10.2 Configuration of External Interrupt Vectors

As defined in [Section 5.8, “Interrupts” on page 2:114](#), external interrupts originate from one of four sources:

- From external sources, e.g. external interrupt controllers or intelligent external I/O devices, or
- From the processor’s LINT0 or LINT1 pins<sup>1</sup> (typically connected to an Intel 8259A compatible interrupt controller), or
- From internal processor sources, e.g. timers or performance monitors, or
- From other processors, e.g. inter-processor interrupts (IPIs).

All interrupts are point-to-point communications. There is no facility for broadcasting of interrupts. The interrupt message protocol used by the processor-to-processor and the external source-to-processor is not defined architecturally, and is not visible to software.

A number of external interrupt control registers (LID,TPR, ITV, PMV, CMCV, LRR0 and LRR1) allow software to directly configure the processor interrupt resources. The Local ID register (LID) establishes a processor’s unique physical interrupt identifier. The Task Priority Register (TPR) allows masking of external interrupts based on vector priority classes. The ITV, PMV, CMCV, LRR0 and LRR1 external interrupt control registers configure the vector number for the processor’s local interrupt sources. Configuration of the external controllers and devices is controller-/device-specific, and is beyond the scope of this document.

## 10.3 External Interrupt Masking

The Itanium architecture provides four mechanisms to prevent external interrupts from being delivered to a processor: a bit in the processor status register (PSR.i), the interrupt vector register (IVR) and the end-of-interrupt (EOI) register, the task priority register (TPR), and the external task priority register (XTPR). The next four sections discuss these mechanisms.

### 10.3.1 PSR.i

When PSR.i is zero, the processor does not accept any external interrupts. However, interrupts continue to be pended by the processor. Software can use PSR.i to temporarily disable taking of external interrupts, e.g. to ensure uninterruptable execution of critical code sections. Since clearing of PSR.i takes effect immediately (refer to the `rsm` instruction page), software is not necessarily required to explicitly serialize clearing of PSR.i (unless another processor resource requires serialization). On

---

1. Processors optionally support two external interrupt pins. Software can query for the presence of LINT pins via the `PAL_PROC_GET_FEATURES` procedure call.



the way out of an uninterruptable code section software is not required to serialize the setting of PSR.i either, unless it is of interest to software to be able to take interrupts in the very next instruction group. A code example for this case is given below:

```
rsm i ;;
// rsm of PSR.i takes effect on the next instruction

// uninterruptable code sequence here

ssm i ;;
// ssm of PSR.i does require data serialization, if we need to ensure
// that external interrupts are enabled at the very next instruction. If
// data serialization is omitted, PSR.i is set to 1 at the latest when
// the next exception is taken.
```

By avoiding the serialization operations on PSR.i the performance of such uninterruptable code sections is improved.

### 10.3.2 IVR Reads and EOI Writes

As described in [Section 10.4](#), IVR reads return the highest priority, pending, unmasked vector, and places this vector “in-service.” Additionally, IVR reads have the side-effect of masking all vectors that have equal or lower priority than one that is returned by the IVR read. Correspondingly, writes to the EOI register unmask all vectors with equal or lower priority than the highest priority “in-service” vector. Due to nesting of higher priority interrupts, it is possible to have multiple vectors in the “in-service” state.

### 10.3.3 Task Priority Register (TPR)

The Task Priority Register (TPR) provides an additional interrupt masking capability. It allows software to mask interrupt “priority classes” of 16 vectors each by specifying the mask priority class in the TPR.mic field. The TPR.mmi field allows masking of all maskable external interrupts (essentially all but NMI).

An example of TPR use is shown in [Section 10.5.2, “TPR and XPTR Usage Example”](#) on page 2:608.

### 10.3.4 External Task Priority Register (XTPR)

The External Task Priority Register (XTPR) is a per-processor resource that can be provided by external bus logic in some Itanium architecture-based platforms. If supported by the platform, XTPR can be used by the operating system to redirect external interrupts to other processors in a multiprocessor system.

The XTPR is updated by performing a 1-byte store to the XTP byte which is located at an offset of 0x1e0008 in the Processor Interrupt Block (see [Section 5.8.4, “Processor Interrupt Block”](#) for details). Since the timing of the modification of the XTP register is not time critical there is no serialization required. Effects of the one byte store operation are platform specific. Typically, it will generate a transaction on the system bus identifying it as an XTP register update transaction, and will indicate which processor generated the transaction as well as the stored data.

An example of XTPR use is included in [Section 10.5.2, “TPR and XPTR Usage Example”](#) on page 2:608.

## 10.4 External Interrupt Delivery

The architectural interrupt model in [Section 5.8](#) defines how each interrupt vector cycles through one of four states:

- *Inactive*: there is no interrupt *pending* on this vector.
- *Pending*: an interrupt has been received by the processor on this vector, but has not been *accepted* by the processor and has not been *acquired* by software. The processor hardware will *accept* the interrupt when this vector's priority level is higher than the highest currently in-service vector, PSR.i is one, and TPR settings do not mask the interrupt. This will cause the processor to transfer control flow to the external interrupt handler. Software can then *acquire* the highest priority, pending, unmasked vector by reading the IVR control register. The IVR read returns the 8-bit vector number in a register and masks all vectors that have equal or lower priority. This vector now enters the In-Service/None Pending state.
- *In-Service/None Pending*: an interrupt has been received by the processor on this vector, and has been acquired by software (by reading the IVR control register), but software has not *completed servicing* this interrupt. In this state, the processor masks all vectors that have equal or lower priority. In this state, the processor can receive and remember a second interrupt on this vector. If this happens, the processor transitions this vector to the "In-Service/One Pending" state. If software *completes the interrupt* service routine (indicated to the processor by writing the EOI register) before another interrupt is received on this vector, then the processor returns this vector to the Inactive state, and all vectors with equal or lower priority are unmasked.
- *In-Service/One Pending*: an interrupt has been received by the processor on this vector, and has been acquired by software (by reading the IVR control register), and software has not completed servicing this interrupt. Additionally, the processor received a second interrupt on this vector, which is now held pending. If additional interrupts on this vector are received by the processor while this vector is in the "In-Service/One Pending" state, those additional interrupts are not distinguishable by the processor hardware. When software completes the interrupt service routine for the original interrupt on this vector (indicated to the processor by writing the EOI register), then the processor returns this interrupt vector to the Pending state for the second interrupt that was received on this vector. Additionally, all vectors with equal or lower priority are unmasked.

It is recommended the following structure for an Itanium architecture-based external interrupt handler:

1. Read and Save TPR, i.e. save Old Task Priority variable (optional).
2. External Interrupt Harvest Loop:
  - a. Read the IVR control register to determine which vector is being delivered. If the returned IVR value is 15, then this is a spurious interrupt and it can be ignored; software can now clear PSR.ic, restore IPSR and IIP and then `rfi` to the interrupted context. If the returned IVR value is not 15, continue with step 2b.
  - b. Raise TPR register to the interrupt class to which the level read out of IVR belongs (optional).

- c. Software must preserve IIP and IPSR prior to re-enabling PSR.ic and PSR.i which will re-enable taking of exceptions and higher priority external interrupts.
- d. Issue a `srlz.d` instruction. This ensures that updated PSR.ic and PSR.i settings are visible, and it also makes sure that the IVR read side effect of masking lower or equal priority interrupts is visible when PSR.i becomes 1.
- e. Dispatch the appropriate interrupt service routine.
- f. Disable external interrupts by clearing PSR.i with an `rsm 0x4000` instruction. This ensures that external interrupts are disabled prior to the EOI write in the next step.
- g. Notify the processor that interrupt handling for this vector is completed by writing to the EOI register. This will unmask any pending lower priority interrupts. If this was a level triggered interrupt, write to the I/O SAPIC EOI register.
- h. Lower TPR register to Old Task Priority (optional).
- i. Issue a `srlz.d` instruction. This ensures that ensure the EOI write from step 2g is reflected in the future IVR read (in step 2a). It also ensures that the TPR update from step 2h unmask any interrupts in the priority classes (including the current task priority level) that were masked by the previous value of TPR.
- j. Return to top of loop (step 2a).

These steps assume that the routine's caller already performed the required state preservation of interruption resources. Therefore the focus of the steps above is to check the IVR to acquire the vector so the operating system can determine what device the interrupt is associated with. The code is setup to loop, servicing interrupts until the spurious interrupt vector (15) is returned. Looping and harvesting outstanding interrupts reduces the time wasted by returning to the previous state just to get interrupted again. The benefit of interrupt harvesting is that the processor pipeline is not unnecessarily flushed and that the interrupted context is only saved/restored once for a sequence of external interrupts. Once the vector is obtained the specific interrupt service routine is called to service the device request. Upon return from the interrupt service routine, an EOI is written and the IVR is checked once again.

If the operating system does not implement priority levels then there is no need to save and restore the task priority level (steps 1, 2b, and 2h are optional). As described in [Section 10.3](#) above, an IVR read automatically masks interrupts at the current in-service level and below until the corresponding EOI is issued. For level triggered interrupts, the programmer must not only inform the processor, but the external interrupt controller that the level triggered interrupt has been serviced.

## 10.5 Interrupt Control Register Usage Examples

The examples in this section provide an overview of using the Itanium external interrupt control registers. Actual and pseudo code fragments are listed to aid in the development of OS code which will utilize these registers. It is up to the operating system and its writer to determine what minimum set of control registers are required to be used.

## 10.5.1 Notation

Preprocessor macros for function ENTRY and END are used in the examples to reduce duplication of code and reduce document space requirements.

```
#define ENTRY(label) \  
    .text; \  
    .align 32;; \  
    .global label; \  
    .proc label; \  
label::  
  
#define END(label) .endp
```

## 10.5.2 TPR and XPTR Usage Example

This code will allow certain interrupts to be masked by increasing/decreasing the task priority register. If you don't want to mask all external interrupts, you can raise the priority level to mask out only the interrupts that have higher priority (and no effect on your current critical section).

We also take the expensive route here by updating not only the processor TPR, but the External Task Priority Register used by the chipset (if supported) as a hint to what processor should receive the next external interrupt.

```
//  
// routine to set the task priority register to mask  
// interrupts at the specific level or below  
//  
// INPUT: SPL level  
//  
TPR_MIC=4  
TPR_MIC_LEN=4  
  
.global external_task_pri_reg// address points to Interrupt Delivery block  
  
ENTRY(set_spl)  
    alloc r18=ar.pfs,1,0,0,0  
    dep.z r22=r32,TPR_MIC,TPR_MIC_LEN  
    movl r19=external_task_pri_reg  
    ;;  
    mov cr.tpr=r22  
    ld8 r20=[r19] // get address of Ext. TASK Priority Register  
    ;;  
    srlz.d // srlz.d only required if want TPR update effective  
immediately  
    stl [r20]=r32 // if supported by platform: update eXternal Task Priority  
(XTP)  
    br.ret.sptk b0  
    ;;  
END(set_spl)
```

### 10.5.3 EOI Usage Example

This example is a typical return from an interrupt service routine to the generic interrupt handler. Interrupts are disabled before returning to the main trap handler in preparation for returning from kernel space.

```
return_from_interrupt:
// disable interrupts here

    rsm 0x4000          // make sure interrupts disabled

// interrupt_eoi# clear the sapic/pic interrupt
sapic_eoi:
    mov cr.eoi=r0      // issue and eoi
    ;;
    srlz.d             // make sure it takes effect

// issue the appropriate EOI sequence to the external interrupt
// controller here.
```

For level trigger interrupts, the OS is required to issue an EOI not only to the processor, but also the external interrupt controller where the interrupt originated. This forces the OS to keep track of whether the vector is associated with a level or an edge trigger interrupt line.

### 10.5.4 IRR Usage Example

Waiting on an interrupt with interrupts disabled.

```
my_interrupt_loop::
//
// check for vector 192 (0xc0) via irr3
//

    mov    r3=cr.irr3
    ;;
    and    r3=0x1,r3
    ;;
    cmp.eq p6,p7=0x1,r3
    (p7)br.cond.sptk.few my_interrupt_loop
    ;;
    mov    r4=cr.ivr      // read the vector
    ;;
    mov    cr.eoi=r0     // clear it
    ;;
```

### 10.5.5 Interval Timer Usage Example

The Itanium architecture provides a 64 bit interval timer for elapsed time notification interrupts. It is similar to the IA-32 Time Stamp Counter (TSC). Programming the Itanium interval timer consists of initializing the ITV (CR 72), ITM (CR 1), and ITC (AR 44).

The Interval Timer Vector (ITV) specifies the external interrupt vector number for the Interval Timer Interrupts. The code examples below show how to clear and initialize the timers vector, match register, and count registers.

The Interval Time Counter (ITC) gets updated at a fixed relation to the processor clock. The ITM, Interval Timer Match, is used to determine when a interval timer interrupt is generated. When the ITC matches the ITM and the timer is unmasked via ITV then an interrupt will be generated.

```
//
// routine to reset the interval timer to zero..
//

ENTRY(em_timer_reinit)
    mov    ar.itc=r0           // reset itimer counter
    br.ret.sptn.few rp
END(em_timer_reinit)

//
// routine to setup the interval timer.
//
// 1) setup the interval timer vector
// 2) initialize the time counter to zero
// 3) initialize the match register
//
// INPUTS: timermatch -- value to initialize ITM register with.
//         vector number -- vector to interrupt with
// OUTPUTS: none
//
ENTRY(enable_minterval)
    alloc  r14=ar.pfs,0x2,0,0,0 // get ready for input parameters
    mov    ar.itc=r0           // initialize counter to zero
    ;;
    mov    cr.itm=r32          // set match register
    ;;
    srlz.d
    mov    cr.itv=r33          // set interval timer vector
    ;;
    srlz.d                     // make sure it goes through
    br.ret.sptk.few rp        // return
    .endp
```

Since the ITC gets updated at a fixed relation to the processor clock, in order to find out the frequency at run time, one can use a firmware call to obtain the input frequency information to the interval time. Using this frequency information the ITM can be set to deliver an interrupt at a specific time interval (i.e. for operating system scheduling purposes). Assuming the frequency information returned by the firmware is in ticks per second, the programmer could use a time-out delta for delivering a timer interrupt every 10 milliseconds as follows:

```
timeout_delta=ticks_per_second/100;
```

where `ticks_per_second` is the frequency value returned by the firmware and `timeout_delta` will be the value added to the ITC for setting the next ITM. Therefore, the ITC is left free running, but the ITM must be updated upon every timer interrupt with its next time out match value, i.e.  $ITM = ITC + \text{timeout\_delta}$ .

The only issue with this setup is if the timer interrupt delivery is delayed beyond the point of the original intended delivery time (i.e.  $ITC > ITM$ ). This could happen if interrupts were disabled or blocked by the operating system/device driver longer than

the time-out value. In this case the ITM has to be adjusted in order for the next ITM to be accurate. The following algorithm could be used to adjust the next ITM before returning from the timer interrupt handler.

```

for (;;) {
    itm_next = itm_next + timeout_delta + (read current ITC - read current ITM);
    if (itm_next < current ITC) {
        /* we missed the next interrupt already, continue */
    } else {
        set_itm(itm_next);
        break;
    }
}

```

where `itm_next` was initialized to `current ITC + timeout_delta`, and `set_itm` in Itanium architecture-based assembly would look like:

```

.global set_itm
.proc set_itm
set_itm:
    alloc r18=ar.pfs,1,0,0,0
    mov cr.itm=r32
    ;;
    srlz.d
    br.ret.sptk b0
    ;;
.endp set_itm

```

## 10.5.6 Resource Utilization Counter Usage Example

The Itanium architecture provides a 64-bit counter to provide information on how many execution cycles a given logical processor is getting. It is similar to the Interval Timer (ITC, AR 44), except that it is clocked only when the logical processor is active. Optimizations such as hardware multi-threading and processor virtualization may cause a logical processor to sometimes be inactive. The Resource Utilization Counter allows for better cycle accounting for logical processors, given these types of optimizations.

RUC should only be written by Virtual Machine Monitors; other Operating Systems should not write to RUC, but should only read it.

