

Intel[®] 64 and IA-32 Architectures Optimization Reference Manual

Order Number: 248966-023
January 2011

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.

UNLESS OTHERWISE AGREED IN WRITING BY INTEL, THE INTEL PRODUCTS ARE NOT DESIGNED NOR INTENDED FOR ANY APPLICATION IN WHICH THE FAILURE OF THE INTEL PRODUCT COULD CREATE A SITUATION WHERE PERSONAL INJURY OR DEATH MAY OCCUR.

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. The information here is subject to change without notice. Do not finalize a design with this information.

The Intel® 64 architecture processors 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.

Hyper-Threading Technology requires a computer system with an Intel® processor supporting Hyper-Threading Technology and an HT Technology enabled chipset, BIOS and operating system. Performance will vary depending on the specific hardware and software you use. For more information, see <http://www.intel.com/technology/hyperthread/index.htm>; including details on which processors support HT Technology.

Intel® Virtualization Technology requires a computer system with an enabled Intel® processor, BIOS, virtual machine monitor (VMM) and for some uses, certain platform software enabled for it. Functionality, performance or other benefits will vary depending on hardware and software configurations. Intel® Virtualization Technology-enabled BIOS and VMM applications are currently in development.

64-bit computing on Intel architecture requires a computer system with a processor, chipset, BIOS, operating system, device drivers and applications enabled for Intel® 64 architecture. Processors will not operate (including 32-bit operation) without an Intel® 64 architecture-enabled BIOS. Performance will vary depending on your hardware and software configurations. Consult with your system vendor for more information.

Intel, Pentium, Intel Atom, Intel Centrino, Intel Centrino Duo, Intel Xeon, Intel NetBurst, Intel Core, Intel Core Solo, Intel Core Duo, Intel Core 2 Duo, Intel Core 2 Extreme, Intel Pentium D, Itanium, Intel SpeedStep, MMX, and VTune 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.

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>

Copyright © 1997-2011 Intel Corporation

CHAPTER 1 INTRODUCTION

1.1	TUNING YOUR APPLICATION.....	1-1
1.2	ABOUT THIS MANUAL.....	1-2
1.3	RELATED INFORMATION.....	1-4

CHAPTER 2 INTEL® 64 AND IA-32 PROCESSOR ARCHITECTURES

2.1	INTEL® MICROARCHITECTURE CODE NAME SANDY BRIDGE	2-2
2.1.1	Intel® microarchitecture code name Sandy Bridge Pipeline Overview	2-3
2.1.2	The Front End	2-5
2.1.2.1	Legacy Decode Pipeline	2-6
2.1.2.2	Decoded ICache	2-9
2.1.2.3	Branch Prediction	2-10
2.1.2.4	Micro-op Queue and the Loop Stream Detector (LSD)	2-10
2.1.3	The Out-of-Order Engine	2-11
2.1.3.1	Renamer	2-12
2.1.3.2	Scheduler	2-13
2.1.4	The Execution Core	2-13
2.1.5	Cache Hierarchy	2-15
2.1.5.1	Load and Store Operation Overview	2-16
2.1.5.2	L1 DCache	2-17
2.1.5.3	Ring Interconnect and Last Level Cache	2-23
2.1.5.4	Data Prefetching	2-23
2.1.6	System Agent	2-25
2.2	INTEL® CORE™ MICROARCHITECTURE AND ENHANCED INTEL CORE MICROARCHITECTURE . 2-26	
2.2.1	Intel® Core™ Microarchitecture Pipeline Overview	2-27
2.2.2	Front End	2-29
2.2.2.1	Branch Prediction Unit	2-30
2.2.2.2	Instruction Fetch Unit	2-30
2.2.2.3	Instruction Queue (IQ)	2-31
2.2.2.4	Instruction Decode	2-32
2.2.2.5	Stack Pointer Tracker	2-32
2.2.2.6	Micro-fusion	2-33
2.2.3	Execution Core	2-33
2.2.3.1	Issue Ports and Execution Units	2-34
2.2.4	Intel® Advanced Memory Access	2-37
2.2.4.1	Loads and Stores	2-38
2.2.4.2	Data Prefetch to L1 caches	2-39
2.2.4.3	Data Prefetch Logic	2-39
2.2.4.4	Store Forwarding	2-40
2.2.4.5	Memory Disambiguation	2-41
2.2.5	Intel® Advanced Smart Cache	2-42
2.2.5.1	Loads	2-43
2.2.5.2	Stores	2-44
2.3	INTEL® MICROARCHITECTURE CODE NAME NEHALEM	2-45

CONTENTS

	PAGE
2.3.1	Microarchitecture Pipeline 2-45
2.3.2	Front End Overview 2-47
2.3.3	Execution Engine 2-49
2.3.3.1	Issue Ports and Execution Units 2-50
2.3.4	Cache and Memory Subsystem 2-51
2.3.5	Load and Store Operation Enhancements 2-52
2.3.5.1	Efficient Handling of Alignment Hazards 2-52
2.3.5.2	Store Forwarding Enhancement 2-53
2.3.6	REP String Enhancement 2-55
2.3.7	Enhancements for System Software 2-56
2.3.8	Efficiency Enhancements for Power Consumption 2-56
2.3.9	Hyper-Threading Technology Support in Intel® microarchitecture code name Nehalem2-57
2.4	INTEL NETBURST® MICROARCHITECTURE 2-57
2.4.1	Design Goals 2-58
2.4.2	Pipeline 2-58
2.4.2.1	Front End 2-60
2.4.2.2	Out-of-order Core 2-60
2.4.2.3	Retirement 2-61
2.4.3	Front End Pipeline Detail 2-61
2.4.3.1	Prefetching 2-61
2.4.3.2	Decoder 2-62
2.4.3.3	Execution Trace Cache 2-62
2.4.3.4	Branch Prediction 2-62
2.4.4	Execution Core Detail 2-63
2.4.4.1	Instruction Latency and Throughput 2-64
2.4.4.2	Execution Units and Issue Ports 2-64
2.4.4.3	Caches 2-66
2.4.4.4	Data Prefetch 2-67
2.4.4.5	Loads and Stores 2-69
2.4.4.6	Store Forwarding 2-70
2.5	INTEL® PENTIUM® M PROCESSOR MICROARCHITECTURE 2-70
2.5.1	Front End 2-71
2.5.2	Data Prefetching 2-72
2.5.3	Out-of-Order Core 2-73
2.5.4	In-Order Retirement 2-73
2.6	MICROARCHITECTURE OF INTEL® CORE™ SOLO AND INTEL® CORE™ DUO PROCESSORS 2-74
2.6.1	Front End 2-74
2.6.2	Data Prefetching 2-75
2.7	INTEL® HYPER-THREADING TECHNOLOGY 2-75
2.7.1	Processor Resources and HT Technology 2-76
2.7.1.1	Replicated Resources 2-77
2.7.1.2	Partitioned Resources 2-77
2.7.1.3	Shared Resources 2-77
2.7.2	Microarchitecture Pipeline and HT Technology 2-78
2.7.3	Front End Pipeline 2-78
2.7.4	Execution Core 2-78
2.7.5	Retirement 2-79
2.8	MULTICORE PROCESSORS 2-79
2.8.1	Microarchitecture Pipeline and MultiCore Processors 2-81
2.8.2	Shared Cache in Intel® Core™ Duo Processors 2-81
2.8.2.1	Load and Store Operations 2-82

	PAGE	
2.9	INTEL® 64 ARCHITECTURE.....	2-83
2.10	SIMD TECHNOLOGY.....	2-83
2.10.1	Summary of SIMD Technologies.....	2-86
2.10.1.1	MMX™ Technology.....	2-86
2.10.1.2	Streaming SIMD Extensions.....	2-86
2.10.1.3	Streaming SIMD Extensions 2.....	2-87
2.10.1.4	Streaming SIMD Extensions 3.....	2-87
2.10.1.5	Supplemental Streaming SIMD Extensions 3.....	2-87
2.10.1.6	SSE4.1.....	2-88
2.10.1.7	SSE4.2.....	2-88
2.10.1.8	AESNI and PCLMULQDQ.....	2-89
2.10.1.9	Intel® Advanced Vector Extensions.....	2-89

CHAPTER 3

GENERAL OPTIMIZATION GUIDELINES

3.1	PERFORMANCE TOOLS.....	3-1
3.1.1	Intel® C++ and Fortran Compilers.....	3-1
3.1.2	General Compiler Recommendations.....	3-2
3.1.3	VTune™ Performance Analyzer.....	3-2
3.2	PROCESSOR PERSPECTIVES.....	3-3
3.2.1	CPUID Dispatch Strategy and Compatible Code Strategy.....	3-4
3.2.2	Transparent Cache-Parameter Strategy.....	3-5
3.2.3	Threading Strategy and Hardware Multithreading Support.....	3-5
3.3	CODING RULES, SUGGESTIONS AND TUNING HINTS.....	3-5
3.4	OPTIMIZING THE FRONT END.....	3-6
3.4.1	Branch Prediction Optimization.....	3-6
3.4.1.1	Eliminating Branches.....	3-7
3.4.1.2	Spin-Wait and Idle Loops.....	3-9
3.4.1.3	Static Prediction.....	3-9
3.4.1.4	Inlining, Calls and Returns.....	3-11
3.4.1.5	Code Alignment.....	3-12
3.4.1.6	Branch Type Selection.....	3-13
3.4.1.7	Loop Unrolling.....	3-15
3.4.1.8	Compiler Support for Branch Prediction.....	3-16
3.4.2	Fetch and Decode Optimization.....	3-17
3.4.2.1	Optimizing for Micro-fusion.....	3-17
3.4.2.2	Optimizing for Macro-fusion.....	3-18
3.4.2.3	Length-Changing Prefixes (LCP).....	3-22
3.4.2.4	Optimizing the Loop Stream Detector (LSD).....	3-24
3.4.2.5	Exploit LSD Micro-op Emission Bandwidth in Intel® microarchitecture code name Sandy Bridge.....	3-25
3.4.2.6	Optimization for Decoded ICache.....	3-26
3.4.2.7	Scheduling Rules for the Pentium M Processor Decoder.....	3-27
3.4.2.8	Other Decoding Guidelines.....	3-28
3.5	OPTIMIZING THE EXECUTION CORE.....	3-28
3.5.1	Instruction Selection.....	3-28
3.5.1.1	Use of the INC and DEC Instructions.....	3-29
3.5.1.2	Integer Divide.....	3-30
3.5.1.3	Using LEA.....	3-30
3.5.1.4	ADC and SBB in Intel® microarchitecture code name Sandy Bridge.....	3-32
3.5.1.5	Bitwise Rotation.....	3-33

	PAGE	
3.5.1.6	Address Calculations	3-34
3.5.1.7	Clearing Registers and Dependency Breaking Idioms	3-34
3.5.1.8	Compares	3-36
3.5.1.9	Using NOPs	3-37
3.5.1.10	Mixing SIMD Data Types	3-38
3.5.1.11	Spill Scheduling	3-38
3.5.2	Avoiding Stalls in Execution Core	3-39
3.5.2.1	ROB Read Port Stalls	3-39
3.5.2.2	Writeback Bus Conflicts	3-40
3.5.2.3	Bypass between Execution Domains	3-41
3.5.2.4	Partial Register Stalls	3-41
3.5.2.5	Partial XMM Register Stalls	3-42
3.5.2.6	Partial Flag Register Stalls	3-44
3.5.2.7	Floating Point/SIMD Operands	3-45
3.5.3	Vectorization	3-46
3.5.4	Optimization of Partially Vectorizable Code	3-47
3.5.4.1	Alternate Packing Techniques	3-49
3.5.4.2	Simplifying Result Passing	3-50
3.5.4.3	Stack Optimization	3-51
3.5.4.4	Tuning Considerations	3-52
3.6	OPTIMIZING MEMORY ACCESSES	3-53
3.6.1	Load and Store Execution Bandwidth	3-54
3.6.1.1	Make Use of Load Bandwidth in Intel® microarchitecture code name Sandy Bridge	3-54
3.6.1.2	L1D Cache Latency in Intel® microarchitecture code name Sandy Bridge	3-56
3.6.1.3	Handling L1D Cache Bank Conflict	3-56
3.6.2	Minimize Register Spills	3-58
3.6.3	Enhance Speculative Execution and Memory Disambiguation	3-59
3.6.4	Alignment	3-60
3.6.5	Store Forwarding	3-62
3.6.5.1	Store-to-Load-Forwarding Restriction on Size and Alignment	3-63
3.6.5.2	Store-forwarding Restriction on Data Availability	3-68
3.6.6	Data Layout Optimizations	3-69
3.6.7	Stack Alignment	3-72
3.6.8	Capacity Limits and Aliasing in Caches	3-72
3.6.8.1	Capacity Limits in Set-Associative Caches	3-73
3.6.8.2	Aliasing Cases in the Pentium® M, Intel® Core™ Solo, Intel® Core™ Duo and Intel® Core™ 2 Duo Processors	3-74
3.6.9	Mixing Code and Data	3-75
3.6.9.1	Self-modifying Code	3-75
3.6.9.2	Position Independent Code	3-76
3.6.10	Write Combining	3-77
3.6.11	Locality Enhancement	3-78
3.6.12	Minimizing Bus Latency	3-79
3.6.13	Non-Temporal Store Bus Traffic	3-80
3.7	PREFETCHING	3-81
3.7.1	Hardware Instruction Fetching and Software Prefetching	3-81
3.7.2	Software and Hardware Prefetching in Prior Microarchitectures	3-82
3.7.3	Hardware Prefetching for First-Level Data Cache	3-83
3.7.4	Hardware Prefetching for Second-Level Cache	3-86
3.7.5	Cacheability Instructions	3-87
3.7.6	REP Prefix and Data Movement	3-87

	PAGE
3.8	FLOATING-POINT CONSIDERATIONS 3-90
3.8.1	Guidelines for Optimizing Floating-point Code 3-90
3.8.2	Floating-point Modes and Exceptions 3-92
3.8.2.1	Floating-point Exceptions 3-92
3.8.2.2	Dealing with floating-point exceptions in x87 FPU code 3-92
3.8.2.3	Floating-point Exceptions in SSE/SSE2/SSE3 Code 3-93
3.8.3	Floating-point Modes 3-94
3.8.3.1	Rounding Mode 3-94
3.8.3.2	Precision 3-96
3.8.3.3	Improving Parallelism and the Use of FXCH 3-97
3.8.4	x87 vs. Scalar SIMD Floating-point Trade-offs 3-97
3.8.4.1	Scalar SSE/SSE2 Performance on Intel® Core™ Solo and Intel® Core™ Duo Processors. 3-98
3.8.4.2	x87 Floating-point Operations with Integer Operands 3-99
3.8.4.3	x87 Floating-point Comparison Instructions 3-99
3.8.4.4	Transcendental Functions 3-99
3.9	MAXIMIZING PCIE PERFORMANCE 3-100

CHAPTER 4

CODING FOR SIMD ARCHITECTURES

4.1	CHECKING FOR PROCESSOR SUPPORT OF SIMD TECHNOLOGIES 4-1
4.1.1	Checking for MMX Technology Support 4-2
4.1.2	Checking for Streaming SIMD Extensions Support 4-2
4.1.3	Checking for Streaming SIMD Extensions 2 Support 4-3
4.1.4	Checking for Streaming SIMD Extensions 3 Support 4-3
4.1.5	Checking for Supplemental Streaming SIMD Extensions 3 Support 4-4
4.1.6	Checking for SSE4.1 Support 4-4
4.1.7	Checking for SSE4.2 Support 4-5
4.1.8	Detection of PCLMULQDQ and AESNI Instructions 4-5
4.1.9	Detection of AVX Instructions 4-6
4.1.10	Detection of VEX-Encoded AES and VPCLMULQDQ 4-8
4.2	CONSIDERATIONS FOR CODE CONVERSION TO SIMD PROGRAMMING 4-9
4.2.1	Identifying Hot Spots 4-11
4.2.2	Determine If Code Benefits by Conversion to SIMD Execution 4-11
4.3	CODING TECHNIQUES 4-12
4.3.1	Coding Methodologies 4-13
4.3.1.1	Assembly 4-14
4.3.1.2	Intrinsics 4-15
4.3.1.3	Classes 4-16
4.3.1.4	Automatic Vectorization 4-17
4.4	STACK AND DATA ALIGNMENT 4-18
4.4.1	Alignment and Contiguity of Data Access Patterns 4-18
4.4.1.1	Using Padding to Align Data 4-19
4.4.1.2	Using Arrays to Make Data Contiguous 4-19
4.4.2	Stack Alignment For 128-bit SIMD Technologies 4-20
4.4.3	Data Alignment for MMX Technology 4-21
4.4.4	Data Alignment for 128-bit data 4-21
4.4.4.1	Compiler-Supported Alignment 4-21
4.5	IMPROVING MEMORY UTILIZATION 4-23
4.5.1	Data Structure Layout 4-23
4.5.2	Strip-Mining 4-27

	PAGE	
4.5.3	Loop Blocking.....	4-28
4.6	INSTRUCTION SELECTION	4-30
4.6.1	SIMD Optimizations and Microarchitectures.....	4-32
4.7	TUNING THE FINAL APPLICATION.....	4-32

CHAPTER 5 OPTIMIZING FOR SIMD INTEGER APPLICATIONS

5.1	GENERAL RULES ON SIMD INTEGER CODE	5-2
5.2	USING SIMD INTEGER WITH X87 FLOATING-POINT	5-2
5.2.1	Using the EMMS Instruction.....	5-3
5.2.2	Guidelines for Using EMMS Instruction.....	5-3
5.3	DATA ALIGNMENT	5-4
5.4	DATA MOVEMENT CODING TECHNIQUES	5-6
5.4.1	Unsigned Unpack	5-6
5.4.2	Signed Unpack	5-7
5.4.3	Interleaved Pack with Saturation	5-8
5.4.4	Interleaved Pack without Saturation	5-10
5.4.5	Non-Interleaved Unpack.....	5-10
5.4.6	Extract Data Element	5-12
5.4.7	Insert Data Element	5-13
5.4.8	Non-Unit Stride Data Movement	5-14
5.4.9	Move Byte Mask to Integer	5-15
5.4.10	Packed Shuffle Word for 64-bit Registers	5-16
5.4.11	Packed Shuffle Word for 128-bit Registers	5-17
5.4.12	Shuffle Bytes	5-18
5.4.13	Conditional Data Movement	5-18
5.4.14	Unpacking/interleaving 64-bit Data in 128-bit Registers	5-18
5.4.15	Data Movement.....	5-19
5.4.16	Conversion Instructions	5-19
5.5	GENERATING CONSTANTS.....	5-19
5.6	BUILDING BLOCKS.....	5-20
5.6.1	Absolute Difference of Unsigned Numbers	5-20
5.6.2	Absolute Difference of Signed Numbers.....	5-21
5.6.3	Absolute Value	5-21
5.6.4	Pixel Format Conversion	5-22
5.6.5	Endian Conversion	5-24
5.6.6	Clipping to an Arbitrary Range [High, Low].....	5-25
5.6.6.1	Highly Efficient Clipping	5-26
5.6.6.2	Clipping to an Arbitrary Unsigned Range [High, Low]	5-27
5.6.7	Packed Max/Min of Byte, Word and Dword	5-28
5.6.8	Packed Multiply Integers	5-28
5.6.9	Packed Sum of Absolute Differences.....	5-28
5.6.10	MPSADBW and PHMINPOSUW	5-29
5.6.11	Packed Average (Byte/Word)	5-29
5.6.12	Complex Multiply by a Constant.....	5-30
5.6.13	Packed 64-bit Add/Subtract	5-30
5.6.14	128-bit Shifts.....	5-31
5.6.15	PTEST and Conditional Branch	5-31
5.6.16	Vectorization of Heterogeneous Computations across Loop Iterations	5-32
5.6.17	Vectorization of Control Flows in Nested Loops.....	5-33
5.7	MEMORY OPTIMIZATIONS	5-35

	PAGE	
5.7.1	Partial Memory Accesses	5-36
5.7.1.1	Supplemental Techniques for Avoiding Cache Line Splits.....	5-38
5.7.2	Increasing Bandwidth of Memory Fills and Video Fills.....	5-39
5.7.2.1	Increasing Memory Bandwidth Using the MOVDQ Instruction	5-39
5.7.2.2	Increasing Memory Bandwidth by Loading and Storing to and from the Same DRAM Page.....	5-40
5.7.2.3	Increasing UC and WC Store Bandwidth by Using Aligned Stores	5-40
5.7.3	Reverse Memory Copy.....	5-40
5.8	CONVERTING FROM 64-BIT TO 128-BIT SIMD INTEGERS.....	5-43
5.8.1	SIMD Optimizations and Microarchitectures	5-44
5.8.1.1	Packed SSE2 Integer versus MMX Instructions.....	5-44
5.8.1.2	Work-around for False Dependency Issue.....	5-45
5.9	TUNING PARTIALLY VECTORIZABLE CODE	5-46
5.10	PARALLEL MODE AES ENCRYPTION AND DECRYPTION	5-49
5.10.1	AES Counter Mode of Operation.....	5-50
5.10.2	AES Key Expansion Alternative.....	5-59

CHAPTER 6 OPTIMIZING FOR SIMD FLOATING-POINT APPLICATIONS

6.1	GENERAL RULES FOR SIMD FLOATING-POINT CODE	6-1
6.2	PLANNING CONSIDERATIONS.....	6-1
6.3	USING SIMD FLOATING-POINT WITH X87 FLOATING-POINT	6-2
6.4	SCALAR FLOATING-POINT CODE.....	6-2
6.5	DATA ALIGNMENT	6-3
6.5.1	Data Arrangement.....	6-3
6.5.1.1	Vertical versus Horizontal Computation.....	6-3
6.5.1.2	Data Swizzling	6-6
6.5.1.3	Data Deswizzling	6-9
6.5.1.4	Horizontal ADD Using SSE.....	6-10
6.5.2	Use of CVTTPS2PI/CVTTSS2SI Instructions	6-13
6.5.3	Flush-to-Zero and Denormals-are-Zero Modes	6-13
6.6	SIMD OPTIMIZATIONS AND MICROARCHITECTURES.....	6-14
6.6.1	SIMD Floating-point Programming Using SSE3.....	6-14
6.6.1.1	SSE3 and Complex Arithmetics.....	6-15
6.6.1.2	Packed Floating-Point Performance in Intel Core Duo Processor.....	6-18
6.6.2	Dot Product and Horizontal SIMD Instructions.....	6-18
6.6.3	Vector Normalization	6-21
6.6.4	Using Horizontal SIMD Instruction Sets and Data Layout.....	6-23
6.6.4.1	SOA and Vector Matrix Multiplication	6-26

CHAPTER 7 OPTIMIZING CACHE USAGE

7.1	GENERAL PREFETCH CODING GUIDELINES.....	7-1
7.2	HARDWARE PREFETCHING OF DATA	7-3
7.3	PREFETCH AND CACHEABILITY INSTRUCTIONS.....	7-4
7.4	PREFETCH	7-4
7.4.1	Software Data Prefetch	7-4
7.4.2	Prefetch Instructions - Pentium® 4 Processor Implementation	7-5
7.4.3	Prefetch and Load Instructions.....	7-6
7.5	CACHEABILITY CONTROL	7-7

- 7.5.1 The Non-temporal Store Instructions. 7-7
 - 7.5.1.1 Fencing. 7-7
 - 7.5.1.2 Streaming Non-temporal Stores 7-7
 - 7.5.1.3 Memory Type and Non-temporal Stores 7-8
 - 7.5.1.4 Write-Combining 7-8
- 7.5.2 Streaming Store Usage Models 7-9
 - 7.5.2.1 Coherent Requests 7-9
 - 7.5.2.2 Non-coherent requests. 7-9
- 7.5.3 Streaming Store Instruction Descriptions. 7-10
- 7.5.4 The Streaming Load Instruction. 7-10
- 7.5.5 FENCE Instructions. 7-11
 - 7.5.5.1 SFENCE Instruction. 7-11
 - 7.5.5.2 LFENCE Instruction. 7-11
 - 7.5.5.3 MFENCE Instruction 7-12
- 7.5.6 CLFLUSH Instruction 7-12
- 7.6 MEMORY OPTIMIZATION USING PREFETCH 7-13
 - 7.6.1 Software-Controlled Prefetch 7-13
 - 7.6.2 Hardware Prefetch 7-13
 - 7.6.3 Example of Effective Latency Reduction
with Hardware Prefetch. 7-14
 - 7.6.4 Example of Latency Hiding with S/W Prefetch Instruction 7-16
 - 7.6.5 Software Prefetching Usage Checklist 7-17
 - 7.6.6 Software Prefetch Scheduling Distance 7-18
 - 7.6.7 Software Prefetch Concatenation. 7-19
 - 7.6.8 Minimize Number of Software Prefetches 7-20
 - 7.6.9 Mix Software Prefetch with Computation Instructions 7-22
 - 7.6.10 Software Prefetch and Cache Blocking Techniques 7-23
 - 7.6.11 Hardware Prefetching and Cache Blocking Techniques 7-27
 - 7.6.12 Single-pass versus Multi-pass Execution. 7-28
- 7.7 MEMORY OPTIMIZATION USING NON-TEMPORAL STORES 7-31
 - 7.7.1 Non-temporal Stores and Software Write-Combining. 7-31
 - 7.7.2 Cache Management 7-32
 - 7.7.2.1 Video Encoder 7-32
 - 7.7.2.2 Video Decoder 7-32
 - 7.7.2.3 Conclusions from Video Encoder and Decoder Implementation. 7-33
 - 7.7.2.4 Optimizing Memory Copy Routines. 7-33
 - 7.7.2.5 TLB Priming 7-34
 - 7.7.2.6 Using the 8-byte Streaming Stores and Software Prefetch. 7-35
 - 7.7.2.7 Using 16-byte Streaming Stores and Hardware Prefetch. 7-35
 - 7.7.2.8 Performance Comparisons of Memory Copy Routines 7-37
 - 7.7.3 Deterministic Cache Parameters 7-38
 - 7.7.3.1 Cache Sharing Using Deterministic Cache Parameters 7-40
 - 7.7.3.2 Cache Sharing in Single-Core or Multicore 7-40
 - 7.7.3.3 Determine Prefetch Stride 7-40

**CHAPTER 8
MULTICORE AND HYPER-THREADING TECHNOLOGY**

- 8.1 PERFORMANCE AND USAGE MODELS 8-1
 - 8.1.1 Multithreading 8-2
 - 8.1.2 Multitasking Environment 8-3
- 8.2 PROGRAMMING MODELS AND MULTITHREADING 8-4

8.2.1	Parallel Programming Models.....	8-5
8.2.1.1	Domain Decomposition.....	8-5
8.2.2	Functional Decomposition.....	8-5
8.2.3	Specialized Programming Models.....	8-6
8.2.3.1	Producer-Consumer Threading Models.....	8-7
8.2.4	Tools for Creating Multithreaded Applications.....	8-10
8.2.4.1	Programming with OpenMP Directives.....	8-10
8.2.4.2	Automatic Parallelization of Code.....	8-10
8.2.4.3	Supporting Development Tools.....	8-11
8.2.4.4	Intel® Thread Checker.....	8-11
8.2.4.5	Intel® Thread Profiler.....	8-11
8.2.4.6	Intel® Threading Building Block.....	8-11
8.3	OPTIMIZATION GUIDELINES.....	8-11
8.3.1	Key Practices of Thread Synchronization.....	8-12
8.3.2	Key Practices of System Bus Optimization.....	8-12
8.3.3	Key Practices of Memory Optimization.....	8-12
8.3.4	Key Practices of Front-end Optimization.....	8-13
8.3.5	Key Practices of Execution Resource Optimization.....	8-13
8.3.6	Generality and Performance Impact.....	8-14
8.4	THREAD SYNCHRONIZATION.....	8-14
8.4.1	Choice of Synchronization Primitives.....	8-15
8.4.2	Synchronization for Short Periods.....	8-16
8.4.3	Optimization with Spin-Locks.....	8-18
8.4.4	Synchronization for Longer Periods.....	8-18
8.4.4.1	Avoid Coding Pitfalls in Thread Synchronization.....	8-19
8.4.5	Prevent Sharing of Modified Data and False-Sharing.....	8-21
8.4.6	Placement of Shared Synchronization Variable.....	8-21
8.5	SYSTEM BUS OPTIMIZATION.....	8-23
8.5.1	Conserve Bus Bandwidth.....	8-23
8.5.2	Understand the Bus and Cache Interactions.....	8-24
8.5.3	Avoid Excessive Software Prefetches.....	8-25
8.5.4	Improve Effective Latency of Cache Misses.....	8-25
8.5.5	Use Full Write Transactions to Achieve Higher Data Rate.....	8-26
8.6	MEMORY OPTIMIZATION.....	8-26
8.6.1	Cache Blocking Technique.....	8-27
8.6.2	Shared-Memory Optimization.....	8-27
8.6.2.1	Minimize Sharing of Data between Physical Processors.....	8-27
8.6.2.2	Batched Producer-Consumer Model.....	8-28
8.6.3	Eliminate 64-KByte Aliased Data Accesses.....	8-29
8.7	FRONT-END OPTIMIZATION.....	8-30
8.7.1	Avoid Excessive Loop Unrolling.....	8-30
8.8	AFFINITIES AND MANAGING SHARED PLATFORM RESOURCES.....	8-30
8.8.1	Topology Enumeration of Shared Resources.....	8-32
8.8.2	Non-Uniform Memory Access.....	8-32
8.9	OPTIMIZATION OF OTHER SHARED RESOURCES.....	8-35
8.9.1	Expanded Opportunity for HT Optimization.....	8-35

CHAPTER 9 64-BIT MODE CODING GUIDELINES

9.1	INTRODUCTION.....	9-1
9.2	CODING RULES AFFECTING 64-BIT MODE.....	9-1

	PAGE	
9.2.1	Use Legacy 32-Bit Instructions When Data Size Is 32 Bits	9-1
9.2.2	Use Extra Registers to Reduce Register Pressure	9-2
9.2.3	Effective Use of 64-Bit by 64-Bit Multiplies	9-2
9.2.4	Replace 128-bit Integer Division with 128-bit Multiplies	9-3
9.2.5	Sign Extension to Full 64-Bits	9-6
9.3	ALTERNATE CODING RULES FOR 64-BIT MODE	9-7
9.3.1	Use 64-Bit Registers Instead of Two 32-Bit Registers for 64-Bit Arithmetic Result	9-7
9.3.2	CVTSI2SS and CVTSI2SD	9-8
9.3.3	Using Software Prefetch	9-8

CHAPTER 10 SSE4.2 AND SIMD PROGRAMMING FOR TEXT-PROCESSING/LEXING/PARSING

10.1	SSE4.2 STRING AND TEXT INSTRUCTIONS	10-1
10.1.1	CRC32	10-5
10.2	USING SSE4.2 STRING AND TEXT INSTRUCTIONS	10-7
10.2.1	Unaligned Memory Access and Buffer Size Management	10-7
10.2.2	Unaligned Memory Access and String Library	10-8
10.3	SSE4.2 APPLICATION CODING GUIDELINE AND EXAMPLES	10-8
10.3.1	Null Character Identification (Strlen equivalent)	10-8
10.3.2	White-Space-Like Character Identification	10-12
10.3.3	Substring Searches	10-16
10.3.4	String Token Extraction and Case Handling	10-25
10.3.5	Unicode Processing and PCMPxSTRY	10-30
10.3.6	Replacement String Library Function Using SSE4.2	10-37
10.4	SSE4.2 ENABLED NUMERICAL AND LEXICAL COMPUTATION	10-39