## 10.5.7 Local Redirection Example

The Local Redirection Registers (LRR0-1) serves to steer external signal-based interrupts that are directly connected to the processor. LRR0 and LRR1 control the external interrupt signals (pins) referred to as Local Interrupt 0 (LINT0) and Local Interrupt 1 (LINT1) respectively. The example below shows how to mask interrupt delivery on LINT0.

```

movl r18=(1<<16)
;;
mov cr.lrr0=r18
;;
srlz.d // srlz.d is required after LRR write to ensure write effect

```

**Note:** LINT0 and LINT1 pins are not required to be supported. Writes to LRR0-1 control registers would have not effect, and reads from LRR0-1 control registers would return 0.

## 10.5.8 Inter-processor Interrupts Layout and Example

A processor generates an inter-processor interrupt (IPI) by storing a 64-bit interrupt command to an 8-byte aligned address in the Interrupt delivery region of the Processor Interrupt block. The address being stored to determines what target processor receives the IPI. The example below is an example of sending an interrupt to a specific processor based on the destination ID passed in. The destination ID consists of the Local interrupt ID and the Extended interrupt ID.

Writing to improperly aligned addresses in the delivery region or failure to store less than 64 bits can result in an invalid operation fault. The access must be uncacheable in order to generate an IPI.

```
//
// send_ipi_physical (dest_id, vector)
//
// inputs:      processor destination ID vector to send
//              (Local ID (8 bits << 8) | EID ( 8 bits))
//
//
//
.global ipi_block          // pointer to processor I/O block

IPI_DEST_EID=0x4

ENTRY(send_ipi_physical)
    alloc r19=ar.pfs,2,0,0,0
    movl r17=ipi_block;;
    ld8 r17=[r17]          // get pointer to processor block
    shl r21=r32,IPI_DEST_EID;;
    add r20=r21,r17;;      // point to proper processor
    st8.rel [r20]=r33      // send the IPI
    br.ret.sptk b0;;

END(send_ipi_physical)
```

## 10.5.9 INTA Example

External interrupt controllers, that are compatible with the Intel 8259A interrupt controller can not issue interrupt messages, so the vector number is not available at the time of the interrupt request. When an interrupt is accepted the software must check to see if it came from an external controller by the vector number (via IVR) to see if it is the ExtINT vector.



Once the software determines it is an ExtINT, it must obtain the actual vector by doing an uncached 1-byte load from the INTA byte located in the upper half of the processor interrupt block, offset 0x1e0000 from the base.

```
EXTINT=r0
INTA_PHYS_ADDRESS=0x80000000fefe0000
inta_address=r31

    movl inta_address=INTA_PHYS_ADDRESS
    ;;
    srlz.d          // make sure everything is up to date
    mov r14 = cr.ivr // read ivr
    ;;
    srlz.d          // serialize before the EOI is written...
    ;;
    cmp.ne p1,p2 = EXTINT,r14 ;;
    (p1)br.cond.sptk process_interrupt
    ;;

//
// A single byte load from the INTA address should cause
// the processor to emit the INTA cycle on the processor
// system bus. Any Intel 8259A compatible external interrupt
// controller must respond with the actual interrupt
// vector number as the data to be loaded.
//
//
    ldl r17 = [inta_address] // get the real vector..
    ;;
// vector obtained

process_interrupt:
```

§



I/O devices can be accessed from Itanium architecture-based programs using regular loads and stores to uncacheable space. While cacheable Itanium memory references may be reordered by the processor, uncacheable I/O references are always presented to the platform in program order. This “sequentiality” of uncacheable references is discussed in [Section 2.2.2, “Memory Attributes” on page 2:524](#) and in more detail in [Section 4.4.7, “Sequentiality Attribute and Ordering” on page 2:82](#).

Additionally, uncacheable memory pages are defined to be “non-speculative” which causes all data and control speculative loads to uncacheable pages to defer. Control speculative loads to uncacheable memory return a NaT/NaTVal to their target register. Data speculative loads to uncacheable memory return zero to their target register. For details, refer to [Section 4.4.6, “Speculation Attributes” on page 2:79](#).

When configuring chipset registers or setting up device registers, it is sometimes required to know when a memory transaction has been completed. Completion means the processor received acknowledgment that the transaction finished successfully in the platform, and that all its side-effects have occurred and will be visible to the next memory operation (issued by the same processor). To ensure completion of prior accesses on the platform, the Itanium architecture provides the `mf.a` instruction. Unlike the `mf` instruction that waits for **visibility** of prior operations, the `mf.a` waits for **completion** of prior operations on the platform. More details in [Section 11.1](#).

To fully leverage the large set of existing platform infrastructure and I/O devices, the architecture also supports the IA-32 platform I/O port space. The Itanium instruction set does not provide IN and OUT instructions, but they can be emulated. The I/O port space can be mapped into user-space, and IA-32 applications can use IN and OUT instructions to directly communicate with the I/O port space. More details in [Section 11.2](#).

The Itanium architecture provides a high-performance, high-bandwidth uncacheable memory attribute that supports write-coalescing. This allows the processor to burst writes to uncacheable locations at much higher bandwidth. The Itanium architecture does **not** guarantee the FIFO delivery of write-coalescing stores. More details in [Section 4.4.5, “Coalescing Attribute” on page 2:78](#).

## 11.1 Memory Acceptance Fence (mf.a)

An `mf` instruction ensures that all cache coherent agents have observed all prior memory operations made by the processor issuing the `mf`. However, it does **not** ensure that those operations have completed, in the Itanium architecture parlance it does not ensure that they have been “accepted” by the external platform. For instance, a load may have been made visible to all processors by snooping their caches, but the data return may still be in progress. Such a load would be visible, but not complete.

The `mf.a` instruction on the other hand ensures that all prior data memory references made by the processor issuing the `mf.a` have been “accepted” by the external platform. However by itself the `mf.a` does not guarantee that all cache coherent agents have observed all prior memory operations. For instance, an uncacheable store to a chipset register may have completed on the system bus, however, that does not entail that all prior cacheable transactions (from the processor issuing the store) have been observed by all other processors in the coherence domain.

If software needs to ensure that all prior memory operations have been accepted by the platform **and** have been observed by all cache coherent agents, both an `mf.a` and an `mf` instruction must be issued. The `mf.a` must be issued first, and the `mf` must be issued second. For more details on memory ordering between cache coherent agents please refer to [Chapter 2, “MP Coherence and Synchronization.”](#)

Typically `mf.a` is used to configure a system’s I/O space, e.g. to setup chipset registers that affect all subsequent memory operations. Specifically, the `mf.a` instruction restrains further data accesses from initiating on the external platform interface until:

1. All previous sequential (i.e. non write-coalescing uncacheable) loads have been returned data, and
2. All previous stores have been “accepted” by the platform. Typically acceptance is indicated by a bus-specific signals/phase, e.g. completion of response phase on the system bus.

Architecturally, the definition of “acceptance” is platform dependent. The next section discusses the usage of the `mf.a` instruction in the context of the I/O port space.

## 11.2 I/O Port Space

IA-32 processors support two I/O models: memory mapped I/O and the 64KB I/O port space. To support IA-32 platforms, the Itanium architecture allows operating systems to map the 64KB I/O port space into the 64-bit virtual address space. This allows Itanium architecture-based operating systems to see all I/O devices as a single unified memory mapped I/O model, and permits “normal” Itanium load and store instructions as well as IA-32 IN and OUT instructions to directly access the I/O port space.

As described in [Section 10.7, “I/O Port Space Model” on page 2:267](#), Itanium architecture-based operating systems can map the physical 64KB I/O port space into a spread-out 64MB block of virtual address space. The virtual base address of the I/O port space (IOBase) is maintained by the operating system in kernel register KR0. When the processor issues Itanium load and stores accesses to the I/O port space, a port’s virtual address is computed as:

```
port_virtual_address = IOBase | (port{15:2}<<12) | port{11:0}
```

For Itanium loads and stores, this address computation places four 1-byte ports on each 4KB page and expands the space to 64MB, with the ports being at a relative offset specified by `port{11:0}` within each 4KB virtual page. When executing an IA-32 IN or OUT instruction a processor based on the Itanium architecture automatically converts the IA-32 address to the appropriate expanded I/O port space address.

As a result of the spreading-out of the I/O ports into individual 4KB pages, Itanium architecture-based operating system code can control IA-32 IN, OUT instruction and IA-32 or Itanium load/store accessibility to blocks of 4 virtual I/O ports using the TLBs. This allows Itanium architecture-based operating systems to securely map devices that inhabit the I/O port space to different Itanium architecture-based device drivers or to user-space Itanium architecture-based applications.

Itanium architecture-based operating systems must ensure that the I/O port space is always mapped as uncacheable memory, and that Itanium architecture-based software only issues aligned 1, 2 or 4 byte references to I/O port space, otherwise device behavior is undefined.

When porting an IA-32 device driver to the Itanium architecture it can be useful to emulate the behavior of IA-32 IN and OUT instructions. The following code examples should be used for this purpose, since they enforce the strict memory ordering and platform acceptance requirements that IA-32 IN and OUT instructions are subject to. The following Itanium architecture-based assembly code outb (out byte) and inb (in byte) examples assume that the io\_port\_base is the virtual address mapping pointer set up by the IA\_64 operating system. An mf.a instruction is used to verify acceptance by the platform before returning to the calling routine. Interrupts would be expected to be disabled if these routines are called from user mode. This is for possible issues with process migration after servicing an interrupt.

```
//
// void outb(unsigned char *io_port,unsigned char byte)
//
//Output a byte to an I/O port.
//
ENTRY(outb)
    base_addr = r16
    port_addr = r17
    port_offset = r18
    mask = r19

    alloc    r13 = ar.pfs, 2, 0, 0, 0        // 2 in, 0 local, 0 out, 0 rot
    movl     base_addr = io_port_base
    extr.u   port_offset = in0, 2, 14
    mov      mask = 0xffff
    ;;
    ld8      port_addr = [base_addr]
    shl      port_offset = port_offset, 12
    and      in0 = mask, in0
    ;;
    add      port_offset = port_offset, in0
    ;;
    mf
    add      port_addr = port_addr, port_offset
    ;;
    st1.rel [port_addr] = in1
    mf.a
    mf
    br.ret.spnt.few rp
END(outb)

//
// unsigned char inb(unsigned char *io_port)
//
// Input a byte from an I/O port.
//
ENTRY(inb)
    base_addr = r16
    port_addr = r17
    port_offset = r18
```

```

mask = r19

alloc   r13 = ar.pfs, 2, 0, 0, 0           // 2 in, 0 local, 0 out, 0 rot
movl   base_addr = io_port_base
extr.u port_offset = in0, 2, 14
mov    mask = 0xfff
;;
ld8    port_addr = [base_addr]
shl    port_offset = port_offset, 12
and    in0 = mask, in0
;;
add    port_offset = port_offset, in0
;;
mf     port_addr = port_addr, port_offset
;;
ldl.acq r8 = [port_addr]
mf.a
mf
br.ret.spnt.few rp
END(inb)

```

§

Processors based on the Itanium architecture include a minimum of four performance counters which can be programmed to count processor events. These event counts can be used to analyze both hardware and software performance. Performance counters can be configured to generate a counter overflow interrupt. This interrupt can be used for event- or time-based profiling. For hot-spot analysis of running code, performance monitor interrupts can be used to create a profile of frequently occurring instruction pointers (IP). Another common use of event counts is to compute processor performance metrics such as cycles per instructions (CPI), the current branch, cache or TLB miss rates, etc.

The Itanium architecture provides architected support for context switching of performance monitors by an Itanium architecture-based operating system. If supported by the operating system, this allows performance counter events to be broken down per thread or per process which is important for effective performance tuning of Itanium architecture-based applications.

The remainder of this chapter reviews the architected performance monitoring mechanisms. It also discusses the Itanium architecture-based operating system support needed for two monitoring usage models: per process/thread and system-wide event monitoring.

## 12.1 Architected Performance Monitoring Mechanisms

As defined in [Section 7.2, “Performance Monitoring” on page 2:155](#), processors based on the Itanium architecture provide a minimum of four generic performance counter pairs (PMC/PMD[4..7]). The performance monitor control (PMC) registers are used to select the event to be counted, and to define under what conditions the event should qualify for being counted (for details refer to [Section 7.2.1, “Generic Performance Counter Registers” on page 2:156](#)). The performance monitor data (PMD) registers contain the event count or data.

The PMC/PMD registers can only be written by privileged software (PSR.cpl must be zero). A counter can be configured as a “privileged” counter or a “user-level” counter by setting of the PMC[i].pm bit. Privileged counters can only read at privilege level 0, while user-level counters can be read by user mode code (unless the operating system has explicitly disabled the user-level monitor reads using PSR.sp).

Once the PMC/PMD registers have been configured, counting is enabled and disabled by setting bits in the PSR. User-level counters can be controlled at user-level using the rum and sum instructions to toggle PSR.up. Privileged counters are controlled by privileged software using the rsm, ssm, mov from/to PSR instructions to toggle PSR.pp. Counting for all counters is further controlled by the PMC[0] freeze bit. When PMC[0].fr is 0, all counters are disabled. When PMC[0].fr is 1, counting is enabled based on PMC[i].pm, PSR.pp and PSR.up. For more details on controlling of the performance monitors please refer to [Section 7.2.1, “Generic Performance Counter Registers” on page 2:156](#).

The PAL firmware provides information about the performance monitor registers that are implemented on the processor through the PAL\_PERF\_MON\_INFO PAL call. Information provided by the PAL includes bit masks which indicate which PMC/PMD registers are implemented on this processor model, as well as the implemented number of generic PMC/PMD pairs, and the counter width of the generic counters.

## 12.2 Operating System Support

The monitoring mechanisms discussed in the previous section support two performance monitoring usage models that need support from an Itanium architecture-based operating system.

- Per Thread/Process Event Monitoring

To monitor processor events per thread the operating system needs to save and restore performance monitor state at thread/process context switches. This save/restore of PMC and PMD registers only needs to be done for monitored threads. The effect of the save/restore is that when a monitored thread is running, PMD reads will reflect events for the monitored thread/process only. [Section 7.2.4.2, "Performance Monitor Context Switch"](#) defines the steps required for per-thread context switch of performance monitors. It is worth noting that the PMC/PMD masks returned from PAL\_PERF\_MON\_INFO indicate which PMC/PMD registers are implemented. The context switch routine can use the mask to save/restore implemented monitors without knowing the function of the monitors.

- System Wide Event Monitoring

To monitor processor events system wide (across all processes and the operating system kernel itself), a monitor must be enabled continuously across all contexts. This can be achieved by configuring a privileged monitor (PMC.pm=1), and by ensuring that PSR.pp and DCR.pp remain set for the duration of the monitor session. Since the operating system typically reloads PSR and possibly DCR on context switch, this requires the operating system to set PSR.pp and DCR.pp for all contexts that are active during the monitoring session. One way to accomplish this is to have code in the context switch routine to always set PSR.pp and DCR.pp when system wide monitoring is in effect. Another technique is to set the initial state for all new threads/processes to PSR.pp=1, PSR.up=0, PSR.sp=0 and DCR.pp=1. Setting the per thread PSR and DCR in this way ensures that privileged monitors will be enabled across all contexts. When system wide monitoring is in effect, PSR.pp, DCR.pp as well as the PMC and PMD registers should not be altered by the context switch routine.

To support both per thread and system wide monitoring, the operating system needs to be aware which type of monitoring is being performed at any given moment. If per thread/process monitoring is active, then the operating system must save/restore monitor state for monitored threads. If system wide monitoring is active, then the operating system must ensure that PSR.pp and DCR.pp remain set.