CHAPTER 11 OPTIMIZATIONS FOR INTEL AVX

11.1	INTEL AVX INTRINSICS CODING	11-2
11.1.1	Intel AVX Assembly Coding	11-5
11.2	NON-DESTRUCTIVE SOURCE (NDS)	11-7
11.3	MIXING AVX CODE WITH SSE CODE	11-8
11.3.1	Mixing Intel AVX and Intel SSE in Function Calls	11-10
11.4	128-BIT LANE OPERATION AND AVX	11-11
11.4.1	Programming With the Lane Concept	11-11
11.4.2	Strided Load Technique	11-12
11.4.3	The Register Overlap Technique	11-15
11.5	DATA GATHER AND SCATTER	11-17
11.5.1	Data Gather	11-17
11.5.2	Data Scatter	11-20
11.6	DATA ALIGNMENT FOR INTEL AVX	11-21
11.6.1	Align Data to 32 Bytes	11-22
11.6.2	Consider 16-Byte Memory Access when Memory is Unaligned	11-23
11.6.3	Prefer Aligned Stores Over Aligned Loads	11-25
11.7	L1D CACHE LINE REPLACEMENTS	11-25
11.8	4K ALIASING	11-26
11.9	CONDITIONAL SIMD PACKED LOADS AND STORES	11-26
11.9.1	Conditional Loops	11-28
11.10	MIXING INTEGER AND FLOATING-POINT CODE	11-29

	PAGE
11.11 HANDLING PORT 5 PRESSURE	11-33
11.11.1 Replace Shuffles with Blends	11-33
11.11.2 Design Algorithm With Fewer Shuffles	11-36
11.11.3 Perform Basic Shuffles on Load Ports	11-39
11.12 DIVIDE AND SQUARE ROOT OPERATIONS	11-40
11.12.1 Single-Precision Divide	11-41
11.12.2 Single-Precision Reciprocal Square Root	11-44
11.12.3 Single-Precision Square Root	11-46
11.13 OPTIMIZATION OF ARRAY SUB SUM EXAMPLE	11-48

CHAPTER 12

POWER OPTIMIZATION FOR MOBILE USAGES

12.1 OVERVIEW	12-1
12.2 MOBILE USAGE SCENARIOS	12-2
12.2.1 Intelligent Energy Efficient Software	12-3
12.3 ACPI C-STATES	12-4
12.3.1 Processor-Specific C4 and Deep C4 States	12-5
12.3.2 Processor-Specific Deep C-States and Intel® Turbo Boost Technology	12-6
12.3.3 Processor-Specific Deep C-States for Intel® microarchitecture code name Sandy Bridge 12-7	
12.3.4 Intel® Turbo Boost Technology 2.0	12-8
12.4 GUIDELINES FOR EXTENDING BATTERY LIFE	12-8
12.4.1 Adjust Performance to Meet Quality of Features	12-9
12.4.2 Reducing Amount of Work	12-10
12.4.3 Platform-Level Optimizations	12-10
12.4.4 Handling Sleep State Transitions	12-11
12.4.5 Using Enhanced Intel SpeedStep® Technology	12-11
12.4.6 Enabling Intel® Enhanced Deeper Sleep	12-13
12.4.7 Multicore Considerations	12-14
12.4.7.1 Enhanced Intel SpeedStep® Technology	12-14
12.4.7.2 Thread Migration Considerations	12-14
12.4.7.3 Multicore Considerations for C-States	12-15
12.5 TUNING SOFTWARE FOR INTELLIGENT POWER CONSUMPTION	12-16
12.5.1 Reduction of Active Cycles	12-17
12.5.1.1 Multi-threading to reduce Active Cycles	12-17
12.5.1.2 Vectorization	12-18
12.5.2 PAUSE and Sleep(0) Loop Optimization	12-18
12.5.3 Spin-Wait Loops	12-20
12.5.4 Using Event Driven Service Instead of Polling in Code	12-20
12.5.5 Reducing Interrupt Rate	12-20
12.5.6 Reducing Privileged Time	12-21
12.5.7 Setting Context Awareness in the Code	12-22
12.5.8 Saving Energy by Optimizing for Performance	12-23
12.6 PROCESSOR SPECIFIC POWER MANAGEMENT OPTIMIZATION FOR SYSTEM SOFTWARE 12-24	
12.6.1 Power Management Recommendation of Processor-Specific Inactive State Configurations	12-24
12.6.1.1 Balancing Power Management and Responsiveness of Inactive To Active State Transitions	12-25

CHAPTER 13**INTEL® ATOM™ MICROARCHITECTURE AND SOFTWARE OPTIMIZATION**

13.1	OVERVIEW	13-1
13.2	INTEL® ATOM™ MICROARCHITECTURE	13-1
13.2.1	Hyper-Threading Technology Support in Intel® Atom™ Microarchitecture	13-3
13.3	CODING RECOMMENDATIONS FOR INTEL® ATOM™ MICROARCHITECTURE	13-4
13.3.1	Optimization for Front End of Intel® Atom™ Microarchitecture	13-4
13.3.2	Optimizing the Execution Core	13-6
13.3.2.1	Integer Instruction Selection	13-6
13.3.2.2	Address Generation	13-7
13.3.2.3	Integer Multiply	13-8
13.3.2.4	Integer Shift Instructions	13-9
13.3.2.5	Partial Register Access	13-9
13.3.2.6	FP/SIMD Instruction Selection	13-9
13.3.3	Optimizing Memory Access	13-12
13.3.3.1	Store Forwarding	13-12
13.3.3.2	First-level Data Cache	13-13
13.3.3.3	Segment Base	13-13
13.3.3.4	String Moves	13-14
13.3.3.5	Parameter Passing	13-15
13.3.3.6	Function Calls	13-15
13.3.3.7	Optimization of Multiply/Add Dependent Chains	13-15
13.3.3.8	Position Independent Code	13-17
13.4	INSTRUCTION LATENCY	13-18

APPENDIX A**APPLICATION PERFORMANCE****TOOLS**

A.1	COMPILERS	A-2
A.1.1	Recommended Optimization Settings for Intel® 64 and IA-32 Processors	A-2
A.1.2	Vectorization and Loop Optimization	A-5
A.1.2.1	Multithreading with OpenMP*	A-5
A.1.2.2	Automatic Multithreading	A-5
A.1.3	Inline Expansion of Library Functions (/Oi, /Oi-)	A-6
A.1.4	Floating-point Arithmetic Precision (/Op, /Op-, /Qprec, /Qprec_div, /Qpc, /Qlong_double).	A-6
A.1.5	Rounding Control Option (/Qrcr, /Qrcd)	A-6
A.1.6	Interprocedural and Profile-Guided Optimizations	A-6
A.1.6.1	Interprocedural Optimization (IPO)	A-6
A.1.6.2	Profile-Guided Optimization (PGO)	A-6
A.1.7	Auto-Generation of Vectorized Code	A-7
A.2	INTEL® VTUNE™ PERFORMANCE ANALYZER	A-11
A.2.1	Sampling	A-11
A.2.1.1	Time-based Sampling	A-12
A.2.1.2	Event-based Sampling	A-12
A.2.1.3	Workload Characterization	A-12
A.2.2	Call Graph	A-12
A.2.3	Counter Monitor	A-13
A.3	INTEL® PERFORMANCE LIBRARIES	A-13
A.3.1	Benefits Summary	A-14

	PAGE
A.3.2	Optimizations with the Intel® Performance Libraries A-14
A.4	INTEL® THREADING ANALYSIS TOOLS..... A-15
A.4.1	Intel® Thread Checker 3.0 A-15
A.4.2	Intel® Thread Profiler 3.0..... A-15
A.4.3	Intel® Threading Building Blocks 1.0. A-16
A.5	INTEL® CLUSTER TOOLS..... A-17
A.5.1	Intel® MPI Library 3.1 A-17
A.5.2	Intel® Trace Analyzer and Collector 7.1..... A-17
A.5.3	Intel® MPI Benchmarks 3.1 A-17
A.5.4	Benefits Summary A-18
A.5.4.1	Multiple usability improvements..... A-18
A.5.4.2	Improved application performance A-18
A.5.4.3	Extended interoperability..... A-18
A.6	INTEL® XML PRODUCTS A-18
A.6.1	Intel® XML Software Suite 1.0..... A-18
A.6.1.1	Intel® XSLT Accelerator A-18
A.6.1.2	Intel® XPath Accelerator A-19
A.6.1.3	Intel® XML Schema Accelerator..... A-19
A.6.1.4	Intel® XML Parsing Accelerator A-19
A.6.2	Intel® SOA Security Toolkit 1.0 Beta for Axis2..... A-19
A.6.2.1	High Performance A-20
A.6.2.2	Standards Compliant..... A-20
A.6.2.3	Easy Integration..... A-20
A.6.3	Intel® XSLT Accelerator 1.1 for Java® Environments on Linux® and Windows® Operating Systems A-20
A.6.3.1	High Performance Transformations A-20
A.6.3.2	Large XML File Transformations..... A-20
A.6.3.3	Standards Compliant..... A-21
A.6.3.4	Thread-Safe A-21
A.7	INTEL® SOFTWARE COLLEGE A-21

APPENDIX B USING PERFORMANCE MONITORING EVENTS

B.1	INTEL® XEON® PROCESSOR 5500 SERIES..... B-1
B.2	PERFORMANCE ANALYSIS TECHNIQUES FOR INTEL® XEON® PROCESSOR 5500 SERIES B-2
B.2.1	Cycle Accounting and Uop Flow Analysis B-3
B.2.1.1	Cycle Drill Down and Branch Mispredictions B-5
B.2.1.2	Basic Block Drill Down B-8
B.2.2	Stall Cycle Decomposition and Core Memory Accesses B-9
B.2.2.1	Measuring Costs of Microarchitectural Conditions B-10
B.2.3	Core PMU Precise Events B-11
B.2.3.1	Precise Memory Access Events..... B-12
B.2.3.2	Load Latency Event B-14
B.2.3.3	Precise Execution Events B-16
B.2.3.4	Last Branch Record (LBR)..... B-18
B.2.3.5	Measuring Core Memory Access Latency..... B-21
B.2.3.6	Measuring Per-Core Bandwidth..... B-24
B.2.3.7	Miscellaneous L1 and L2 Events for Cache Misses B-25
B.2.3.8	TLB Misses..... B-25
B.2.3.9	L1 Data Cache..... B-26
B.2.4	Front End Monitoring Events B-27

CONTENTS

	PAGE
B.2.4.1	Branch Mispredictions..... B-27
B.2.4.2	Front End Code Generation Metrics B-27
B.2.5	Uncore Performance Monitoring Events..... B-28
B.2.5.1	Global Queue Occupancy B-28
B.2.5.2	Global Queue Port Events B-31
B.2.5.3	Global Queue Snoop Events B-31
B.2.5.4	L3 Events B-32
B.2.6	Intel QuickPath Interconnect Home Logic (QHL)..... B-32
B.2.7	Measuring Bandwidth From the Uncore B-39
B.3	PERFORMANCE TUNING TECHNIQUES FOR INTEL® MICROARCHITECTURE CODE NAME SANDY BRIDGE..... B-39
B.3.1	Correlating Performance Bottleneck to Source Location..... B-40
B.3.2	Locating Stalls in the Microarchitecture Pipeline B-40
B.3.3	Back End Stalls..... B-42
B.3.4	Memory Sub-System Stalls B-43
B.3.4.1	Accounting for Load Latency B-43
B.3.4.2	Cache-line Replacement Analysis B-46
B.3.4.3	Lock Contention Analysis..... B-47
B.3.4.4	Other Memory Access Issues B-48
B.3.5	Execution Stalls..... B-51
B.3.5.1	Longer Instruction Latencies B-51
B.3.5.2	Assists B-52
B.3.6	Bad Speculation..... B-52
B.3.6.1	Branch Mispredicts..... B-53
B.3.7	Front End Stalls..... B-53
B.3.7.1	Understanding the Micro-op Delivery Rate B-53
B.3.7.2	Understanding the Sources of the Micro-op Queue..... B-55
B.3.7.3	The Decoded ICache..... B-57
B.3.7.4	Issues in the Legacy Decode Pipeline..... B-58
B.3.7.5	Instruction Cache B-59
B.4	USING PERFORMANCE EVENTS OF INTEL® CORE™ SOLO AND INTEL® CORE™ DUO PROCESSORS B-59
B.4.1	Understanding the Results in a Performance Counter..... B-59
B.4.2	Ratio Interpretation..... B-60
B.4.3	Notes on Selected Events B-61
B.5	DRILL-DOWN TECHNIQUES FOR PERFORMANCE ANALYSIS B-62
B.5.1	Cycle Composition at Issue Port..... B-64
B.5.2	Cycle Composition of OOO Execution..... B-65
B.5.3	Drill-Down on Performance Stalls B-66
B.6	EVENT RATIOS FOR INTEL CORE MICROARCHITECTURE B-67
B.6.1	Clocks Per Instructions Retired Ratio (CPI)..... B-68
B.6.2	Front-end Ratios..... B-68
B.6.2.1	Code Locality B-68
B.6.2.2	Branching and Front-end B-69
B.6.2.3	Stack Pointer Tracker B-69
B.6.2.4	Macro-fusion B-69
B.6.2.5	Length Changing Prefix (LCP) Stalls..... B-69
B.6.2.6	Self Modifying Code Detection..... B-70
B.6.3	Branch Prediction Ratios B-70
B.6.3.1	Branch Mispredictions..... B-70
B.6.3.2	Virtual Tables and Indirect Calls..... B-70
B.6.3.3	Mispredicted Returns B-71

	PAGE
B.6.4	Execution RatiosB-71
B.6.4.1	Resource StallsB-71
B.6.4.2	ROB Read Port StallsB-71
B.6.4.3	Partial Register StallsB-71
B.6.4.4	Partial Flag StallsB-72
B.6.4.5	Bypass Between Execution DomainsB-72
B.6.4.6	Floating Point Performance RatiosB-72
B.6.5	Memory Sub-System - Access Conflicts RatiosB-73
B.6.5.1	Loads Blocked by the L1 Data CacheB-73
B.6.5.2	4K Aliasing and Store Forwarding Block DetectionB-73
B.6.5.3	Load Block by Preceding StoresB-73
B.6.5.4	Memory DisambiguationB-74
B.6.5.5	Load Operation Address TranslationB-74
B.6.6	Memory Sub-System - Cache Misses RatiosB-74
B.6.6.1	Locating Cache Misses in the CodeB-74
B.6.6.2	L1 Data Cache MissesB-75
B.6.6.3	L2 Cache MissesB-75
B.6.7	Memory Sub-system - PrefetchingB-75
B.6.7.1	L1 Data PrefetchingB-75
B.6.7.2	L2 Hardware PrefetchingB-75
B.6.7.3	Software PrefetchingB-76
B.6.8	Memory Sub-system - TLB Miss RatiosB-76
B.6.9	Memory Sub-system - Core InteractionB-77
B.6.9.1	Modified Data SharingB-77
B.6.9.2	Fast Synchronization PenaltyB-77
B.6.9.3	Simultaneous Extensive Stores and Load MissesB-77
B.6.10	Memory Sub-system - Bus CharacterizationB-78
B.6.10.1	Bus UtilizationB-78
B.6.10.2	Modified Cache Lines EvictionB-78

APPENDIX C INSTRUCTION LATENCY AND THROUGHPUT

C.1	OVERVIEW C-1
C.2	DEFINITIONS C-2
C.3	LATENCY AND THROUGHPUT C-3
C.3.1	Latency and Throughput with Register Operands C-3
C.3.2	Table Footnotes C-33
C.3.3	Instructions with Memory Operands C-35

APPENDIX D STACK ALIGNMENT

D.4	STACK FRAMES D-1
D.4.1	Aligned ESP-Based Stack Frames D-3
D.4.2	Aligned EBP-Based Stack Frames D-4
D.4.3	Stack Frame Optimizations D-6
D.5	INLINED ASSEMBLY AND EBX D-7

APPENDIX E
SUMMARY OF RULES AND SUGGESTIONS

E.1 ASSEMBLY/COMPILER CODING RULES E-1
E.2 USER/SOURCE CODING RULES E-8
E.3 TUNING SUGGESTIONS E-11
E.4 SSE4.2 CODING RULES E-12
E.5 ASSEMBLY/COMPILER CODING RULES FOR THE INTEL® ATOM™ PROCESSOR..... E-12

EXAMPLES

Example 3-1. Assembly Code with an Unpredictable Branch.....	3-8
Example 3-2. Code Optimization to Eliminate Branches	3-8
Example 3-4. Use of PAUSE Instruction	3-9
Example 3-3. Eliminating Branch with CMOV Instruction	3-9
Example 3-5. Pentium 4 Processor Static Branch Prediction Algorithm.....	3-10
Example 3-6. Static Taken Prediction	3-11
Example 3-7. Static Not-Taken Prediction.....	3-11
Example 3-8. Indirect Branch With Two Favored Targets	3-14
Example 3-9. A Peeling Technique to Reduce Indirect Branch Misprediction.....	3-15
Example 3-10. Loop Unrolling.....	3-16
Example 3-11. Macro-fusion, Unsigned Iteration Count.....	3-20
Example 3-12. Macro-fusion, If Statement.....	3-20
Example 3-13. Macro-fusion, Signed Variable	3-21
Example 3-14. Macro-fusion, Signed Comparison.....	3-22
Example 3-15. Additional Macro-fusion Benefit in Intel microarchitecture code name Sandy Bridge. 3-22	
Example 3-16. Avoiding False LCP Delays with 0xF7 Group Instructions	3-24
Example 3-17. Unrolling Loops in LSD to Optimize Emission Bandwidth	3-25
Example 3-18. Independent Two-Operand LEA Example	3-30
Example 3-19. Alternative to Three-Operand LEA	3-31
Example 3-20. Examples of 512-bit Additions	3-32
Example 3-21. Clearing Register to Break Dependency While Negating Array Elements.....	3-35
Example 3-22. Spill Scheduling Code	3-38
Example 3-23. Avoiding Partial Register Stalls in Integer Code.....	3-42
Example 3-24. Avoiding Partial Register Stalls in SIMD Code.....	3-43
Example 3-25. Avoiding Partial Flag Register Stalls	3-44
Example 3-26. Partial Flag Register Accesses in Intel microarchitecture code name Sandy Bridge 3-45	
Example 3-27. Reference Code Template for Partially Vectorizable Program.....	3-48
Example 3-28. Three Alternate Packing Methods for Avoiding Store Forwarding Difficulty.....	3-50
Example 3-29. Using Four Registers to Reduce Memory Spills and Simplify Result Passing	3-50
Example 3-30. Stack Optimization Technique to Simplify Parameter Passing	3-51
Example 3-31. Base Line Code Sequence to Estimate Loop Overhead.....	3-52
Example 3-32. Optimize for Load Port Bandwidth in Intel microarchitecture code name Sandy Bridge 3-55	
Example 3-33. Index versus Pointers in Pointer-Chasing Code.....	3-56
Example 3-34. Example of Bank Conflicts in L1D Cache and Remedy.....	3-57
Example 3-35. Using XMM Register in Lieu of Memory for Register Spills	3-59
Example 3-36. Loads Blocked by Stores of Unknown Address	3-60
Example 3-37. Code That Causes Cache Line Split.....	3-61
Example 3-38. Situations Showing Small Loads After Large Store	3-65
Example 3-39. Non-forwarding Example of Large Load After Small Store	3-65
Example 3-40. A Non-forwarding Situation in Compiler Generated Code.....	3-66
Example 3-41. Two Ways to Avoid Non-forwarding Situation in Example 3-40	3-66
Example 3-42. Large and Small Load Stalls	3-66

CONTENTS

	PAGE
Example 3-43. Loop-carried Dependence Chain.....	3-69
Example 3-44. Rearranging a Data Structure.....	3-69
Example 3-45. Decomposing an Array.....	3-70
Example 3-46. Dynamic Stack Alignment.....	3-72
Example 3-47. Aliasing Between Loads and Stores Across Loop Iterations.....	3-75
Example 3-48. Instruction Pointer Query Techniques.....	3-76
Example 3-49. Using Non-temporal Stores and 64-byte Bus Write Transactions.....	3-80
Example 3-50. On-temporal Stores and Partial Bus Write Transactions.....	3-81
Example 3-51. Using DCU Hardware Prefetch.....	3-84
Example 3-52. Avoid Causing DCU Hardware Prefetch to Fetch Un-needed Lines.....	3-85
Example 3-53. Technique For Using L1 Hardware Prefetch.....	3-86
Example 3-54. REP STOSD with Arbitrary Count Size and 4-Byte-Aligned Destination.....	3-89
Example 3-55. Algorithm to Avoid Changing Rounding Mode.....	3-95
Example 4-1. Identification of MMX Technology with CPUID.....	4-2
Example 4-2. Identification of SSE with CPUID.....	4-3
Example 4-3. Identification of SSE2 with cpuid.....	4-3
Example 4-4. Identification of SSE3 with CPUID.....	4-4
Example 4-5. Identification of SSSE3 with cpuid.....	4-4
Example 4-6. Identification of SSE4.1 with cpuid.....	4-5
Example 4-7. Identification of SSE4.2 with cpuid.....	4-5
Example 4-8. Detection of AESNI Instructions.....	4-6
Example 4-9. Detection of PCLMULQDQ Instruction.....	4-6
Example 4-10. Detection of AVX Instruction.....	4-7
Example 4-11. Detection of VEX-Encoded AESNI Instructions.....	4-8
Example 4-12. Detection of VEX-Encoded AESNI Instructions.....	4-8
Example 4-13. Simple Four-Iteration Loop.....	4-14
Example 4-14. Streaming SIMD Extensions Using Inlined Assembly Encoding.....	4-15
Example 4-15. Simple Four-Iteration Loop Coded with Intrinsics.....	4-16
Example 4-16. C++ Code Using the Vector Classes.....	4-17
Example 4-17. Automatic Vectorization for a Simple Loop.....	4-18
Example 4-18. C Algorithm for 64-bit Data Alignment.....	4-21
Example 4-20. SoA Data Structure.....	4-24
Example 4-21. AoS and SoA Code Samples.....	4-24
Example 4-19. AoS Data Structure.....	4-24
Example 4-22. Hybrid SoA Data Structure.....	4-26
Example 4-23. Pseudo-code Before Strip Mining.....	4-27
Example 4-24. Strip Mined Code.....	4-28
Example 4-25. Loop Blocking.....	4-29
Example 4-26. Emulation of Conditional Moves.....	4-31
Example 5-1. Resetting Register Between __m64 and FP Data Types Code.....	5-4
Example 5-2. FIR Processing Example in C language Code.....	5-5
Example 5-3. SSE2 and SSSE3 Implementation of FIR Processing Code.....	5-5
Example 5-5. Signed Unpack Code.....	5-7
Example 5-4. Zero Extend 16-bit Values into 32 Bits Using Unsigned Unpack Instructions Code.....	5-7
Example 5-6. Interleaved Pack with Saturation Code.....	5-9
Example 5-7. Interleaved Pack without Saturation Code.....	5-10

Example 5-8. Unpacking Two Packed-word Sources in Non-interleaved Way Code	5-12
Example 5-9. PEXTRW Instruction Code	5-13
Example 5-11. Repeated PINSRW Instruction Code	5-14
Example 5-10. PINSRW Instruction Code	5-14
Example 5-12. Non-Unit Stride Load/Store Using SSE4.1 Instructions	5-15
Example 5-13. Scatter and Gather Operations Using SSE4.1 Instructions	5-15
Example 5-14. PMOVMSKB Instruction Code	5-16
Example 5-15. Broadcast a Word Across XMM, Using 2 SSE2 Instructions	5-17
Example 5-16. Swap/Reverse words in an XMM, Using 3 SSE2 Instructions	5-18
Example 5-17. Generating Constants	5-19
Example 5-18. Absolute Difference of Two Unsigned Numbers	5-21
Example 5-19. Absolute Difference of Signed Numbers	5-21
Example 5-20. Computing Absolute Value	5-22
Example 5-21. Basic C Implementation of RGBA to BGRA Conversion	5-22
Example 5-22. Color Pixel Format Conversion Using SSE2	5-23
Example 5-23. Color Pixel Format Conversion Using SSSE3	5-24
Example 5-24. Big-Endian to Little-Endian Conversion	5-25
Example 5-25. Clipping to a Signed Range of Words [High, Low]	5-26
Example 5-26. Clipping to an Arbitrary Signed Range [High, Low]	5-26
Example 5-28. Clipping to an Arbitrary Unsigned Range [High, Low]	5-27
Example 5-27. Simplified Clipping to an Arbitrary Signed Range	5-27
Example 5-29. Complex Multiply by a Constant	5-30
Example 5-30. Using PTEST to Separate Vectorizable and non-Vectorizable Loop Iterations ..	5-31
Example 5-31. Using PTEST and Variable BLEND to Vectorize Heterogeneous Loops	5-32
Example 5-32. Baseline C Code for Mandelbrot Set Map Evaluation	5-33
Example 5-33. Vectorized Mandelbrot Set Map Evaluation Using SSE4.1 Intrinsics	5-34
Example 5-34. A Large Load after a Series of Small Stores (Penalty)	5-36
Example 5-36. A Series of Small Loads After a Large Store	5-37
Example 5-37. Eliminating Delay for a Series of Small Loads after a Large Store	5-37
Example 5-35. Accessing Data Without Delay	5-37
Example 5-38. An Example of Video Processing with Cache Line Splits	5-38
Example 5-39. Video Processing Using LDDQU to Avoid Cache Line Splits	5-39
Example 5-40. Un-optimized Reverse Memory Copy in C	5-41
Example 5-41. Using PSHUFB to Reverse Byte Ordering 16 Bytes at a Time	5-42
Example 5-42. PMOVX/PMOVZX Work-around to Avoid False Dependency	5-45
Example 5-43. Table Look-up Operations in C Code	5-46
Example 5-44. Shift Techniques on Non-Vectorizable Table Look-up	5-47
Example 5-45. PEXTRD Techniques on Non-Vectorizable Table Look-up	5-48
Example 5-46. Pseudo-Code Flow of AES Counter Mode Operation	5-50
Example 5-47. AES128-CTR Implementation with Eight Block in Parallel	5-51
Example 5-48. AES128 Key Expansion	5-60
Example 6-1. Pseudocode for Horizontal (xyz, AoS) Computation	6-6
Example 6-2. Pseudocode for Vertical (xxxx, yyyy, zzzz, SoA) Computation	6-6
Example 6-3. Swizzling Data Using SHUFPS, MOVLHPS, MOVHLPS	6-7
Example 6-4. Swizzling Data Using UNPCKxxx Instructions	6-8
Example 6-5. Deswizzling Single-Precision SIMD Data	6-9
Example 6-6. Deswizzling Data Using SIMD Integer Instructions	6-10

Example 6-7. Horizontal Add Using MOVHLPS/MOVLHPS 6-12

Example 6-8. Horizontal Add Using Intrinsics with MOVHLPS/MOVLHPS 6-12

Example 6-9. Multiplication of Two Pair of Single-precision Complex Number 6-15

Example 6-10. Division of Two Pair of Single-precision Complex Numbers 6-16

Example 6-11. Double-Precision Complex Multiplication of Two Pairs 6-17

Example 6-12. Double-Precision Complex Multiplication Using Scalar SSE2 6-17

Example 6-13. Dot Product of Vector Length 4 Using SSE/SSE2 6-19

Example 6-14. Dot Product of Vector Length 4 Using SSE3 6-19

Example 6-15. Dot Product of Vector Length 4 Using SSE4.1 6-19

Example 6-16. Unrolled Implementation of Four Dot Products 6-20

Example 6-17. Normalization of an Array of Vectors 6-21

Example 6-18. Normalize (x, y, z) Components of an Array of Vectors Using SSE2 6-22

Example 6-19. Normalize (x, y, z) Components of an Array of Vectors Using SSE4.1 6-23

Example 6-20. Data Organization in Memory for AOS Vector-Matrix Multiplication 6-24

Example 6-21. AOS Vector-Matrix Multiplication with HADDPS 6-24

Example 6-22. AOS Vector-Matrix Multiplication with DPPS 6-25

Example 6-23. Data Organization in Memory for SOA Vector-Matrix Multiplication 6-26

Example 6-24. Vector-Matrix Multiplication with Native SOA Data Layout 6-27

Example 7-1. Pseudo-code Using CLFLUSH 7-13

Example 7-2. Populating an Array for Circular Pointer Chasing with Constant Stride 7-15

Example 7-3. Prefetch Scheduling Distance 7-18

Example 7-4. Using Prefetch Concatenation 7-20

Example 7-5. Concatenation and Unrolling the Last Iteration of Inner Loop 7-20

Example 7-6. Data Access of a 3D Geometry Engine without Strip-mining 7-26

Example 7-7. Data Access of a 3D Geometry Engine with Strip-mining 7-26

Example 7-8. Using HW Prefetch to Improve Read-Once Memory Traffic 7-28

Example 7-9. Basic Algorithm of a Simple Memory Copy 7-33

Example 7-10. A Memory Copy Routine Using Software Prefetch 7-34

Example 7-11. Memory Copy Using Hardware Prefetch and Bus Segmentation 7-36

Example 8-1. Serial Execution of Producer and Consumer Work Items 8-6

Example 8-2. Basic Structure of Implementing Producer Consumer Threads 8-7

Example 8-3. Thread Function for an Interlaced Producer Consumer Model 8-9

Example 8-4. Spin-wait Loop and PAUSE Instructions 8-17

Example 8-5. Coding Pitfall using Spin Wait Loop 8-20

Example 8-6. Placement of Synchronization and Regular Variables 8-22

Example 8-7. Declaring Synchronization Variables without Sharing a Cache Line 8-22