The preferred approach for performance monitoring is for Itanium architecture-based operating systems to provide a set of kernel mode services that allow performance monitoring software to be implemented in a loadable device driver. Such a loadable device driver can support various usage monitoring models, can be adapted to



model-specific processor monitoring capabilities, and is a well-defined isolated and easily replaceable software component. The following operating system services allow a kernel mode device driver to take full advantage of the performance monitors:

- Allocation/Free Performance monitors – operating system should delegate management of the performance monitor resources to device driver.
- Process create/terminate notification – operating system should notify driver on process create/terminate.
- Thread create/terminate notification – operating system should notify driver on thread create/terminate.
- Context switch notification – operating system should notify driver on thread and process context switch. The driver will perform the required save/restore depending on the currently active usage model.
- Performance counter overflow interrupt – operating system should notify driver when a performance monitor overflow interrupt occurs.
- Get Current Process Identifier – returns a unique identifier for the current process or address space. This should be callable in any context, e.g. by an interrupt handler.
- Get Current Thread Identifier – returns a unique identifier for the current thread of execution. This should be callable in any context, e.g. by an interrupt handler.

One of the challenges when doing instruction pointer (IP) profiling is to relate the current IP to an executable binary module and to an instruction within that module. If appropriate symbol information is available, the IP can be mapped to a line of source code.

To support this IP to module mapping, it is recommended that the OS provide services to enumerate all kernel and user mode modules in memory, and to allow a kernel mode driver to be notified of each module load. The following services are recommended:

- Enumerate kernel mode modules – provides information each kernel mode module currently loaded in memory.
- Enumerate threads/processes – provides a list of current threads/processes. The list should include the unique identifier for each thread/process.
- Enumerate all user mode modules – provides information on each user mode module that is currently loaded in memory (all processes).
- Enumerate modules for a process – provides information on each user mode module that is currently loaded in memory for the selected process.
- Module load notification – OS should notify a driver when the OS loads a kernel or user mode module into memory for execution. The notification should occur before the module begins execution.

In the above services for module enumeration and load notification, the module information provided for a module should include module name, load address, size in bytes, section number (if a section of a module is loaded non-contiguously), and a process/thread identifier that identifies the process into which the module is loaded.

## §



Itanium-based systems make use of several firmware components: Processor Abstraction Layer (PAL), System Abstraction Layer (SAL), Unified Extensible Firmware Interface (UEFI) and Advanced Configuration and Power Interface (ACPI).

The PAL and SAL components work together to handle the reset abort event. The reset abort handling performs processor and system initialization for operating system (OS) boot and provides an API to the operating system loader. The PAL and SAL firmware layers work together to handle machine check aborts (MCA), initialization events (INIT), and platform management interrupt (PMI) handling. All firmware components also provide runtime procedure calls to abstract processor and platform functions that may vary across implementations.

This chapter will provide an overview of the firmware components and how the firmware components interact with each other as well as with the operating system. For the full architecture specifications of the PAL firmware please refer to [Chapter 11, "Processor Abstraction Layer."](#) For full architecture specifications on SAL, UEFI and ACPI firmware components please refer to [Section 1.2, "Related Documents" on page 2:505.](#)

The PAL layer is developed by Intel Corporation and delivered with the processor. The SAL, UEFI and ACPI firmware is developed by the platform manufacturer and provide a means of supporting value added platform features from different vendors.

The interaction of the various functional firmware blocks with the processor, platform and operating system is shown in [Figure 13-1, "Firmware Model" on page 2:624.](#)

## 13.1 Processor Boot Flow Overview

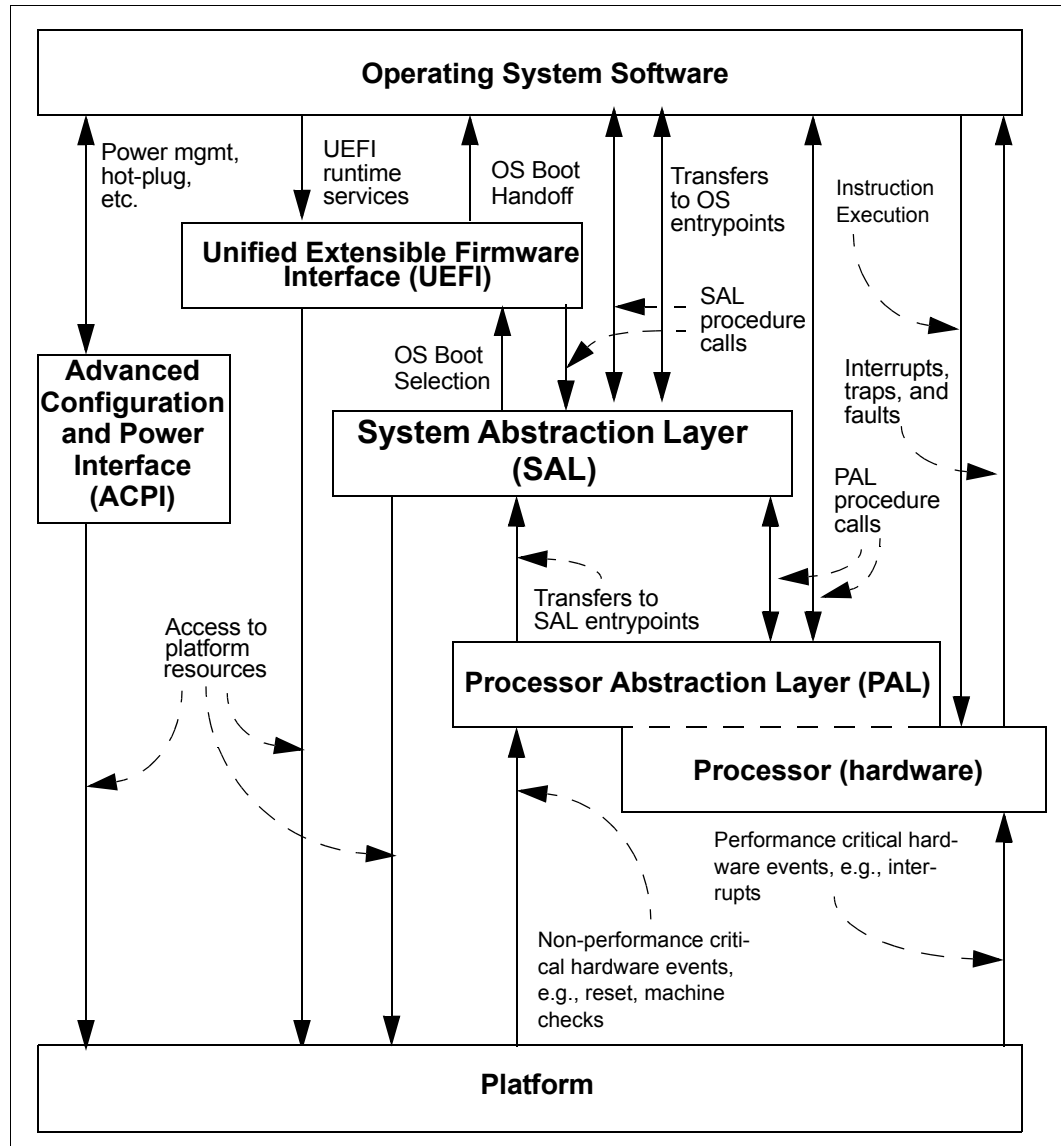
### 13.1.1 Firmware Boot Flow

Upon detection of a reset event on a processor based on the Itanium architecture, execution begins at an architected entry point inside of PAL. This PAL code will verify the integrity of the PAL code and may perform some basic processor testing. PAL will then branch to an entry point within the SAL firmware. This first branch to SAL is to determine if a firmware update is needed requiring re-programming of the firmware code. If no firmware update is needed SAL will branch back to PAL.

PAL now performs additional processor testing and initialization. These first processor tests are performed without platform memory. PAL indicates the outcome of the testing and branches to an entry point within SAL firmware for the second time. SAL will now begin platform testing and initialization. The exact division of work between SAL and UEFI from that point on is platform implementation dependent. It is required that the SAL runtime services, the UEFI boot and runtime services, and the ACPI tables and control methods be exposed to the operating systems for correct operation.

The order of steps within the UEFI/SAL firmware is platform implementation dependent and may vary. In general, the UEFI/SAL firmware selects a Bootstrap processor (BSP) in multiprocessor (MP) configurations early in the boot sequence. Next, UEFI/SAL will find and initialize memory and invoke PAL procedures to conduct additional processor tests to ensure the health of the processors. UEFI/SAL then initializes the system fabric and platform devices.

**Figure 13-1. Firmware Model**



The UEFI firmware may incorporate a Boot Manager. The UEFI firmware specification [UEFI] enables booting from a variety of mass storage devices such as hard disk, CD, DVD as well as remote boot via a network. At a minimum, one of the mass storage devices contains an UEFI system partition.

The UEFI Boot Manager displays the list of operating system choices and permits the user to select the operating system for booting. To support this functionality, the OS setup program stores the boot paths of the OS loaders and boot options in non-volatile storage managed by the UEFI firmware. The UEFI reserves the environment variables `Boot####` (#### represents values 0000 to 0xFFFF) for this purpose. The OS setup program must also store the OS loader binary images within the UEFI System Partition. The UEFI Boot Manager will also allow the user to add boot options, delete boot options, launch an UEFI application, and set the auto-boot time out value.

The UEFI System Partition also contains UEFI drivers that may be loaded by the UEFI firmware prior to transfer of control to an OS loader. The floating-point software assist (FPSWA) library is included in a UEFI runtime driver. The FPSWA library may be invoked by the OS during floating-point exception faults and traps. Please see [Section 8.1.1, “Software Assistance Exceptions \(Faults and Traps\)” on page 2:587](#) for more information on the usage of this library.

If the user elects to boot an Itanium architecture-based operating system, the UEFI loads the appropriate OS loader from the UEFI System Partition and passes control to it. The OS loader will load other files including the OS kernel from an OS partition using the UEFI boot services which provides an API interface to the OS loader.

The OS loader can obtain information about the memory map usage of the firmware by making the UEFI procedure call `GetMemoryMap()`. This procedure provides information related to the size and attributes of the memory regions currently used by firmware.

The OS loader will then jump to the OS kernel that takes control of the system. Until this point, system firmware retained control of key system resources such as the Interrupt Vector Table and provided the necessary interrupt, trap and fault handlers.

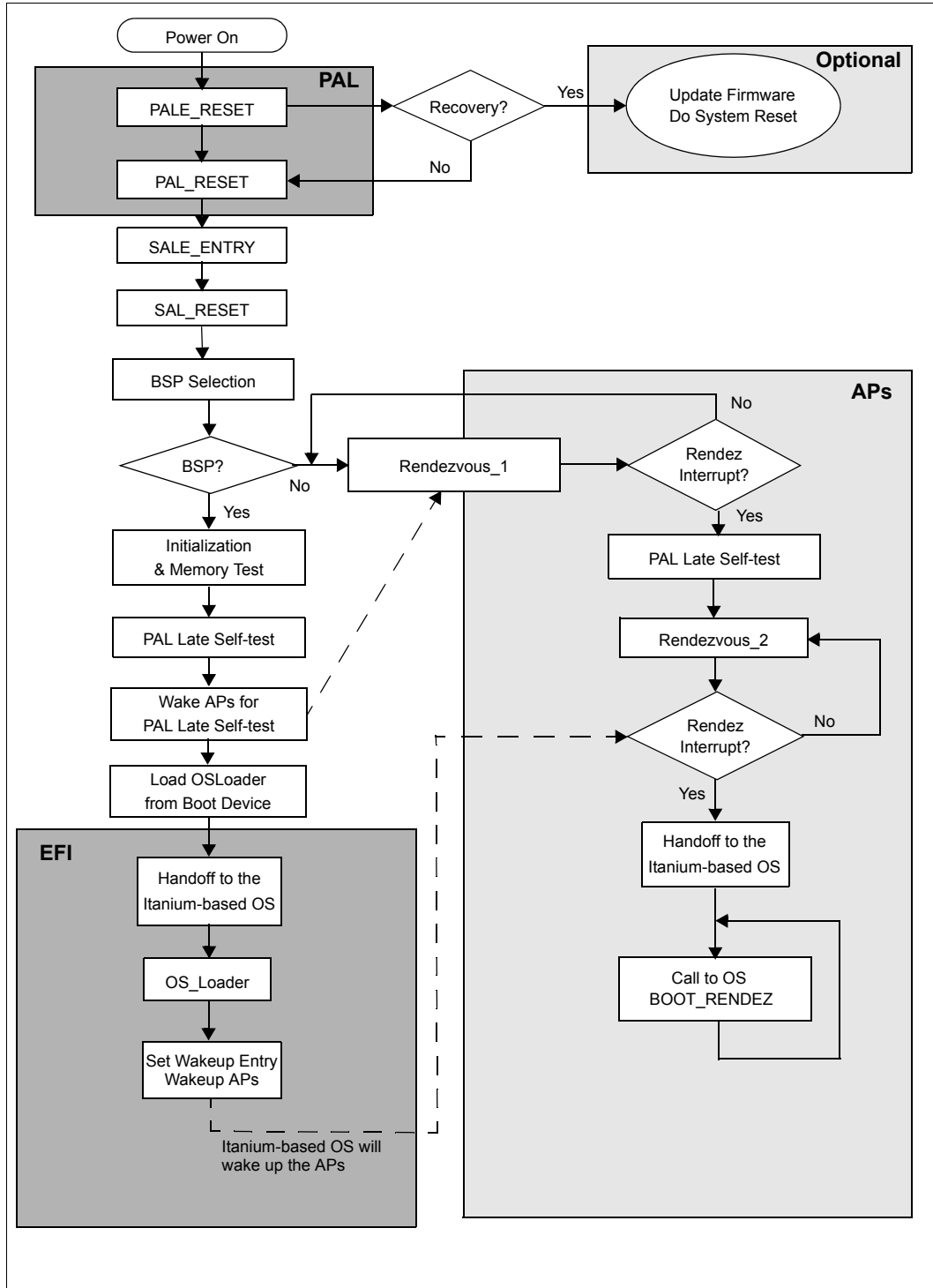
[Figure 13-2, “Control Flow of Boot Process in a Multiprocessor Configuration” on page 2:626](#) depicts the booting steps in a MP configuration.

## 13.1.2 Operating System Boot Steps

The firmware will initialize the processor(s) and platform to a specific state before handing off to the operating system boot loader. The boot loader is then responsible for copying the operating system from some storage medium into memory for running. Once this is done the operating system will need to initialize some key registers before entering into a higher level language code such as C. This section will describe code that an OS will need to execute in order to initialize system registers for preparing an OS to run in virtual mode and handle interrupts. [Appendix A, “Code Examples”](#) provides the Itanium architecture-based sample assembly code described in this section.

Assuming the specific operating system boot loader hands off to the OS kernel in physical mode, the operating system should first disable interrupts and interrupt collection via the PSR. This is done to avoid taking external interrupts from timers, etc and also prepares for writing specific system registers that require PSR.ic to be 0 when written.

**Figure 13-2. Control Flow of Boot Process in a Multiprocessor Configuration**



Next the operating system startup code invalidates the ALAT via the `invala` instruction. The `invala` in complete form will invalidate all entries in the ALAT.

The register stack should be invalidated. This can be done by setting the Register Stack Configuration Register (RSC) to zero followed by a loadrs instruction. Setting the RSC to zero will put the register stack in enforced lazy mode and set the RSC.loadrs, load distance to tear point, to zero. The loadrs will invalidate all stacked registers outside current frame.

The region registers and protection key registers are then initialized with operating system implementation dependent values. For example, the OS will initialize the region register with a preferred page size. It would also disable the VHPT until it was ready for it. In the example, all region registers are initialized with an 8-KB page size.

An OS must setup a kernel stack pointer and backing store pointer for the register stack. The stack pointer (GR12) is set to the OS kernel stack area with scratch space to cover calling conventions. AR.RSC must be set to enforced lazy mode before writing to the bspstore register. Initializing the bspstore has effects on all three RSE pointers (BSP, BSPSTORE, and RSE.BspLoad).

In order for the operating systems to handle interruptions, the operating system interrupt vector table base address must be set up. The size of the vector table is 32K bytes and is 32K byte aligned. Setting the location of the table is accomplished by moving the address into CR.IVA.