Example 8-8. Batched Implementation of the Producer Consumer Threads 8-29

Example 8-9. Parallel Memory Initialization Technique Using OpenMP and NUMA 8-34

Example 9-1. Compute 64-bit Quotient and Remainder with 64-bit Divisor 9-4

Example 9-2. Quotient and Remainder of 128-bit Dividend with 64-bit Divisor 9-5

Example 10-1. A Hash Function Examples 10-5

Example 10-2. Hash Function Using CRC32 10-6

Example 10-3. Strlen() Using General-Purpose Instructions 10-9

Example 10-4. Sub-optimal PCMPISTR1 Implementation of EOS handling 10-11

Example 10-5. Strlen() Using PCMPISTR1 without Loop-Carry Dependency 10-12

Example 10-6. WordCnt() Using C and Byte-Scanning Technique 10-13

Example 10-7. WordCnt() Using PCMPISTRM 10-15

	PAGE
Example 10-8. KMP Substring Search in C	10-17
Example 10-9. Brute-Force Substring Search Using PCMPISTRI Intrinsic	10-19
Example 10-10. Substring Search Using PCMPISTRI and KMP Overlap Table	10-22
Example 10-11. I Equivalent Strtok_s() Using PCMPISTRI Intrinsic	10-26
Example 10-12. I Equivalent Strupr() Using PCMPISTRM Intrinsic	10-29
Example 10-13. UTF16 VerStrlen() Using C and Table Lookup Technique	10-31
Example 10-14. Assembly Listings of UTF16 VerStrlen() Using PCMPISTRI	10-32
Example 10-15. Intrinsic Listings of UTF16 VerStrlen() Using PCMPISTRI	10-35
Example 10-16. Replacement String Library Strcmp Using SSE4.2	10-38
Example 10-17. High-level flow of Character Subset Validation for String Conversion.	10-40
Example 10-18. Intrinsic Listings of atol() Replacement Using PCMPISTRI	10-41
Example 10-19. Auxiliary Routines and Data Constants Used in sse4i_atol() listing.	10-44
Example 11-1. Cartesian Coordinate Transformation with Intrinsics.	11-3
Example 11-2. Cartesian Coordinate Transformation with Assembly	11-5
Example 11-3. Direct Polynomial Calculation.	11-7
Example 11-4. Function Calls and AVX/SSE transitions.	11-10
Example 11-5. AoS to SoA Conversion of Complex Numbers in C Code	11-13
Example 11-6. Aos to SoA Conversion of Complex Numbers Using AVX	11-14
Example 11-7. Register Overlap Method for Median of 3 Numbers.	11-16
Example 11-8. Data Gather - AVX versus Scalar Code	11-18
Example 11-9. Scatter Operation Using AVX	11-20
Example 11-10. SAXPY using Intel AVX	11-22
Example 11-11. Using 16-Byte Memory Operations for Unaligned 32-Byte Memory Operation	11-24
Example 11-12. SAXPY Implementations for Unaligned Data Addresses	11-24
Example 11-13. Loop with Conditional Expression	11-28
Example 11-14. Handling Loop Conditional with VMASKMOV	11-28
Example 11-15. Three-Tap Filter in C Code	11-30
Example 11-16. Three-Tap Filter with 128-bit Mixed Integer and FP SIMD.	11-30
Example 11-17. 256-bit AVX Three-Tap Filter Code with VSHUFPS	11-31
Example 11-18. Three-Tap Filter Code with Mixed 256-bit AVX and 128-bit AVX Code	11-32
Example 11-19. 8x8 Matrix Transpose - Replace Shuffles with Blends	11-34
Example 11-20. 8x8 Matrix Transpose Using VINSRTPS	11-37
Example 11-21. Port 5 versus Load Port Shuffles	11-40
Example 11-22. Divide Using DIVPS for 24-bit Accuracy	11-42
Example 11-23. Divide Using RCPPS 11-bit Approximation	11-42
Example 11-24. Divide Using RCPPS and Newton-Raphson Iteration	11-43
Example 11-25. Reciprocal Square Root Using DIVPS+SQRTPS for 24-bit Accuracy	11-44
Example 11-26. Reciprocal Square Root Using RCPPS 11-bit Approximation	11-45
Example 11-27. Reciprocal Square Root Using RCPPS and Newton-Raphson Iteration.	11-45
Example 11-28. Square Root Using SQRTPS for 24-bit Accuracy	11-46
Example 11-29. Square Root Using RCPPS 11-bit Approximation	11-47
Example 11-30. Square Root Using RCPPS and One Taylor Series Expansion	11-47
Example 11-31. Array Sub Sums Algorithm.	11-50
Example 12-1. Unoptimized Sleep Loop	12-19
Example 12-2. Power Consumption Friendly Sleep Loop Using PAUSE	12-19
Example 13-1. Instruction Pairing and Alignment to Optimize Decode Throughput on Intel® Atom™ Microarchitecture 13-5	

CONTENTS

	PAGE
Example 13-2. Alternative to Prevent AGU and Execution Unit Dependency	13-8
Example 13-3. Pipeling Instruction Execution in Integer Computation.	13-9
Example 13-4. Memory Copy of 64-byte.....	13-14
Example 13-5. Examples of Dependent Multiply and Add Computation	13-16
Example 13-6. Instruction Pointer Query Techniques.....	13-17
Example 13-8. Auto-Generated Code of Storing Absolutes.....	A-8
Example 13-9. Changes Signs	A-8
Example 13-7. Storing Absolute Values.....	A-8
Example 13-11.Data Conversion	A-9
Example 13-10.Auto-Generated Code of Sign Conversion	A-9
Example 13-13.Un-aligned Data Operation	A-10
Example 13-12.Auto-Generated Code of Data Conversion	A-10
Example 13-14.Auto-Generated Code to Avoid Unaligned Loads.....	A-11
Example D-1. Aligned esp-Based Stack Frame.....	D-3
Example D-2. Aligned ebp-based Stack Frames.....	D-5

FIGURES

Figure 2-1.	Intel microarchitecture code name Sandy Bridge Pipeline Functionality	2-4
Figure 2-2.	Intel Core Microarchitecture Pipeline Functionality	2-28
Figure 2-3.	Execution Core of Intel Core Microarchitecture	2-36
Figure 2-4.	Store-Forwarding Enhancements in Enhanced Intel Core Microarchitecture	2-41
Figure 2-5.	Intel Advanced Smart Cache Architecture	2-42
Figure 2-6.	Intel microarchitecture code name Nehalem Pipeline Functionality	2-46
Figure 2-7.	Front End of Intel microarchitecture code name Nehalem	2-47
Figure 2-8.	Store-Forwarding Scenarios of 16-Byte Store Operations	2-54
Figure 2-9.	Store-Forwarding Enhancement in Intel microarchitecture code name Nehalem	2-55
Figure 2-10.	The Intel NetBurst Microarchitecture	2-59
Figure 2-11.	Execution Units and Ports in Out-Of-Order Core	2-65
Figure 2-12.	The Intel Pentium M Processor Microarchitecture	2-71
Figure 2-13.	Hyper-Threading Technology on an SMP	2-76
Figure 2-14.	Pentium D Processor, Pentium Processor Extreme Edition, Intel Core Duo Processor, Intel Core 2 Duo Processor, and Intel Core 2 Quad Processor 2-80	
Figure 2-15.	Typical SIMD Operations	2-84
Figure 2-16.	SIMD Instruction Register Usage	2-85
Figure 3-1.	Generic Program Flow of Partially Vectorized Code	3-47
Figure 3-2.	Cache Line Split in Accessing Elements in a Array	3-62
Figure 3-3.	Size and Alignment Restrictions in Store Forwarding	3-64
Figure 4-1.	General Procedural Flow of Application Detection of AVX	4-7
Figure 4-2.	Converting to Streaming SIMD Extensions Chart	4-10
Figure 4-3.	Hand-Coded Assembly and High-Level Compiler Performance Trade-offs	4-13
Figure 4-4.	Loop Blocking Access Pattern	4-30
Figure 5-1.	PACKSSDW mm, mm/mm64 Instruction	5-8
Figure 5-2.	Interleaved Pack with Saturation	5-9
Figure 5-4.	Result of Non-Interleaved Unpack High in MM1	5-11
Figure 5-3.	Result of Non-Interleaved Unpack Low in MM0	5-11
Figure 5-5.	PEXTRW Instruction	5-12
Figure 5-6.	PINSRW Instruction	5-13
Figure 5-7.	PMOVSMKB Instruction	5-16
Figure 5-8.	Data Alignment of Loads and Stores in Reverse Memory Copy	5-41
Figure 5-9.	A Technique to Avoid Cacheline Split Loads in Reverse Memory Copy Using Two Aligned Loads	5-43
Figure 6-1.	Homogeneous Operation on Parallel Data Elements	6-4
Figure 6-2.	Horizontal Computation Model	6-4
Figure 6-3.	Dot Product Operation	6-5
Figure 6-4.	Horizontal Add Using MOVHLP/MOVHPS	6-11
Figure 6-5.	Asymmetric Arithmetic Operation of the SSE3 Instruction	6-14
Figure 6-6.	Horizontal Arithmetic Operation of the SSE3 Instruction HADDPD	6-15
Figure 7-1.	Effective Latency Reduction as a Function of Access Stride	7-15
Figure 7-2.	Memory Access Latency and Execution Without Prefetch	7-16
Figure 7-3.	Memory Access Latency and Execution With Prefetch	7-17
Figure 7-4.	Prefetch and Loop Unrolling	7-21

CONTENTS

	PAGE
Figure 7-5. Memory Access Latency and Execution With Prefetch	7-22
Figure 7-6. Spread Prefetch Instructions	7-23
Figure 7-7. Cache Blocking – Temporally Adjacent and Non-adjacent Passes	7-24
Figure 7-8. Examples of Prefetch and Strip-mining for Temporally Adjacent and Non-Adjacent Passes Loops	7-25
Figure 7-9. Single-Pass Vs. Multi-Pass 3D Geometry Engines	7-30
Figure 8-1. Amdahl’s Law and MP Speed-up	8-2
Figure 8-2. Single-threaded Execution of Producer-consumer Threading Model	8-6
Figure 8-3. Execution of Producer-consumer Threading Model on a Multicore Processor	8-7
Figure 8-4. Interlaced Variation of the Producer Consumer Model	8-8
Figure 8-5. Batched Approach of Producer Consumer Model	8-28
Figure 10-1. SSE4.2 String/Text Instruction Immediate Operand Control	10-2
Figure 10-2. Retrace Inefficiency of Byte-Granular, Brute-Force Search	10-17
Figure 10-3. SSE4.2 Speedup of SubString Searches	10-25
Figure 12-1. Performance History and State Transitions	12-2
Figure 12-2. Active Time Versus Halted Time of a Processor	12-4
Figure 12-3. Application of C-states to Idle Time	12-5
Figure 12-4. Profiles of Coarse Task Scheduling and Power Consumption	12-12
Figure 12-5. Thread Migration in a Multicore Processor	12-15
Figure 12-6. Progression to Deeper Sleep	12-15
Figure 12-7. Energy Saving due to Performance Optimization	12-17
Figure 12-8. Energy Saving due to Vectorization	12-18
Figure 12-9. Energy Saving Comparison of Synchronization Primitives	12-22
Figure 12-10. Power Saving Comparison of Power-Source-Aware Frame Rate Configurations .	12-23
Figure 13-1. Intel Atom Microarchitecture Pipeline	13-2
Figure A-1. Intel Thread Profiler Showing Critical Paths of Threaded Execution Timelines	A-16
Figure B-1. System Topology Supported by Intel® Xeon® Processor 5500 Series	B-1
Figure B-2. PMU Specific Event Logic Within the Pipeline	B-4
Figure B-3. LBR Records and Basic Blocks	B-19
Figure B-4. Using LBR Records to Rectify Skewed Sample Distribution	B-20
Figure B-5. RdData Request after LLC Miss to Local Home (Clean Rsp)	B-35
Figure B-6. RdData Request after LLC Miss to Remote Home (Clean Rsp)	B-35
Figure B-8. RdData Request after LLC Miss to Local Home (Hitm Response)	B-36
Figure B-7. RdData Request after LLC Miss to Remote Home (Hitm Response)	B-36
Figure B-9. RdData Request after LLC Miss to Local Home (Hit Response)	B-37
Figure B-10. RdInvOwn Request after LLC Miss to Remote Home (Clean Res)	B-37
Figure B-12. RdInvOwn Request after LLC Miss to Local Home (Hit Res)	B-38
Figure B-11. RdInvOwn Request after LLC Miss to Remote Home (Hitm Res)	B-38
Figure B-13. Performance Events Drill-Down and Software Tuning Feedback Loop	B-63
Figure D-1. Stack Frames Based on Alignment Type	D-2

TABLES

Table 2-1.	Components of the Front End of Intel microarchitecture code name Sandy Bridge	2-5
Table 2-2.	ICache and ITLB of Intel microarchitecture code name Sandy Bridge	2-6
Table 2-3.	Dispatch Port and Execution Stacks	2-14
Table 2-4.	Execution Core Writeback Latency (cycles)	2-15
Table 2-5.	Cache Parameters	2-16
Table 2-6.	Lookup Order and Load Latency	2-16
Table 2-7.	L1 Data Cache Components	2-18
Table 2-8.	Effect of Addressing Modes on Load Latency	2-19
Table 2-9.	DTLB and STLB Parameters	2-20
Table 2-10.	Store Forwarding Conditions (1 and 2 byte stores)	2-20
Table 2-11.	Store Forwarding Conditions (4-16 byte stores)	2-21
Table 2-12.	32-byte Store Forwarding Conditions (0-15 byte alignment)	2-21
Table 2-13.	32-byte Store Forwarding Conditions (16-31 byte alignment)	2-22
Table 2-14.	Components of the Front End	2-29
Table 2-15.	Issue Ports of Intel Core Microarchitecture and Enhanced Intel Core Microarchitecture	2-35
Table 2-16.	Cache Parameters of Processors based on Intel Core Microarchitecture	2-43
Table 2-17.	Characteristics of Load and Store Operations in Intel Core Microarchitecture	2-44
Table 2-18.	Bypass Delay Between Producer and Consumer Micro-ops (cycles)	2-49
Table 2-19.	Issue Ports of Intel microarchitecture code name Nehalem	2-50
Table 2-20.	Cache Parameters of Intel Core i7 Processors	2-51
Table 2-21.	Performance Impact of Address Alignments of MOVDQU from L1	2-53
Table 2-22.	Pentium 4 and Intel Xeon Processor Cache Parameters	2-66
Table 2-23.	Trigger Threshold and CPUID Signatures for Processor Families	2-73
Table 2-24.	Cache Parameters of Pentium M, Intel Core Solo, and Intel Core Duo Processors	2-73
Table 2-25.	Family And Model Designations of Microarchitectures	2-81
Table 2-26.	Characteristics of Load and Store Operations in Intel Core Duo Processors	2-82
Table 3-1.	Macro-Fusible Instructions in Intel microarchitecture code name Sandy Bridge	3-19
Table 3-2.	Store Forwarding Restrictions of Processors Based on Intel Core Microarchitecture	3-67
Table 5-1.	PSHUF Encoding	5-17
Table 6-1.	SoA Form of Representing Vertices Data	6-5
Table 7-1.	Software Prefetching Considerations into Strip-mining Code	7-27
Table 7-2.	Relative Performance of Memory Copy Routines	7-37
Table 7-3.	Deterministic Cache Parameters Leaf	7-39
Table 8-1.	Properties of Synchronization Objects	8-15
Table 8-2.	Design-Time Resource Management Choices	8-31
Table 8-3.	Microarchitectural Resources Comparisons of HT Implementations	8-36
Table 10-1.	SSE4.2 String/Text Instructions Compare Operation on N-elements	10-3
Table 10-2.	SSE4.2 String/Text Instructions Unary Transformation on IntRes1	10-3
Table 10-3.	SSE4.2 String/Text Instructions Output Selection Imm[6]	10-4

	PAGE
Table 10-4.	SSE4.2 String/Text Instructions Element-Pair Comparison Definition. 10-4
Table 10-5.	SSE4.2 String/Text Instructions Eflags Behavior. 10-5
Table 11-1.	Features between 256-bit AVX, 128-bit AVX and Legacy SSE Extensions 11-2
Table 11-2.	State Transitions of Mixing AVX and SSE Code 11-9
Table 11-3.	Single-Precision Divide and Square Root Alternatives 11-41
Table 11-4.	Comparison of Single-Precision Divide Alternatives 11-44
Table 11-5.	Comparison of Single-Precision Reciprocal Square Root Operation 11-46
Table 11-6.	Comparison of Single-Precision Square Root Operation. 11-48
Table 12-1.	ACPI C-State Type Mappings to Processor Specific C-State for Mobile Processors Based on Intel microarchitecture code name Nehalem12-6
Table 12-2.	ACPI C-State Type Mappings to Processor Specific C-State of Intel microarchitecture code name Sandy Bridge12-7
Table 12-3.	C-State Total Processor Exit Latency for Client Systems (Core+ Package Exit Latency) with Slow VR12-24
Table 12-4.	C-State Total Processor Exit Latency for Client Systems (Core+ Package Exit Latency) with Fast VR12-25
Table 12-5.	C-State Core-Only Exit Latency for Client Systems with Slow VR. 12-25
Table 12-6.	POWER_CTL MSR in Next Generation Intel Processor (Intel® microarchitecture code name Sandy Bridge)12-26
Table 13-1.	Instruction Latency/Throughput Summary of Intel® Atom™ Microarchitecture. . 13-10
Table 13-2.	Intel® Atom™ Microarchitecture Instructions Latency Data 13-19
Table A-1.	Recommended IA-32 Processor Optimization Options. A-2
Table A-2.	Recommended Processor Optimization Options for 64-bit Code A-4
Table A-3.	Vectorization Control Switch Options A-5
Table B-1.	Cycle Accounting and Micro-ops Flow Recipe. B-3
Table B-2.	CMask/Inv/Edge/Thread Granularity of Events for Micro-op Flow B-4
Table B-3.	Cycle Accounting of Wasted Work Due to Misprediction. B-6
Table B-4.	Cycle Accounting of Instruction Starvation. B-7
Table B-5.	CMask/Inv/Edge/Thread Granularity of Events for Micro-op Flow B-8
Table B-6.	Approximate Latency of L2 Misses of Intel Xeon Processor 5500. B-11
Table B-7.	Load Latency Event Programming B-14
Table B-8.	Data Source Encoding for Load Latency PEBS Record B-15
Table B-9.	Core PMU Events to Drill Down L2 Misses B-21
Table B-10.	Core PMU Events for Super Queue Operation. B-22
Table B-11.	Core PMU Event to Drill Down OFFCore Responses B-22
Table B-12.	OFFCORE_RSP_0 MSR Programming B-22
Table B-13.	Common Request and Response Types for OFFCORE_RSP_0 MSR. B-23
Table B-14.	Uncore PMU Events for Occupancy Cycles. B-30
Table B-15.	Common QHL Opcode Matching Facility Programming. B-33
Table C-1.	SIMD Instruction Extensions Support by Microarchitectures (CPUID Signature) . . C-4
Table C-2.	256-bit AVX Instructions. C-5
Table C-3.	AESNI and PCLMULQDQ Instructions C-6
Table C-4.	SSE4.2 Instructions C-7
Table C-5.	SSE4.1 Instructions C-7
Table C-6.	Supplemental Streaming SIMD Extension 3 Instructions. C-9
Table C-7.	Streaming SIMD Extension 3 SIMD Floating-point Instructions C-10
Table C-8.	Streaming SIMD Extension 2 128-bit Integer Instructions C-11

Table C-9. Streaming SIMD Extension 2 Double-precision Floating-point Instructions C-16

Table C-10. Streaming SIMD Extension Single-precision Floating-point Instructions C-20

Table C-11. Streaming SIMD Extension 64-bit Integer InstructionsC-24

Table C-12. MMX Technology 64-bit InstructionsC-25

Table C-13. MMX Technology 64-bit InstructionsC-26

Table C-14. x87 Floating-point Instructions.....C-27

Table C-15. General Purpose Instructions.....C-30

CHAPTER 1 INTRODUCTION

The *Intel® 64 and IA-32 Architectures Optimization Reference Manual* describes how to optimize software to take advantage of the performance characteristics of IA-32 and Intel 64 architecture processors. Optimizations described in this manual apply to processors based on the Intel® Core™ microarchitecture, Enhanced Intel® Core™ microarchitecture, Intel® microarchitecture code name Nehalem, Intel® microarchitecture code name Westmere, Intel® microarchitecture code name Sandy Bridge, Intel NetBurst® microarchitecture, the Intel® Core™ Duo, Intel® Core™ Solo, Pentium® M processor families.

The target audience for this manual includes software programmers and compiler writers. This manual assumes that the reader is familiar with the basics of the IA-32 architecture and has access to the *Intel® 64 and IA-32 Architectures Software Developer's Manual* (five volumes). A detailed understanding of Intel 64 and IA-32 processors is often required. In many cases, knowledge of the underlying microarchitectures is required.

The design guidelines that are discussed in this manual for developing high-performance software generally apply to current as well as to future IA-32 and Intel 64 processors. The coding rules and code optimization techniques listed target the Intel Core microarchitecture, the Intel NetBurst microarchitecture and the Pentium M processor microarchitecture. In most cases, coding rules apply to software running in 64-bit mode of Intel 64 architecture, compatibility mode of Intel 64 architecture, and IA-32 modes (IA-32 modes are supported in IA-32 and Intel 64 architectures). Coding rules specific to 64-bit modes are noted separately.

1.1 TUNING YOUR APPLICATION

Tuning an application for high performance on any Intel 64 or IA-32 processor requires understanding and basic skills in:

- Intel 64 and IA-32 architecture
- C and Assembly language
- hot-spot regions in the application that have impact on performance
- optimization capabilities of the compiler
- techniques used to evaluate application performance

The Intel® VTune™ Performance Analyzer can help you analyze and locate hot-spot regions in your applications. On the Intel® Core™ i7, Intel® Core™2 Duo, Intel® Core™ Duo, Intel® Core™ Solo, Pentium® 4, Intel® Xeon® and Pentium® M processors, this tool can monitor an application through a selection of performance moni-

toring events and analyze the performance event data that is gathered during code execution.

This manual also describes information that can be gathered using the performance counters through Pentium 4 processor's performance monitoring events.

1.2 ABOUT THIS MANUAL

The Intel® Xeon® processor 3000, 3200, 5100, 5300, 7200 and 7300 series, Intel® Pentium® dual-core, Intel® Core™2 Duo, Intel® Core™2 Quad, and Intel® Core™2 Extreme processors are based on Intel® Core™ microarchitecture. In this document, references to the Core 2 Duo processor refer to processors based on the Intel® Core™ microarchitecture.

The Intel® Xeon® processor 3100, 3300, 5200, 5400, 7400 series, Intel® Core™2 Quad processor Q8000 series, and Intel® Core™2 Extreme processors QX9000 series are based on 45nm Enhanced Intel® Core™ microarchitecture.

The Intel® Core™ i7 processor and Intel® Xeon® processor 3400, 5500, 7500 series are based on 45 nm Intel® microarchitecture code name Nehalem. Intel® microarchitecture code name Westmere is a 32nm version of Intel® microarchitecture code name Nehalem. Intel® Xeon® processor 5600 series and various Intel Core i7, i5, i3 processors are based on Intel® microarchitecture code name Westmere.

Intel® Core™ i7 processor 2xxx series are based on Intel® microarchitecture code name Sandy Bridge.

In this document, references to the Pentium 4 processor refer to processors based on the Intel NetBurst® microarchitecture. This includes the Intel Pentium 4 processor and many Intel Xeon processors based on Intel NetBurst microarchitecture. Where appropriate, differences are noted (for example, some Intel Xeon processors have third level cache).

The Dual-core Intel® Xeon® processor LV is based on the same architecture as Intel® Core™ Duo and Intel® Core™ Solo processors.

Intel® Atom™ processor is based on Intel® Atom™ microarchitecture.

The following bullets summarize chapters in this manual.

- **Chapter 1: Introduction** — Defines the purpose and outlines the contents of this manual.
- **Chapter 2: Intel® 64 and IA-32 Processor Architectures** — Describes the microarchitecture of recent IA-32 and Intel 64 processor families, and other features relevant to software optimization.
- **Chapter 3: General Optimization Guidelines** — Describes general code development and optimization techniques that apply to all applications designed to take advantage of the common features of the Intel Core microarchitecture, Enhanced Intel Core microarchitecture, Intel NetBurst microarchitecture and Pentium M processor microarchitecture.

- **Chapter 4: Coding for SIMD Architectures** — Describes techniques and concepts for using the SIMD integer and SIMD floating-point instructions provided by the MMX™ technology, Streaming SIMD Extensions, Streaming SIMD Extensions 2, Streaming SIMD Extensions 3, SSSE3, and SSE4.1.
- **Chapter 5: Optimizing for SIMD Integer Applications** — Provides optimization suggestions and common building blocks for applications that use the 128-bit SIMD integer instructions.
- **Chapter 6: Optimizing for SIMD Floating-point Applications** — Provides optimization suggestions and common building blocks for applications that use the single-precision and double-precision SIMD floating-point instructions.
- **Chapter 7: Optimizing Cache Usage** — Describes how to use the PREFETCH instruction, cache control management instructions to optimize cache usage, and the deterministic cache parameters.
- **Chapter 8: Multiprocessor and Hyper-Threading Technology** — Describes guidelines and techniques for optimizing multithreaded applications to achieve optimal performance scaling. Use these when targeting multicore processor, processors supporting Hyper-Threading Technology, or multiprocessor (MP) systems.
- **Chapter 9: 64-Bit Mode Coding Guidelines** — This chapter describes a set of additional coding guidelines for application software written to run in 64-bit mode.
- **Chapter 10: SSE4.2 and SIMD Programming for Text-Processing/Lexing/Parsing**— Describes SIMD techniques of using SSE4.2 along with other instruction extensions to improve text/string processing and lexing/parsing applications.
- **Chapter 11: Power Optimization for Mobile Usages** — This chapter provides background on power saving techniques in mobile processors and makes recommendations that developers can leverage to provide longer battery life.
- **Chapter 12: Intel® Atom™ Processor Architecture and Optimization** — Describes the microarchitecture of processor families based on Intel Atom microarchitecture, and software optimization techniques targeting Intel Atom microarchitecture.
- **Appendix A: Application Performance Tools** — Introduces tools for analyzing and enhancing application performance without having to write assembly code.
- **Appendix B: Intel® Pentium® 4 Processor Performance Metrics** — Provides information that can be gathered using Pentium 4 processor's performance monitoring events. These performance metrics can help programmers determine how effectively an application is using the features of the Intel NetBurst microarchitecture.
- **Appendix C: IA-32 Instruction Latency and Throughput** — Provides latency and throughput data for the IA-32 instructions. Instruction timing data specific to recent processor families are provided.

- **Appendix D: Stack Alignment** — Describes stack alignment conventions and techniques to optimize performance of accessing stack-based data.
- **Appendix E: Summary of Rules and Suggestions** — Summarizes the rules and tuning suggestions referenced in the manual.

1.3 RELATED INFORMATION

For more information on the Intel® architecture, techniques, and the processor architecture terminology, the following are of particular interest:

- *Intel® 64 and IA-32 Architectures Software Developer's Manual* (in five volumes)
- *Intel® Processor Identification with the CPUID Instruction, AP-485*
- *Developing Multi-threaded Applications: A Platform Consistent Approach*
- Intel® C++ Compiler documentation and online help
- Intel® Fortran Compiler documentation and online help
- Intel® VTune™ Performance Analyzer documentation and online help
- *Using Spin-Loops on Intel Pentium 4 Processor and Intel Xeon Processor MP*

More relevant links are:

- Software network link:
<http://softwarecommunity.intel.com/isn/home/>
- Developer centers:
<http://www3.intel.com/cd/ids/developer/asmo-na/eng/dc/index.htm>
- Processor support general link:
<http://www.intel.com/support/processors/>
- Software products and packages:
<http://www3.intel.com/cd/software/products/asmo-na/eng/index.htm>
- Intel 64 and IA-32 processor manuals (printed or PDF downloads):
<http://developer.intel.com/products/processor/manuals/index.htm>
- Intel Multi-Core Technology:
<http://developer.intel.com/technology/multi-core/index.htm>
- Hyper-Threading Technology (HT Technology):
<http://developer.intel.com/technology/hyperthread/>
- SSE4.1 Application Note: Motion Estimation with Intel® Streaming SIMD Extensions 4
<http://softwarecommunity.intel.com/articles/eng/1246.htm>
- SSE4.1 Application Note: Increasing Memory Throughput with Intel® Streaming SIMD Extensions 4

<http://softwarecommunity.intel.com/articles/eng/1248.htm>

- Processor Topology and Cache Topology white paper and reference code
<http://software.intel.com/en-us/articles/intel-64-architecture-processor-topology-enumeration>

INTRODUCTION

CHAPTER 2

INTEL® 64 AND IA-32 PROCESSOR ARCHITECTURES

This chapter gives an overview of features relevant to software optimization for current generations of Intel 64 and IA-32 processors (processors based on Intel microarchitecture code name Sandy Bridge, processors based on the Intel Core microarchitecture, Enhanced Intel Core microarchitecture, Intel microarchitecture code name Nehalem, Intel NetBurst microarchitecture; including Intel Core Solo, Intel Core Duo, and Intel Pentium M processors). These features are:

- Microarchitectures that enable executing instructions with high throughput at high clock rates, a high speed cache hierarchy and high speed system bus
- Multicore architecture available in Intel Core i7, Intel Core 2 Extreme, Intel Core 2 Quad, Intel Core 2 Duo, Intel Core Duo, Intel Pentium D processors, Pentium processor Extreme Edition¹, and Quad-core Intel Xeon, Dual-core Intel Xeon processors
- Hyper-Threading Technology² (HT Technology) support
- Intel 64 architecture on Intel 64 processors
- SIMD instruction extensions: MMX technology, Streaming SIMD Extensions (SSE), Streaming SIMD Extensions 2 (SSE2), Streaming SIMD Extensions 3 (SSE3), Supplemental Streaming SIMD Extensions 3 (SSSE3), SSE4.1, and SSE4.2.
- Intel® Advanced Vector Extensions (Intel® AVX).

The Intel Pentium M processor introduced a power-efficient microarchitecture with balanced performance. Dual-core Intel Xeon processor LV, Intel Core Solo and Intel Core Duo processors incorporate enhanced Pentium M processor microarchitecture. The Intel Core 2, Intel Core 2 Extreme, Intel Core 2 Quad processor family, Intel Xeon processor 3000, 3200, 5100, 5300, 7300 series are based on the high-performance and power-efficient Intel Core microarchitecture. Intel Xeon processor 3100, 3300, 5200, 5400, 7400 series, Intel Core 2 Extreme processor QX9600, QX9700

1. Quad-core platforms require an Intel Xeon processor 3200, 3300, 5300, 5400, 7300 series, an Intel Core 2 Extreme processor QX6000, QX9000 series, or an Intel Core 2 Quad processor, with appropriate chipset, BIOS, and operating system. Six-core platform requires an Intel Xeon processor 7400 series, with appropriate chipset, BIOS, and operating system. Dual-core platform requires an Intel Xeon processor 3000, 3100 series, Intel Xeon processor 5100, 5200, 7100 series, Intel Core 2 Duo, Intel Core 2 Extreme processor X6800, Dual-core Intel Xeon processors, Intel Core Duo, Pentium D processor or Pentium processor Extreme Edition, with appropriate chipset, BIOS, and operating system. Performance varies depending on the hardware and software used.
2. Hyper-Threading Technology requires a computer system with an Intel processor supporting HT Technology and an HT Technology enabled chipset, BIOS and operating system. Performance varies depending on the hardware and software used.

series, Intel Core 2 Quad Q9000 series, Q8000 series are based on the enhanced Intel Core microarchitecture. Intel Core i7 processor is based on Intel microarchitecture code name Nehalem. Intel® Xeon® processor 5600 series and Intel Core i7, i5, i3 processors are based on Intel microarchitecture code name Westmere.

Intel® Core™ i7 processor 2xxx series are based on Intel® microarchitecture code name Sandy Bridge.

Intel Core 2 Extreme QX6700 processor, Intel Core 2 Quad processors, Intel Xeon processors 3200 series, 5300 series are quad-core processors. Intel Pentium 4 processors, Intel Xeon processors, Pentium D processors, and Pentium processor Extreme Editions are based on Intel NetBurst microarchitecture.

2.1 INTEL® MICROARCHITECTURE CODE NAME SANDY BRIDGE

Intel® microarchitecture code name Sandy Bridge builds on the successes of Intel® Core™ microarchitecture and Intel microarchitecture code name Nehalem. It offers the following innovative features:

- Intel Advanced Vector Extensions (Intel AVX)
 - 256-bit floating-point instruction set extensions to the 128-bit Intel Streaming SIMD Extensions, providing up to 2X performance benefits relative to 128-bit code.
 - Non-destructive destination encoding offers more flexible coding techniques.
 - Supports flexible migration and co-existence between 256-bit AVX code, 128-bit AVX code and legacy 128-bit SSE code.
- Enhanced front-end and execution engine
 - New decoded Icache component that improves front-end bandwidth and reduces branch misprediction penalty.
 - Advanced branch prediction.
 - Additional macro-fusion support.
 - Larger dynamic execution window.
 - Multi-precision integer arithmetic enhancements (ADC/SBB, MUL/IMUL).
 - LEA bandwidth improvement.
 - Reduction of general execution stalls (read ports, writeback conflicts, bypass latency, partial stalls).
 - Fast floating-point exception handling.
 - XSAVE/XRSTORE performance improvements and XSAVEOPT new instruction.
- Cache hierarchy improvements for wider data path

- Doubling of bandwidth enabled by two symmetric ports for memory operation.
- Simultaneous handling of more in-flight loads and stores enabled by increased buffers.
- Internal bandwidth of two loads and one store each cycle.
- Improved prefetching.
- High bandwidth low latency LLC architecture.
- High bandwidth ring architecture of on-die interconnect.
- System-on-a-chip support
 - Integrated graphics and media engine in second generation Intel Core processors.
 - Integrated PCIE controller.
 - Integrated memory controller.
- Next generation Intel Turbo Boost Technology
 - Leverage TDP headroom to boost performance of CPU cores and integrated graphic unit.

2.1.1 Intel® microarchitecture code name Sandy Bridge Pipeline Overview

Figure 2-1 depicts the pipeline and major components of a processor core that's based on Intel microarchitecture code name Sandy Bridge. The pipeline consists of

- An in-order issue front end that fetches instructions and decodes them into micro-ops (micro-operations). The front end feeds the next pipeline stages with a continuous stream of micro-ops from the most likely path that the program will execute.
- An out-of-order, superscalar execution engine that dispatches up to six micro-ops to execution, per cycle. The allocate/rename block reorders micro-ops to "dataflow" order so they can execute as soon as their sources are ready and execution resources are available.
- An in-order retirement unit that ensures that the results of execution of the micro-ops, including any exceptions they may have encountered, are visible according to the original program order.

The flow of an instruction in the pipeline can be summarized in the following progression:

1. The Branch Prediction Unit chooses the next block of code to execute from the program. The processor searches for the code in the following resources, in this order:
 - a. Decoded ICache

- b. Instruction Cache, via activating the legacy decode pipeline
- c. L2 cache, last level cache (LLC) and memory, as necessary

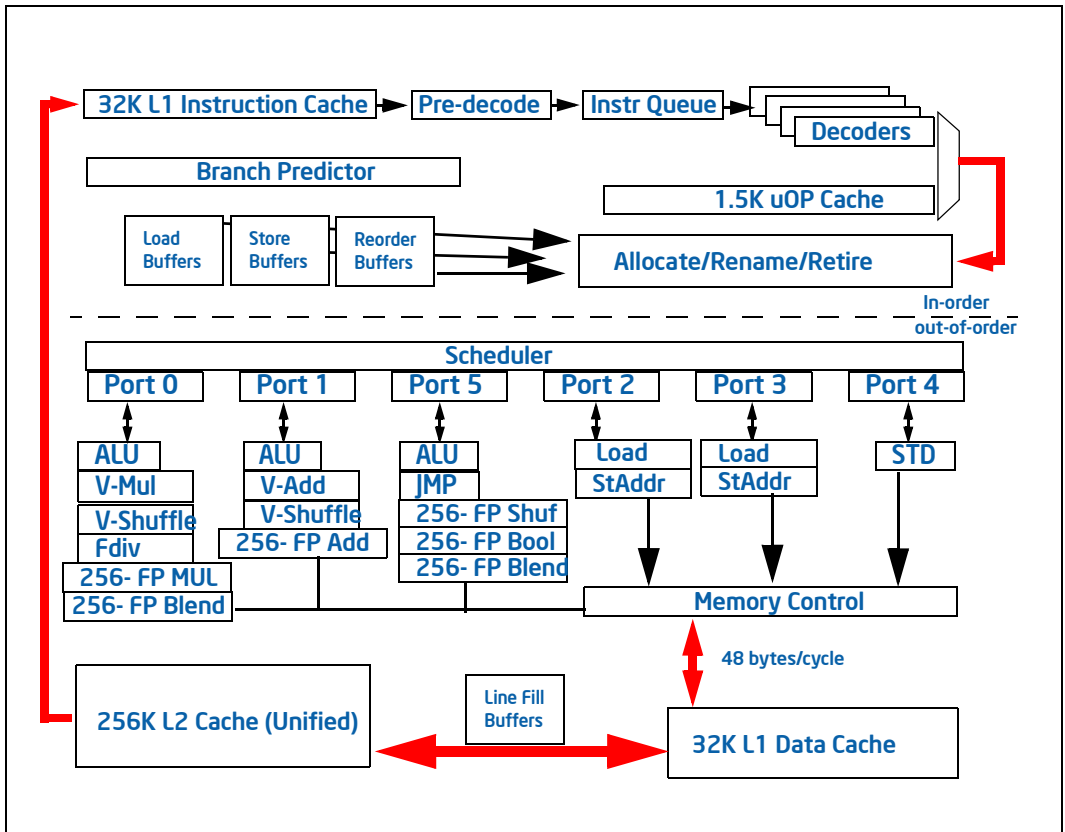


Figure 2-1. Intel microarchitecture code name Sandy Bridge Pipeline Functionality

2. The micro-ops corresponding to this code are sent to the Rename/retirement block. They enter into the scheduler in program order, but execute and are de-allocated from the scheduler according to data-flow order. For simultaneously ready micro-ops, FIFO ordering is nearly always maintained.

Micro-op execution is executed using execution resources arranged in three stacks. The execution units in each stack are associated with the data type of the instruction.

Branch mispredictions are signaled at branch execution. It re-steers the front end which delivers micro-ops from the correct path. The processor can overlap work preceding the branch misprediction with work from the following corrected path.

3. Memory operations are managed and reordered to achieve parallelism and maximum performance. Misses to the L1 data cache go to the L2 cache. The data cache is non-blocking and can handle multiple simultaneous misses.
4. Exceptions (Faults, Traps) are signaled at retirement (or attempted retirement) of the faulting instruction.

Each processor core based on Intel microarchitecture code name Sandy Bridge can support two logical processor if Intel HyperThreading Technology is enabled.

2.1.2 The Front End

This section describes the key characteristics of the front end. Table 2-1 lists the components of the front end, their functions, and the problems they address.

Table 2-1. Components of the Front End of Intel microarchitecture code name Sandy Bridge

Component	Functions	Performance Challenges
Instruction Cache	32-Kbyte backing store of instruction bytes	Fast access to hot code instruction bytes
Legacy Decode Pipeline	Decode instructions to micro-ops, delivered to the micro-op queue and the Decoded ICache.	Provides the same decode latency and bandwidth as prior Intel processors. Decoded ICache warm-up
Decoded ICache	Provide stream of micro-ops to the micro-op queue.	Provides higher micro-op bandwidth at lower latency and lower power than the legacy decode pipeline
MSROM	Complex instruction micro-op flow store, accessible from both Legacy Decode Pipeline and Decoded ICache	
Branch Prediction Unit (BPU)	Determine next block of code to be executed and drive lookup of Decoded ICache and legacy decode pipelines.	Improves performance and energy efficiency through reduced branch mispredictions.
Micro-op queue	Queues micro-ops from the Decoded ICache and the legacy decode pipeline.	Hide front-end bubbles; provide execution micro-ops at a constant rate.

2.1.2.1 Legacy Decode Pipeline

The Legacy Decode Pipeline comprises the instruction translation lookaside buffer (ITLB), the instruction cache (ICache), instruction predecode, and instruction decode units.

Instruction Cache and ITLB

An instruction fetch is a 16-byte aligned lookup through the ITLB and into the instruction cache. The instruction cache can deliver every cycle 16 bytes to the instruction pre-decoder. Table 2-1 compares the ICache and ITLB with prior generation.

Table 2-2. ICache and ITLB of Intel microarchitecture code name Sandy Bridge

Component	Intel microarchitecture code name Sandy Bridge	Intel microarchitecture code name Nehalem
ICache Size	32-Kbyte	32-Kbyte
ICache Ways	8	4
ITLB 4K page entries	128	128
ITLB large page (2M or 4M) entries	8	7

Upon ITLB miss there is a lookup to the Second level TLB (STLB) that is common to the DTLB and the ITLB. The penalty of an ITLB miss and a STLB hit is seven cycles.

Instruction PreDecode

The predecode unit accepts the 16 bytes from the instruction cache and determines the length of the instructions.

The following length changing prefixes (LCPs) imply instruction length that is different from the default length of instructions. Therefore they cause an additional penalty of three cycles per LCP during length decoding. Previous processors incur a six-cycle penalty for each 16-byte chunk that has one or more LCPs in it. Since usually there is no more than one LCP in a 16-byte chunk, in most cases, Intel microarchitecture code name Sandy Bridge introduces an improvement over previous processors.

- Operand Size Override (66H) preceding an instruction with a word/double immediate data. This prefix might appear when the code uses 16 bit data types, unicode processing, and image processing.
- Address Size Override (67H) preceding an instruction with a `modr/m` in real, big real, 16-bit protected or 32-bit protected modes. This prefix may appear in boot code sequences.

- The REX prefix (4xh) in the Intel® 64 instruction set can change the size of two classes of instructions: MOV offset and MOV immediate. Despite this capability, it does not cause an LCP penalty and hence is not considered an LCP.

Instruction Decode

There are four decoding units that decode instruction into micro-ops. The first can decode all IA-32 and Intel 64 instructions up to four micro-ops in size. The remaining three decoding units handle single-micro-op instructions. All four decoding units support the common cases of single micro-op flows including micro-fusion and macro-fusion.

Micro-ops emitted by the decoders are directed to the micro-op queue and to the Decoded ICache. Instructions longer than four micro-ops generate their micro-ops from the MSROM. The MSROM bandwidth is four micro-ops per cycle. Instructions whose micro-ops come from the MSROM can start from either the legacy decode pipeline or from the Decoded ICache

MicroFusion

Micro-fusion fuses multiple micro-ops from the same instruction into a single complex micro-op. The complex micro-op is dispatched in the out-of-order execution core as many times as it would if it were not micro-fused.

Micro-fusion enables you to use memory-to-register operations, also known as the complex instruction set computer (CISC) instruction set, to express the actual program operation without worrying about a loss of decode bandwidth. Micro-fusion improves instruction bandwidth delivered from decode to retirement and saves power.

Coding an instruction sequence by using single-uop instructions will increase the code size, which can decrease fetch bandwidth from the legacy pipeline.

The following are examples of micro-fused micro-ops that can be handled by all decoders.

- All stores to memory, including store immediate. Stores execute internally as two separate functions, store-address and store-data.
- All instructions that combine load and computation operations (load+op), for example:
 - `ADDPS XMM9, QWORD PTR [RSP+40]`
 - `FADD DOUBLE PTR [RDI+RSI*8]`
 - `XOR RAX, QWORD PTR [RBP+32]`
- All instructions of the form "load and jump," for example:
 - `JMP [RDI+200]`
 - `RET`
- `CMP` and `TEST` with immediate operand and memory

An instruction with RIP relative addressing is not micro-fused in the following cases:

- An additional immediate is needed, for example:
 - `CMP [RIP+400], 27`
 - `MOV [RIP+3000], 142`
- The instruction is a control flow instruction with an indirect target specified using RIP-relative addressing, for example:
 - `JMP [RIP+5000000]`

In these cases, an instruction that can not be micro-fused will require decoder 0 to issue two micro-ops, resulting in a slight loss of decode bandwidth.

In 64-bit code, the usage of RIP Relative addressing is common for global data. Since there is no micro-fusion in these cases, performance may be reduced when porting 32-bit code to 64-bit code.

Macro-Fusion

Macro-fusion merges two instructions into a single micro-op. In Intel Core microarchitecture, this hardware optimization is limited to specific conditions specific to the first and second of the macro-fusable instruction pair.

- The first instruction of the macro-fused pair modifies the flags. The following instructions can be macro-fused:
 - In Intel microarchitecture code name Nehalem: `CMP`, `TEST`.
 - In Intel microarchitecture code name Sandy Bridge: `CMP`, `TEST`, `ADD`, `SUB`, `AND`, `INC`, `DEC`
 - These instructions can fuse if
 - The first source / destination operand is a register.
 - The second source operand (if exists) is one of: immediate, register, or non RIP-relative memory.
- The second instruction of the macro-fusable pair is a conditional branch. Table 3-1 describes, for each instruction, what branches it can fuse with.

Macro fusion does not happen if the first instruction ends on byte 63 of a cache line, and the second instruction is a conditional branch that starts at byte 0 of the next cache line.

Since these pairs are common in many types of applications, macro-fusion improves performance even on non-recompiled binaries.

Each macro-fused instruction executes with a single dispatch. This reduces latency and frees execution resources. You also gain increased rename and retire bandwidth, increased virtual storage, and power savings from representing more work in fewer bits.

2.1.2.2 Decoded ICache

The Decoded ICache is essentially an accelerator of the legacy decode pipeline. By storing decoded instructions, the Decoded ICache enables the following features:

- Reduced latency on branch mispredictions
- Increased micro-op delivery bandwidth to the out-of-order engine
- Reduced front end power consumption

The Decoded ICache caches the output of the instruction decoder. The next time the micro-ops are consumed for execution the decoded micro-ops are taken from the Decoded ICache. This enables skipping the fetch and decode stages for these micro-ops and reduces power and latency of the Front End. The Decoded ICache provides average hit rates of above 80% of the micro-ops; furthermore, "hot spots" typically have hit rates close to 100%.

Typical integer programs average less than four bytes per instruction, and the front end is able to race ahead of the back end, filling in a large window for the scheduler to find instruction level parallelism. However, for high performance code with a basic block consisting of many instructions, for example, Intel SSE media algorithms or excessively unrolled loops, the 16 instruction bytes per cycle is occasionally a limitation. The 32-byte orientation of the Decoded ICache helps such code to avoid this limitation.

The Decoded ICache automatically improves performance of programs with temporal and spatial locality. However, to fully utilize the Decoded ICache potential, you might need to understand its internal organization.

The Decoded ICache consists of 32 sets. Each set contains eight Ways. Each Way can hold up to six micro-ops. The Decoded ICache can ideally hold up to 1536 micro-ops.

The following are some of the rules how the Decoded ICache is filled with micro-ops:

- All micro-ops in a Way represent instructions which are statically contiguous in the code and have their EIPs within the same aligned 32-byte region.
- Up to three Ways may be dedicated to the same 32-byte aligned chunk, allowing a total of 18 micro-ops to be cached per 32-byte region of the original IA program.
- A multi micro-op instruction cannot be split across Ways.
- Up to two branches are allowed per Way.
- An instruction which turns on the MSR0M consumes an entire Way.
- A non-conditional branch is the last micro-op in a Way.
- Micro-fused micro-ops (load+op and stores) are kept as one micro-op.
- A pair of macro-fused instructions is kept as one micro-op.
- Instructions with 64-bit immediate require two slots to hold the immediate.

When micro-ops cannot be stored in the Decoded ICache due to these restrictions, they are delivered from the legacy decode pipeline. Once micro-ops are delivered

from the legacy pipeline, fetching micro-ops from the Decoded ICache can resume only after the next branch micro-op. Frequent switches can incur a penalty.

The Decoded ICache is virtually included in the Instruction cache and ITLB. That is, any instruction with micro-ops in the Decoded ICache has its original instruction bytes present in the instruction cache. Instruction cache evictions must also be evicted from the Decoded ICache, which evicts only the necessary lines.

There are cases where the entire Decoded ICache is flushed. One reason for this can be an ITLB entry eviction. Other reasons are not usually visible to the application programmer, as they occur when important controls are changed, for example, mapping in CR3, or feature and mode enabling in CR0 and CR4. There are also cases where the Decoded ICache is disabled, for instance, when the CS base address is NOT set to zero.

2.1.2.3 Branch Prediction

Branch prediction predicts the branch target and enables the processor to begin executing instructions long before the branch true execution path is known. All branches utilize the branch prediction unit (BPU) for prediction. This unit predicts the target address not only based on the EIP of the branch but also based on the execution path through which execution reached this EIP. The BPU can efficiently predict the following branch types:

- Conditional branches
- direct calls and jumps
- indirect calls and jumps
- returns

2.1.2.4 Micro-op Queue and the Loop Stream Detector (LSD)

The micro-op queue decouples the front end and the out-of order engine. It stays between the micro-op generation and the renamer as shown in Figure 2-1. This queue helps to hide bubbles which are introduced between the various sources of micro-ops in the front end and ensures that four micro-ops are delivered for execution, each cycle.

The micro-op queue provides post-decode functionality for certain instructions types. In particular, loads combined with computational operations and all stores, when used with indexed addressing, are represented as a single micro-op in the decoder or Decoded ICache. In the micro-op queue they are fragmented into two micro-ops through a process called un-lamination, one does the load and the other does the operation. A typical example is the following "load plus operation" instruction:

```
ADD                                RAX, [RBP+RSI] ; rax := rax + LD( RBP+RSI )
```

Similarly, the following store instruction has three register sources and is broken into "generate store address" and "generate store data" sub-components.

MOV [ESP+ECX*4+12345678], AL

The additional micro-ops generated by unlamination use the rename and retirement bandwidth. However, it has an overall power benefit. For code that is dominated by indexed addressing (as often happens with array processing), recoding algorithms to use base (or base+displacement) addressing can sometimes improve performance by keeping the load plus operation and store instructions fused.

The Loop Stream Detector (LSD)

The Loop Stream Detector was introduced in Intel® Core microarchitectures. The LSD detects small loops that fit in the micro-op queue and locks them down. The loop streams from the micro-op queue, with no more fetching, decoding, or reading micro-ops from any of the caches, until a branch miss-prediction inevitably ends it.

The loops with the following attributes qualify for LSD/micro-op queue replay:

- Up to eight chunk fetches of 32-instruction-bytes
- Up to 28 micro-ops (~28 instructions)
- All micro-ops are also resident in the Decoded ICache
- Can contain no more than eight taken branches and none of them can be a CALL or RET
- Cannot have mismatched stack operations. For example, more PUSH than POP instructions.

Many calculation-intensive loops, searches and software string moves match these characteristics.

Use the loop cache functionality opportunistically. For high performance code, loop unrolling is generally preferable for performance even when it overflows the LSD capability.

2.1.3 The Out-of-Order Engine

The Out-of-Order engine provides improved performance over prior generations with excellent power characteristics. It detects dependency chains and sends them to execution out-of-order while maintaining the correct data flow. When a dependency chain is waiting for a resource, such as a second-level data cache line, it sends micro-ops from another chain to the execution core. This increases the overall rate of instructions executed per cycle (IPC).

The out-of-order engine consists of two blocks, shown in Figures 2-1: Core Functional Diagram, the Rename/retirement block, and the Scheduler.

The Out-of-Order engine contains the following major components:

Renamer. The Renamer component moves micro-ops from the front end to the execution core. It eliminates false dependencies among micro-ops, thereby enabling out-of-order execution of micro-ops.

Scheduler. The Scheduler component queues micro-ops until all source operands are ready. Schedules and dispatches ready micro-ops to the available execution units in as close to a first in first out (FIFO) order as possible.

Retirement. The Retirement component retires instructions and micro-ops in order and handles faults and exceptions.

2.1.3.1 Renamer

The Renamer is the bridge between the in-order part in Figures 2-1, and the dataflow world of the Scheduler. It moves up to four micro-ops every cycle from the micro-op queue to the out-of-order engine. Although the renamer can send up to 4 micro-ops (unfused, micro-fused, or macro-fused) per cycle, this is equivalent to the issue port can dispatch six micro-ops per cycle. In this process, the out-of-order core carries out the following steps:

- Renames architectural sources and destinations of the micro-ops to micro-architectural sources and destinations.
- Allocates resources to the micro-ops. For example, load or store buffers.
- Binds the micro-op to an appropriate dispatch port.

Some micro-ops can execute to completion during rename and are removed from the pipeline at that point, effectively costing no execution bandwidth. These include:

- Zero idioms (dependency breaking idioms)
- NOP
- VZEROUPPER
- FXCHG

The renamer can allocate two branches each cycle, compared to one branch each cycle in the previous microarchitecture. This can eliminate some bubbles in execution.

Micro-fused load and store operations that use an index register are decomposed to two micro-ops, hence consume two out of the four slots the Renamer can use every cycle.

Dependency Breaking Idioms

Instruction parallelism can be improved by using common instructions to clear register contents to zero. The renamer can detect them on the zero evaluation of the destination register.

Use one of these dependency breaking idioms to clear a register when possible.

- XOR REG,REG
- SUB REG,REG
- PXOR/VPXOR XMMREG,XMMREG
- PSUBB/W/D/Q XMMREG,XMMREG

- VPSUBB/W/D/Q XMMREG, XMMREG
- XORPS/PD XMMREG, XMMREG
- VXORPS/PD YMMREG, YMMREG

Since zero idioms are detected and removed by the renamer, they have no execution latency.

There is another dependency breaking idiom - the "ones idiom".

- CMPEQ XMM1, XMM1; "ones idiom" set all elements to all "ones"

In this case, the micro-op must execute, however, since it is known that regardless of the input data the output data is always "all ones" the micro-op dependency upon its sources does not exist as with the zero idiom and it can execute as soon as it finds a free execution port.

2.1.3.2 Scheduler

The scheduler controls the dispatch of micro-ops onto their execution ports. In order to do this, it must identify which micro-ops are ready and where its sources come from: a register file entry, or a bypass directly from an execution unit. Depending on the availability of dispatch ports and writeback buses, and the priority of ready micro-ops, the scheduler selects which micro-ops are dispatched every cycle.

2.1.4 The Execution Core

The execution core is superscalar and can process instructions out of order. The execution core optimizes overall performance by handling the most common operations efficiently, while minimizing potential delays.

The out-of-order execution core improves execution unit organization over prior generation in the following ways:

- Reduction in read port stalls
- Reduction in writeback conflicts and delays
- Reduction in power

The out-of-order core consist of three execution stacks, where each stack encapsulates a certain type of data. The execution core contains the following execution stacks:

- General purpose integer
- SIMD integer and floating point
- X87

The execution core also contains connections to and from the cache hierarchy. The loaded data is fetched from the caches and written back into one of the stacks.

The scheduler can dispatch up to six micro-ops every cycle, one on each port. The following table summarizes which operations can be dispatched on which port.

Table 2-3. Dispatch Port and Execution Stacks

	Port 0	Port 1	Port 2	Port 3	Port 4	Port 5
Integer	ALU, Shift	ALU, Fast LEA, Slow LEA, MUL	Load_Ad dr, Store_ad dr	Load_Ad dr Store_ad dr	Store_dat a	ALU, Shift, Branch, Fast LEA
SSE-Int, AVX-Int, MMX	Mul, Shift, STTNI, Int- Div, 128b-Mov	ALU, Shuf, Blend, 128b-Mov			Store_dat a	ALU, Shuf, Shift, Blend, 128b-Mov
SSE-FP, AVX- FP_low	Mul, Div, Blend, 256b-Mov	Add, CVT			Store_dat a	Shuf, Blend, 256b-Mov
X87, AVX- FP_High	Mul, Div, Blend, 256b-Mov	Add, CVT			Store_dat a	Shuf, Blend, 256b-Mov

After execution, the data is written back on a writeback bus corresponding to the dispatch port and the data type of the result. Micro-ops that are dispatched on the same port but have different latencies may need the write back bus at the same cycle. In these cases the execution of one of the micro-ops is delayed until the write-back bus is available. For example, MULPS (five cycles) and BLENDPS (one cycle) may collide if both are ready for execution on port 0: first the MULPS and four cycles later the BLENDPS. Intel microarchitecture code name Sandy Bridge eliminates such collisions as long as the micro-ops write the results to different stacks. For example, integer ADD (one cycle) can be dispatched four cycles after MULPS (five cycles) since the integer ADD uses the integer stack while the MULPS uses the FP stack.

When a source of a micro-op executed in one stack comes from a micro-op executed in another stack, a one- or two-cycle delay can occur. The delay occurs also for transitions between Intel SSE integer and Intel SSE floating-point operations. In some of the cases the data transition is done using a micro-op that is added to the instruction flow. The following table describes how data, written back after execution, can bypass to micro-op execution in the following cycles

Table 2-4. Execution Core Writeback Latency (cycles)

	Integer	SSE-Int, AVX-Int, MMX	SSE-FP, AVX-FP_low	X87, AVX-FP_High
Integer	0	micro-op (port 0)	micro-op (port 0)	micro-op (port 0) + 1 cycle
SSE-Int, AVX-Int, MMX	micro-op (port 5) or micro-op (port 5) +1 cycle	0	1 cycle delay	0
SSE-FP, AVX-FP_low	micro-op (port 5) or micro-op (port 5) +1 cycle	1 cycle delay	0	micro-op (port 5) +1 cycle
X87, AVX-FP_High	micro-op (port 5) +1 cycle	0	micro-op (port 5) +1 cycle	0
Load	0	1 cycle delay	1 cycle delay	2 cycle delay

2.1.5 Cache Hierarchy

The cache hierarchy contains a first level instruction cache, a first level data cache (L1 DCache) and a second level (L2) cache, in each core. The L1D cache may be shared by two logical processors if the processor support Intel HyperThreading Technology. The L2 cache is shared by instructions and data. All cores in a physical processor package connect to a shared last level cache (LLC) via a ring connection.

The caches use the services of the Instruction Translation Lookaside Buffer (ITLB), Data Translation Lookaside Buffer (DTLB) and Shared Translation Lookaside Buffer (STLB) to translate linear addresses to physical address. Data coherency in all cache levels is maintained using the MESI protocol. For more information, see the Intel® 64 IA-32 Architectures Software Developer's Manual, Volume 3. Cache hierarchy details can be obtained at run-time using the CPUID instruction. see the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A*.

Table 2-5. Cache Parameters

Level	Capacity	Associativity (ways)	Line Size (bytes)	Write Update Policy	Inclusive
L1 Data	32 KB	8	64	Writeback	-
Instruction	32 KB	8	N/A	N/A	-
L2 (Unified)	256 KB	8	64	Writeback	No
Third Level (LLC)	0.5, 1, 1.5, or 2 MB	Varies with cache size	64	Writeback	Yes

2.1.5.1 Load and Store Operation Overview

This section provides an overview of the load and store operations.

Loads

When an instruction reads data from a memory location that has write-back (WB) type, the processor looks for it in the caches and memory. Table 2-6 shows the access lookup order and best case latency. The actual latency can vary depending on the cache queue occupancy, LLC ring occupancy, memory components, and their parameters.

Table 2-6. Lookup Order and Load Latency

Level	Latency (cycles)	Bandwidth (per core per cycle)
L1 Data	4 ¹	2 x 16 bytes
L2 (Unified)	12	1 x 32 bytes
Third Level (LLC)	26-31	1 x 32 bytes
L2 and L1 DCache in other cores if applicable	43- clean hit; 60 - dirty hit	

NOTES:

1. Subject to execution core bypass restriction shown in Table 2-4.

The LLC is inclusive of all cache levels above it - data contained in the core caches must also reside in the LLC. Each cache line in the LLC holds an indication of the cores that may have this line in their L2 and L1 caches. If there is an indication in the LLC that other cores may hold the line of interest and its state might have to modify, there is a lookup into the L1 DCache and L2 of these cores too. The lookup is called "clean" if it does not require fetching data from the other core caches. The lookup is called "dirty" if modified data has to be fetched from the other core caches and transferred to the loading core.

The latencies shown above are the best-case scenarios. Sometimes a modified cache line has to be evicted to make space for a new cache line. The modified cache line is evicted in parallel to bringing the new data and does not require additional latency. However, when data is written back to memory, the eviction uses cache bandwidth and possibly memory bandwidth as well. Therefore, when multiple cache misses require the eviction of modified lines within a short time, there is an overall degradation in cache response time. Memory access latencies vary based on occupancy of the memory controller queues, DRAM configuration, DDR parameters, and DDR paging behavior (if the requested page is a page-hit, page-miss or page-empty).

Stores

When an instruction writes data to a memory location that has a write back memory type, the processor first ensures that it has the line containing this memory location in its L1 DCache, in Exclusive or Modified MESI state. If the cache line is not there, in the right state, the processor fetches it from the next levels of the memory hierarchy using a Read for Ownership request. The processor looks for the cache line in the following locations, in the specified order:

1. L1 DCache
2. L2
3. Last Level Cache
4. L2 and L1 DCache in other cores, if applicable
5. Memory

Once the cache line is in the L1 DCache, the new data is written to it, and the line is marked as Modified.

Reading for ownership and storing the data happens after instruction retirement and follows the order of store instruction retirement. Therefore, the store latency usually does not affect the store instruction itself. However, several sequential stores that miss the L1 DCache may have cumulative latency that can affect performance. As long as the store does not complete, its entry remains occupied in the store buffer. When the store buffer becomes full, new micro-ops cannot enter the execution pipe and execution might stall.

2.1.5.2 L1 DCache

The L1 DCache is the first level data cache. It manages all load and store requests from all types through its internal data structures. The L1 DCache:

- enables loads and stores to issue speculatively and out of order
- ensures that retired loads and stores have the correct data upon retirement
- ensures that loads and stores follow the memory ordering rules of the IA-32 and Intel 64 instruction set architecture.

Table 2-7. L1 Data Cache Components

Component	Intel microarchitecture code name Sandy Bridge	Intel microarchitecture code name Nehalem
Data Cache Unit (DCU)	32KB, 8 ways	32KB, 8 ways
Load buffers	64 entries	48 entries
Store buffers	36 entries	32 entries
Line fill buffers (LFB)	10 entries	10 entries

The DCU is organized as 32 KBytes, eight-way set associative. Cache line size is 64-bytes arranged in eight banks.

Internally, accesses are up to 16 bytes, with 256-bit Intel AVX instructions utilizing two 16-byte accesses. Two load operations and one store operation can be handled each cycle.