Operating systems setup system address translations for the kernel text and data by using the translation insertion format described in [Section 4.1.1.5, "Translation Insertion Format" on page 2:53](#). A combination of a general register, Interruption TLB Insertion Register (ITIR), and the Interruption Faulting Address register (IFA) are used to insert entries into the TLB. To void TLB faults on specific text and data areas the operating system can lock critical virtual memory translations in the TLB by use of Translation Register (TR) section of the TLB. The entries are placed into a TR via the Insert Translation Register (itr) instruction. The translation will remain unless the software issues the Purge Translation (ptr) instruction. Other important areas might be locked also, such as entries for memory mapped I/O, etc.

After the initial translations have been entered, the OS can make final preparations for enabling virtual addressing. The OS needs to set several important bits in the IPSR, such as data address translation (dt), register stack translation (rt), instruction address translation (it), enabling interruption collection (ic), and setting the specific register bank (bn).

The Default Control Register (DCR) specifies the default parameters for PSR values on interruption, some additional global controls, and whether speculative load faults can be deferred. The example defers all speculation faults. Also, if the operating system is utilizing the performance monitors then the DCR.pp bit should be set so that on interruption the PSR.pp bit will be set.

The global pointer (GR1) should point to the global data area. It must be setup properly before using higher level languages such as C. The startup code should also set the following registers to zero, the Interruption Function State (CR.IFS, to set frame marker to zero), and AR.RNAT (to make sure no NaT bits are set before OS kernel begins using the RSE).

Before enabling virtual addressing, the Interruption Instruction Bundle Pointer (IIP) is set to point a virtual address. This is done so when the return from interruption instruction (*rfi*) is executed the instruction fetched will have a virtual address. The *rfi* will switch modes based on IPSR values which are moved into the PSR. The IIP value becomes the new IP.

## 13.2 Runtime Procedure Calls

The PAL, SAL, and UEFI firmware components provide entry points as runtime interfaces to the OS. These runtime interfaces allow the OS to obtain information about the processor and platform as well as perform implementation-specific functions on the processor and platform.

The calling conventions for these runtime procedures are documented in the respective firmware architecture specifications. For PAL and SAL, the first input argument to the procedure call specifies the index of the procedure within the list of supported procedures for each firmware layer.

### 13.2.1 PAL Procedure Calls

PAL procedure calls are classified into two types: static and stacked. The static calls are intended for boot-time use before main memory is available or in error recovery situations where memory or the RSE may not be reliable. All parameters will be passed in the general registers GR28 to GR31 of Bank 1. The stacked registers (GR32 to GR127) will not be used for these calls. The static calls can be called at both boot-time and runtime.

Stacked register calls are intended for use after memory has been made available. The stacked registers are used for parameter passing and local variable allocation. These calls also allow memory pointers may be passed as arguments. These calls can be made at boot-time after memory has been tested and initialized as well as runtime.

For a listing of all the PAL procedures and their classification please see [Section 11.10.1, "PAL Procedure Summary" on page 2:354](#).

All PAL calls are re-entrant and can be executed simultaneously on multiple processors.

#### 13.2.1.1 Making a Static PAL Call

Since the static PAL calls do not use stacked registers, these calls are made as a pure jump with branch register B0 containing the address of the bundle to which control will return. The following code example describes how to make a static PAL call:



```

GetFeaturesCall:

mov r14 = ip // Get the ip of the current bundle
movl r28 = PAL_PROC_GET_FEATURES// Index of the PAL procedure
movl r4 = AddressOfPALProc;;; Address of the PAL proc entry point
ld8 r4 = [r4];; Read address from local pointer
mov b5 = r4 // Move address into a branch register

// Compute the return address in a position independent manner

addl r14 = (BackHome - GetFeaturesCall),r14;;
mov b0 = r14 // b0 is the return link
mov r29 = r0 // Initialize rest of input arguments
mov r30 = r0 // to zero as required by the
mov r31 = r0 // architecture.

br.sptk b5;; // Make the PAL call.

// PAL will return here when the call is completed

BackHome:

```

The sample code is position independent and functions in both physical and virtual addressing modes. Since the return address is evaluated by using the runtime instruction pointer (IP value), it will run from any address. This attribute is important for any relocatable code.

The address of the PAL procedure entry point is passed to SAL at the hand-off from PAL to SAL during reset. SAL will pass this information on to the OS during OS boot as well.

### 13.2.1.2 Making a Stacked PAL Call

A stacked PAL call uses the stacked registers for argument passing and local variable allocation. The stacked PAL calls conform to the calling conventions document [SWC], with the exception that general register GR28 must also contain the function index input argument. The following code example describes how to make a stacked PAL call.

```

movl r4 = AddressOfPALProc;;// Address of the PAL proc entry point
ld8 r4 = [r4];;// Read address from local pointer
mov b5 = r4 // Move address into a branch register

// Make the PAL_HALT_INFO procedure call. PAL_HALT_INFO uses stacked
register
// convention and parameters are passed with in0-in3

mov r28 = PAL_HALT_INFO;;// Index of the PAL procedure
mov out0 = r28// r28 and in0 must both contain the
// index value for stacked PAL calls.
mov out1 = ScratchMem_Pointer// Pointer to the memory argument
mov out2 = 0x0// Write zero to unused input arguments
mov out3 = 0x0

br.call.sptk.few b0 = b5;;// PAL stacked call

// PAL will return here when the call is completed

```

### 13.2.1.3 PAL Procedure Calls and Performance

PAL procedure calls are designed for a number of different functions varying from boot-time usage before platform memory is available to processor-specific functions used during runtime by the OS. PAL runtime procedure calls made by the OS are designed to be flexible with minimal overhead. The following features aid in this goal:

- PAL procedure calls are relocatable. This feature is useful for platforms that have PAL stored in non-volatile storage, such as flash. During OS boot the PAL procedures are copied into RAM which will reduce the memory latency.
- A number of PAL procedure calls are defined to be called in both physical and virtual addressing. This allows the caller to make the call in its currently executing addressing mode, thus reducing the need to switch between physical and virtual addressing.

### 13.2.2 SAL Procedure Calls

All SAL procedure calls use the stacked register calling convention. SAL follows the floating-point register conventions specified in the calling conventions document [SWC], with the exception that SAL does not use the floating-point registers FR32 to FR127. This exception eliminates the need for the OS to save these registers across SAL procedure calls.

SAL procedures are non re-entrant. The OS is required to enforce single threaded access to the SAL procedures except for the following procedures:

- SAL\_MC\_RENDEZ, SAL\_CACHE\_INIT, SAL\_CACHE\_FLUSH

### 13.2.3 UEFI Procedure Calls

UEFI procedure calls are classified into the following two categories: boot services and runtime services. The UEFI boot services execute in physical addressing mode only. The runtime services can execute in either physical or virtual addressing mode. The UEFI boot services are only available during the boot process and are terminated by a call to

the `EfiExitBootServices()` procedure. After this call, UEFI boot services may no longer be invoked by the OS. The UEFI runtime services execute in physical mode until the OS invokes the `EfiSetVirtualAddress()` function to switch the UEFI to virtual mode. After this point, the UEFI runtime services may be invoked in virtual mode only. For full information on all the UEFI boot and runtime services please refer to the UEFI specification [UEFI].

## 13.2.4 ACPI Control Methods

Advanced Configuration and Power Interface (ACPI) firmware provides a method of reporting system resources (up to the boundary of the box) to the operating systems. ACPI uses tables to describe system information, features, and methods for controlling those features. The ACPI tables list devices on the system board, devices that cannot be detected by bus walks, and devices which require the OS for power or temperature management. The ACPI control methods use a pseudo-code language called AML (ACPI Machine Language). AML is a tokenized language. The OS contains and uses an AML interpreter that interprets and executes these methods stored in the ACPI tables.

## 13.2.5 Physical and Virtual Addressing Mode Considerations

All of the PAL procedures can be called in the physical addressing mode. A subset of PAL calls can be made using the virtual addressing mode. For PAL calls that can be invoked using virtual addressing mode, it is the responsibility of the caller to map these PAL procedures with an ITR as well as either a DTR or DTC. If the caller chooses to map the PAL procedures using a DTC it must be able to handle TLB faults that could occur. See [Section 11.10.1, "PAL Procedure Summary"](#) for a summary of all PAL procedures and the calling conventions.

The SAL and UEFI firmware layers have been designed to operate in virtual addressing mode. UEFI provides an interface to the OS loader that describes the physical memory addresses used by firmware and indicates whether the virtual address of such areas need to be registered by the OS with UEFI. The UEFI Specification [UEFI] also provides the interfaces for the OS to register the virtual address mappings. In a MP configuration, the virtual addresses registered by the OS must be valid globally on all the processors in the system.

The SAL runtime services may be called either in virtual or physical addressing mode. SAL procedures that execute during machine check, INIT, and PMI handling must be invoked in physical addressing mode.

The parameters passed to the firmware runtime services must be consistent with the addressing environment, i.e. `PSR.dt`, `PSR.rt` setting. Additionally, the global pointer (`gp`) register [SWC] must contain the physical or virtual address for use by the firmware.

### 13.2.5.1 SAL Procedures that Invoke PAL Procedures

Some of the SAL runtime services, e.g. `SAL_CACHE_FLUSH`, will need to invoke PAL procedures. While invoking these SAL procedures in virtual mode, the OS must provide the appropriate translation resources required by PAL (i.e. ITR and DTC covering the PAL code area).

In general, if SAL needs to invoke a PAL procedure, it will do so in the same addressing mode in which it was called by the OS (i.e. without changing the PSR.dt, PSR.rt, and PSR.it bits). If a particular PAL procedure can only be invoked in physical mode, SAL will turn off translations and then invoke the PAL procedure. SAL will then restore translations before returning to the caller. The PAL\_CACHE\_INIT procedure invoked by the SAL\_CACHE\_INIT is an example of a procedure that would require such an addressing mode transition.

## 13.3 Event Handling in Firmware

The PAL and SAL firmware layers are responsible for handling three events. These events are the machine check abort (MCA), the initialization event (INIT) and the platform management interrupt (PMI). When the processor detects these events it will pass control to PAL for handling. The following sections describe the high level overview of the firmware handling of these events.

### 13.3.1 Machine Check Abort (MCA) Flows

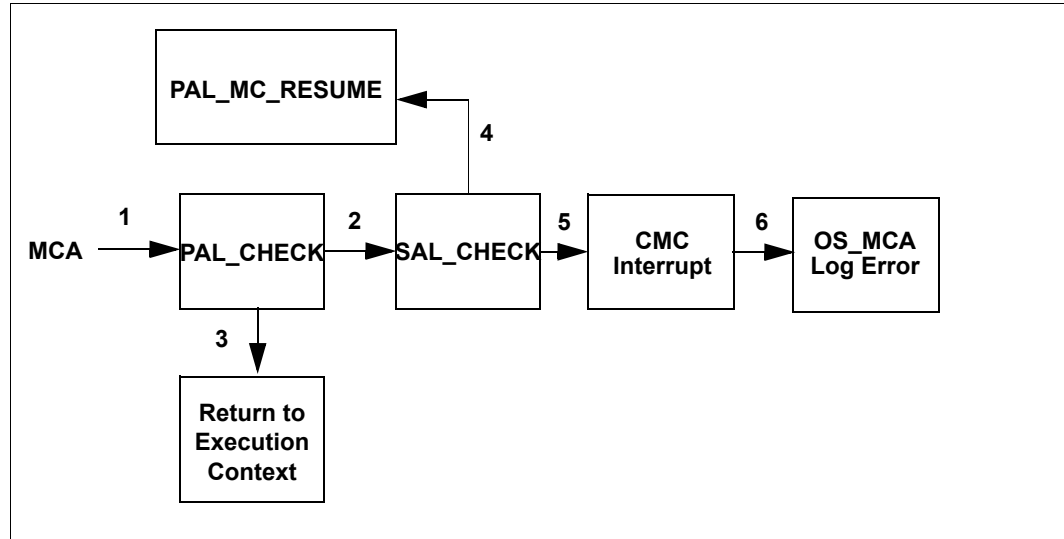
In order to have a highly reliable and fault tolerant computing environment a great deal of coordination and cooperation between the system entities (i.e. the processor, platform, and system software) is required. The PAL firmware, the SAL firmware, and the operating system all work together to meet this goal. This section will provide an overview of the machine check abort handling.

When the processor detects an error, control is transferred to the PAL\_CHECK entrypoint. PAL\_CHECK will perform error analysis and processor error correction where possible. Subsequently, PAL either returns to the interrupted context or hands off control to the SAL\_CHECK component. The level of recovery provided by PAL\_CHECK is implementation dependant and is beyond the scope of this specification. SAL\_CHECK will perform error logging and platform error correction where possible. Errors that are corrected by PAL and SAL firmware are logged and control is transferred back to the interrupted process/context. For corrected errors, no OS intervention is required for error handling, but the OS is notified of the event for logging purposes through a low priority asynchronous corrected machine check interrupt (CMCI). See [Section 5.8.3.8, "Corrected Machine Check Vector \(CMCV – CR74\)"](#) for more information on the CMCI. If the error was not corrected by firmware, SAL hands off control to the OS\_MCA handler.

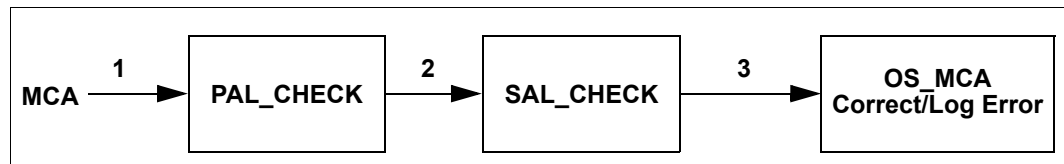
Within the firmware the entire machine check is handled with virtual address translations disabled. However, the OS machine check handler may optionally enable virtual addressing and execute most of MCA handler in virtual mode.

[Figure 13-3](#) and [Figure 13-4](#) depict an overview of Itanium machine check processing. The control flows are slightly different for corrected and uncorrected machine checks.

**Figure 13-3. Correctable Machine Check Code Flow**



**Figure 13-4. Uncorrectable Machine Check Code Flow**



For multiprocessor systems, machine checks are classified as local and global. A global MCA implies a system wide broadcast by hardware of an error condition. During a global MCA condition, all the processors in the system will be notified of the MCA, detected by one or more system components, and each of the processors in the system will start processing the MCA in their respective handlers. The SAL firmware and OS layers will coordinate the handling of the error among the processors.

A local MCA has a scope of influence that is limited to the particular processor which encountered the error. This local MCA will not be broadcast to other processors in the system and will be handled on an individual processor basis. At any point in time, more than one processor in the system may experience a local MCA and handle it without notifying other processors in the system.

The next sections will provide an overview of the responsibilities that the PAL, SAL and OS have for handling machine checks. These sections are not an exhaustive description of the functionality of the handlers but provides a high level description of how the MCA handling is split among the different components.

### **13.3.1.1 Machine Check Handling in PAL**

All machine check abort events are first handled in the PAL firmware layer. The following provides a brief description of some of the functions of the PAL machine check handler:

- Correct processor errors if possible.

- Attempt to contain the error by requesting a rendezvous for all processors in the system if needed.
- Hand off control to SAL for further processing, such as error logging.
- Return processor error log information upon request by SAL.
- Return to the interrupted context by restoring the state of the processor.
- Notify the OS about corrected machine check conditions through the CMC interrupt.

### 13.3.1.2 Machine Check Handling in SAL

Before SAL is ready to handle machine checks, it must register with PAL an uncacheable memory buffer that PAL can use to save away processor state. This area is known as the min-state save area. If a machine check occurs before this memory location has been registered, return to the interrupted context is not possible and the machine check is not recoverable.

The following provides a description of some of the functions of the SAL machine check handler.

- Attempt to rendezvous the other processors in the system on a PAL request.
- Process MCA handling after handoff from PAL.
- Retrieve processor error log information via PAL procedure calls and store this information for logging purposes.
- Issue a PAL clear log request to clear the processor error logs, which enables further logging.
- Log platform state for MCA and retain it until it is retrieved by the OS.
- Attempt to correct processor machine check errors which are not corrected by PAL.
- Attempt to correct platform machine check errors.
- Branch to the OS MCA handler for uncorrected errors or optionally reset the system.
- Return to the interrupted context via a PAL procedure call.