The L1 DCache maintains requests which cannot be serviced immediately to completion. Some reasons for requests that are delayed: cache misses, unaligned access that splits across cache lines, data not ready to be forwarded from a preceding store, loads experiencing bank collisions, and load block due to cache line replacement.

The L1 DCache can maintain up to 64 load micro-ops from allocation until retirement. It can maintain up to 36 store operations from allocation until the store value is committed to the cache, or written to the line fill buffers (LFB) in the case of non-temporal stores.

The L1 DCache can handle multiple outstanding cache misses and continue to service incoming stores and loads. Up to 10 requests of missing cache lines can be managed simultaneously using the LFB.

The L1 DCache is a write-back write-allocate cache. Stores that hit in the DCU do not update the lower levels of the memory hierarchy. Stores that miss the DCU allocate a cache line.

Loads

The L1 DCache architecture can service two loads per cycle, each of which can be up to 16 bytes. Up to 32 loads can be maintained at different stages of progress, from their allocation in the out of order engine until the loaded value is returned to the execution core.

Loads can:

- Read data before preceding stores when the load address and store address ranges are known not to conflict.
- Be carried out speculatively, before preceding branches are resolved.
- Take cache misses out of order and in an overlapped manner.

Loads cannot:

- Speculatively take any sort of fault or trap.
- Speculatively access uncacheable memory.

The common load latency is five cycles. When using a simple addressing mode, base plus offset that is smaller than 2048, the load latency can be four cycles. This technique is especially useful for pointer-chasing code. However, overall latency varies depending on the target register data type due to stack bypass. See Section 2.1.4 for more information.

The following table lists overall load latencies. These latencies assume the common case of flat segment, that is, segment base address is zero. If segment base is not zero, load latency increases

Table 2-8. Effect of Addressing Modes on Load Latency

Data Type/Addressing Mode	Base + Offset > 2048; Base + Index [+ Offset]	Base + Offset < 2048
Integer	5	4
MMX, SSE, 128-bit AVX	6	5
X87	7	6
256-bit AVX	7	7

Stores

Stores to memory are executed in two phases:

- Execution phase. Fills the store buffers with linear and physical address and data. Once store address and data are known, the store data can be forwarded to the following load operations that need it.
- Completion phase. After the store retires, the L1 DCache moves its data from the store buffers to the DCU, up to 16 bytes per cycle.

Address Translation

The DTLB can perform three linear to physical address translations every cycle, two for load addresses and one for a store address. If the address is missing in the DTLB, the processor looks for it in the STLB, which holds data and instruction address translations. The penalty of a DTLB miss that hits the STLB is seven cycles. Large page support include 1G byte pages, in addition to 4K and 2M/4M pages.

The DTLB and STLB are four way set associative. The following table specifies the number of entries in the DTLB and STLB.

Table 2-9. DTLB and STLB Parameters

TLB	Page Size	Entries
DTLB	4KB	64
	2MB/4MB	32
	1GB	4
STLB	4KB	512

Store Forwarding

If a load follows a store and reloads the data that the store writes to memory, the data can forward directly from the store operation to the load. This process, called store to load forwarding, saves cycles by enabling the load to obtain the data directly from the store operation instead of through memory. You can take advantage of store forwarding to quickly move complex structures without losing the ability to forward the subfields. The memory control unit can handle store forwarding situations with less restrictions compared to previous micro-architectures.

The following rules must be met to enable store to load forwarding:

- The store must be the last store to that address, prior to the load.
- The store must contain all data being loaded.
- The load is from a write-back memory type and neither the load nor the store are non-temporal accesses.

Stores cannot forward to loads in the following cases:

- Four byte and eight byte loads that cross eight byte boundary, relative to the preceding 16- or 32-byte store.
- Any load that crosses a 16-byte boundary of a 32-byte store.

Table 2-10 to Table 2-13 detail the store to load forwarding behavior. For a given store size, all the loads that may overlap are shown and specified by 'F'. Forwarding from 32 byte store is similar to forwarding from each of the 16 byte halves of the store. Cases that cannot forward are shown as 'N'.

Table 2-10. Store Forwarding Conditions (1 and 2 byte stores)

Store Size	Load Size	Load Alignment																
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1	1	F																
2	1	F	F															
	2	F	N															

Table 2-11. Store Forwarding Conditions (4-16 byte stores)

		Load Alignment															
Store Size	Load Size	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
4	1	F	F	F	F												
	2	F	F	F	N												
	4	F	N	N	N												
8	1	F	F	F	F	F	F	F	F								
	2	F	F	F	F	F	F	F	N								
	4	F	F	F	F	F	N	N	N								
	8	F	N	N	N	N	N	N	N								
16	1	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F
	2	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	N
	4	F	F	F	F	F	N	N	N	F	F	F	F	F	N	N	N
	8	F	N	N	N	N	N	N	N	F	N	N	N	N	N	N	N
	16	F	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N

Table 2-12. 32-byte Store Forwarding Conditions (0-15 byte alignment)

		Load Alignment															
Store Size	Load Size	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
32	1	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F
	2	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	N
	4	F	F	F	F	F	N	N	N	F	F	F	F	F	N	N	N
	8	F	N	N	N	N	N	N	N	F	N	N	N	N	N	N	N
	16	F	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N
	32	F	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N

Table 2-13. 32-byte Store Forwarding Conditions (16-31 byte alignment)

Store Size	Load Size	Load Alignment															
		16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
32	1	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F
	2	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	N
	4	F	F	F	F	F	N	N	N	F	F	F	F	F	N	N	N
	8	F	N	N	N	N	N	N	N	F	N	N	N	N	N	N	N
	16	F	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N
	32	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N	N

Memory Disambiguation

A load operation may depend on a preceding store. Many microarchitectures block loads until all preceding store addresses are known. The memory disambiguator predicts which loads will not depend on any previous stores. When the disambiguator predicts that a load does not have such a dependency, the load takes its data from the L1 data cache even when the store address is unknown. This hides the load latency. Eventually, the prediction is verified. If an actual conflict is detected, the load and all succeeding instructions are re-executed.

The following loads are not disambiguated. The execution of these loads is stalled until addresses of all previous stores are known.

- Loads that cross the 16-byte boundary
- 32-byte Intel AVX loads that are not 32-byte aligned.

The memory disambiguator always assumes dependency between loads and earlier stores that have the same address bits 0:11.

Bank Conflict

Since 16-byte loads can cover up to three banks, and two loads can happen every cycle, it is possible that six of the eight banks may be accessed per cycle, for loads. A bank conflict happens when two load accesses need the same bank (their address has the same 2-4 bit value) in different sets, at the same time. When a bank conflict occurs, one of the load accesses is recycled internally.

In many cases two loads access exactly the same bank in the same cache line, as may happen when popping operands off the stack, or any sequential accesses. In these cases, conflict does not occur and the loads are serviced simultaneously.

2.1.5.3 Ring Interconnect and Last Level Cache

The system-on-a-chip design provides a high bandwidth bi-directional ring bus to connect between the IA cores and various sub-systems in the uncore. In the second generation Intel Core processor 2xxx series, the uncore subsystem include a system agent, the graphics unit (GT) and the last level cache (LLC).

The LLC consists of multiple cache slices. The number of slices is equal to the number of IA cores. Each slice has logic portion and data array portion. The logic portion handles data coherency, memory ordering, access to the data array portion, LLC misses and writeback to memory, and more. The data array portion stores cache lines. Each slice contains a full cache port that can supply 32 bytes/cycle.

The physical addresses of data kept in the LLC data arrays are distributed among the cache slices by a hash function, such that addresses are uniformly distributed. The data array in a cache block may have 4/8/12/16 ways corresponding to 0.5M/1M/1.5M/2M block size. However, due to the address distribution among the cache blocks from the software point of view, this does not appear as a normal N-way cache.

From the processor cores and the GT view, the LLC act as one shared cache with multiple ports and bandwidth that scales with the number of cores. The LLC hit latency, ranging between 26-31 cycles, depends on the core location relative to the LLC block, and how far the request needs to travel on the ring.

The number of cache-slices increases with the number of cores, therefore the ring and LLC are not likely to be a bandwidth limiter to core operation.

The GT sits on the same ring interconnect, and uses the LLC for its data operations as well. In this respect it is very similar to an IA core. Therefore, high bandwidth graphic applications using cache bandwidth and significant cache footprint, can interfere, to some extent, with core operations.

All the traffic that cannot be satisfied by the LLC, such as LLC misses, dirty line write-back, non-cacheable operations, and MMIO/IO operations, still travels through the cache-slice logic portion and the ring, to the system agent.

2.1.5.4 Data Prefetching

Data can be speculatively loaded to the L1 DCache using software prefetching, hardware prefetching, or any combination of the two.

You can use the four Streaming SIMD Extensions (SSE) prefetch instructions to enable software-controlled prefetching. These instructions are hints to bring a cache line of data into the desired levels of the cache hierarchy. The software-controlled prefetch is intended for prefetching data, but not for prefetching code.

The rest of this section describes the various hardware prefetching mechanisms provided by Intel microarchitecture code name Sandy Bridge and their improvement over previous processors. The goal of the prefetchers is to automatically predict which data the program is about to consume. If this data is not close-by to the execu-

tion core or inner cache, the prefetchers bring it from the next levels of cache hierarchy and memory. Prefetching has the following effects:

- Improves performance if data is arranged sequentially in the order used in the program.
- May cause slight performance degradation due to bandwidth issues, if access patterns are sparse instead of local.
- On rare occasions, if the algorithm's working set is tuned to occupy most of the cache and unneeded prefetches evict lines required by the program, hardware prefetcher may cause severe performance degradation due to cache capacity of L1.

Data Prefetch to L1 Data Cache

Data prefetching is triggered by load operations when the following conditions are met:

- Load is from writeback memory type.
- The prefetched data is within the same 4K byte page as the load instruction that triggered it.
- No fence is in progress in the pipeline.
- Not many other load misses are in progress.
- There is not a continuous stream of stores.

Two hardware prefetchers load data to the L1 DCache:

- **Data cache unit (DCU) prefetcher.** This prefetcher, also known as the streaming prefetcher, is triggered by an ascending access to very recently loaded data. The processor assumes that this access is part of a streaming algorithm and automatically fetches the next line.
- **Instruction pointer (IP)-based stride prefetcher.** This prefetcher keeps track of individual load instructions. If a load instruction is detected to have a regular stride, then a prefetch is sent to the next address which is the sum of the current address and the stride. This prefetcher can prefetch forward or backward and can detect strides of up to 2K bytes.

Data Prefetch to the L2 and Last Level Cache

The following two hardware prefetchers fetched data from memory to the L2 cache and last level cache:

Spatial Prefetcher: This prefetcher strives to complete every cache line fetched to the L2 cache with the pair line that completes it to a 128-byte aligned chunk.

Streamer: This prefetcher monitors read requests from the L1 cache for ascending and descending sequences of addresses. Monitored read requests include L1 DCache requests initiated by load and store operations and by the hardware prefetchers, and L1 ICache requests for code fetch. When a forward or backward stream of requests is detected, the anticipated cache lines are prefetched. Prefetched cache lines must be in the same 4K page.

The streamer and spatial prefetcher prefetch the data to the last level cache. Typically data is brought also to the L2 unless the L2 cache is heavily loaded with missing demand requests.

Enhancement to the streamer includes the following features:

- The streamer may issue two prefetch requests on every L2 lookup. The streamer can run up to 20 lines ahead of the load request.
- Adjusts dynamically to the number of outstanding requests per core. If there are not many outstanding requests, the streamer prefetches further ahead. If there are many outstanding requests it prefetches to the LLC only and less far ahead.
- When cache lines are far ahead, it prefetches to the last level cache only and not to the L2. This method avoids replacement of useful cache lines in the L2 cache.
- Detects and maintains up to 32 streams of data accesses. For each 4K byte page, you can maintain one forward and one backward stream can be maintained.

2.1.6 System Agent

The system agent contains the following components:

- An arbiter that handles all accesses from the ring domain and from I/O (PCIe* and DMI) and routes the accesses to the right place.
- PCIe controllers connect to external PCIe devices. The PCIe controllers have different configuration possibilities that varies with product segment specifics: x16+x4, x8+x8+x4, x8+x4+x4+x4.
- DMI controller connects to the PCH chipset.
- Integrated display engine, Flexible Display Interconnect, and Display Port, for the internal graphic operations.
- Memory controller

All main memory traffic is routed from the arbiter to the memory controller. The memory controller in the second generation Intel Core processor 2xxx series support two channels of DDR, with data rates of 1066MHz, 1333MHz and 1600MHz, and 8 bytes per cycle, depending on the unit type, system configuration and DRAMs. Addresses are distributed between memory channels based on a local hash function that attempts to balance the load between the channels in order to achieve maximum bandwidth and minimum hotspot collisions.

For best performance, populate both channels with equal amounts of memory, preferably the exact same types of DIMMs. In addition, using more ranks for the same amount of memory, results in somewhat better memory bandwidth, since more DRAM pages can be open simultaneously. For best performance, populate the system with the highest supported speed DRAM (1333MHz or 1600MHz data rates, depending on the max supported frequency) with the best DRAM timings.

The two channels have separate resources and handle memory requests independently. The memory controller contains a high-performance out-of-order scheduler

that attempts to maximize memory bandwidth while minimizing latency. Each memory channel contains a 32 cache-line write-data-buffer. Writes to the memory controller are considered completed when they are written to the write-data-buffer. The write-data-buffer is flushed out to main memory at a later time, not impacting write latency.

Partial writes are not handled efficiently on the memory controller and may result in read-modify-write operations on the DDR channel if the partial-writes do not complete a full cache-line in time. Software should avoid creating partial write transactions whenever possible and consider alternative, such as buffering the partial writes into full cache line writes.

The memory controller also supports high-priority isochronous requests (such as USB isochronous, and Display isochronous requests). High bandwidth of memory requests from the integrated display engine takes up some of the memory bandwidth and impacts core access latency to some degree

2.2 INTEL® CORE™ MICROARCHITECTURE AND ENHANCED INTEL CORE MICROARCHITECTURE

Intel Core microarchitecture introduces the following features that enable high performance and power-efficient performance for single-threaded as well as multi-threaded workloads:

- **Intel® Wide Dynamic Execution** enables each processor core to fetch, dispatch, execute with high bandwidths and retire up to four instructions per cycle. Features include:
 - Fourteen-stage efficient pipeline
 - Three arithmetic logical units
 - Four decoders to decode up to five instruction per cycle
 - Macro-fusion and micro-fusion to improve front-end throughput
 - Peak issue rate of dispatching up to six μ ops per cycle
 - Peak retirement bandwidth of up to four μ ops per cycle
 - Advanced branch prediction
 - Stack pointer tracker to improve efficiency of executing function/procedure entries and exits
- **Intel® Advanced Smart Cache** delivers higher bandwidth from the second level cache to the core, optimal performance and flexibility for single-threaded and multi-threaded applications. Features include:
 - Optimized for multicore and single-threaded execution environments
 - 256 bit internal data path to improve bandwidth from L2 to first-level data cache

- Unified, shared second-level cache of 4 Mbyte, 16 way (or 2 MByte, 8 way)
- **Intel® Smart Memory Access** prefetches data from memory in response to data access patterns and reduces cache-miss exposure of out-of-order execution. Features include:
 - Hardware prefetchers to reduce effective latency of second-level cache misses
 - Hardware prefetchers to reduce effective latency of first-level data cache misses
 - Memory disambiguation to improve efficiency of speculative execution engine
- **Intel® Advanced Digital Media Boost** improves most 128-bit SIMD instructions with single-cycle throughput and floating-point operations. Features include:
 - Single-cycle throughput of most 128-bit SIMD instructions (except 128-bit shuffle, pack, unpack operations)
 - Up to eight floating-point operations per cycle
 - Three issue ports available to dispatching SIMD instructions for execution.

The Enhanced Intel Core microarchitecture supports all of the features of Intel Core microarchitecture and provides a comprehensive set of enhancements.

- **Intel® Wide Dynamic Execution** includes several enhancements:
 - A radix-16 divider replacing previous radix-4 based divider to speedup long-latency operations such as divisions and square roots.
 - Improved system primitives to speedup long-latency operations such as RDTSC, STI, CLI, and VM exit transitions.
- **Intel® Advanced Smart Cache** provides up to 6 MBytes of second-level cache shared between two processor cores (quad-core processors have up to 12 MBytes of L2); up to 24 way/set associativity.
- **Intel® Smart Memory Access** supports high-speed system bus up 1600 MHz and provides more efficient handling of memory operations such as split cache line load and store-to-load forwarding situations.
- **Intel® Advanced Digital Media Boost** provides 128-bit shuffler unit to speedup shuffle, pack, unpack operations; adds support for 47 SSE4.1 instructions.

In the sub-sections of 2.1.x, most of the descriptions on Intel Core microarchitecture also applies to Enhanced Intel Core microarchitecture. Differences between them are note explicitly.

2.2.1 Intel® Core™ Microarchitecture Pipeline Overview

The pipeline of the Intel Core microarchitecture contains:

- An in-order issue front end that fetches instruction streams from memory, with four instruction decoders to supply decoded instruction (μ ops) to the out-of-order execution core.
- An out-of-order superscalar execution core that can issue up to six μ ops per cycle (see Table 2-2) and reorder μ ops to execute as soon as sources are ready and execution resources are available.
- An in-order retirement unit that ensures the results of execution of μ ops are processed and architectural states are updated according to the original program order.

Intel Core 2 Extreme processor X6800, Intel Core 2 Duo processors and Intel Xeon processor 3000, 5100 series implement two processor cores based on the Intel Core microarchitecture. Intel Core 2 Extreme quad-core processor, Intel Core 2 Quad processors and Intel Xeon processor 3200 series, 5300 series implement four processor cores. Each physical package of these quad-core processors contains two processor dies, each die containing two processor cores. The functionality of the subsystems in each core are depicted in Figure 2-1.

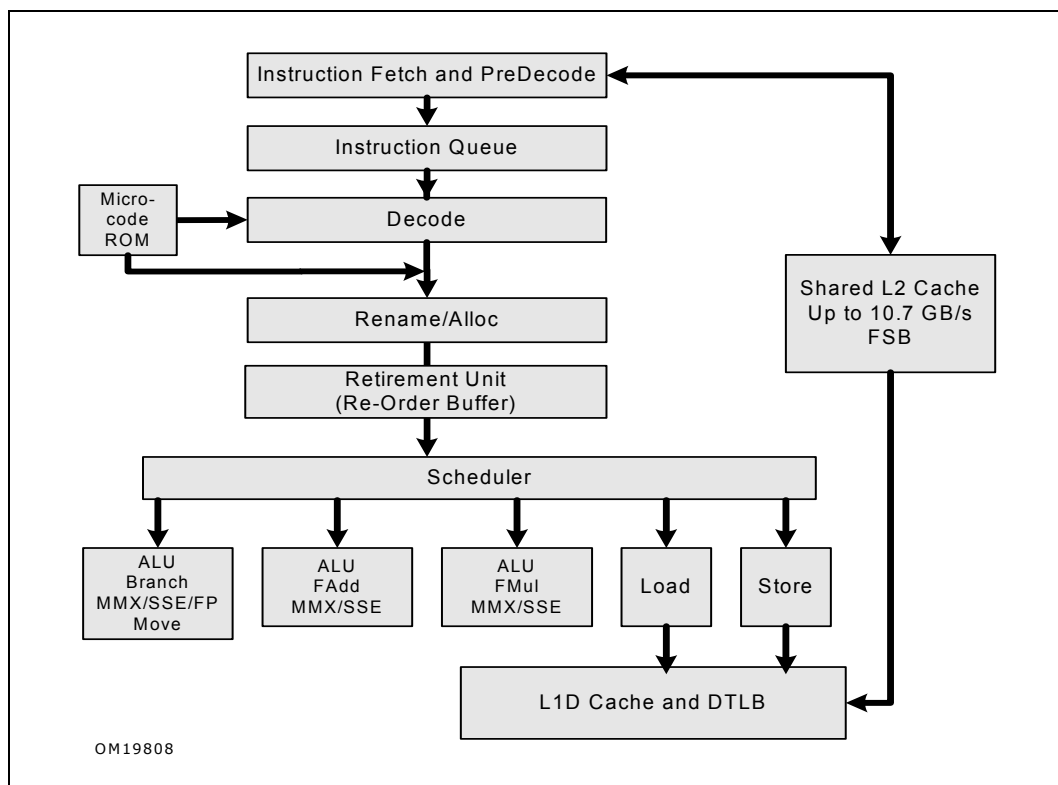


Figure 2-2. Intel Core Microarchitecture Pipeline Functionality

2.2.2 Front End

The front ends needs to supply decoded instructions (μ ops) and sustain the stream to a six-issue wide out-of-order engine. The components of the front end, their functions, and the performance challenges to microarchitectural design are described in Table 2-1.

Table 2-14. Components of the Front End

Component	Functions	Performance Challenges
Branch Prediction Unit (BPU)	<ul style="list-style-type: none"> Helps the instruction fetch unit fetch the most likely instruction to be executed by predicting the various branch types: conditional, indirect, direct, call, and return. Uses dedicated hardware for each type. 	<ul style="list-style-type: none"> Enables speculative execution. Improves speculative execution efficiency by reducing the amount of code in the “non-architected path”¹ to be fetched into the pipeline.
Instruction Fetch Unit	<ul style="list-style-type: none"> Prefetches instructions that are likely to be executed Caches frequently-used instructions Predecodes and buffers instructions, maintaining a constant bandwidth despite irregularities in the instruction stream 	<ul style="list-style-type: none"> Variable length instruction format causes unevenness (bubbles) in decode bandwidth. Taken branches and misaligned targets causes disruptions in the overall bandwidth delivered by the fetch unit.
Instruction Queue and Decode Unit	<ul style="list-style-type: none"> Decodes up to four instructions, or up to five with macro-fusion Stack pointer tracker algorithm for efficient procedure entry and exit Implements the Macro-Fusion feature, providing higher performance and efficiency The Instruction Queue is also used as a loop cache, enabling some loops to be executed with both higher bandwidth and lower power 	<ul style="list-style-type: none"> Varying amounts of work per instruction requires expansion into variable numbers of μops. Prefix adds a dimension of decoding complexity. Length Changing Prefix (LCP) can cause front end bubbles.

NOTES:

- Code paths that the processor thought it should execute but then found out it should go in another path and therefore reverted from its initial intention.

2.2.2.1 Branch Prediction Unit

Branch prediction enables the processor to begin executing instructions long before the branch outcome is decided. All branches utilize the BPU for prediction. The BPU contains the following features:

- 16-entry Return Stack Buffer (RSB). It enables the BPU to accurately predict RET instructions.
- Front end queuing of BPU lookups. The BPU makes branch predictions for 32 bytes at a time, twice the width of the fetch engine. This enables taken branches to be predicted with no penalty.

Even though this BPU mechanism generally eliminates the penalty for taken branches, software should still regard taken branches as consuming more resources than do not-taken branches.

The BPU makes the following types of predictions:

- Direct Calls and Jumps. Targets are read as a target array, without regarding the taken or not-taken prediction.
- Indirect Calls and Jumps. These may either be predicted as having a monotonic target or as having targets that vary in accordance with recent program behavior.
- Conditional branches. Predicts the branch target and whether or not the branch will be taken.

For information about optimizing software for the BPU, see Section 3.4, “Optimizing the Front End”.

2.2.2.2 Instruction Fetch Unit

The instruction fetch unit comprises the instruction translation lookaside buffer (ITLB), an instruction prefetcher, the instruction cache and the predecode logic of the instruction queue (IQ).

Instruction Cache and ITLB

An instruction fetch is a 16-byte aligned lookup through the ITLB into the instruction cache and instruction prefetch buffers. A hit in the instruction cache causes 16 bytes to be delivered to the instruction predecoder. Typical programs average slightly less than 4 bytes per instruction, depending on the code being executed. Since most instructions can be decoded by all decoders, an entire fetch can often be consumed by the decoders in one cycle.

A misaligned target reduces the number of instruction bytes by the amount of offset into the 16 byte fetch quantity. A taken branch reduces the number of instruction bytes delivered to the decoders since the bytes after the taken branch are not decoded. Branches are taken approximately every 10 instructions in typical integer code, which translates into a “partial” instruction fetch every 3 or 4 cycles.

Due to stalls in the rest of the machine, front end starvation does not usually cause performance degradation. For extremely fast code with larger instructions (such as SSE2 integer media kernels), it may be beneficial to use targeted alignment to prevent instruction starvation.

Instruction PreDecode

The predecode unit accepts the sixteen bytes from the instruction cache or prefetch buffers and carries out the following tasks:

- Determine the length of the instructions.
- Decode all prefixes associated with instructions.
- Mark various properties of instructions for the decoders (for example, “is branch.”).

The predecode unit can write up to six instructions per cycle into the instruction queue. If a fetch contains more than six instructions, the predecoder continues to decode up to six instructions per cycle until all instructions in the fetch are written to the instruction queue. Subsequent fetches can only enter predecoding after the current fetch completes.

For a fetch of seven instructions, the predecoder decodes the first six in one cycle, and then only one in the next cycle. This process would support decoding 3.5 instructions per cycle. Even if the instruction per cycle (IPC) rate is not fully optimized, it is higher than the performance seen in most applications. In general, software usually does not have to take any extra measures to prevent instruction starvation.

The following instruction prefixes cause problems during length decoding. These prefixes can dynamically change the length of instructions and are known as length changing prefixes (LCPs):

- Operand Size Override (66H) preceding an instruction with a word immediate data
- Address Size Override (67H) preceding an instruction with a mod R/M in real, 16-bit protected or 32-bit protected modes

When the predecoder encounters an LCP in the fetch line, it must use a slower length decoding algorithm. With the slower length decoding algorithm, the predecoder decodes the fetch in 6 cycles, instead of the usual 1 cycle.

Normal queuing within the processor pipeline usually cannot hide LCP penalties.

The REX prefix (4xh) in the Intel 64 architecture instruction set can change the size of two classes of instruction: MOV offset and MOV immediate. Nevertheless, it does not cause an LCP penalty and hence is not considered an LCP.

2.2.2.3 Instruction Queue (IQ)

The instruction queue is 18 instructions deep. It sits between the instruction predecode unit and the instruction decoders. It sends up to five instructions per cycle, and

supports one macro-fusion per cycle. It also serves as a loop cache for loops smaller than 18 instructions. The loop cache operates as described below.

A Loop Stream Detector (LSD) resides in the BPU. The LSD attempts to detect loops which are candidates for streaming from the instruction queue (IQ). When such a loop is detected, the instruction bytes are locked down and the loop is allowed to stream from the IQ until a misprediction ends it. When the loop plays back from the IQ, it provides higher bandwidth at reduced power (since much of the rest of the front end pipeline is shut off).

The LSD provides the following benefits:

- No loss of bandwidth due to taken branches
- No loss of bandwidth due to misaligned instructions
- No LCP penalties, as the pre-decode stage has already been passed
- Reduced front end power consumption, because the instruction cache, BPU and predecode unit can be idle

Software should use the loop cache functionality opportunistically. Loop unrolling and other code optimizations may make the loop too big to fit into the LSD. For high performance code, loop unrolling is generally preferable for performance even when it overflows the loop cache capability.

2.2.2.4 Instruction Decode

The Intel Core microarchitecture contains four instruction decoders. The first, Decoder 0, can decode Intel 64 and IA-32 instructions up to 4 μ ops in size. Three other decoders handles single- μ op instructions. The microsequencer can provide up to 3 μ ops per cycle, and helps decode instructions larger than 4 μ ops.

All decoders support the common cases of single μ op flows, including: micro-fusion, stack pointer tracking and macro-fusion. Thus, the three simple decoders are not limited to decoding single- μ op instructions. Packing instructions into a 4-1-1-1 template is not necessary and not recommended.

Macro-fusion merges two instructions into a single μ op. Intel Core microarchitecture is capable of one macro-fusion per cycle in 32-bit operation (including compatibility sub-mode of the Intel 64 architecture), but not in 64-bit mode because code that uses longer instructions (length in bytes) more often is less likely to take advantage of hardware support for macro-fusion.

2.2.2.5 Stack Pointer Tracker

The Intel 64 and IA-32 architectures have several commonly used instructions for parameter passing and procedure entry and exit: PUSH, POP, CALL, LEAVE and RET. These instructions implicitly update the stack pointer register (RSP), maintaining a combined control and parameter stack without software intervention. These instructions are typically implemented by several μ ops in previous microarchitectures.

The Stack Pointer Tracker moves all these implicit RSP updates to logic contained in the decoders themselves. The feature provides the following benefits:

- Improves decode bandwidth, as PUSH, POP and RET are single μop instructions in Intel Core microarchitecture.
- Conserves execution bandwidth as the RSP updates do not compete for execution resources.
- Improves parallelism in the out of order execution engine as the implicit serial dependencies between μops are removed.
- Improves power efficiency as the RSP updates are carried out on small, dedicated hardware.

2.2.2.6 Micro-fusion

Micro-fusion fuses multiple μops from the same instruction into a single complex μop . The complex μop is dispatched in the out-of-order execution core. Micro-fusion provides the following performance advantages:

- Improves instruction bandwidth delivered from decode to retirement.
- Reduces power consumption as the complex μop represents more work in a smaller format (in terms of bit density), reducing overall “bit-toggling” in the machine for a given amount of work and virtually increasing the amount of storage in the out-of-order execution engine.

Many instructions provide register flavors and memory flavors. The flavor involving a memory operand will decode into a longer flow of μops than the register version. Micro-fusion enables software to use memory to register operations to express the actual program behavior without worrying about a loss of decode bandwidth.

2.2.3 Execution Core

The execution core of the Intel Core microarchitecture is superscalar and can process instructions out of order. When a dependency chain causes the machine to wait for a resource (such as a second-level data cache line), the execution core executes other instructions. This increases the overall rate of instructions executed per cycle (IPC).

The execution core contains the following three major components:

- **Renamer** — Moves μops from the front end to the execution core. Architectural registers are renamed to a larger set of microarchitectural registers. Renaming eliminates false dependencies known as read-after-read and write-after-read hazards.
- **Reorder buffer (ROB)** — Holds μops in various stages of completion, buffers completed μops , updates the architectural state in order, and manages ordering of exceptions. The ROB has 96 entries to handle instructions in flight.

- **Reservation station (RS)** — Queues μ ops until all source operands are ready, schedules and dispatches ready μ ops to the available execution units. The RS has 32 entries.

The initial stages of the out of order core move the μ ops from the front end to the ROB and RS. In this process, the out of order core carries out the following steps:

- Allocates resources to μ ops (for example: these resources could be load or store buffers).
- Binds the μ op to an appropriate issue port.
- Renames sources and destinations of μ ops, enabling out of order execution.
- Provides data to the μ op when the data is either an immediate value or a register value that has already been calculated.

The following list describes various types of common operations and how the core executes them efficiently:

- **Micro-ops with single-cycle latency** — Most μ ops with single-cycle latency can be executed by multiple execution units, enabling multiple streams of dependent operations to be executed quickly.
- **Frequently-used μ ops with longer latency** — These μ ops have pipelined execution units so that multiple μ ops of these types may be executing in different parts of the pipeline simultaneously.
- **Operations with data-dependent latencies** — Some operations, such as division, have data dependent latencies. Integer division parses the operands to perform the calculation only on significant portions of the operands, thereby speeding up common cases of dividing by small numbers.
- **Floating point operations with fixed latency for operands that meet certain restrictions** — Operands that do not fit these restrictions are considered exceptional cases and are executed with higher latency and reduced throughput. The lower-throughput cases do not affect latency and throughput for more common cases.
- **Memory operands with variable latency, even in the case of an L1 cache hit** — Loads that are not known to be safe from forwarding may wait until a store-address is resolved before executing. The memory order buffer (MOB) accepts and processes all memory operations. See Section 2.1.5 for more information about the MOB.

2.2.3.1 Issue Ports and Execution Units

The scheduler can dispatch up to six μ ops per cycle through the issue ports. The issue ports of Intel Core microarchitecture and Enhanced Intel Core microarchitecture are depicted in Table 2-2, the former is denoted by its CPUID signature of DisplayFamily_DisplayModel value of 06_0FH, the latter denoted by the corresponding signature value of 06_17H. The table provides latency and throughput data of common integer and floating-point (FP) operations for each issue port in cycles.

Table 2-15. Issue Ports of Intel Core Microarchitecture and Enhanced Intel Core Microarchitecture

Executable operations	Latency, Throughput		Comment ¹
	Signature = 06_0FH	Signature = 06_17H	
Integer ALU Integer SIMD ALU FP/SIMD/SSE2 Move and Logic	1, 1 1, 1 1, 1	1, 1 1, 1 1, 1	Includes 64-bit mode integer MUL; Issue port 0; Writeback port 0;
Single-precision (SP) FP MUL Double-precision FP MUL	4, 1 5, 1	4, 1 5, 1	Issue port 0; Writeback port 0
FP MUL (X87) FP Shuffle DIV/SQRT	5, 2 1, 1	5, 2 1, 1	Issue port 0; Writeback port 0 FP shuffle does not handle QW shuffle.
Integer ALU Integer SIMD ALU FP/SIMD/SSE2 Move and Logic	1, 1 1, 1 1, 1	1, 1 1, 1 1, 1	Excludes 64-bit mode integer MUL; Issue port 1; Writeback port 1;
FP ADD QW Shuffle	3, 1 1, 1 ²	3, 1 1, 1 ³	Issue port 1; Writeback port 1;
Integer loads FP loads	3, 1 4, 1	3, 1 4, 1	Issue port 2; Writeback port 2;
Store address ⁴	3, 1	3, 1	Issue port 3;
Store data ⁵ .			Issue Port 4;
Integer ALU Integer SIMD ALU FP/SIMD/SSE2 Move and Logic	1, 1 1, 1 1, 1	1, 1 1, 1 1, 1	Issue port 5; Writeback port 5;
QW shuffles 128-bit Shuffle/Pack/Unpack	1, 1 ² 2-4, 2-4 ⁶	1, 1 ³ 1-3, 1 ⁷	Issue port 5; Writeback port 5;

NOTES:

- Mixing operations of different latencies that use the same port can result in writeback bus conflicts; this can reduce overall throughput
- 128-bit instructions executes with longer latency and reduced throughput
- Uses 128-bit shuffle unit in port 5.
- Prepares the store forwarding and store retirement logic with the address of the data being stored.
- Prepares the store forwarding and store retirement logic with the data being stored

6. Varies with instructions; 128-bit instructions are executed using QW shuffle units
7. Varies with instructions, 128-bit shuffle unit replaces QW shuffle units in Intel Core microarchitecture.

In each cycle, the RS can dispatch up to six μ ops. Each cycle, up to 4 results may be written back to the RS and ROB, to be used as early as the next cycle by the RS. This high execution bandwidth enables execution bursts to keep up with the functional expansion of the micro-fused μ ops that are decoded and retired.

The execution core contains the following three execution stacks:

- SIMD integer
- regular integer
- x87/SIMD floating point

The execution core also contains connections to and from the memory cluster. See Figure 2-2.

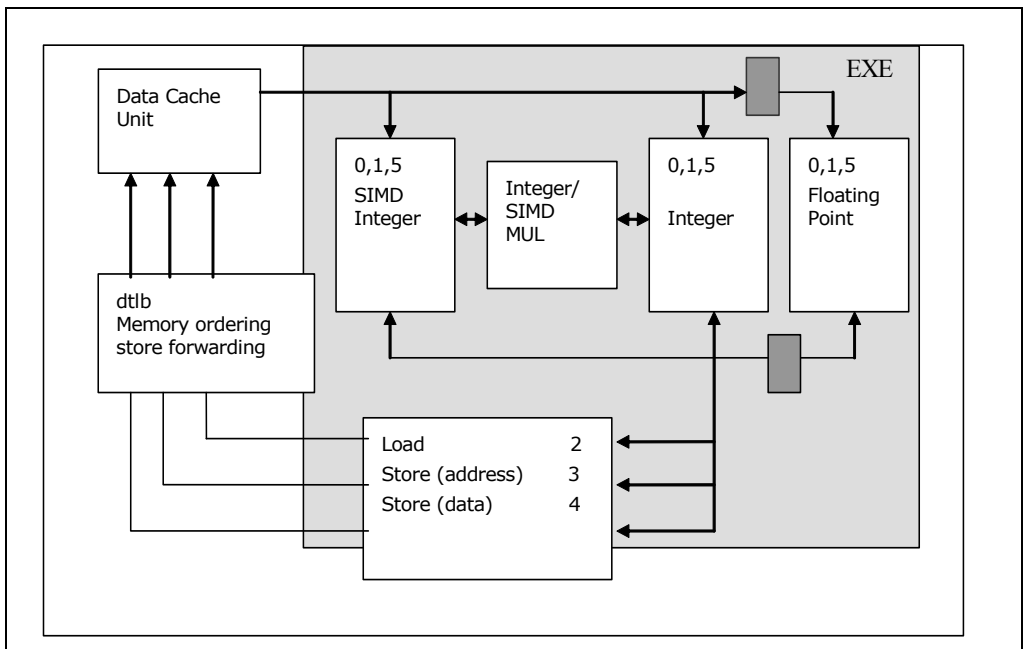


Figure 2-3. Execution Core of Intel Core Microarchitecture

Notice that the two dark squares inside the execution block (in grey color) and appear in the path connecting the integer and SIMD integer stacks to the floating point stack. This delay shows up as an extra cycle called a bypass delay. Data from

the L1 cache has one extra cycle of latency to the floating point unit. The dark-colored squares in Figure 2-2 represent the extra cycle of latency.

2.2.4 Intel® Advanced Memory Access

The Intel Core microarchitecture contains an instruction cache and a first-level data cache in each core. The two cores share a 2 or 4-MByte L2 cache. All caches are writeback and non-inclusive. Each core contains:

- **L1 data cache, known as the data cache unit (DCU)** — The DCU can handle multiple outstanding cache misses and continue to service incoming stores and loads. It supports maintaining cache coherency. The DCU has the following specifications:
 - 32-KBytes size
 - 8-way set associative
 - 64-bytes line size
- **Data translation lookaside buffer (DTLB)** — The DTLB in Intel Core microarchitecture implements two levels of hierarchy. Each level of the DTLB have multiple entries and can support either 4-KByte pages or large pages. The entries of the inner level (DTLB0) is used for loads. The entries in the outer level (DTLB1) support store operations and loads that missed DTLB0. All entries are 4-way associative. Here is a list of entries in each DTLB:
 - DTLB1 for large pages: 32 entries
 - DTLB1 for 4-KByte pages: 256 entries
 - DTLB0 for large pages: 16 entries
 - DTLB0 for 4-KByte pages: 16 entries

An DTLB0 miss and DTLB1 hit causes a penalty of 2 cycles. Software only pays this penalty if the DTLB0 is used in some dispatch cases. The delays associated with a miss to the DTLB1 and PMH are largely non-blocking due to the design of Intel Smart Memory Access.
- **Page miss handler (PMH)**
- **A memory ordering buffer (MOB)** — Which:
 - enables loads and stores to issue speculatively and out of order
 - ensures retired loads and stores have the correct data upon retirement
 - ensures loads and stores follow memory ordering rules of the Intel 64 and IA-32 architectures.

The memory cluster of the Intel Core microarchitecture uses the following to speed up memory operations:

- 128-bit load and store operations
- data prefetching to L1 caches

- data prefetch logic for prefetching to the L2 cache
- store forwarding
- memory disambiguation
- 8 fill buffer entries
- 20 store buffer entries
- out of order execution of memory operations
- pipelined read-for-ownership operation (RFO)

For information on optimizing software for the memory cluster, see Section 3.6, “Optimizing Memory Accesses.”

2.2.4.1 Loads and Stores

The Intel Core microarchitecture can execute up to one 128-bit load and up to one 128-bit store per cycle, each to different memory locations. The microarchitecture enables execution of memory operations out of order with respect to other instructions and with respect to other memory operations.

Loads can:

- issue before preceding stores when the load address and store address are known not to conflict
- be carried out speculatively, before preceding branches are resolved
- take cache misses out of order and in an overlapped manner
- issue before preceding stores, speculating that the store is not going to be to a conflicting address

Loads cannot:

- speculatively take any sort of fault or trap
- speculatively access the uncacheable memory type

Faulting or uncacheable loads are detected and wait until retirement, when they update the programmer visible state. x87 and floating point SIMD loads add 1 additional clock latency.

Stores to memory are executed in two phases:

- **Execution phase** — Prepares the store buffers with address and data for store forwarding. Consumes dispatch ports, which are ports 3 and 4.
- **Completion phase** — The store is retired to programmer-visible memory. It may compete for cache banks with executing loads. Store retirement is maintained as a background task by the memory order buffer, moving the data from the store buffers to the L1 cache.

2.2.4.2 Data Prefetch to L1 caches

Intel Core microarchitecture provides two hardware prefetchers to speed up data accessed by a program by prefetching to the L1 data cache:

- **Data cache unit (DCU) prefetcher** — This prefetcher, also known as the streaming prefetcher, is triggered by an ascending access to very recently loaded data. The processor assumes that this access is part of a streaming algorithm and automatically fetches the next line.
- **Instruction pointer (IP)- based strided prefetcher** — This prefetcher keeps track of individual load instructions. If a load instruction is detected to have a regular stride, then a prefetch is sent to the next address which is the sum of the current address and the stride. This prefetcher can prefetch forward or backward and can detect strides of up to half of a 4KB-page, or 2 KBytes.

Data prefetching works on loads only when the following conditions are met:

- Load is from writeback memory type.
- Prefetch request is within the page boundary of 4 Kbytes.
- No fence or lock is in progress in the pipeline.
- Not many other load misses are in progress.
- The bus is not very busy.
- There is not a continuous stream of stores.

DCU Prefetching has the following effects:

- Improves performance if data in large structures is arranged sequentially in the order used in the program.
- May cause slight performance degradation due to bandwidth issues if access patterns are sparse instead of local.
- On rare occasions, if the algorithm's working set is tuned to occupy most of the cache and unneeded prefetches evict lines required by the program, hardware prefetcher may cause severe performance degradation due to cache capacity of L1.

In contrast to hardware prefetchers relying on hardware to anticipate data traffic, software prefetch instructions relies on the programmer to anticipate cache miss traffic, software prefetch act as hints to bring a cache line of data into the desired levels of the cache hierarchy. The software-controlled prefetch is intended for prefetching data, but not for prefetching code.

2.2.4.3 Data Prefetch Logic

Data prefetch logic (DPL) prefetches data to the second-level (L2) cache based on past request patterns of the DCU from the L2. The DPL maintains two independent arrays to store addresses from the DCU: one for upstreams (12 entries) and one for down streams (4 entries). The DPL tracks accesses to one 4K byte page in each

entry. If an accessed page is not in any of these arrays, then an array entry is allocated.

The DPL monitors DCU reads for incremental sequences of requests, known as streams. Once the DPL detects the second access of a stream, it prefetches the next cache line. For example, when the DCU requests the cache lines A and A+1, the DPL assumes the DCU will need cache line A+2 in the near future. If the DCU then reads A+2, the DPL prefetches cache line A+3. The DPL works similarly for “downward” loops.

The Intel Pentium M processor introduced DPL. The Intel Core microarchitecture added the following features to DPL:

- The DPL can detect more complicated streams, such as when the stream skips cache lines. DPL may issue 2 prefetch requests on every L2 lookup. The DPL in the Intel Core microarchitecture can run up to 8 lines ahead from the load request.
- DPL in the Intel Core microarchitecture adjusts dynamically to bus bandwidth and the number of requests. DPL prefetches far ahead if the bus is not busy, and less far ahead if the bus is busy.
- DPL adjusts to various applications and system configurations.

Entries for the two cores are handled separately.

2.2.4.4 Store Forwarding

If a load follows a store and reloads the data that the store writes to memory, the Intel Core microarchitecture can forward the data directly from the store to the load. This process, called store to load forwarding, saves cycles by enabling the load to obtain the data directly from the store operation instead of through memory.

The following rules must be met for store to load forwarding to occur:

- The store must be the last store to that address prior to the load.
- The store must be equal or greater in size than the size of data being loaded.
- The load cannot cross a cache line boundary.
- The load cannot cross an 8-Byte boundary. 16-Byte loads are an exception to this rule.
- The load must be aligned to the start of the store address, except for the following exceptions:
 - An aligned 64-bit store may forward either of its 32-bit halves
 - An aligned 128-bit store may forward any of its 32-bit quarters
 - An aligned 128-bit store may forward either of its 64-bit halves

Software can use the exceptions to the last rule to move complex structures without losing the ability to forward the subfields.

In Enhanced Intel Core microarchitecture, the alignment restrictions to permit store forwarding to proceed have been relaxed. Enhanced Intel Core microarchitecture permits store-forwarding to proceed in several situations that the succeeding load is not aligned to the preceding store. Table 2-3 shows six situations (in gradient-filled background) of store-forwarding that are permitted in Enhanced Intel Core microarchitecture but not in Intel Core microarchitecture. The cases with backward slash background depicts store-forwarding that can proceed in both Intel Core microarchitecture and Enhanced Intel Core microarchitecture.

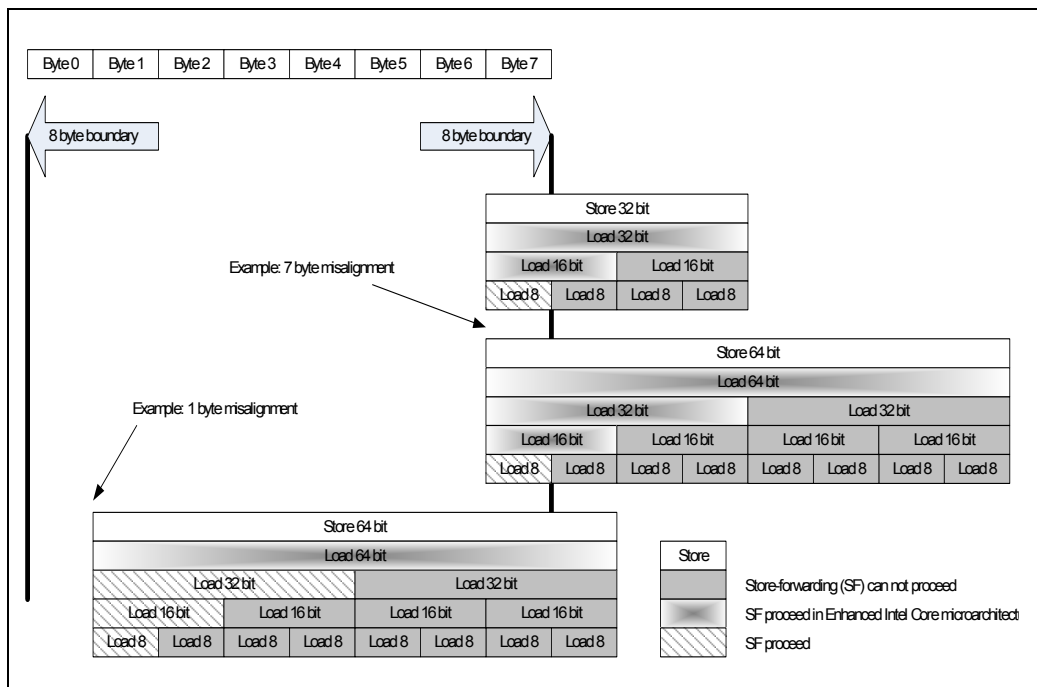


Figure 2-4. Store-Forwarding Enhancements in Enhanced Intel Core Microarchitecture

2.2.4.5 Memory Disambiguation

A load instruction μop may depend on a preceding store. Many microarchitectures block loads until all preceding store address are known.

The memory disambiguator predicts which loads will not depend on any previous stores. When the disambiguator predicts that a load does not have such a dependency, the load takes its data from the L1 data cache.

Eventually, the prediction is verified. If an actual conflict is detected, the load and all succeeding instructions are re-executed.

2.2.5 Intel® Advanced Smart Cache

The Intel Core microarchitecture optimized a number of features for two processor cores on a single die. The two cores share a second-level cache and a bus interface unit, collectively known as Intel Advanced Smart Cache. This section describes the components of Intel Advanced Smart Cache. Figure 2-4 illustrates the architecture of the Intel Advanced Smart Cache.

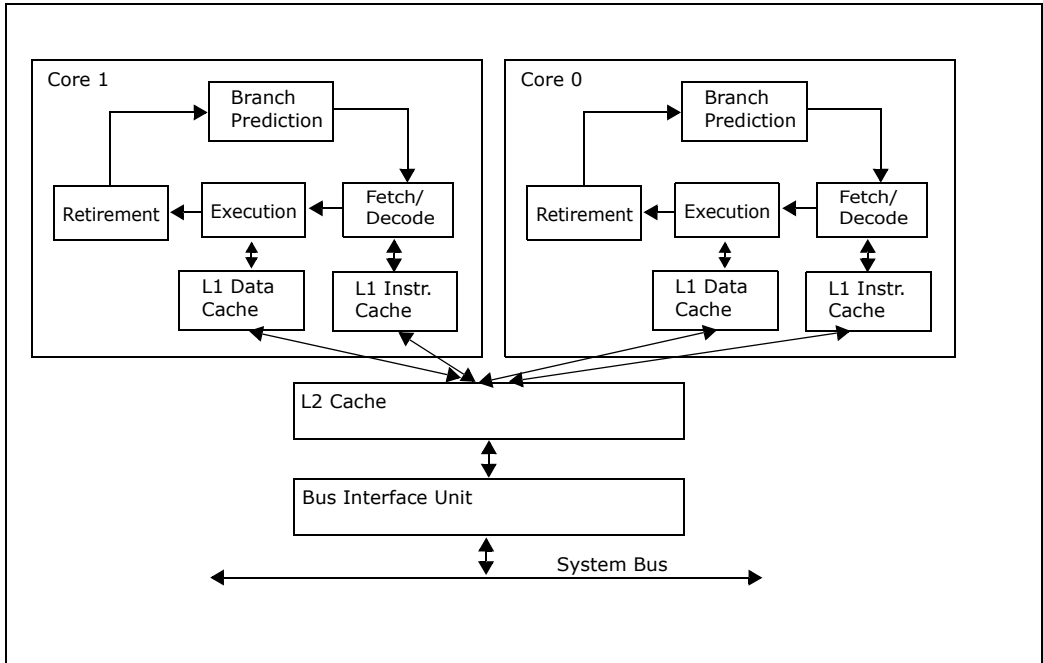


Figure 2-5. Intel Advanced Smart Cache Architecture

Table 2-3 details the parameters of caches in the Intel Core microarchitecture. For information on enumerating the cache hierarchy identification using the deterministic cache parameter leaf of CPUID instruction, see the *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 2A*.

Table 2-16. Cache Parameters of Processors based on Intel Core Microarchitecture

Level	Capacity	Associativity (ways)	Line Size (bytes)	Access Latency (clocks)	Access Throughput (clocks)	Write Update Policy
First Level	32 KB	8	64	3	1	Writeback
Instruction	32 KB	8	N/A	N/A	N/A	N/A
Second Level (Shared L2) ¹	2, 4 MB	8 or 16	64	14 ²	2	Writeback
Second Level (Shared L2) ³	3, 6MB	12 or 24	64	15 ²	2	Writeback
Third Level ⁴	8, 12, 16 MB	16	64	~110	12	Writeback

NOTES:

1. Intel Core microarchitecture (CUID signature DisplayFamily = 06H, DisplayModel = 0FH).
2. Software-visible latency will vary depending on access patterns and other factors.
3. Enhanced Intel Core microarchitecture (CUID signature DisplayFamily = 06H, DisplayModel = 17H or 1DH).
4. Enhanced Intel Core microarchitecture (CUID signature DisplayFamily = 06H, DisplayModel = 1DH).

2.2.5.1 Loads

When an instruction reads data from a memory location that has write-back (WB) type, the processor looks for the cache line that contains this data in the caches and memory in the following order:

1. DCU of the initiating core
2. DCU of the other core and second-level cache
3. System memory

The cache line is taken from the DCU of the other core only if it is modified, ignoring the cache line availability or state in the L2 cache.

Table 2-4 shows the characteristics of fetching the first four bytes of different localities from the memory cluster. The latency column provides an estimate of access latency. However, the actual latency can vary depending on the load of cache, memory components, and their parameters.

Table 2-17. Characteristics of Load and Store Operations in Intel Core Microarchitecture

Data Locality	Load		Store	
	Latency	Throughput	Latency	Throughput
DCU	3	1	2	1
DCU of the other core in modified state	14 + 5.5 bus cycles	14 + 5.5 bus cycles	14 + 5.5 bus cycles	
2nd-level cache	14	3	14	3
Memory	14 + 5.5 bus cycles + memory	Depends on bus read protocol	14 + 5.5 bus cycles + memory	Depends on bus write protocol

Sometimes a modified cache line has to be evicted to make space for a new cache line. The modified cache line is evicted in parallel to bringing the new data and does not require additional latency. However, when data is written back to memory, the eviction uses cache bandwidth and possibly bus bandwidth as well. Therefore, when multiple cache misses require the eviction of modified lines within a short time, there is an overall degradation in cache response time.

2.2.5.2 Stores

When an instruction writes data to a memory location that has WB memory type, the processor first ensures that the line is in Exclusive or Modified state in its own DCU. The processor looks for the cache line in the following locations, in the specified order:

1. DCU of initiating core
2. DCU of the other core and L2 cache
3. System memory

The cache line is taken from the DCU of the other core only if it is modified, ignoring the cache line availability or state in the L2 cache. After reading for ownership is completed, the data is written to the first-level data cache and the line is marked as modified.

Reading for ownership and storing the data happens after instruction retirement and follows the order of retirement. Therefore, the store latency does not effect the store instruction itself. However, several sequential stores may have cumulative latency that can affect performance. Table 2-4 presents store latencies depending on the initial cache line location.

2.3 INTEL® MICROARCHITECTURE CODE NAME NEHALEM

Intel microarchitecture code name Nehalem provides the foundation for many innovative features of Intel Core i7 processors and Intel Xeon processor 3400, 5500, and 7500 series. It builds on the success of 45nm enhanced Intel Core microarchitecture and provides the following feature enhancements:

- **Enhanced processor core**
 - Improved branch prediction and recovery from misprediction.
 - Enhanced loop streaming to improve front end performance and reduce power consumption.
 - Deeper buffering in out-of-order engine to extract parallelism.
 - Enhanced execution units to provide acceleration in CRC, string/text processing and data shuffling.
- **Hyper-Threading Technology**
 - Provides two hardware threads (logical processors) per core.
 - Takes advantage of 4-wide execution engine, large L3, and massive memory bandwidth.
- **Smart Memory Access**
 - Integrated memory controller provides low-latency access to system memory and scalable memory bandwidth
 - New cache hierarchy organization with shared, inclusive L3 to reduce snoop traffic
 - Two level TLBs and increased TLB size.
 - Fast unaligned memory access.
- **Dedicated Power management Innovations**
 - Integrated microcontroller with optimized embedded firmware to manage power consumption.
 - Embedded real-time sensors for temperature, current, and power.
 - Integrated power gate to turn off/on per-core power consumption
 - Versatility to reduce power consumption of memory, link subsystems.

Intel microarchitecture code name Westmere is a 32nm version of Intel microarchitecture code name Nehalem. All of the features of latter also apply to the former.

2.3.1 Microarchitecture Pipeline

Intel microarchitecture code name Nehalem continues the four-wide microarchitecture pipeline pioneered by the 65nm Intel Core Microarchitecture. Figure 2-6 illustrates the basic components of the pipeline of Intel microarchitecture code name

Nehalem as implemented in Intel Core i7 processor, only two of the four cores are sketched in the Figure 2-6 pipeline diagram.

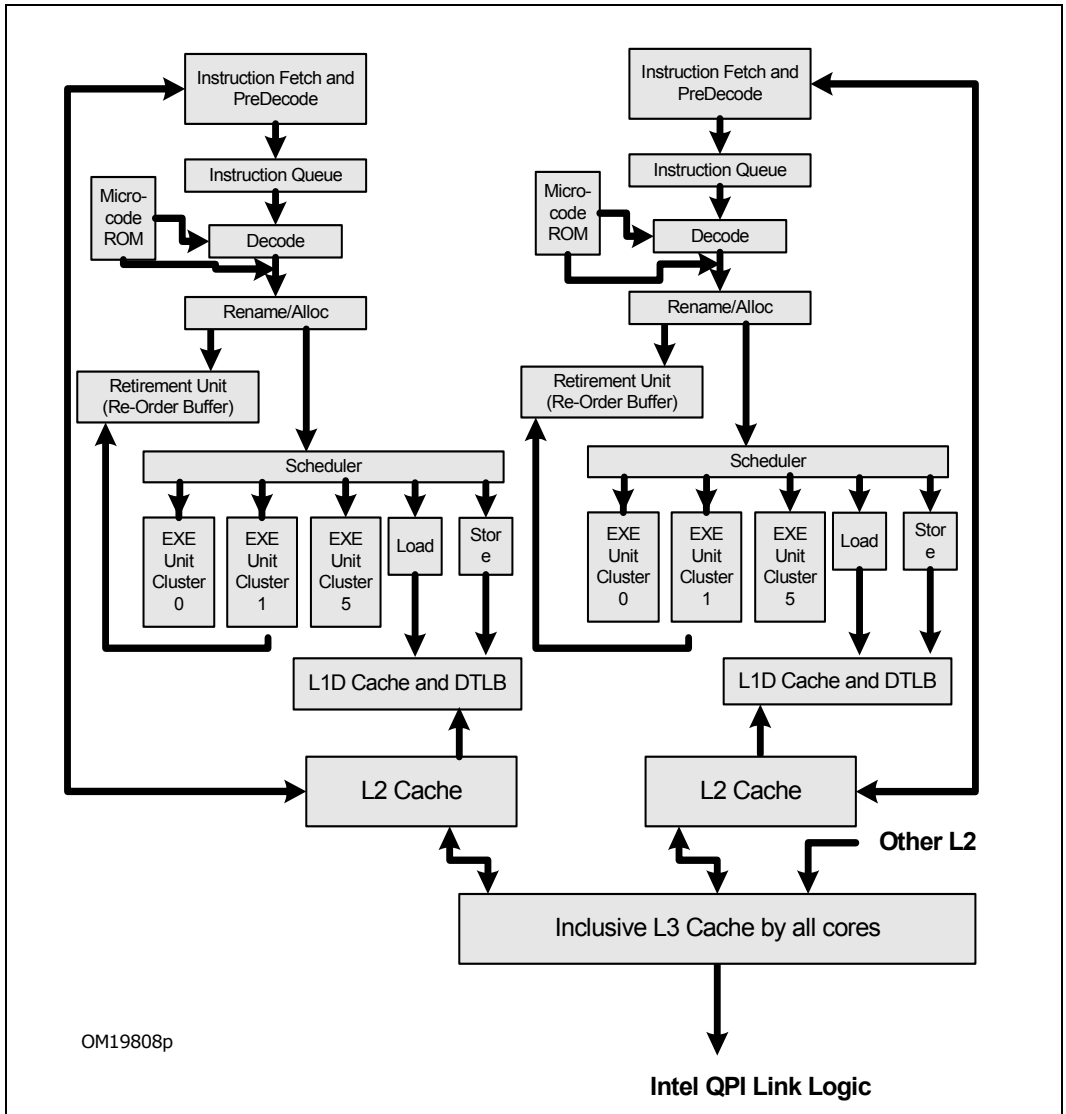


Figure 2-6. Intel microarchitecture code name Nehalem Pipeline Functionality

The length of the pipeline in Intel microarchitecture code name Nehalem is two cycles longer than its predecessor in 45nm Intel Core 2 processor family, as measured by

branch misprediction delay. The front end can decode up to 4 instructions in one cycle and supports two hardware threads by decoding the instruction streams between two logical processors in alternate cycles. The front end includes enhancement in branch handling, loop detection, MSROM throughput, etc. These are discussed in subsequent sections.

The scheduler (or reservation station) can dispatch up to six micro-ops in one cycle through six issue ports (five issue ports are shown in Figure 2-6; store operation involves separate ports for store address and store data but is depicted as one in the diagram).

The out-of-order engine has many execution units that are arranged in three execution clusters shown in Figure 2-6. It can retire four micro-ops in one cycle, same as its predecessor.

2.3.2 Front End Overview

Figure 2-7 depicts the key components of the front end of the microarchitecture. The instruction fetch unit (IFU) can fetch up to 16 bytes of aligned instruction bytes each cycle from the instruction cache to the instruction length decoder (ILD). The instruction queue (IQ) buffers the ILD-processed instructions and can deliver up to four instructions in one cycle to the instruction decoder.

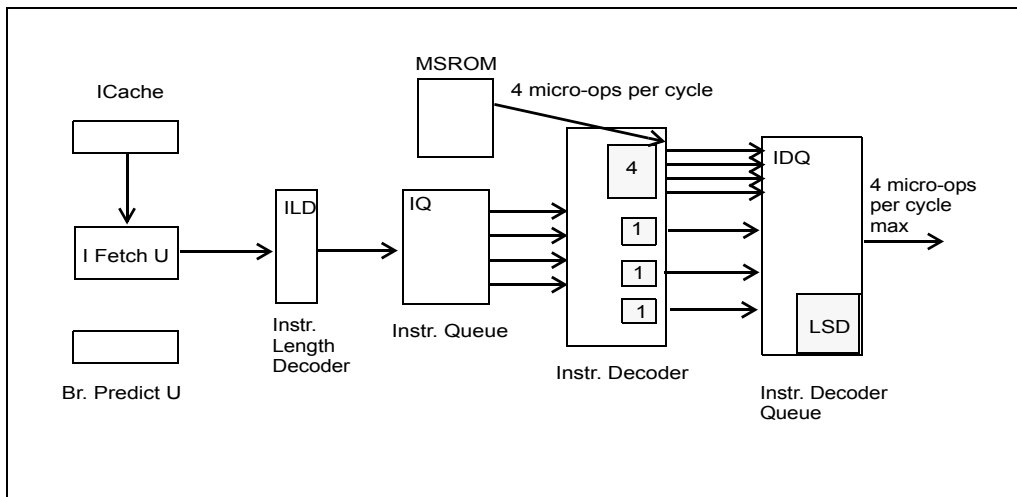


Figure 2-7. Front End of Intel microarchitecture code name Nehalem

The instruction decoder has three decoder units that can decode one simple instruction per cycle per unit. The other decoder unit can decode one instruction every cycle, either simple instruction or complex instruction made up of several micro-ops.