### 13.3.1.3 Machine Check Abort Handling in OS

Before the OS kernel is ready to handle machine checks, it must register the address of the OS\_MCA entry point and the GP [SWC] value for the OS\_MCA handler with SAL. If the OS does not register its entry point, the occurrence of a machine check will cause a system reset. In MP configurations, the OS must also register with SAL:

- A rendezvous interrupt vector which SAL firmware can use to rendezvous the processors.
- The mechanism that the OS will employ to wake up the processors at the end of machine check processing.

When the OS registers the OS\_MCA entry point with SAL, it also supplies the length of the code (or at least the length of the first level OS\_MCA handler). SAL computes and saves the checksum of this code area. Prior to entering OS\_MCA, SAL ensures that the OS\_MCA vector is valid by verifying the checksum of the OS\_MCA code. Hence, the OS\_MCA code must not contain any self modifying code.

When an uncorrected machine check event occurs, SAL will invoke the OS\_MCA handler. The functionality of this handler is dependent on the OS. At a minimum, it must call a SAL procedure to retrieve the error logging and state information and then call another SAL procedure to release these resources for future error logging and state save.

When the OS\_MCA code completes, it decides whether or not to return to the interrupted context. The OS must take into account the state information retrieved from the SAL with respect to the continuability of the processor and system. Thus, even if the OS could correct the error, if PAL or SAL reports that it did not capture the entire processor context, resumption of the interrupted context will not be possible.

The OS must also determine from values stored by PAL in the min-state save area whether the machine check occurred while operating with PSR.ic set to 0 and whether the processor supports recovery for this case. Please refer to [Section 11.3.1.1, "Resources Required for Machine Check and Initialization Event Recovery"](#) for more information on processor recovery under this condition.

To provide better software error handling, some operating systems build mechanisms to identify whether machine checks occurred during execution of the OS kernel code or in the application context. One technique to achieve this is to call the PAL\_MC\_DRAIN procedure when an application makes a system call to the OS. This procedure completes all outstanding transactions within the processor and reports any pending machine checks. This technique impacts system call and interrupt handling performance significantly, but will improve system reliability by allowing the OS to recover from more errors than if this mechanism was not included.

## 13.3.2 INIT Flows

INIT is an initialization event generated by the platform or by software through an inter-processor interrupt message. The INIT can be due to a platform INIT event or due to a failed rendezvous on an application processor.

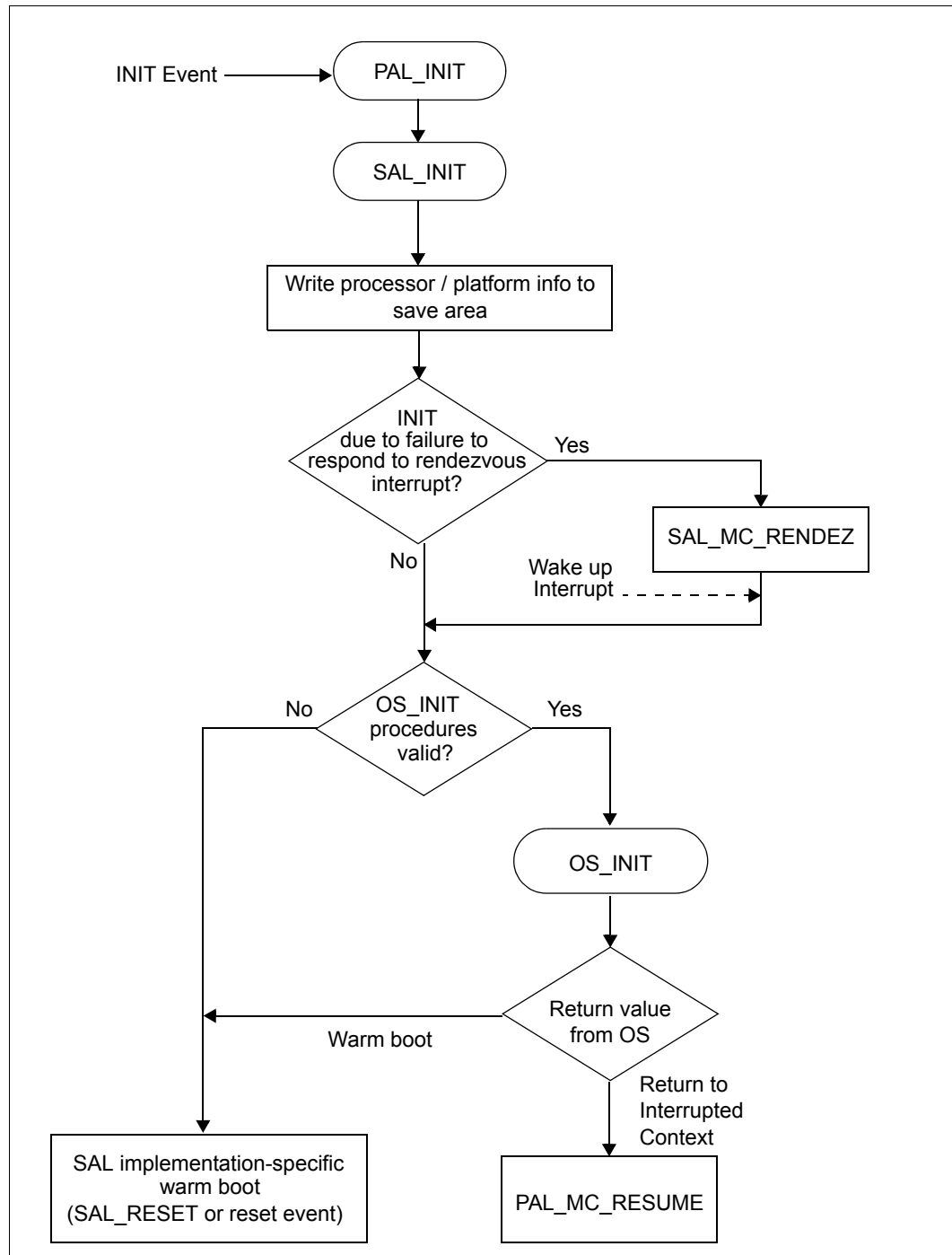
The INIT event will pass control to the PAL firmware INIT handler. The PAL INIT handler saves processor state to the registered min-state save area and sets up the architected hand off state before branching to SAL. See [Section 11.5, "Platform Management Interrupt \(PMI\)"](#) for more information on the PAL INIT handling.

The SAL INIT handler logs processor state and platform state information and then calls the OS\_INIT handler if one is registered. The OS\_INIT handler gains control in physical mode but may switch to virtual mode if necessary. The OS may choose to implement a crash dump or an interactive debugger within the OS\_INIT handler.

The OS must register the OS\_INIT entry point with SAL, otherwise the occurrence of an INIT event will cause a system reset. At the end of OS\_INIT handling, the OS must return to SAL with the appropriate exit status.

[Figure 13-5](#) illustrates the flow of control during INIT processing.

Figure 13-5. INIT Flow





### 13.3.3 PMI Flows

Processors based on the Itanium architecture implement the Platform Management Interrupt (PMI) to enable platform developers to provide high level system functions, such as power management and security, in a manner that is transparent not only to the application software but also to the operating system.

When the processor detects a PMI event it will transfer control to the registered PAL PMI entrypoint. PAL will set up the hand off state which includes the vector information for the PMI and hand off control to the registered SAL PMI handler. To reduce the PMI overhead time, the PAL PMI handler will not save any processor architectural state to memory. Please see [Section 11.5, "Platform Management Interrupt \(PMI\)"](#) for more information on PAL PMI handling.

The SAL PMI handler may choose to save some additional register state to SAL allocated memory to handle the specific platform event that generated the PMI.

The OS will not see the PMI events generated by the platform. The platform developer can use PMI interrupts to provide features to differentiate their platform.

PMI handling was designed to be executed with minimal overhead. The SAL firmware code copies the PAL and SAL PMI handlers to RAM during system reset and registers these entry-points with the processor. This code is then run with the cacheable memory attribute to improve performance.

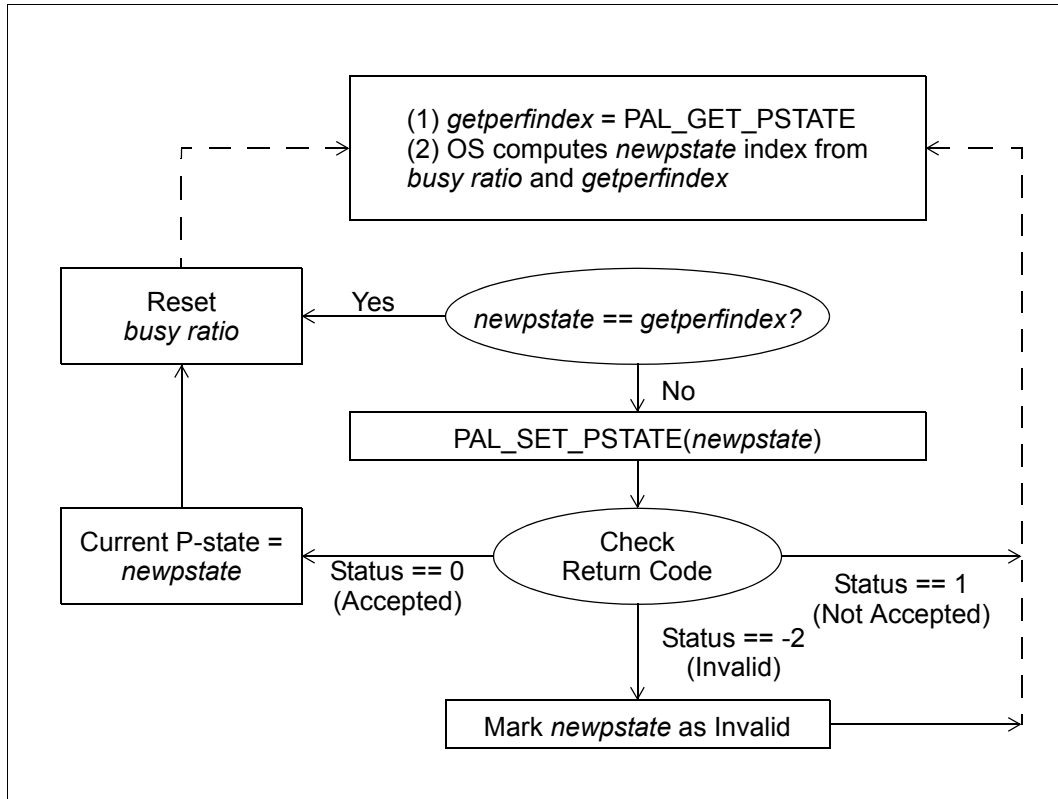
Depending on the implementation and the platform, there may be no special hardware protection of the PMI code's memory area in RAM, and the protection of this code space may be through the OS memory management's paging mechanism. SAL sets the correct attributes for this memory space and passes this information to the OS through the Memory Descriptor Table from `EfiGetMemoryMap()` [UEFI].

### 13.3.4 P-state Feedback Mechanism Flow Diagram

The example flowchart shown below illustrates how the caller can utilize the `PAL_SET_PSTATE` and the `PAL_GET_PSTATE` procedures to manage system utilization and power consumption, for a processor implementation that belongs to either a hardware-coordinated dependency domain or a hardware-independent dependency domain. At the beginning of the loop, `PAL_GET_PSTATE` gives the performance characteristics of the processor over the last time period. It is assumed that the caller maintains an internal count for determining the busy ratio of the logical processor (busy ratio can be defined as the percentage of time the processor was busy executing instructions and not idle). The caller then seeks to adjust the P-state for the next time period to match the busy ratio from the previous time period. For example, if the busy ratio for a given period was 100%, and the *performance\_index* returned by `PAL_GET_PSTATE` was 60, then this indicates that the P-state for the next time period should be P0 (which has performance index of 100). The caller would then call the `PAL_SET_PSTATE` procedure to transition the processor to the P0 state. In essence, if the busy ratio is greater than the *performance\_index* returned by `PAL_GET_PSTATE`, the caller responds to the increased demand requirement of the workload by transitioning the processor to a higher-performance P-state. Alternatively, if the busy ratio is lower

than the *performance\_index* returned by PAL\_GET\_PSTATE, the caller responds by transitioning the processor to a lower performance P-state, which consumes less power and operates at reduced performance.

**Figure 13-6. Flowchart Showing P-state Feedback Policy**



Such an adaptive policy implemented by the caller to dynamically respond to system workload characteristics using P-states allows for efficient power utilization – the processor consumes additional power by operating at a higher performance level only when the current workload requires it to do so.

§

## A.1 OS Boot Flow Sample Code

The sample code given below is an example of setting up operating system register state to prepare the processor for running in virtual mode as described in [Section 13.1.2, “Operating System Boot Steps”](#) on page 2:625.

```
// This code will perform the following steps:
//1.Initialize PSR with interrupt disabled (bit 13)
//2.Invalidate ALAT via invala instruction
//3.Invalidate register stack
//4.Set region registers rr[r0] - rr[r7] to RID=0, PS=8K, E=0.
//5.Disable the VHPT
//6.Initialize protection key registers
//7.Initialize SP
//8.Initialize BSP
//9.Enable register stack engine.
//10.Setup IVA
//11.Setup virtual->physical address translation
//12.Setup GP.

.file"start.s"

// globals

        .global main
        .type main, @function        // C function we will return to

        .global __GLOB_DATA_PTR     // External pointer to Global Data area
        .global IVT_BASE            // External pointer to IVT_BASE

        .text

// This is the entry point where primary boot loader
// passes control.

pstart::

        mov     psr.l = r0           // Initialize psr.l
        ;;
        invala                                // Invalidate ALAT
        mov     ar.rsc = r0           // Invalidate register stack
        ;;
        loadrs

// Initialize Region Registers

        mov     r2 = (13 << 2)       // 8K page size
        mov     r3 = r0
        mov     r4 = 61
        ;;

Loader_RRLoop:
        shl     r10 = r3, r4
        ;;
        mov     rr[r10] = r2
        add     r3 = 1, r3
        ;;
        cmp4.geu p6, p7 = 8, r3
```

```

(p6)br.cond.sptk.few.clr Loader_RRLoop
;;

// Disable the VHPT walker and set up the minimum size for it (32K) by writing
// to the page table address register (cr.pta)

    mov r2 = (15<<2)
        ;;
    mov cr.pta = r2

// Initialize the protection key registers for kernel

    mov r2 = (1<< 0)
    mov r3 = r0
    ;;
    mov pkr[r3] = r2 // validate pkr[zero]
    ;;
    mov r2 = r0
    ;;

pkr_loop:
    add r3=r3,r0, 1 // start with index 1
    ;;
    cmp.gtu p6,p7 = 8,r3
    ;;
(p6)mov pkr[r3] = r2
(p6)br.cond.sptk.few.clr pkr_loop // loop until 8

// Setup kernel stack pointer (r12)

    movl    sp = kstack + (64*1024) // 64K stack
    ;;

// Set up the scratch area on stack

    add    sp = - 32, sp

// Setup the Register stack backing store
//
// 1st deal with Register Stack Configuration register
//
// NOTE: the RSC mode must be enforced lazy (00) to write to bspstore
//
// mode: = enforced lazy
// be = little endian

    mov ar.rsc = r0
    ;;

//Now have to setup the RSE backing store pointer
//
//NOTE: initializing the bspstore has effects on all 3 RSE pointers
// (BSP, BSPSTORE, and RSE.BspLoad)

    movl r2 = kstack + ((96 + (96/63))*8)
    ;;
    mov ar.bspstore = r2

// Need to setup base address for interrupt vector table...

    movl r3 = IVT_BASE
    ;;
    mov cr.iva = r3

// Setup system address translation for the kernel

```

```

//
//      The Translation Insertion Format looks like the following...
//
//      Below is the register interface to insert entries into the TLB
//
//1) A general register contains an address, attributes, and permissions
//2) ITIR: additional info such as protection key page size info
//3) IFA: specifies the virtual page number for instruction and data
// TLB inserts
//
//Registers used:
//-----
//      | 63 53 | 52 | 51 50 | 49 12 | 11 9| 8 7 | 6 | 5 |4 1| 0 |
//GR   |  ig | ed |  rv |  ppn | ar | pl | d | a | ma | p |
//
// ITIR | rv {63:32} | key {31:8} | ps {7:2} | rv {1:0}|
//
//IFA   | vpn {63:12}| ignored {11:0} |
//
//RR[vrn] | reserved{63:32} | rid {31:8}| ignored {7:2} | rv{1} | ignored {0}|
//
//
//where
//ig = ignored bits
//rv= reserved bits
//p = present bit
//ma = memory attribute
//a = accessed bit
//d = dirty bit
//pl= privilege level
//ar= access rights
//ppn= physical page number
//ed= exception deferral
//ps= page size of mapping (2**ps)
//vpn= virtual page number
//
// Setup virtual page number
//
// NOTE:The virtual page number depends on a translation's
//page size.
//
// Add entry for TEXT section

    movl r2 = 0x0
    ;;
    mov  cr.ifa = r2

//setup ITIR (Interruption TLB Insertion Register)

    movl r3=( ( 24 << 2 ) | ( 0 << 8 ) ) // set page size to 16 MB
    ;;
    mov  cr.itir = r3

//now setup the general register to use with itr (insert translation
//register), use physical page of zero

    movl r10 = ((1 << 52) | ( 0x00000000 << 12 ) | ( 3 << 9 ) | ( 0 << 7 ) | \
                (1 <<6 ) | ( 1 << 5 ) | ( 1 << 0 ))
    mov  r11 = r0
    ;;
    itr.i itr[r11] = r10      // Insert translation register

//Entry for OS Data section

    add  r11 = 1, r11          // skip to tr next index

```

```

    movl r2 = 0x0                // use vpn 0
    ;;
    mov cr.ifa = r2

//Setup ITIR (Interruption TLB Insertion Register)

    movl r3 = ( ( 24 << 2 ) | ( 0 << 8 ) ) // 16 MB
    ;;
    mov cr.itir = r3

//Now setup the general register to use with itr (insert translation
//register)

    movl r10 = ((1 << 52 ) | (0x0 << 12 ) | (3 << 9 ) | (0 << 7) | \
                (1 << 6) | ( 1 << 5 ) | (1 << 0))

    ;;
    itr.d dtr[r11] = r10        // Insert translation register
    ;;

//It is now time to set the appropriate bits in the PSR (processor
//status register)

    movl r3 = ((1 << 44) | (1 << 36) | (1 << 38) | (1 << 27) | (1 << 17) | \
                (1 << 15) | (1 << 14) | (1 << 13))

    ;;
    mov cr.ipsr = r3

//Initialize DCR to defer all speculation faults

    movl r2 = 0x7f00
    ;;
    mov cr.dcr = r2

// Initialize the global pointer (gp = r1)

    movl gp = __GLOB_DATA_PTR

// Clear out ifs

    mov cr.ifs=r0

// Need to do a "rfi" in order to synchronize above instructions and set
// "it" and "ed" bits in the PSR.

    movl r3 = main              // Setup for main, C code
    ;;
    mov cr.iip = r3            // Setup iip to hit main
    ;;
    rfi
    ;;

// Setup kernel stack

.data
.globalkstack
.align 16
kstack:
.skip(64*1024)

```

## §

# ***Index***





# INDEX FOR VOLUMES 1, 2, 3 AND 4

## A

- AAA Instruction 4:21
- AAD Instruction 4:22
- AAM Instruction 4:23
- AAS Instruction 4:24
- Aborts 2:95, 2:538
- ACPI 2:631
  - P-states 2:315, 2:637
- Acquire Semantics 2:507
- ADC Instruction 4:25, 4:26
- ADD Instruction 4:27, 4:28
- add Instruction 3:14
- addp4 Instruction 3:15
- ADDPS Instruction 4:486
- Address Space Model 2:561
- ADDSS Instruction 4:487
- Advanced Load 1:153, 1:154
- Advanced Load Address Table (ALAT) 1:64
- Advanced Load Check 1:154
- ALAT (Advanced Load Address Table) 1:64
  - Coherency 2:554
  - Data Speculation 1:17
- alloc Instruction 3:16
- AND Instruction 4:29, 4:30
- and Instruction 3:18
- andcm Instruction 3:19
- ANDNPS Instruction 4:488
- ANDPS Instruction 4:489
- Application Architecture Guide 1:1
- Application Memory Addressing Model 1:36
- Application Register (AR) 1:23, 1:28, 1:140
- AR (Application Register) 1:28, 1:140
- Arithmetic Instructions 1:51
- ARPL Instruction 4:31, 4:32

## B

- Backing Store 2:133
- Banked General Registers 2:42
- Bit Field and Shift Instructions 1:52
- Bit Strings 1:84
- Boot Sequence 2:13
- BOUND Instruction 4:33
- BR (Branch Register) 1:26, 1:140
- br Instruction 3:20
  - br.ia 1:112, 2:596
- Branch Hints 1:78, 1:176
- Branch Instructions 1:74, 1:145
- Branch Register (BR) 1:19, 1:26, 1:140
- break Instruction 2:556, 3:29
- Break Instruction Fault 2:151
- brl Instruction 3:30
- brp Instruction 3:32
- BSF Instruction 4:35
- BSP (RSE Backing Store Pointer Register) 1:29
- BSPSTORE (RSE Backing Store Pointer for Memory

- Stores Register) 1:30

- BSR Instruction 4:37
- bsw Instruction 3:34
- BSWAP Instruction 4:39
- BT Instruction 4:40
- BTC Instruction 4:42
- BTR Instruction 4:44
- BTS Instruction 4:46
- Bundle Format 1:38
- Bundles 1:38, 1:141
- Byte Ordering 1:36

## C

- CALL Instruction 4:48
- CBW Instruction 4:57
- CCV (Compare and Exchange Value Register) 1:30
- CDQ Instruction 4:85
- CFM (Current Frame Marker) 1:27
- Character Strings 1:83
- Check Code 1:161
- Check Load 1:154
- chk Instruction 3:35
- CLC Instruction 4:59
- CLD Instruction 4:60
- CLI Instruction 4:61
- clrrrb Instruction 3:37
- CLTS Instruction 4:63
- clz Instruction 3:38
- CMC (Corrected Machine Check) 2:350
- CMC Instruction 4:64
- CMCV (Corrected Machine Check Vector) 2:126
- CMP Instruction 4:69
- cmp Instruction 3:39
- cmp4 Instruction 3:43
- CMPPS Instruction 4:490
- CMPS Instruction 4:71
- CMPSB Instruction 4:71
- CMPSD Instruction 4:71
- CMPSW Instruction 4:493
- CMPSW Instruction 4:71
- CMPXCHG Instruction 4:74
- cmpxchg Instruction 2:508, 3:46
- CMPXCHG8B Instruction 4:76
- Coalescing Attribute 2:78
- COMISS Instruction 4:496
- Compare and Exchange Value Register (CCV) 1:30
- Compare and Store Data Register (CSD) 1:30
- Compare Types 1:55
- Context Management 2:549
- Context Switching 2:557
  - Operating System Kernel 2:558
  - User-Level 2:557
- Control Dependencies 1:148
- Control Registers 2:29
- Control Speculation 1:16, 1:60, 1:142, 1:151,

1:155, 2:579  
 Control Speculative Load 1:156  
 Corrected Error 2:350  
 Corrected Machine Check Vector (CMCV) 2:126  
 cover Instruction 3:48  
 CPUID (Processor Identification Register) 1:34  
 CPUID Instruction 4:78  
 Cross-modifying Code 2:533  
 CSD (Compare and Store Data Register) 1:30  
 Current Frame Marker (CFM) 1:27  
 CVTPI2PS Instruction 4:498  
 CVTPS2PI Instruction 4:500  
 CVTSI2SS Instruction 4:502  
 CVTSS2SI Instruction 4:503  
 CVTTPS2PI Instruction 4:504  
 CVTTSS2SI Instruction 4:506  
 CWD Instruction 4:85  
 CWDE Instruction 4:57, 4:86  
 czx Instruction 3:49

**D**

DAA Instruction 4:87  
 DAS Instruction 4:88  
 Data Arrangement 1:81  
 Data Breakpoint Register (DBR) 2:151, 2:152  
 Data Debug Faults 2:152  
 Data Dependencies 1:149, 1:150, 3:371  
 Data Poisoning 2:302  
 Data Prefetch Hint 1:148  
 Data Serialization 2:18  
 Data Speculation 1:17, 1:63, 1:143, 1:151, 2:579  
 Data Speculative Load 1:154  
 DBR (Data Breakpoint Register) 2:151, 2:152  
 DCR (Default Control Register) 2:31  
 Debugging 2:151  
 DEC Instruction 4:89  
 Default Control Register (DCR) 2:31  
 Dekker's Algorithm 2:529  
 dep Instruction 3:51  
 DIV Instruction 4:91  
 DIVPS Instruction 4:507  
 DIVSS Instruction 4:508

**E**

EC (Epilog Count Register) 1:33  
 EFLAG (IA-32 EFLAG Register) 1:123  
 EMMS Instruction 4:400  
 End of External Interrupt Register (EOI) 2:124  
 Endian 1:36  
 ENTER Instruction 4:94  
 EOI (End of External Interrupt Register) 2:124  
 epc Instruction 2:555, 3:53  
 Epilog Count Register (EC) 1:33  
 Explicit Prefetch 1:70  
 External Controller Interrupts 2:96

External Interrupt 2:96, 2:538  
 External Interrupt Control Registers (CR64-81)  
 2:42  
 External Interrupt Request Registers (IRR0-3)  
 2:125  
 External Interrupt Vector Register (IVR) 2:123  
 External Task Priority Cycle (XTP) 2:130  
 External Task Priority Register (XTPR) 2:605  
 ExtINT (External Controller Interrupt) 2:96  
 extr Instruction 3:54

**F**

F2XM1 Instruction 4:97  
 FABS Instruction 4:99  
 fabs Instruction 3:55  
 FADD Instruction 4:100  
 fadd Instruction 3:56  
 FADDP Instruction 4:100  
 famax Instruction 3:57  
 famin Instruction 3:58  
 fand Instruction 3:59  
 fandcm Instruction 3:60  
 Fatal Error 2:350  
 Fault Handlers 2:583  
 Faults 2:96, 2:537  
 FBLD Instruction 4:103  
 FBSTP Instruction 4:105  
 fc Instruction 3:61  
 fchkf Instruction 3:63  
 FCHS Instruction 4:108  
 fclass Instruction 3:64  
 FCLEX Instruction 4:109  
 fclrf Instruction 3:66  
 FCMOI Instruction 4:115  
 FCMOVcc Instruction 4:110  
 fcmp Instruction 3:67  
 FCOM Instruction 4:112  
 FCOMIP Instruction 4:115  
 FCOMP Instruction 4:112  
 FCOMPP Instruction 4:112  
 FCOS Instruction 4:118  
 FCR (IA-32 Floating-point Control Register) 1:126  
 fcvt Instruction  
   fcvt.fx 3:70  
   fcvt.xf 3:72  
   fcvt.xuf 3:73  
 FDECSTP Instruction 4:120  
 FDIV Instruction 4:121  
 FDIVP Instruction 4:121  
 FDIVR Instruction 4:124  
 FDIVRP Instruction 4:124  
 Fence Semantics 2:508  
 fetchadd Instruction 2:508, 3:74  
 FFREE Instruction 4:127  
 FIADD Instruction 4:100

- FICOM Instruction 4:128
- FICOMP Instruction 4:128
- FIDIV Instruction 4:121
- FIDIVR Instruction 4:124
- FILD Instruction 4:130
- FIMUL Instruction 4:145
- FINCSTP Instruction 4:132
- Firmware 1:7, 2:623
- Firmware Address Space 2:283
- Firmware Entrypoint 2:281, 2:350
- Firmware Interface Table (FIT) 2:287
- FIST Instruction 4:134
- FISTP Instruction 4:134
- FISUB Instruction 4:182, 4:183
- FISUBR Instruction 4:185
- FIT (Firmware Interface Table) 2:287
- FLD Instruction 4:137
- FLD1 Instruction 4:139
- FLDCW Instruction 4:141
- FLDENV Instruction 4:143
- FLDL2E Instruction 4:139
- FLDL2T Instruction 4:139
- FLDLG2 Instruction 4:139
- FLDLN2 Instruction 4:139
- FLDPI Instruction 4:139
- FLDZ Instruction 4:139
- Floating-point Architecture 1:19, 1:85, 1:205
- Floating-point Exception Fault 1:102
- Floating-point Instructions 1:91
- Floating-point Register (FR) 1:139
- Floating-point Software Assistance Exception Handler (FPSWA) 2:587
- Floating-point Status Register (FPSR) 1:31, 1:88
- flushrs Instruction 3:76
- fma Instruction 1:210, 3:77
- fmax Instruction 3:79
- fmerge Instruction 3:80
- fmin Instruction 3:82
- fmix Instruction 3:83
- fmpy Instruction 3:85
- fms Instruction 3:86
- FMUL Instruction 4:145
- FMULP Instruction 4:145
- FNCLEX Instruction 4:109
- fneg Instruction 3:88
- fnegabs Instruction 3:89
- FNINIT Instruction 4:133
- fnma Instruction 3:90
- fnmpy Instruction 3:92
- FNOP Instruction 4:148
- fnorm Instruction 3:93
- FNSAVE Instruction 4:162
- FNSTCW Instruction 4:176
- FNSTENV Instruction 4:178
- FNSTSW Instruction 4:180
- for Instruction 3:94
- fpabs Instruction 3:95
- fpack Instruction 3:96
- fpamax Instruction 3:97
- fpamin Instruction 3:99
- FPATAN Instruction 4:149
- fpcmp Instruction 3:101
- fpcvt Instruction 3:104
- fpma Instruction 3:107
- fpmax Instruction 3:109
- fpmerge Instruction 3:111
- fpmmin Instruction 3:113
- fpmpy Instruction 3:115
- fpms Instruction 3:116
- fpneg Instruction 3:118
- fpnegabs Instruction 3:119
- fpnma Instruction 3:120
- fpnmpy Instruction 3:122
- fprcpa Instruction 3:123
- FPREM Instruction 4:151
- FPREM1 Instruction 4:154
- fprsqta Instruction 3:126
- FPSR (Floating-point Status Register) 1:31, 1:88
- FPSWA (Floating-point Software Assistance Handler) 2:587
- FPTAN Instruction 4:157
- FR (Floating-point Register) 1:139
- frcpa Instruction 3:128
- FRNDINT Instruction 4:159
- frsqta Instruction 3:131
- FRSTOR Instruction 4:160
- FSAVE Instruction 4:162
- FSCALE Instruction 4:165
- fselect Instruction 3:134
- fsetc Instruction 3:135
- FSIN Instruction 4:167
- FSINCOS Instruction 4:169
- FSQRT Instruction 4:171
- FSR (IA-32 Floating-point Status Register) 1:126
- FST Instruction 4:173
- FSTCW Instruction 4:176
- FSTENV Instruction 4:178
- FSTP Instruction 4:173
- FSTSW Instruction 4:180
- FSUB Instruction 4:182, 4:183
- fsub Instruction 3:136
- FSUBP Instruction 4:182, 4:183
- FSUBR Instruction 4:185
- FSUBRP Instruction 4:185
- fswap Instruction 3:137
- fsxt Instruction 3:139
- FTST Instruction 4:188
- FUCOM Instruction 4:190
- FUCOMI Instruction 4:115
- FUCOMIP Instruction 4:115
- FUCOMP Instruction 4:190
- FUCOMPP Instruction 4:190

FWAIT Instruction 4:386  
 fwb Instruction 3:141  
 FXAM Instruction 4:193  
 FXCH Instruction 4:195  
 fxor Instruction 3:142  
 FXRSTOR Instruction 4:509  
 FXSAVE Instruction 4:512, 4:515  
 FXTRACT Instruction 4:197  
 FYL2X Instruction 4:199  
 FYL2XP1 Instruction 4:201

**G**

General Register (GR) 1:25, 1:139  
 getf Instruction 3:143  
 GR (General Register) 1:139