Instructions made up of more than four micro-ops are delivered from the MSROM. Up to four micro-ops can be delivered each cycle to the instruction decoder queue (IDQ).

The loop stream detector is located inside the IDQ to improve power consumption and front end efficiency for loops with a short sequence of instructions.

The instruction decoder supports micro-fusion to improve front end throughput, increase the effective size of queues in the scheduler and re-order buffer (ROB). The rules for micro-fusion are similar to those of Intel Core microarchitecture.

The instruction queue also supports macro-fusion to combine adjacent instructions into one micro-ops where possible. In previous generations of Intel Core microarchitecture, macro-fusion support for CMP/Jcc sequence is limited to the CF and ZF flag, and macrofusion is not supported in 64-bit mode.

In Intel microarchitecture code name Nehalem, macro-fusion is supported in 64-bit mode, and the following instruction sequences are supported:

- CMP or TEST can be fused when comparing (unchanged):
 - REG-REG. For example: `CMP EAX,ECX; JZ label`
 - REG-IMM. For example: `CMP EAX,0x80; JZ label`
 - REG-MEM. For example: `CMP EAX,[ECX]; JZ label`
 - MEM-REG. For example: `CMP [EAX],ECX; JZ label`
- TEST can fused with all conditional jumps (unchanged).
- CMP can be fused with the following conditional jumps. These conditional jumps check carry flag (CF) or zero flag (ZF). The list of macro-fusion-capable conditional jumps are (unchanged):
 - JA or JNBE
 - JAE or JNB or JNC
 - JE or JZ
 - JNA or JBE
 - JNAE or JC or JB
 - JNE or JNZ
- CMP can be fused with the following conditional jumps in Intel microarchitecture code name Nehalem, (this is an enhancement):

- JL or JNGE
- JGE or JNL
- JLE or JNG
- JG or JNLE

The hardware improves branch handling in several ways. Branch target buffer has increased to increase the accuracy of branch predictions. Renaming is supported with return stack buffer to reduce mispredictions of return instructions in the code.

Furthermore, hardware enhancement improves the handling of branch misprediction by expediting resource reclamation so that the front end would not be waiting to decode instructions in an architected code path (the code path in which instructions

will reach retirement) while resources were allocated to executing mispredicted code path. Instead, new micro-ops stream can start forward progress as soon as the front end decodes the instructions in the architected code path.

2.3.3 Execution Engine

The IDQ (Figure 2-7) delivers micro-op stream to the allocation/renaming stage (Figure 2-6) of the pipeline. The out-of-order engine supports up to 128 micro-ops in flight. Each micro-ops must be allocated with the following resources: an entry in the re-order buffer (ROB), an entry in the reservation station (RS), and a load/store buffer if a memory access is required.

The allocator also renames the register file entry of each micro-op in flight. The input data associated with a micro-op are generally either read from the ROB or from the retired register file.

The RS is expanded to 36 entry deep (compared to 32 entries in previous generation). It can dispatch up to six micro-ops in one cycle if the micro-ops are ready to execute. The RS dispatch a micro-op through an issue port to a specific execution cluster, each cluster may contain a collection of integer/FP/SIMD execution units.

The result from the execution unit executing a micro-op is written back to the register file, or forwarded through a bypass network to a micro-op in-flight that needs the result. Intel microarchitecture code name Nehalem can support write back throughput of one register file write per cycle per port. The bypass network consists of three domains of integer/FP/SIMD. Forwarding the result within the same bypass domain from a producer micro-op to a consumer micro is done efficiently in hardware without delay. Forwarding the result across different bypass domains may be subject to additional bypass delays. The bypass delays may be visible to software in addition to the latency and throughput characteristics of individual execution units. The bypass delays between a producer micro-op and a consumer micro-op across different bypass domains are shown in Table 2-18.

Table 2-18. Bypass Delay Between Producer and Consumer Micro-ops (cycles)

	FP	Integer	SIMD
FP	0	2	2
Integer	2	0	1
SIMD	2	1	0

2.3.3.1 Issue Ports and Execution Units

Table 2-19 summarizes the key characteristics of the issue ports and the execution unit latency/throughputs for common operations in the microarchitecture.

Table 2-19. Issue Ports of Intel microarchitecture code name Nehalem

Port	Executable operations	Latency	Throughput	Domain	Comment
Port 0	Integer ALU	1	1	Integer	
	Integer Shift	1	1		
Port 0	Integer SIMD ALU	1	1	SIMD	
	Integer SIMD Shuffle	1	1		
Port 0	Single-precision (SP) FP MUL	4	1	FP	
	Double-precision FP MUL	5	1		
	FP MUL (X87)	5	1		
	FP/SIMD/SSE2 Move and Logic	1	1		
	FP Shuffle DIV/SQRT	1	1		
Port 1	Integer ALU	1	1	Integer	
	Integer LEA	1	1		
	Integer Mul	3	1		
Port 1	Integer SIMD MUL	1	1	SIMD	
	Integer SIMD Shift	1	1		
	PSAD	3	1		
	StringCompare				
Port 1	FP ADD	3	1	FP	
Port 2	Integer loads	4	1	Integer	
Port 3	Store address	5	1	Integer	
Port 4	Store data			Integer	
Port 5	Integer ALU	1	1	Integer	
	Integer Shift	1	1		
	Jmp	1	1		
Port 5	Integer SIMD ALU	1	1	SIMD	
	Integer SIMD Shuffle	1	1		

Table 2-19. Issue Ports of Intel microarchitecture code name Nehalem

Port	Executable operations	Latency	Throughput	Domain	Comment
Port 5	FP/SIMD/SSE2 Move and Logic	1	1	FP	

2.3.4 Cache and Memory Subsystem

Intel microarchitecture code name Nehalem contains an instruction cache, a first-level data cache and a second-level unified cache in each core (see Figure 2-6). Each physical processor may contain several processor cores and a shared collection of sub-systems that are referred to as “uncore”. Specifically in Intel Core i7 processor, the uncore provides a unified third-level cache shared by all cores in the physical processor, Intel QuickPath Interconnect links and associated logic. The L1 and L2 caches are writeback and non-inclusive.

The shared L3 cache is writeback and inclusive, such that a cache line that exists in either L1 data cache, L1 instruction cache, unified L2 cache also exists in L3. The L3 is designed to use the inclusive nature to minimize snoop traffic between processor cores. Table 2-20 lists characteristics of the cache hierarchy. The latency of L3 access may vary as a function of the frequency ratio between the processor and the uncore sub-system.

Table 2-20. Cache Parameters of Intel Core i7 Processors

Level	Capacity	Associativity (ways)	Line Size (bytes)	Access Latency (clocks)	Access Throughput (clocks)	Write Update Policy
First Level Data	32 KB	8	64	4	1	Writeback
Instruction	32 KB	4	N/A	N/A	N/A	N/A
Second Level	256KB	8	64	10 ¹	Varies	Writeback
Third Level (Shared L3) ²	8MB	16	64	35-40+ ²	Varies	Writeback

NOTES:

1. Software-visible latency will vary depending on access patterns and other factors.
2. Minimal L3 latency is 35 cycles if the frequency ratio between core and uncore is unity.

The Intel microarchitecture code name Nehalem implements two levels of translation lookaside buffer (TLB). The first level consists of separate TLBs for data and code. DTLB0 handles address translation for data accesses, it provides 64 entries to support 4KB pages and 32 entries for large pages. The ITLB provides 64 entries (per thread) for 4KB pages and 7 entries (per thread) for large pages.

The second level TLB (STLB) handles both code and data accesses for 4KB pages. It support 4KB page translation operation that missed DTLB0 or ITLB. All entries are 4-way associative. Here is a list of entries in each DTLB:

- STLB for 4-KByte pages: 512 entries (services both data and instruction look-ups)
- DTLB0 for large pages: 32 entries
- DTLB0 for 4-KByte pages: 64 entries

An DTLB0 miss and STLB hit causes a penalty of 7cycles. Software only pays this penalty if the DTLB0 is used in some dispatch cases. The delays associated with a miss to the STLB and PMH are largely non-blocking.

2.3.5 Load and Store Operation Enhancements

The memory cluster of Intel microarchitecture code name Nehalem provides the following enhancements to speed up memory operations:

- Peak issue rate of one 128-bit load and one 128-bit store operation per cycle
- Deeper buffers for load and store operations: 48 load buffers, 32 store buffers and 10 fill buffers
- Fast unaligned memory access and robust handling of memory alignment hazards
- Improved store-forwarding for aligned and non-aligned scenarios
- Store forwarding for most address alignments

2.3.5.1 Efficient Handling of Alignment Hazards

The cache and memory subsystems handles a significant percentage of instructions in every workload. Different address alignment scenarios will produce varying performance impact for memory and cache operations. For example, 1-cycle throughput of L1 (see Table 2-21) generally applies to naturally-aligned loads from L1 cache. But using unaligned load instructions (e.g. MOVUPS, MOVUPD, MOVDQU, etc.) to access data from L1 will experience varying amount of delays depending on specific microarchitectures and alignment scenarios.

Table 2-21. Performance Impact of Address Alignments of MOVDQU from L1

Throughput (cycle)	Intel Core i7 Processor	45 nm Intel Core Microarchitecture	65 nm Intel Core Microarchitecture
Alignment Scenario	06_1AH	06_17H	06_0FH
16B aligned	1	2	2

Table 2-21. Performance Impact of Address Alignments of MOVDQU from L1

Throughput (cycle)	Intel Core i7 Processor	45 nm Intel Core Microarchitecture	65 nm Intel Core Microarchitecture
Alignment Scenario	06_1AH	06_17H	06_0FH
Not-16B aligned, not cache split	1	~2	~2
Split cache line boundary	~4.5	~20	~20

Table 2-21 lists approximate throughput of issuing MOVDQU instructions with different address alignment scenarios to load data from the L1 cache. If a 16-byte load spans across cache line boundary, previous microarchitecture generations will experience significant software-visible delays.

Intel microarchitecture code name Nehalem provides hardware enhancements to reduce the delays of handling different address alignment scenarios including cache line splits.

2.3.5.2 Store Forwarding Enhancement

When a load follows a store and reloads the data that the store writes to memory, the microarchitecture can forward the data directly from the store to the load in many cases. This situation, called store to load forwarding, saves several cycles by enabling the load to obtain the data directly from the store operation instead of through the memory system.

Several general rules must be met for store to load forwarding to proceed without delay:

- The store must be the last store to that address prior to the load.
- The store must be equal or greater in size than the size of data being loaded.
- The load data must be completely contained in the preceding store.

Specific address alignment and data sizes between the store and load operations will determine whether a store-forward situation may proceed with data forwarding or experience a delay via the cache/memory sub-system. The 45 nm Enhanced Intel Core microarchitecture offers more flexible address alignment and data sizes requirement than previous microarchitectures. Intel microarchitecture code name Nehalem offers additional enhancement with allowing more situations to forward data expeditiously.

The store-forwarding situations for with respect to store operations of 16 bytes are illustrated in Figure .

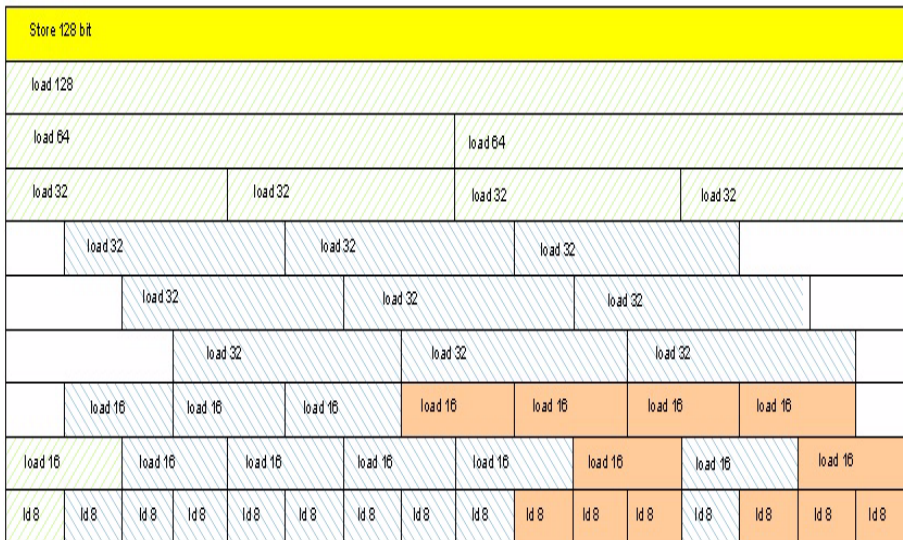
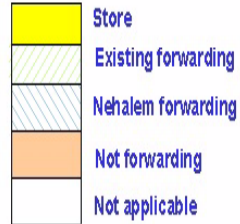
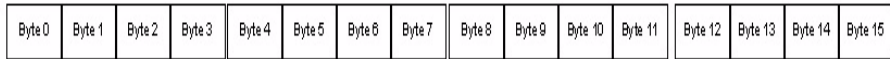


Figure 2-8. Store-Forwarding Scenarios of 16-Byte Store Operations

Intel microarchitecture code name Nehalem allows store-to-load forwarding to proceed regardless of store address alignment (The white space in the diagram does not correspond to an applicable store-to-load scenario). Figure 2-9 illustrates situations for store operation of 8 bytes or less.

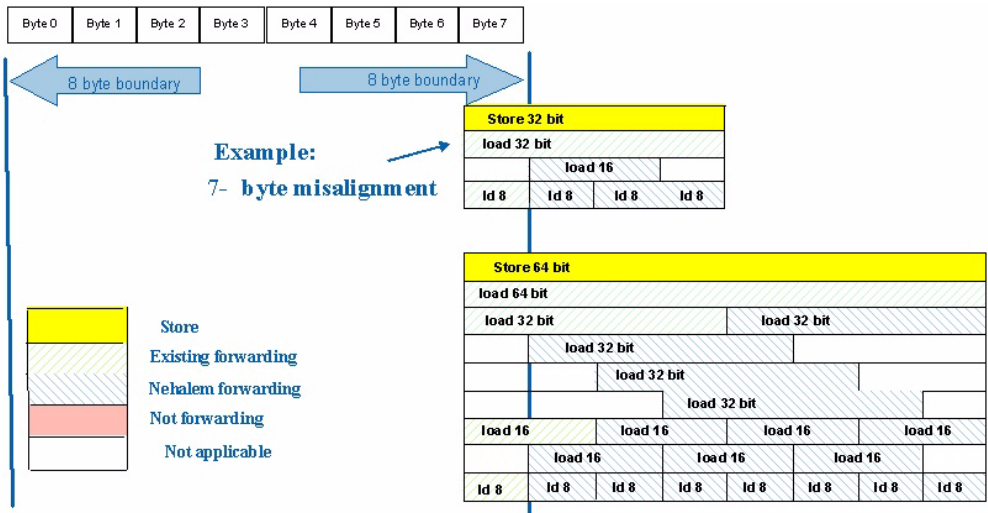


Figure 2-9. Store-Forwarding Enhancement in Intel microarchitecture code name Nehalem

2.3.6 REP String Enhancement

REP prefix in conjunction with MOVSW/STOS instruction and a count value in ECX are frequently used to implement library functions such as memcpy()/memset(). These are referred to as "REP string" instructions. Each iteration of these instruction can copy/write constant a value in byte/word/dword/qword granularity. The performance characteristics of using REP string can be attributed to two components: startup overhead and data transfer throughput.

The two components of performance characteristics of REP String varies further depending on granularity, alignment, and/or count values. Generally, MOVSB is used to handle very small chunks of data. Therefore, processor implementation of REP MOVSB is optimized to handle ECX < 4. Using REP MOVSB with ECX > 3 will achieve low data throughput due to not only byte-granular data transfer but also additional startup overhead. The latency for MOVSB, is 9 cycles if ECX < 4; otherwise REP MOVSB with ECX > 9 have a 50-cycle startup cost.

For REP string of larger granularity data transfer, as ECX value increases, the startup overhead of REP String exhibit step-wise increase:

- Short string (ECX ≤ 12): the latency of REP MOVSW/MOVSD/MOVSQ is about 20 cycles,
- Fast string (ECX ≥ 76: excluding REP MOVSB): the processor implementation provides hardware optimization by moving as many pieces of data in 16 bytes as

possible. The latency of REP string latency will vary if one of the 16-byte data transfer spans across cache line boundary:

- Split-free: the latency consists of a startup cost of about 40 cycles and each 64 bytes of data adds 4 cycles,
 - Cache splits: the latency consists of a startup cost of about 35 cycles and each 64 bytes of data adds 6 cycles.
- Intermediate string lengths: the latency of REP MOVSW/MOVSD/MOVSQ has a startup cost of about 15 cycles plus one cycle for each iteration of the data movement in word/dword/qword.

Intel microarchitecture code name Nehalem improves the performance of REP strings significantly over previous microarchitectures in several ways:

- Startup overhead have been reduced in most cases relative to previous microarchitecture,
- Data transfer throughput are improved over previous generation
- In order for REP string to operate in “fast string” mode, previous microarchitectures requires address alignment. In Intel microarchitecture code name Nehalem, REP string can operate in “fast string” mode even if address is not aligned to 16 bytes.

2.3.7 Enhancements for System Software

In addition to microarchitectural enhancements that can benefit both application-level and system-level software, Intel microarchitecture code name Nehalem enhances several operations that primarily benefit system software.

Lock primitives: Synchronization primitives using the Lock prefix (e.g. XCHG, CMPXCHG8B) executes with significantly reduced latency than previous microarchitectures.

VMM overhead improvements: VMX transitions between a Virtual Machine (VM) and its supervisor (the VMM) can take thousands of cycle each time on previous microarchitectures. The latency of VMX transitions has been reduced in processors based on Intel microarchitecture code name Nehalem.

2.3.8 Efficiency Enhancements for Power Consumption

Intel microarchitecture code name Nehalem is not only designed for high performance and power-efficient performance under wide range of loading situations, it also features enhancement for low power consumption while the system idles. Intel microarchitecture code name Nehalem supports processor-specific C6 states, which have the lowest leakage power consumption that OS can manage through ACPI and OS power management mechanisms.

2.3.9 Hyper-Threading Technology Support in Intel® microarchitecture code name Nehalem

Intel microarchitecture code name Nehalem supports Hyper-Threading Technology (HT). Its implementation of HT provides two logical processors sharing most execution/cache resources in each core. The HT implementation in Intel microarchitecture code name Nehalem differs from previous generations of HT implementations using Intel NetBurst microarchitecture in several areas:

- Intel microarchitecture code name Nehalem provides four-wide execution engine, more functional execution units coupled to three issue ports capable of issuing computational operations.
- Intel microarchitecture code name Nehalem supports integrated memory controller that can provide peak memory bandwidth of up to 25.6 GB/sec in Intel Core i7 processor.
- Deeper buffering and enhanced resource sharing/partition policies:
 - Replicated resource for HT operation: register state, renamed return stack buffer, large-page ITLB
 - Partitioned resources for HT operation: load buffers, store buffers, re-order buffers, small-page ITLB are statically allocated between two logical processors.
 - Competitively-shared resource during HT operation: the reservation station, cache hierarchy, fill buffers, both DTLB0 and STLB.
 - Alternating during HT operation: front-end operation generally alternates between two logical processors to ensure fairness.
 - HT unaware resources: execution units.

2.4 INTEL NETBURST® MICROARCHITECTURE

The Pentium 4 processor, Pentium 4 processor Extreme Edition supporting Hyper-Threading Technology, Pentium D processor, and Pentium processor Extreme Edition implement the Intel NetBurst microarchitecture. Intel Xeon processors that implement Intel NetBurst microarchitecture can be identified using CPUID (family encoding 0FH).

This section describes the features of the Intel NetBurst microarchitecture and its operation common to the above processors. It provides the technical background required to understand optimization recommendations and the coding rules discussed in the rest of this manual. For implementation details, including instruction latencies, see Appendix C, "Instruction Latency and Throughput."

Intel NetBurst microarchitecture is designed to achieve high performance for integer and floating-point computations at high clock rates. It supports the following features:

- hyper-pipelined technology that enables high clock rates
- high-performance, quad-pumped bus interface to the Intel NetBurst microarchitecture system bus
- rapid execution engine to reduce the latency of basic integer instructions
- out-of-order speculative execution to enable parallelism
- superscalar issue to enable parallelism
- hardware register renaming to avoid register name space limitations
- cache line sizes of 64 bytes
- hardware prefetch

2.4.1 Design Goals

The design goals of Intel NetBurst microarchitecture are:

- To execute legacy IA-32 applications and applications based on single-instruction, multiple-data (SIMD) technology at high throughput
- To operate at high clock rates and to scale to higher performance and clock rates in the future

Design advances of the Intel NetBurst microarchitecture include:

- A deeply pipelined design that allows for high clock rates (with different parts of the chip running at different clock rates).
- A pipeline that optimizes for the common case of frequently executed instructions; the most frequently-executed instructions in common circumstances (such as a cache hit) are decoded efficiently and executed with short latencies.
- Employment of techniques to hide stall penalties; Among these are parallel execution, buffering, and speculation. The microarchitecture executes instructions dynamically and out-of-order, so the time it takes to execute each individual instruction is not always deterministic.

Chapter 3, “General Optimization Guidelines,” lists optimizations to use and situations to avoid. The chapter also gives a sense of relative priority. Because most optimizations are implementation dependent, the chapter does not quantify expected benefits and penalties.

The following sections provide more information about key features of the Intel NetBurst microarchitecture.

2.4.2 Pipeline

The pipeline of the Intel NetBurst microarchitecture contains:

- an in-order issue front end
- an out-of-order superscalar execution core

- an in-order retirement unit

The front end supplies instructions in program order to the out-of-order core. It fetches and decodes instructions. The decoded instructions are translated into μ ops. The front end's primary job is to feed a continuous stream of μ ops to the execution core in original program order.

The out-of-order core aggressively reorders μ ops so that μ ops whose inputs are ready (and have execution resources available) can execute as soon as possible. The core can issue multiple μ ops per cycle.

The retirement section ensures that the results of execution are processed according to original program order and that the proper architectural states are updated.

Figure 2-5 illustrates a diagram of the major functional blocks associated with the Intel NetBurst microarchitecture pipeline. The following subsections provide an overview for each.

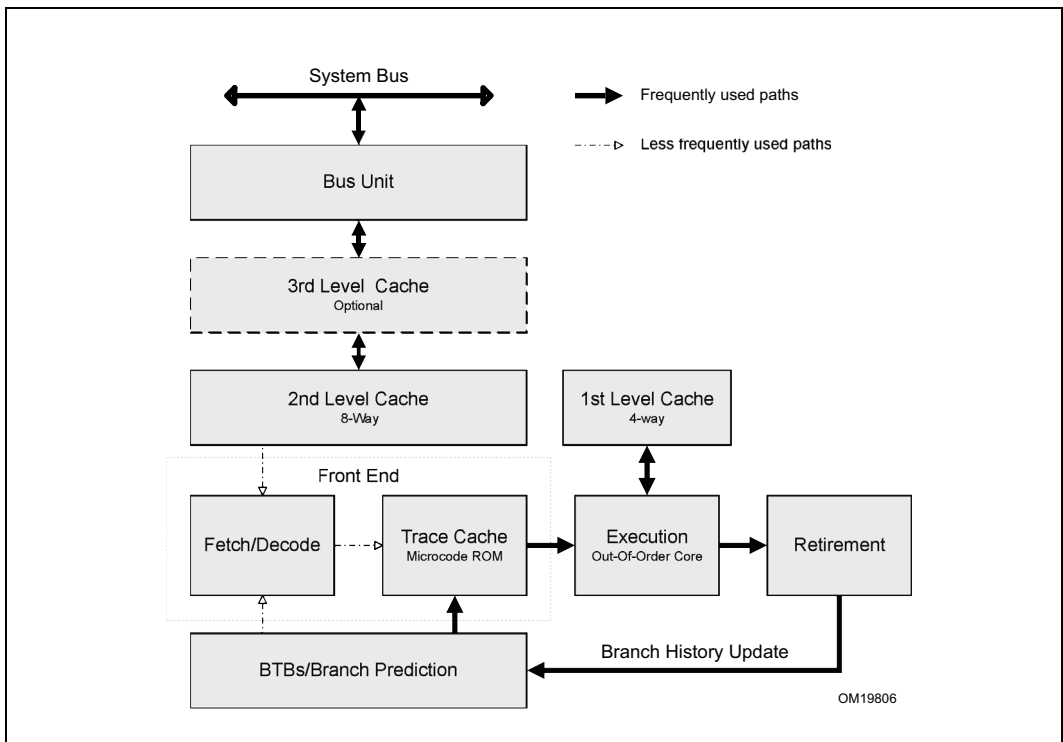


Figure 2-10. The Intel NetBurst Microarchitecture

2.4.2.1 Front End

The front end of the Intel NetBurst microarchitecture consists of two parts:

- fetch/decode unit
- execution trace cache

It performs the following functions:

- prefetches instructions that are likely to be executed
- fetches required instructions that have not been prefetched
- decodes instructions into μ ops
- generates microcode for complex instructions and special-purpose code
- delivers decoded instructions from the execution trace cache
- predicts branches using advanced algorithms

The front end is designed to address two problems that are sources of delay:

- time required to decode instructions fetched from the target
- wasted decode bandwidth due to branches or a branch target in the middle of a cache line

Instructions are fetched and decoded by a translation engine. The translation engine then builds decoded instructions into μ op sequences called traces. Next, traces are then stored in the execution trace cache.

The execution trace cache stores μ ops in the path of program execution flow, where the results of branches in the code are integrated into the same cache line. This increases the instruction flow from the cache and makes better use of the overall cache storage space since the cache no longer stores instructions that are branched over and never executed.

The trace cache can deliver up to 3 μ ops per clock to the core.

The execution trace cache and the translation engine have cooperating branch prediction hardware. Branch targets are predicted based on their linear address using branch prediction logic and fetched as soon as possible. Branch targets are fetched from the execution trace cache if they are cached, otherwise they are fetched from the memory hierarchy. The translation engine's branch prediction information is used to form traces along the most likely paths.

2.4.2.2 Out-of-order Core

The core's ability to execute instructions out of order is a key factor in enabling parallelism. This feature enables the processor to reorder instructions so that if one μ op is delayed while waiting for data or a contended resource, other μ ops that appear later in the program order may proceed. This implies that when one portion of the pipeline experiences a delay, the delay may be covered by other operations executing in parallel or by the execution of μ ops queued up in a buffer.

The core is designed to facilitate parallel execution. It can dispatch up to six μ ops per cycle through the issue ports (Figure 2-6). Note that six μ ops per cycle exceeds the trace cache and retirement μ op bandwidth. The higher bandwidth in the core allows for peak bursts of greater than three μ ops and to achieve higher issue rates by allowing greater flexibility in issuing μ ops to different execution ports.

Most core execution units can start executing a new μ op every cycle, so several instructions can be in flight at one time in each pipeline. A number of arithmetic logical unit (ALU) instructions can start at two per cycle; many floating-point instructions start one every two cycles. Finally, μ ops can begin execution out of program order, as soon as their data inputs are ready and resources are available.

2.4.2.3 Retirement

The retirement section receives the results of the executed μ ops from the execution core and processes the results so that the architectural state is updated according to the original program order. For semantically correct execution, the results of Intel 64 and IA-32 instructions must be committed in original program order before they are retired. Exceptions may be raised as instructions are retired. For this reason, exceptions cannot occur speculatively.

When a μ op completes and writes its result to the destination, it is retired. Up to three μ ops may be retired per cycle. The reorder buffer (ROB) is the unit in the processor which buffers completed μ ops, updates the architectural state and manages the ordering of exceptions.

The retirement section also keeps track of branches and sends updated branch target information to the branch target buffer (BTB). This updates branch history. Figure 2-10 illustrates the paths that are most frequently executing inside the Intel NetBurst microarchitecture: an execution loop that interacts with multilevel cache hierarchy and the system bus.

The following sections describe in more detail the operation of the front end and the execution core. This information provides the background for using the optimization techniques and instruction latency data documented in this manual.

2.4.3 Front End Pipeline Detail

The following information about the front end operation is be useful for tuning software with respect to prefetching, branch prediction, and execution trace cache operations.

2.4.3.1 Prefetching

The Intel NetBurst microarchitecture supports three prefetching mechanisms:

- a hardware instruction fetcher that automatically prefetches instructions

- a hardware mechanism that automatically fetches data and instructions into the unified second-level cache
- a mechanism fetches data only and includes two distinct components: (1) a hardware mechanism to fetch the adjacent cache line within a 128-byte sector that contains the data needed due to a cache line miss, this is also referred to as adjacent cache line prefetch (2) a software controlled mechanism that fetches data into the caches using the prefetch instructions.

The hardware instruction fetcher reads instructions along the path predicted by the branch target buffer (BTB) into instruction streaming buffers. Data is read in 32-byte chunks starting at the target address. The second and third mechanisms are described later.

2.4.3.2 Decoder

The front end of the Intel NetBurst microarchitecture has a single decoder that decodes instructions at the maximum rate of one instruction per clock. Some complex instructions must enlist the help of the microcode ROM. The decoder operation is connected to the execution trace cache.

2.4.3.3 Execution Trace Cache

The execution trace cache (TC) is the primary instruction cache in the Intel NetBurst microarchitecture. The TC stores decoded instructions (μ ops).

In the Pentium 4 processor implementation, TC can hold up to 12-Kbyte μ ops and can deliver up to three μ ops per cycle. TC does not hold all of the μ ops that need to be executed in the execution core. In some situations, the execution core may need to execute a microcode flow instead of the μ op traces that are stored in the trace cache.

The Pentium 4 processor is optimized so that most frequently-executed instructions come from the trace cache while only a few instructions involve the microcode ROM.

2.4.3.4 Branch Prediction

Branch prediction is important to the performance of a deeply pipelined processor. It enables the processor to begin executing instructions long before the branch outcome is certain. Branch delay is the penalty that is incurred in the absence of correct prediction. For Pentium 4 and Intel Xeon processors, the branch delay for a correctly predicted instruction can be as few as zero clock cycles. The branch delay for a mispredicted branch can be many cycles, usually equivalent to the pipeline depth.

Branch prediction in the Intel NetBurst microarchitecture predicts near branches (conditional calls, unconditional calls, returns and indirect branches). It does not predict far transfers (far calls, irets and software interrupts).