**H**

hint Instruction 3:145  
 HLT Instruction 4:203

**I**

I/O Architecture 2:615

## IA-32

IA-32 Application Execution 1:109  
 IA-32 Applications 2:239, 2:595  
 IA-32 Architecture 1:7, 1:21  
 IA-32 Current Privilege Level (PSR.cpl) 2:243  
 IA-32 EFLAG Register 1:123, 2:243  
 IA-32 Exception  
   Alignment Check Fault 2:229  
   Code Breakpoint Fault 2:215  
   Data Breakpoint, Single Step, Taken  
     Branch Trap 2:216  
   Device Not Available Fault 2:221  
   Divide Fault 2:214  
   Double Fault 2:222  
   General Protection Fault 2:226  
   INT 3 Trap 2:217  
   Invalid Opcode Fault 2:220  
   Invalid TSS Fault 2:223  
   Machine Check 2:230  
   Overflow Trap 2:218  
   Page Fault 2:227  
   Pending Floating-point Error 2:228  
   Segment Not Present Fault 2:224  
   SSE Numeric Error Fault 2:231  
   Stack Fault 2:225  
 IA-32 Execution Layer 1:109  
 IA-32 Floating-point Control Registers 1:126  
 IA-32 Instruction Reference 4:11  
 IA-32 Instruction Set 2:253  
 IA-32 Intel® MMX™ Technology 1:129  
 IA-32 Intercept  
   Gate Intercept Trap 2:235  
   Instruction Intercept Fault 2:233

Locked Data Reference Fault 2:237  
 System Flag Trap 2:236

## IA-32 Interrupt

Software Trap 2:232

IA-32 Interruption 2:111

IA-32 Interruption Vector Definitions 2:213

IA-32 Interruption Vector Descriptions 2:213

IA-32 Memory Ordering 2:265

IA-32 Physical Memory References 2:262

IA-32 SSE Extensions 1:20, 1:130

IA-32 System Registers 2:246

IA-32 System Segment Registers 2:241

IA-32 Trap Code 2:213

IA-32 Virtual Memory References 2:261

IBR (Index Breakpoint Register) 2:151, 2:152

IDIV Instruction 4:204

IFA (Interruption Faulting Address) 2:541

IFS (Interruption Function State) 2:541

IHA (Interruption Hash Address) 2:41, 2:541

IIB0 (Interruption Instruction Bundle 0) 2:541

IIB1 (Interruption Instruction Bundle 1) 2:541

IIM (Interruption Immediate) 2:541

IIP (Interruption Instruction Pointer) 2:541

IIPA (Interruption Instruction Previous Address)  
 2:541

Implicit Prefetch 1:70

IMUL Instruction 4:207

IN Instruction 4:210

INC Instruction 4:212

In-flight Resources 2:19

INIT (Initialization Event) 2:96, 2:306, 2:635

Initialization Event (INIT) 2:96

INS Instruction 4:214

INSB Instruction 4:214

INSD Instruction 4:214

Instruction Breakpoint Register (IBR) 2:151,  
 2:152

Instruction Debug Faults 2:151

Instruction Dependencies 1:148

Instruction Encoding 1:38

Instruction Formats 3:293

SSE 4:483

Instruction Group 1:40

Instruction Level Parallelism 1:15

Instruction Pointer (IP) 1:27, 1:140

Instruction Scheduling 1:148, 1:150, 1:164

Instruction Serialization 2:18

Instruction Set Architecture (ISA) 1:7

Instruction Set Modes 1:110

Instruction Set Transition 1:14

Instruction Set Transitions 2:239, 2:596

Instruction Slot Mapping 1:38

Instruction Slots 1:38

INSW Instruction 4:214

INT (External Interrupt) 2:96

INT3 Instruction 4:217

- INTA (Interrupt Acknowledge) 2:130
  - Inter-processor Interrupt (IPI) 2:127
  - Interrupt Acknowledge Cycle 2:130
  - Interrupt Control Registers (CR16-27) 2:36
  - Interrupt Handler 2:537
  - Interrupt Handling 2:543
  - Interrupt Hash Address 2:41
  - Interrupt Instruction Bundle Registers (IIB0-1) 2:42
  - Interrupt Processor Status Register (IPSR) 2:36
  - Interrupt Register State 2:540
  - Interrupt Registers 2:538
  - Interrupt Status Register (ISR) 2:36
  - Interrupt Vector 2:165
    - Alternate Data TLB 2:178
    - Alternate Instruction TLB 2:177
    - Break Instruction 2:185
    - Data Access Rights 2:191
    - Data Access-Bit 2:184
    - Data Key Miss 2:181
    - Data Nested TLB 2:179
    - Data TLB 2:176
    - Debug 2:200
    - Dirty-Bit 2:182
    - Disabled FP-Register 2:195
    - External Interrupt 2:186
    - Floating-point Fault 2:203
    - Floating-point Trap 2:204
    - General Exception 2:192
    - IA-32 Exception 2:210
    - IA-32 Intercept 2:211
    - IA-32 Interrupt 2:212
    - Instruction Access Rights 2:190
    - Instruction Access-Bit 2:183
    - Instruction Key Miss 2:180
    - Instruction TLB 2:175
    - Key Permission 2:189
    - Lower-Privilege Transfer Trap 2:205
    - NaT Consumption 2:196
    - Page Not Present 2:188
    - Single Step Trap 2:208
    - Speculation 2:198
    - Taken Branch Trap 2:207
    - Unaligned Reference 2:201
    - Unsupported Data Reference 2:202
    - Virtual External Interrupt 2:187
    - Virtualization 2:209
  - Interrupt Vector Address 2:35, 2:538
  - Interrupt Vector Table 2:538
  - Interruptions 2:95, 2:537
  - Interrupts 2:96, 2:114
    - External Interrupt Architecture 2:603
  - Interval Time Counter (ITC) 1:31
  - Interval Timer Match Register (ITM) 2:32
  - Interval Timer Offset (ITO) 2:34
  - Interval Timer Vector (ITV) 2:125
  - INTn Instruction 4:217
  - INTO Instruction 4:217
  - invala Instruction 3:146
  - INVD instructions 4:228
  - INVLPG Instruction 4:230
  - IP (Instruction Pointer) 1:27, 1:140
  - IPI (Inter-processor Interrupt) 2:127
  - IPSR (Interrupt Processor Status Register) 2:36, 2:541
  - IRET Instruction 4:231
  - IRETD Instruction 4:231
  - IRR (External Interrupt Request Registers) 2:125
  - ISR (Interrupt Status Register) 2:36, 2:165, 2:541
  - Itanium Architecture 1:7
  - Itanium Instruction Set 1:21
  - Itanium System Architecture 1:20
  - Itanium System Environment 1:7, 1:21
  - ITC (Interval Time Counter) 1:31, 2:32
  - itc Instruction 3:147
  - ITIR (Interrupt TLB Insertion Register) 2:541
  - ITM (Interval Time Match Register) 2:32
  - ITO (Interval Timer Offset) 2:34
  - itr Instruction 3:149
  - ITV (Interval Timer Vector) 2:125
  - IVA (Interrupt Vector Address) 2:35, 2:538
  - IVA-based interruptions 2:95, 2:537
  - IVR (External Interrupt Vector Register) 2:123
- ## J
- Jcc Instruction 4:239
  - JMP Instruction 4:243
  - JMPE Instruction 1:111, 2:597, 4:249
- ## K
- Kernel Register (KR) 1:29
  - KR (Kernel Register) 1:29
- ## L
- LAHF Instruction 4:251
  - Lamport's Algorithm 2:530
  - LAR Instruction 4:252
  - Large Constants 1:53
  - LC (Loop Count Register) 1:33
  - ld Instruction 3:151
  - ldf Instruction 3:157
  - ldfp Instruction 3:161
  - LDMXCSR Instruction 4:516
  - LDS Instruction 4:255
  - LEA Instruction 4:258
  - LEAVE Instruction 4:260
  - LES Instruction 4:255
  - lfetch Instruction 3:164
  - LFS Instruction 4:255
  - LGDT Instruction 4:264

LGS Instruction 4:255  
 LIDT Instruction 4:264  
 LLDT Instruction 4:267  
 LMSW Instruction 4:270  
 Load Instructions 1:58  
 loadrs Instruction 3:167  
 Loads from Memory 1:147  
 Local Redirection Registers (LRR0-1) 2:126  
 Locality Hints 1:70  
 LOCK Instruction 4:272  
 LODS Instruction 4:274  
 LODSB Instruction 4:274  
 LODSD Instruction 4:274  
 LODSW Instruction 4:274  
 Logical Instructions 1:51  
 Loop Count Register (LC) 1:33  
 LOOP Instruction 4:276  
 Loop Optimization 1:160, 1:181  
 LOOPcc Instruction 4:276  
 Lower Privilege Transfer Trap 2:151  
 LRR (Local Redirection Registers) 2:126  
 LSL Instruction 4:278  
 LSS Instruction 4:255  
 LTR Instruction 4:282

**M**

Machine Check (MC) 2:95, 2:296, 2:351  
 Machine Check Abort (MCA) 2:632  
 MASKMOVQ Instruction 4:576  
 MAXPS Instruction 4:519  
 MAXSS Instruction 4:521  
 MC (Machine Check) 2:351  
 MCA (Machine Check Abort) 2:95, 2:296, 2:632  
 Memory 1:36
 

- Cacheable Page 2:77
- Memory Access 1:142
- Memory Access Ordering 1:73
- Memory Attribute Transition 2:88
- Memory Attributes 2:75, 2:524
- Memory Consistency 1:72
- Memory Fences 2:510
- Memory Instructions 1:57
- Memory Management 2:561
- Memory Ordering 2:507, 2:510
  - IA-32 2:525
- Memory Reference 1:147
- Memory Regions 2:561
- Memory Synchronization 2:526

 mf Instruction 2:510, 2:526, 3:168
 

- mf.a 2:615

 MINPS Instruction 4:523  
 MINSS Instruction 4:525  
 mix Instruction 3:169  
 MMX technology 1:20  
 MOV Instruction 4:284  
 mov Instruction 3:172

MOVAPS Instruction 4:527  
 MOVD Instruction 4:401  
 MOVHLPS Instruction 4:529  
 MOVHPS Instruction 4:530  
 movl Instruction 3:187  
 MOVLHPS Instruction 4:532  
 MOVLPS Instruction 4:533  
 MOVMSKPS Instruction 4:535  
 MOVNTPS Instruction 4:578  
 MOVNTQ Instruction 4:579  
 MOVQ Instruction 4:403  
 MOVS Instruction 4:292  
 MOVSB Instruction 4:292  
 MOVSD Instruction 4:292  
 MOVSS Instruction 4:536  
 MOVSW Instruction 4:292  
 MOVSX Instruction 4:294  
 MOVUPS Instruction 4:538  
 MOVZX Instruction 4:295  
 MP Coherence 2:507  
 mpy4 Instruction 3:188  
 mpyshl4 Instruction 3:189  
 MUL Instruction 4:297  
 MULPS Instruction 4:540  
 MULSS Instruction 4:541  
 Multimedia Instructions 1:79  
 Multimedia Support 1:20  
 Multi-threading 1:177  
 Multiway Branches 1:173  
 mux Instruction 3:190

**N**

NaT (Not a Thing) 1:155  
 NaTPage (Not a Thing Attribute) 2:86  
 NaTVal (Not a Thing Value) 1:26  
 NEG Instruction 4:299  
 NMI (Non-Maskable Interrupt) 2:96  
 Non-Maskable Interrupt (NMI) 2:96  
 NOP Instruction 4:301  
 nop Instruction 3:193  
 Not A Thing (NaT) 1:155  
 Not a Thing Attribute (NaTPage) 2:86  
 Not a Thing Value (NatVal) 1:26  
 NOT Instruction 4:302

**O**

OLR (On Line Replacement) 2:351  
 Operating Environments 1:14  
 Operating System - See OS (Operating System)  
 OR Instruction 4:304  
 or Instruction 3:194  
 ORPS Instruction 4:542  
 OS (Operating System)
 

- Boot Flow Sample Code 2:639
- Boot Sequence 2:625
- FPSWA handler 2:587



- Illegal Dependency Fault 2:584
  - Long Branch Emulation 2:585
  - Multiple Address Spaces 1:20, 2:562
  - OS\_BOOT Entrypoint 2:283
  - OS\_INIT Entrypoint 2:283
  - OS\_MCA Entrypoint 2:283
  - OS\_RENDEZ Entrypoint 2:283
  - Performance Monitoring Support 2:620
  - Single Address Space 1:20, 2:565
  - Unaligned Reference Handler 2:583
  - Unsupported Data Reference Handler 2:584
  - OUT Instruction 4:306
  - OUTS Instruction 4:308
  - OUTSB Instruction 4:308
  - OUTSD Instruction 4:308
  - OUTSW Instruction 4:308
- P**
- pack Instruction 3:195
  - PACKSSDW Instruction 4:405
  - PACKSSWB Instruction 4:405
  - PACKUSWB Instruction 4:408
  - padd Instruction 3:197
  - PADDB Instruction 4:410
  - PADDD Instruction 4:410
  - PADDSB Instruction 4:413
  - PADDSW Instruction 4:413
  - PADDUSB Instruction 4:416
  - PADDUSW Instruction 4:416
  - PADDW Instruction 4:410
  - Page Access Rights 2:56
  - Page Sizes 2:57
  - Page Table Address 2:35
  - PAL (Processor Abstraction Layer) 1:7, 1:21, 2:279, 2:351
    - PAL Entrypoints 2:282
    - PAL Initialization 2:306
    - PAL Intercepts 2:351
    - PAL Intercepts in Virtual Environment 2:332
    - PAL Procedure Calls 2:628
    - PAL Procedures 2:353
    - PAL Self-test Control Word 2:295
    - PAL Virtualization 2:324
    - PAL Virtualization Optimizations 2:335
    - PAL Virtualization Services 2:486
    - PAL Virtualization Disables 2:346
    - PAL\_A 2:283
    - PAL\_B 2:283
    - PAL\_BRAND\_INFO 2:366
    - PAL\_BUS\_GET\_FEATURES 2:367
    - PAL\_BUS\_SET\_FEATURES 2:369
    - PAL\_CACHE\_FLUSH 2:370
    - PAL\_CACHE\_INFO 2:374
    - PAL\_CACHE\_INIT 2:376
    - PAL\_CACHE\_LINE\_INIT 2:377
    - PAL\_CACHE\_PROT\_INFO 2:378
    - PAL\_CACHE\_READ 2:380
    - PAL\_CACHE\_SHARED\_INFO 2:382
    - PAL\_CACHE\_SUMMARY 2:384
    - PAL\_CACHE\_WRITE 2:385
    - PAL\_COPY\_INFO 2:388
    - PAL\_COPY\_PAL 2:389
    - PAL\_DEBUG\_INFO 2:390
    - PAL\_FIXED\_ADDR 2:391
    - PAL\_FREQ\_BASE 2:392
    - PAL\_FREQ\_RATIOS 2:393
    - PAL\_GET\_HW\_POLICY 2:394
    - PAL\_GET\_PSTATE 2:320, 2:396, 2:637
    - PAL\_HALT 2:314
    - PAL\_HALT\_INFO 2:401
    - PAL\_HALT\_LIGHT 2:314, 2:403
    - PAL\_LOGICAL\_TO\_PHYSICAL 2:404
    - PAL\_MC\_CLEAR\_LOG 2:407
    - PAL\_MC\_DRAIN 2:408
    - PAL\_MC\_DYNAMIC\_STATE 2:409
    - PAL\_MC\_ERROR\_INFO 2:410
    - PAL\_MC\_ERROR\_INJECT 2:421
    - PAL\_MC\_EXPECTED 2:434
    - PAL\_MC\_HW\_TRACKING 2:432
    - PAL\_MC\_RESUME 2:436
    - PAL\_MEM\_ATTRIB 2:437
    - PAL\_MEMORY\_BUFFER 2:438
    - PAL\_PERF\_MON\_INFO 2:440
    - PAL\_PLATFORM\_ADDR 2:442
    - PAL\_PMI\_ENTRYPOINT 2:443
    - PAL\_PREFETCH\_VISIBILITY 2:444
    - PAL\_PROC\_GET\_FEATURES 2:446
    - PAL\_PROC\_SET\_FEATURES 2:450
    - PAL\_PSTATE\_INFO 2:319, 2:451
    - PAL\_PTCE\_INFO 2:453
    - PAL\_REGISTER\_INFO 2:454
    - PAL\_RSE\_INFO 2:455
    - PAL\_SET\_HW\_POLICY 2:456
    - PAL\_SET\_PSTATE 2:319, 2:458, 2:637
    - PAL\_SHUTDOWN 2:460
    - PAL\_TEST\_INFO 2:461
    - PAL\_TEST\_PROC 2:462
    - PAL\_VERSION 2:465
    - PAL\_VM\_INFO 2:466
    - PAL\_VM\_PAGE\_SIZE 2:467
    - PAL\_VM\_SUMMARY 2:468
    - PAL\_VM\_TR\_READ 2:470
    - PAL\_VP\_CREATE 2:471
    - PAL\_VP\_ENV\_INFO 2:473
    - PAL\_VP\_EXIT\_ENV 2:475
    - PAL\_VP\_INFO 2:476
    - PAL\_VP\_INIT\_ENV 2:478
    - PAL\_VP\_REGISTER 2:481
    - PAL\_VP\_RESTORE 2:483
    - PAL\_VP\_SAVE 2:484
    - PAL\_VP\_TERMINATE 2:485
    - PAL\_VPS\_RESTORE 2:499

- PAL\_VPS\_RESUME\_HANDLER 2:492
- PAL\_VPS\_RESUME\_NORMAL 2:489
- PAL\_VPS\_SAVE 2:500
- PAL\_VPS\_SET\_PENDING\_INTERRUPT 2:495
- PAL\_VPS\_SYNC\_READ 2:493
- PAL\_VPS\_SYNC\_WRITE 2:494
- PAL\_VPS\_THASH 2:497
- PAL\_VPS\_TTAG 2:498
- PAL-based Interruptions 2:95, 2:537
- PALE\_CHECK 2:282, 2:296
- PALE\_INIT 2:282, 2:306
- PALE\_PMI 2:282, 2:310
- PALE\_RESET 2:282, 2:289
- PAND Instruction 4:419
- PANDN Instruction 4:421
- Parallel Arithmetic 1:79
- Parallel Compares 1:172
- Parallel Shifts 1:81
- pavg Instruction 3:201
- PAVGB Instruction 4:563
- pavgsub Instruction 3:204
- PAVGW Instruction 4:563
- pcmp Instruction 3:206
- PCMPEQB Instruction 4:423
- PCMPEQD Instruction 4:423
- PCMPEQW Instruction 4:423
- PCMPGTB Instruction 4:426
- PCMPGTD Instruction 4:426
- PCMPGTW Instruction 4:426
- Performance Monitor Data Register (PMD) 1:33
- Performance Monitor Events 2:162
- Performance Monitoring 2:155, 2:619
- Performance Monitoring Vector 2:126
- PEXTRW Instruction 4:565
- PFS (Previous Function State Register) 1:32
- Physical Addressing 2:73
- PIB (Processor Interrupt Block) 2:127
- PINSRW Instruction 4:566
- PKR (Protection Key Register) 2:564
- Platform Management Interrupt (PMI) 2:96, 2:310, 2:538, 2:637
- PMADDWD Instruction 4:429
- pmax Instruction 3:209
- PMAXSW Instruction 4:567
- PMAXUB Instruction 4:568
- PMC (Performance Monitor Configuration) 2:155
- PMD (Performance Monitor Data Register) 1:33
- PMD (Performance Monitor Data) 2:155
- PMI (Platform Management Interrupt) 2:96, 2:310, 2:538, 2:637
- pmin Instruction 3:211
- PMINSW Instruction 4:569
- PMINUB Instruction 4:570
- PMOVMKB Instruction 4:571
- pmpy Instruction 3:213
- pmpyshr Instruction 3:214
- PMULHUW Instruction 4:572
- PMULHW Instruction 4:431
- PMULLW Instruction 4:433
- PMV (Performance Monitoring Vector) 2:126
- POP Instruction 4:311
- POPA Instruction 4:315
- POPAD Instruction 4:315
- popcnt Instruction 3:216
- POPF Instruction 4:317
- POPFD Instruction 4:317
- POR Instruction 4:435
- Power Management 2:313
- Power-on Event 2:351
- PR (Predicate Register) 1:26, 1:140
- Predicate Register (PR) 1:26, 1:140
- Predication 1:17, 1:54, 1:143, 1:163, 1:164
- Prefetch Hints 1:176
- PREFETCH Instruction 4:580
- Preserved Values 2:351
- Previous Function State (PFS) 1:32
- Privilege Level Transfer 1:84
- Privilege Levels 2:17
- probe Instruction 3:217
- Procedure Calls 2:549
- Processor Abstraction Layer - See PAL (Processor Abstraction Layer)
- Processor Abstraction Layer (PAL) 2:279
- Processor Boot Flow 2:623
- Processor Identification Registers (CPUID) 1:34
- Processor Interrupt Block (PIB) 2:127
- Processor Min-state Save Area 2:302
- Processor Reset 2:95
- Processor State Parameter (PSP) 2:299, 2:308
- Processor Status Register (PSR) 2:23
- Programmed I/O 2:534
- Protection Keys 2:59, 2:564
- psad Instruction 3:220
- PSADBW Instruction 4:573
- Pseudo-Code Functions 3:281
- pshl Instruction 3:222
- pshladd Instruction 3:223
- pshr Instruction 3:224
- pshradd Instruction 3:226
- PSHUFW Instruction 4:575
- PSLLD Instruction 4:437
- PSLLQ Instruction 4:437
- PSLLW Instruction 4:437
- PSP (Processor State Parameter) 2:308
- PSR (Processor Status Register) 2:23
- PSRAD Instruction 4:440
- PSRAW Instruction 4:440
- PSRLD Instruction 4:443
- PSRLQ Instruction 4:443
- PSRLW Instruction 4:443
- psub Instruction 3:227
- PSUBB Instruction 4:446



PSUBD Instruction 4:446  
 PSUBSB Instruction 4:449  
 PSUBSW Instruction 4:449  
 PSUBUSB Instruction 4:452  
 PSUBUSW Instruction 4:452  
 PSUBW Instruction 4:446  
 PTA (Page Table Address Register) 2:35  
 ptc Instruction  
     ptc.e 2:569, 3:230  
     ptc.g 2:570, 3:231  
     ptc.ga 2:570, 3:231  
     ptc.l 2:568, 3:233  
 ptr Instruction 3:234  
 PUNPCKHBW Instruction 4:455  
 PUNPCKHDQ Instruction 4:455  
 PUNPCKHWD Instruction 4:455  
 PUNPCKLBW Instruction 4:458  
 PUNPCKLDQ Instruction 4:458  
 PUNPCKLWD Instruction 4:458  
 PUSH Instruction 4:320  
 PUSHA Instruction 4:323  
 PUSHAD Instruction 4:323  
 PUSHF Instruction 4:325  
 PUSHFD Instruction 4:325  
 PXOR Instruction 4:461

## R

RAW Dependency 1:149  
 RCL Instruction 4:327  
 RCPPS Instruction 4:543  
 RCPSS Instruction 4:545  
 RCR Instruction 4:327  
 RDMSR Instruction 4:331  
 RDPMC Instruction 4:333  
 RDTSC Instruction 4:335  
 Read-after-write Dependency 1:149  
 Recoverable Error 2:351  
 Recovery Code 1:153, 1:154, 1:156  
 Region Identifier (RID) 2:561  
 Region Register (RR) 2:58, 2:561  
 Register File Transfers 1:82  
 Register Rotation 1:19, 1:185  
 Register Spill and Fill 1:62  
 Register Stack 1:18, 1:47  
 Register Stack Configuration Register (RSC) 1:29  
 Register Stack Engine (RSE) 1:144, 2:133  
 Register State 2:549  
 Release Semantics 2:507  
 Rendezvous 2:301  
 REP Instruction 4:337  
 REPE Instruction 4:337  
 REPNE Instruction 4:337  
 REPNZ Instruction 4:337  
 REPZ Instruction 4:337  
 Reserved Variables 2:351  
 Reset Event 2:95, 2:351  
 Resource Utilization Counter (RUC) 1:31, 2:33  
 RET Instruction 4:340  
 rfi Instruction 2:543, 3:236  
 RID (Region Identifier) 2:561  
 RNAT(RSE NaT Collection Register) 1:30  
 ROL Instruction 4:327  
 ROR Instruction 4:327  
 Rotating Registers 1:145  
 RR (Region Register) 2:58, 2:561  
 RSC (Register Stack Configuration Register) 1:29  
 RSE (Register Stack Engine) 2:133  
 RSE Backing Store Pointer (BSP) 1:29  
 RSE Backing Store Pointer for Memory Stores (BSPSTORE) 1:30  
 RSE NaT Collection Register (RNAT) 1:30  
 RSM Instruction 4:346  
 rsm Instruction 3:239  
 RSQRTPS Instruction 4:547  
 RSQRTSS Instruction 4:548  
 RUC (Resource Utilization Counter) 1:31, 2:33  
 rum Instruction 3:241

## S

SAHF Instruction 4:347  
 SAL (System Abstraction Layer) 1:7, 1:21, 2:352, 2:630  
     SAL\_B 2:283  
     SALE\_ENTRY 2:282, 2:291, 2:305  
     SALE\_PMI 2:282, 2:310  
 SAL Instruction 4:348  
 SAR Instruction 4:348  
 SBB Instruction 4:352  
 SCAS Instruction 4:354  
 SCASB Instruction 4:354  
 SCASD Instruction 4:354  
 SCASW Instruction 4:354  
 Scratch Register 2:352  
 Self Test State Parameter 2:293  
 Self-modifying Code 2:532  
 Semaphore Instructions 1:59  
 Semaphores 2:508  
 Serialization 2:17, 2:537  
 SETcc Instruction 4:356  
 setf Instruction 3:242  
 SFENCE Instruction 4:581  
 SGDT Instruction 4:359  
 SHL Instruction 4:348  
 shl Instruction 3:244  
 shladd Instruction 3:245  
 shladdp4 Instruction 3:246  
 SHLD Instruction 4:362  
 SHR Instruction 4:348  
 shr Instruction 3:247  
 SHRD Instruction 4:364  
 shrp Instruction 3:248  
 SHUFPS Instruction 4:549

SIDT Instruction 4:359  
 Single Step Trap 2:151  
 SLDT Instruction 4:367  
 SMSW Instruction 4:369  
 Software Pipelining 1:19, 1:75, 1:145, 1:181  
 Speculation 1:16, 1:142, 1:151  
     Control Speculation 1:16  
     Data Speculation 1:17  
     Recovery Code 1:17, 2:580  
     Speculation Check 1:156  
 SQRTPS Instruction 4:551  
 SQRTPSS Instruction 4:552  
 srlz Instruction 3:249  
 SSE Instructions 4:463  
 ssm Instruction 3:250  
 st Instruction 3:251  
 Stacked Calling Convention 2:352  
 Stacked General Registers 2:550  
 Stacked Registers 1:144  
 Static Calling Convention 2:352  
 Static General Registers 2:550  
 STC Instruction 4:371  
 STD Instruction 4:372  
 stf Instruction 3:254  
 STI Instruction 4:373  
 STMXCSR Instruction 4:553  
 Stops 1:38  
 Store Instructions 1:59  
 Stores to Memory 1:147  
 STOS Instruction 4:376  
 STOSB Instruction 4:376  
 STOSD Instruction 4:376  
 STOSW Instruction 4:376  
 STR Instruction 4:378  
 SUB Instruction 4:379  
 sub Instruction 3:256  
 SUBPS Instruction 4:554  
 SUBSS Instruction 4:555  
 sum Instruction 3:257  
 sxt Instruction 3:258  
 sync Instruction 3:259  
     sync.i 2:526  
 System Abstraction Layer - See SAL (System Abstraction Layer)  
 System Architecture 1:20  
 System Environment 2:13  
 System Programmer's Guide 2:501  
 System State 2:20

**T**

tak Instruction 3:260  
 Taken Branch trap 2:151  
 Task Priority Register (TPR) 2:123, 2:605  
 tbit Instruction 3:261  
 TC (Translation Cache) 2:49, 2:567

Template Field Encoding 1:38  
 Templates 1:141  
 TEST Instruction 4:381  
 tf Instruction 3:263  
 thash Instruction 3:265  
 TLB (Translation Lookaside Buffer) 2:47, 2:565  
 tnat Instruction 3:266  
 tpa Instruction 3:268  
 TPR (Task Priority Register) 2:123, 2:605  
 TR (Translation Register) 2:48, 2:566  
 Translation Cache (TC) 2:49, 2:567  
     purge 2:568  
 Translation Instructions 2:60  
 Translation Lookaside Buffer (TLB) 2:47, 2:565  
 Translation Register (TR) 2:48, 2:566  
 Traps 2:96, 2:537  
 ttag Instruction 3:269

**U**

UCOMISS Instruction 4:556  
 UD2 Instruction 4:383  
 UEFI (Unified Extensible Firmware Interface) 2:630  
 UM (User Mask Register) 1:33  
 UNAT (User NaT Collection Register) 1:31, 1:156  
 Uncacheable Page 2:77  
 Unchanged Register 2:352  
 Unordered Semantics 2:507  
 unpack Instruction 3:270  
 UNPCKHPS Instruction 4:558  
 UNPCKLPS Instruction 4:560  
 User Mask (UM) 1:33  
 User NaT Collection Register (UNAT) 1:31, 1:156

**V**

VERR Instruction 4:384  
 VERW Instruction 4:384  
 VHPT (Virtual Hash Page Table) 2:61, 2:571  
 VHPT Translation Vector 2:173  
 Virtual Addressing 2:45  
 Virtual Hash Page Table (VHPT) 2:61, 2:571  
 Virtual Machine Monitor (VMM) 2:352  
 Virtual Processor Descriptor (VPD) 2:325, 2:352  
 Virtual Processor State 2:352  
 Virtual Processor Status Register (VPSR) 2:327  
 Virtual Region Number (VRN) 2:561  
 Virtualization 2:44, 2:324  
 Virtualization Acceleration Control (vac) 2:329  
 Virtualization Disable Control (vdc) 2:329  
 VMM (Virtual Machine Monitor) 2:352  
 vmw Instruction 3:273  
 VPD (Virtual Processor Descriptor) 2:325, 2:352  
 VPSR (Virtual Processor Status Register) 2:327  
 VRN (Virtual Region Number) 2:561

**W**

WAIT Instruction 4:386  
WAR Dependency 1:149  
WAW Dependency 1:149  
WBINVD Instruction 4:387  
Write-after-read Dependency 1:149  
Write-after-write Dependency 1:149  
WRMSR Instruction 4:389

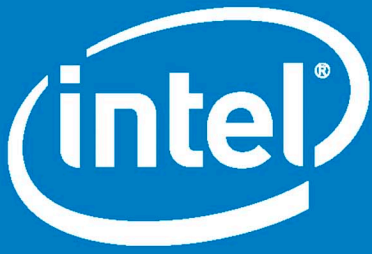
**X**

XADD Instruction 4:391  
XCHG Instruction 4:393  
xchg Instruction 2:508, 3:274  
XLAT Instruction 4:395  
XLATB Instruction 4:395  
xma Instruction 3:276  
xmpy Instruction 3:278  
XOR Instruction 4:397  
xor Instruction 3:279  
XORPS Instruction 4:562  
XTP (External Task Priority Cycle) 2:130  
XTPR (External Task Priority Register) 2:605

**Z**

zxt Instruction 3:280





Copyright ©1999-2010 Intel Corporation. All rights reserved.  
Intel, the Intel logo, Intel Inside, and Itanium are trademarks or  
registered trademarks of Intel Corporation or its subsidiaries  
in the United States and other countries.

Other names and brands may be claimed as the property of others.  
0510/FL/DS/NOD/RRD/2K 245318-006US