Mechanisms have been implemented to aid in predicting branches accurately and to reduce the cost of taken branches. These include:

- ability to dynamically predict the direction and target of branches based on an instruction's linear address, using the branch target buffer (BTB)
- if no dynamic prediction is available or if it is invalid, the ability to statically predict the outcome based on the offset of the target: a backward branch is predicted to be taken, a forward branch is predicted to be not taken
- ability to predict return addresses using the 16-entry return address stack
- ability to build a trace of instructions across predicted taken branches to avoid branch penalties

The Static Predictor. Once a branch instruction is decoded, the direction of the branch (forward or backward) is known. If there was no valid entry in the BTB for the branch, the static predictor makes a prediction based on the direction of the branch. The static prediction mechanism predicts backward conditional branches (those with negative displacement, such as loop-closing branches) as taken. Forward branches are predicted not taken.

To take advantage of the forward-not-taken and backward-taken static predictions, code should be arranged so that the likely target of the branch immediately follows forward branches (see also Section 3.4.1, "Branch Prediction Optimization").

Branch Target Buffer. Once branch history is available, the Pentium 4 processor can predict the branch outcome even before the branch instruction is decoded. The processor uses a branch history table and a branch target buffer (collectively called the BTB) to predict the direction and target of branches based on an instruction's linear address. Once the branch is retired, the BTB is updated with the target address.

Return Stack. Returns are always taken; but since a procedure may be invoked from several call sites, a single predicted target does not suffice. The Pentium 4 processor has a Return Stack that can predict return addresses for a series of procedure calls. This increases the benefit of unrolling loops containing function calls. It also mitigates the need to put certain procedures inline since the return penalty portion of the procedure call overhead is reduced.

Even if the direction and target address of the branch are correctly predicted, a taken branch may reduce available parallelism in a typical processor (since the decode bandwidth is wasted for instructions which immediately follow the branch and precede the target, if the branch does not end the line and target does not begin the line). The branch predictor allows a branch and its target to coexist in a single trace cache line, maximizing instruction delivery from the front end.

2.4.4 Execution Core Detail

The execution core is designed to optimize overall performance by handling common cases most efficiently. The hardware is designed to execute frequent operations in a

common context as fast as possible, at the expense of infrequent operations using rare contexts.

Some parts of the core may speculate that a common condition holds to allow faster execution. If it does not, the machine may stall. An example of this pertains to store-to-load forwarding (see “Store Forwarding” in this chapter). If a load is predicted to be dependent on a store, it gets its data from that store and tentatively proceeds. If the load turned out not to depend on the store, the load is delayed until the real data has been loaded from memory, then it proceeds.

2.4.4.1 Instruction Latency and Throughput

The superscalar out-of-order core contains hardware resources that can execute multiple μ ops in parallel. The core’s ability to make use of available parallelism of execution units can be enhanced by software’s ability to:

- Select instructions that can be decoded in less than 4 μ ops and/or have short latencies
- Order instructions to preserve available parallelism by minimizing long dependence chains and covering long instruction latencies
- Order instructions so that their operands are ready and their corresponding issue ports and execution units are free when they reach the scheduler

This subsection describes port restrictions, result latencies, and issue latencies (also referred to as throughput). These concepts form the basis to assist software for ordering instructions to increase parallelism. The order that μ ops are presented to the core of the processor is further affected by the machine’s scheduling resources.

It is the execution core that reacts to an ever-changing machine state, reordering μ ops for faster execution or delaying them because of dependence and resource constraints. The ordering of instructions in software is more of a suggestion to the hardware.

Appendix C, “Instruction Latency and Throughput,” lists some of the more-commonly-used Intel 64 and IA-32 instructions with their latency, their issue throughput, and associated execution units (where relevant). Some execution units are not pipelined (meaning that μ ops cannot be dispatched in consecutive cycles and the throughput is less than one per cycle). The number of μ ops associated with each instruction provides a basis for selecting instructions to generate. All μ ops executed out of the microcode ROM involve extra overhead.

2.4.4.2 Execution Units and Issue Ports

At each cycle, the core may dispatch μ ops to one or more of four issue ports. At the microarchitecture level, store operations are further divided into two parts: store data and store address operations. The four ports through which μ ops are dispatched to execution units and to load and store operations are shown in Figure 2-6. Some ports can dispatch two μ ops per clock. Those execution units are marked Double Speed.

Port 0. In the first half of the cycle, port 0 can dispatch either one floating-point move μ op (a floating-point stack move, floating-point exchange or floating-point store data) or one arithmetic logical unit (ALU) μ op (arithmetic, logic, branch or store data). In the second half of the cycle, it can dispatch one similar ALU μ op.

Port 1. In the first half of the cycle, port 1 can dispatch either one floating-point execution (all floating-point operations except moves, all SIMD operations) μ op or one normal-speed integer (multiply, shift and rotate) μ op or one ALU (arithmetic) μ op. In the second half of the cycle, it can dispatch one similar ALU μ op.

Port 2. This port supports the dispatch of one load operation per cycle.

Port 3. This port supports the dispatch of one store address operation per cycle.

The total issue bandwidth can range from zero to six μ ops per cycle. Each pipeline contains several execution units. The μ ops are dispatched to the pipeline that corresponds to the correct type of operation. For example, an integer arithmetic logic unit and the floating-point execution units (adder, multiplier, and divider) can share a pipeline.

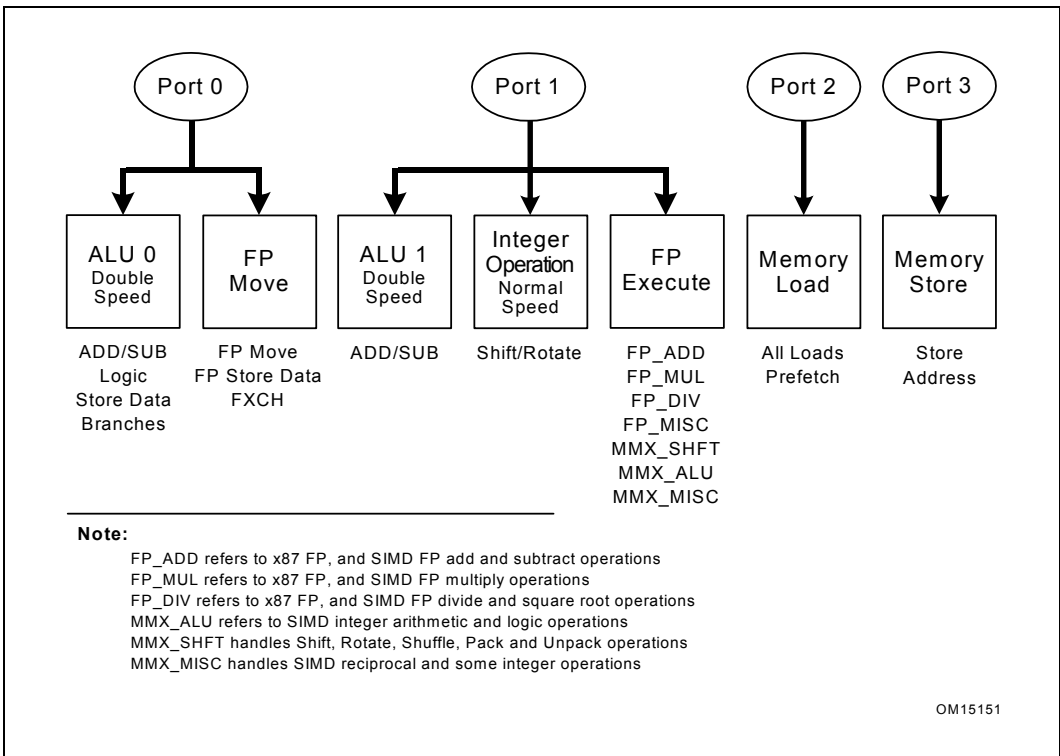


Figure 2-11. Execution Units and Ports in Out-Of-Order Core

2.4.4.3 Caches

The Intel NetBurst microarchitecture supports up to three levels of on-chip cache. At least two levels of on-chip cache are implemented in processors based on the Intel NetBurst microarchitecture. The Intel Xeon processor MP and selected Pentium and Intel Xeon processors may also contain a third-level cache.

The first level cache (nearest to the execution core) contains separate caches for instructions and data. These include the first-level data cache and the trace cache (an advanced first-level instruction cache). All other caches are shared between instructions and data.

Levels in the cache hierarchy are not inclusive. The fact that a line is in level *i* does not imply that it is also in level *i*+1. All caches use a pseudo-LRU (least recently used) replacement algorithm.

Table 2-5 provides parameters for all cache levels for Pentium and Intel Xeon Processors with CPUID model encoding equals 0, 1, 2 or 3.

Table 2-22. Pentium 4 and Intel Xeon Processor Cache Parameters

Level (Model)	Capacity	Associativity (ways)	Line Size (bytes)	Access Latency, Integer/ floating-point (clocks)	Write Update Policy
First (Model 0, 1, 2)	8 KB	4	64	2/9	write through
First (Model 3)	16 KB	8	64	4/12	write through
TC (All models)	12K μ ops	8	N/A	N/A	N/A
Second (Model 0, 1, 2)	256 KB or 512 KB ¹	8	64 ²	7/7	write back
Second (Model 3, 4)	1 MB	8	64 ²	18/18	write back
Second (Model 3, 4, 6)	2 MB	8	64 ²	20/20	write back
Third (Model 0, 1, 2)	0, 512 KB, 1 MB or 2 MB	8	64 ²	14/14	write back

NOTES:

1. Pentium 4 and Intel Xeon processors with CPUID model encoding value of 2 have a second level cache of 512 KB.
2. Each read due to a cache miss fetches a sector, consisting of two adjacent cache lines; a write operation is 64 bytes.

On processors without a third level cache, the second-level cache miss initiates a transaction across the system bus interface to the memory sub-system. On processors with a third level cache, the third-level cache miss initiates a transaction across the system bus. A bus write transaction writes 64 bytes to cacheable memory, or separate 8-byte chunks if the destination is not cacheable. A bus read transaction from cacheable memory fetches two cache lines of data.

The system bus interface supports using a scalable bus clock and achieves an effective speed that quadruples the speed of the scalable bus clock. It takes on the order of 12 processor cycles to get to the bus and back within the processor, and 6-12 bus cycles to access memory if there is no bus congestion. Each bus cycle equals several processor cycles. The ratio of processor clock speed to the scalable bus clock speed is referred to as bus ratio. For example, one bus cycle for a 100 MHz bus is equal to 15 processor cycles on a 1.50 GHz processor. Since the speed of the bus is implementation-dependent, consult the specifications of a given system for further details.

2.4.4.4 Data Prefetch

The Pentium 4 processor and other processors based on the NetBurst microarchitecture have two type of mechanisms for prefetching data: software prefetch instructions and hardware-based prefetch mechanisms.

Software controlled prefetch is enabled using the four prefetch instructions (PREFETCHh) introduced with SSE. The software prefetch is not intended for prefetching code. Using it can incur significant penalties on a multiprocessor system if code is shared.

Software prefetch can provide benefits in selected situations. These situations include when:

- the pattern of memory access operations in software allows the programmer to hide memory latency
- a reasonable choice can be made about how many cache lines to fetch ahead of the line being execute
- a choice can be made about the type of prefetch to use

SSE prefetch instructions have different behaviors, depending on cache levels updated and the processor implementation. For instance, a processor may implement the non-temporal prefetch by returning data to the cache level closest to the processor core. This approach has the following effect:

- minimizes disturbance of temporal data in other cache levels
- avoids the need to access off-chip caches, which can increase the realized bandwidth compared to a normal load-miss, which returns data to all cache levels

Situations that are less likely to benefit from software prefetch are:

- For cases that are already bandwidth bound, prefetching tends to increase bandwidth demands.

- Prefetching far ahead can cause eviction of cached data from the caches prior to the data being used in execution.
- Not prefetching far enough can reduce the ability to overlap memory and execution latencies.

Software prefetches are treated by the processor as a hint to initiate a request to fetch data from the memory system, and consume resources in the processor and the use of too many prefetches can limit their effectiveness. Examples of this include prefetching data in a loop for a reference outside the loop and prefetching in a basic block that is frequently executed, but which seldom precedes the reference for which the prefetch is targeted.

See: Chapter 7, “Optimizing Cache Usage.”

Automatic hardware prefetch is a feature in the Pentium 4 processor. It brings cache lines into the unified second-level cache based on prior reference patterns.

Software prefetching has the following characteristics:

- handles irregular access patterns, which do not trigger the hardware prefetcher
- handles prefetching of short arrays and avoids hardware prefetching start-up delay before initiating the fetches
- must be added to new code; so it does not benefit existing applications

Hardware prefetching for Pentium 4 processor has the following characteristics:

- works with existing applications
- does not require extensive study of prefetch instructions
- requires regular access patterns
- avoids instruction and issue port bandwidth overhead
- has a start-up penalty before the hardware prefetcher triggers and begins initiating fetches

The hardware prefetcher can handle multiple streams in either the forward or backward directions. The start-up delay and fetch-ahead has a larger effect for short arrays when hardware prefetching generates a request for data beyond the end of an array (not actually utilized). The hardware penalty diminishes if it is amortized over longer arrays.

Hardware prefetching is triggered after two successive cache misses in the last level cache and requires these cache misses to satisfy a condition that the linear address distance between these cache misses is within a threshold value. The threshold value depends on the processor implementation (see Table 2-6). However, hardware prefetching will not cross 4-KByte page boundaries. As a result, hardware prefetching can be very effective when dealing with cache miss patterns that have small strides and that are significantly less than half the threshold distance to trigger hardware prefetching. On the other hand, hardware prefetching will not benefit cache miss patterns that have frequent DTLB misses or have access strides that cause successive cache misses that are spatially apart by more than the trigger threshold distance.

Software can proactively control data access pattern to favor smaller access strides (e.g., stride that is less than half of the trigger threshold distance) over larger access strides (stride that is greater than the trigger threshold distance), this can achieve additional benefit of improved temporal locality and reducing cache misses in the last level cache significantly.

Software optimization of a data access pattern should emphasize tuning for hardware prefetch first to favor greater proportions of smaller-stride data accesses in the workload; before attempting to provide hints to the processor by employing software prefetch instructions.

2.4.4.5 Loads and Stores

The Pentium 4 processor employs the following techniques to speed up the execution of memory operations:

- speculative execution of loads
- reordering of loads with respect to loads and stores
- multiple outstanding misses
- buffering of writes
- forwarding of data from stores to dependent loads

Performance may be enhanced by not exceeding the memory issue bandwidth and buffer resources provided by the processor. Up to one load and one store may be issued for each cycle from a memory port reservation station. In order to be dispatched to a reservation station, there must be a buffer entry available for each memory operation. There are 48 load buffers and 24 store buffers³. These buffers hold the μop and address information until the operation is completed, retired, and deallocated.

The Pentium 4 processor is designed to enable the execution of memory operations out of order with respect to other instructions and with respect to each other. Loads can be carried out speculatively, that is, before all preceding branches are resolved. However, speculative loads cannot cause page faults.

Reordering loads with respect to each other can prevent a load miss from stalling later loads. Reordering loads with respect to other loads and stores to different addresses can enable more parallelism, allowing the machine to execute operations as soon as their inputs are ready. Writes to memory are always carried out in program order to maintain program correctness.

A cache miss for a load does not prevent other loads from issuing and completing. The Pentium 4 processor supports up to four (or eight for Pentium 4 processor with CPUID signature corresponding to family 15, model 3) outstanding load misses that can be serviced either by on-chip caches or by memory.

3. Pentium 4 processors with CPUID model encoding equal to 3 have more than 24 store buffers.

Store buffers improve performance by allowing the processor to continue executing instructions without having to wait until a write to memory and/or cache is complete. Writes are generally not on the critical path for dependence chains, so it is often beneficial to delay writes for more efficient use of memory-access bus cycles.

2.4.4.6 Store Forwarding

Loads can be moved before stores that occurred earlier in the program if they are not predicted to load from the same linear address. If they do read from the same linear address, they have to wait for the store data to become available. However, with store forwarding, they do not have to wait for the store to write to the memory hierarchy and retire. The data from the store can be forwarded directly to the load, as long as the following conditions are met:

- **Sequence** — Data to be forwarded to the load has been generated by a programmatically-earlier store which has already executed.
- **Size** — Bytes loaded must be a subset of (including a proper subset, that is, the same) bytes stored.
- **Alignment** — The store cannot wrap around a cache line boundary, and the linear address of the load must be the same as that of the store.

2.5 INTEL® PENTIUM® M PROCESSOR MICROARCHITECTURE

Like the Intel NetBurst microarchitecture, the pipeline of the Intel Pentium M processor microarchitecture contains three sections:

- in-order issue front end
- out-of-order superscalar execution core
- in-order retirement unit

Intel Pentium M processor microarchitecture supports a high-speed system bus (up to 533 MHz) with 64-byte line size. Most coding recommendations that apply to the Intel NetBurst microarchitecture also apply to the Intel Pentium M processor.

The Intel Pentium M processor microarchitecture is designed for lower power consumption. There are other specific areas of the Pentium M processor microarchitecture that differ from the Intel NetBurst microarchitecture. They are described next. A block diagram of the Intel Pentium M processor is shown in Figure 2-7.

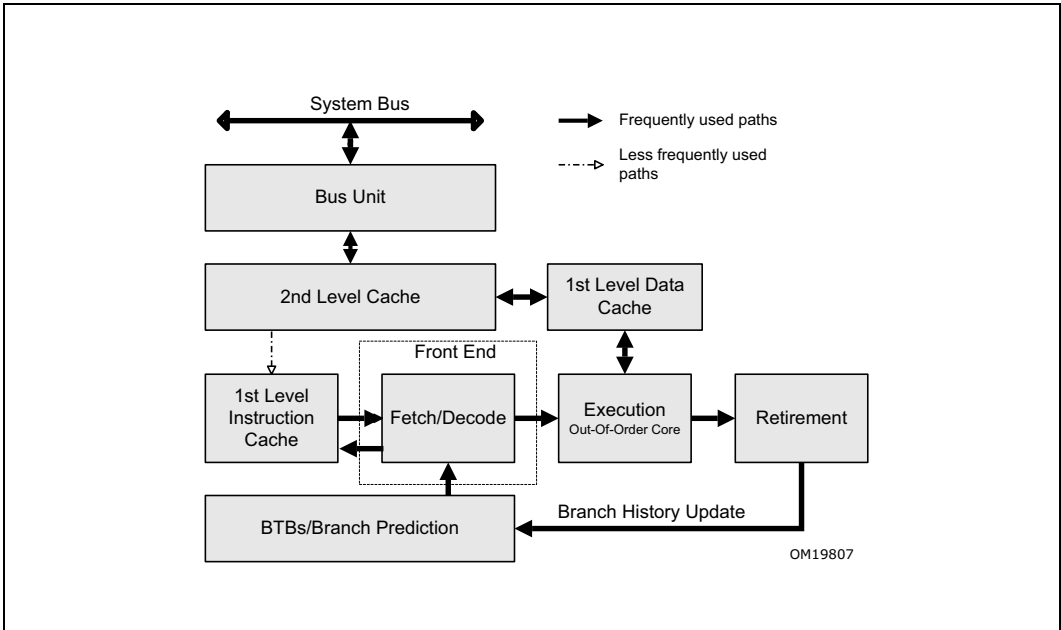


Figure 2-12. The Intel Pentium M Processor Microarchitecture

2.5.1 Front End

The Intel Pentium M processor uses a pipeline depth that enables high performance and low power consumption. It's shorter than that of the Intel NetBurst microarchitecture.

The Intel Pentium M processor front end consists of two parts:

- fetch/decode unit
- instruction cache

The fetch and decode unit includes a hardware instruction prefetcher and three decoders that enable parallelism. It also provides a 32-KByte instruction cache that stores un-decoded binary instructions.

The instruction prefetcher fetches instructions in a linear fashion from memory if the target instructions are not already in the instruction cache. The prefetcher is designed to fetch efficiently from an aligned 16-byte block. If the modulo 16 remainder of a branch target address is 14, only two useful instruction bytes are fetched in the first cycle. The rest of the instruction bytes are fetched in subsequent cycles.

The three decoders decode instructions and break them down into μ ops. In each clock cycle, the first decoder is capable of decoding an instruction with four or fewer

μops. The remaining two decoders each decode a one μop instruction in each clock cycle.

The front end can issue multiple μops per cycle, in original program order, to the out-of-order core.

The Intel Pentium M processor incorporates sophisticated branch prediction hardware to support the out-of-order core. The branch prediction hardware includes dynamic prediction, and branch target buffers.

The Intel Pentium M processor has enhanced dynamic branch prediction hardware. Branch target buffers (BTB) predict the direction and target of branches based on an instruction's address.

The Pentium M Processor includes two techniques to reduce the execution time of certain operations:

- **ESP folding** — This eliminates the ESP manipulation μops in stack-related instructions such as PUSH, POP, CALL and RET. It increases decode rename and retirement throughput. ESP folding also increases execution bandwidth by eliminating μops which would have required execution resources.
- **Micro-ops (μops) fusion** — Some of the most frequent pairs of μops derived from the same instruction can be fused into a single μops. The following categories of fused μops have been implemented in the Pentium M processor:
 - “Store address” and “store data” μops are fused into a single “Store” μop. This holds for all types of store operations, including integer, floating-point, MMX technology, and Streaming SIMD Extensions (SSE and SSE2) operations.
 - A load μop in most cases can be fused with a successive execution μop. This holds for integer, floating-point and MMX technology loads and for most kinds of successive execution operations. Note that SSE Loads can not be fused.

2.5.2 Data Prefetching

The Intel Pentium M processor supports three prefetching mechanisms:

- The first mechanism is a hardware instruction fetcher and is described in the previous section.
- The second mechanism automatically fetches data into the second-level cache. The implementation of automatic hardware prefetching in Pentium M processor family is basically similar to those described for NetBurst microarchitecture. The trigger threshold distance for each relevant processor models is shown in Table 2-6. The third mechanism is a software mechanism that fetches data into the caches using the prefetch instructions.

Table 2-23. Trigger Threshold and CPUID Signatures for Processor Families

Trigger Threshold Distance (Bytes)	Extended Model ID	Extended Family ID	Family ID	Model ID
512	0	0	15	3, 4, 6
256	0	0	15	0, 1, 2
256	0	0	6	9, 13, 14

Data is fetched 64 bytes at a time; the instruction and data translation lookaside buffers support 128 entries. See Table 2-7 for processor cache parameters.

Table 2-24. Cache Parameters of Pentium M, Intel Core Solo, and Intel Core Duo Processors

Level	Capacity	Associativity (ways)	Line Size (bytes)	Access Latency (clocks)	Write Update Policy
First	32 KByte	8	64	3	Writeback
Instruction	32 KByte	8	N/A	N/A	N/A
Second (mode 9)	1 MByte	8	64	9	Writeback
Second (model 13)	2 MByte	8	64	10	Writeback
Second (model 14)	2 MByte	8	64	14	Writeback

2.5.3 Out-of-Order Core

The processor core dynamically executes μ ops independent of program order. The core is designed to facilitate parallel execution by employing many buffers, issue ports, and parallel execution units.

The out-of-order core buffers μ ops in a Reservation Station (RS) until their operands are ready and resources are available. Each cycle, the core may dispatch up to five μ ops through the issue ports.

2.5.4 In-Order Retirement

The retirement unit in the Pentium M processor buffers completed μ ops in the reorder buffer (ROB). The ROB updates the architectural state in order. Up to three μ ops may be retired per cycle.

2.6 MICROARCHITECTURE OF INTEL® CORE™ SOLO AND INTEL® CORE™ DUO PROCESSORS

Intel Core Solo and Intel Core Duo processors incorporate an microarchitecture that is similar to the Pentium M processor microarchitecture, but provides additional enhancements for performance and power efficiency. Enhancements include:

- **Intel Smart Cache** — This second level cache is shared between two cores in an Intel Core Duo processor to minimize bus traffic between two cores accessing a single-copy of cached data. It allows an Intel Core Solo processor (or when one of the two cores in an Intel Core Duo processor is idle) to access its full capacity.
- **Stream SIMD Extensions 3** — These extensions are supported in Intel Core Solo and Intel Core Duo processors.
- **Decoder improvement** — Improvement in decoder and μ op fusion allows the front end to see most instructions as single μ op instructions. This increases the throughput of the three decoders in the front end.
- **Improved execution core** — Throughput of SIMD instructions is improved and the out-of-order engine is more robust in handling sequences of frequently-used instructions. Enhanced internal buffering and prefetch mechanisms also improve data bandwidth for execution.
- **Power-optimized bus** — The system bus is optimized for power efficiency; increased bus speed supports 667 MHz.
- **Data Prefetch** — Intel Core Solo and Intel Core Duo processors implement improved hardware prefetch mechanisms: one mechanism can look ahead and prefetch data into L1 from L2. These processors also provide enhanced hardware prefetchers similar to those of the Pentium M processor (see Table 2-6).

2.6.1 Front End

Execution of SIMD instructions on Intel Core Solo and Intel Core Duo processors are improved over Pentium M processors by the following enhancements:

- **Micro-op fusion** — Scalar SIMD operations on register and memory have single μ op flows comparable to X87 flows. Many packed instructions are fused to reduce its μ op flow from four to two μ ops.
- **Eliminating decoder restrictions** — Intel Core Solo and Intel Core Duo processors improve decoder throughput with micro-fusion and macro-fusion, so that many more SSE and SSE2 instructions can be decoded without restriction. On Pentium M processors, many single μ op SSE and SSE2 instructions must be decoded by the main decoder.
- **Improved packed SIMD instruction decoding** — On Intel Core Solo and Intel Core Duo processors, decoding of most packed SSE instructions is done by all three decoders. As a result the front end can process up to three packed SSE instructions every cycle. There are some exceptions to the above; some shuffle/unpack/shift operations are not fused and require the main decoder.

2.6.2 Data Prefetching

Intel Core Solo and Intel Core Duo processors provide hardware mechanisms to prefetch data from memory to the second-level cache. There are two techniques:

1. One mechanism activates after the data access pattern experiences two cache-reference misses within a trigger-distance threshold (see Table 2-6). This mechanism is similar to that of the Pentium M processor, but can track 16 forward data streams and 4 backward streams.
2. The second mechanism fetches an adjacent cache line of data after experiencing a cache miss. This effectively simulates the prefetching capabilities of 128-byte sectors (similar to the sectoring of two adjacent 64-byte cache lines available in Pentium 4 processors).

Hardware prefetch requests are queued up in the bus system at lower priority than normal cache-miss requests. If bus queue is in high demand, hardware prefetch requests may be ignored or cancelled to service bus traffic required by demand cache-misses and other bus transactions. Hardware prefetch mechanisms are enhanced over that of Pentium M processor by:

- Data stores that are not in the second-level cache generate read for ownership requests. These requests are treated as loads and can trigger a prefetch stream.
- Software prefetch instructions are treated as loads, they can also trigger a prefetch stream.

2.7 INTEL® HYPER-THREADING TECHNOLOGY

Intel® Hyper-Threading Technology (HT Technology) is supported by specific members of the Intel Pentium 4 and Xeon processor families. The technology enables software to take advantage of task-level, or thread-level parallelism by providing multiple logical processors within a physical processor package. In its first implementation in Intel Xeon processor, Hyper-Threading Technology makes a single physical processor appear as two logical processors.

The two logical processors each have a complete set of architectural registers while sharing one single physical processor's resources. By maintaining the architecture state of two processors, an HT Technology capable processor looks like two processors to software, including operating system and application code.

By sharing resources needed for peak demands between two logical processors, HT Technology is well suited for multiprocessor systems to provide an additional performance boost in throughput when compared to traditional MP systems.

Figure 2-8 shows a typical bus-based symmetric multiprocessor (SMP) based on processors supporting HT Technology. Each logical processor can execute a software thread, allowing a maximum of two software threads to execute simultaneously on one physical processor. The two software threads execute simultaneously, meaning that in the same clock cycle an "add" operation from logical processor 0 and another

“add” operation and load from logical processor 1 can be executed simultaneously by the execution engine.

In the first implementation of HT Technology, the physical execution resources are shared and the architecture state is duplicated for each logical processor. This minimizes the die area cost of implementing HT Technology while still achieving performance gains for multithreaded applications or multitasking workloads.

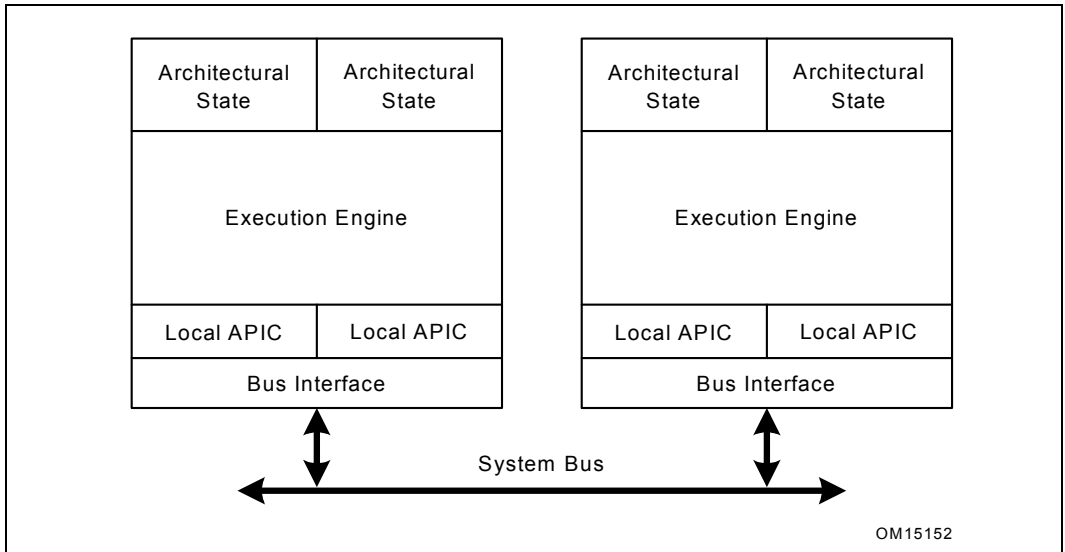


Figure 2-13. Hyper-Threading Technology on an SMP

The performance potential due to HT Technology is due to:

- The fact that operating systems and user programs can schedule processes or threads to execute simultaneously on the logical processors in each physical processor
- The ability to use on-chip execution resources at a higher level than when only a single thread is consuming the execution resources; higher level of resource utilization can lead to higher system throughput

2.7.1 Processor Resources and HT Technology

The majority of microarchitecture resources in a physical processor are shared between the logical processors. Only a few small data structures were replicated for each logical processor. This section describes how resources are shared, partitioned or replicated.

2.7.1.1 Replicated Resources

The architectural state is replicated for each logical processor. The architecture state consists of registers that are used by the operating system and application code to control program behavior and store data for computations. This state includes the eight general-purpose registers, the control registers, machine state registers, debug registers, and others. There are a few exceptions, most notably the memory type range registers (MTRRs) and the performance monitoring resources. For a complete list of the architecture state and exceptions, see the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volumes 3A & 3B*.

Other resources such as instruction pointers and register renaming tables were replicated to simultaneously track execution and state changes of the two logical processors. The return stack predictor is replicated to improve branch prediction of return instructions.

In addition, a few buffers (for example, the 2-entry instruction streaming buffers) were replicated to reduce complexity.

2.7.1.2 Partitioned Resources

Several buffers are shared by limiting the use of each logical processor to half the entries. These are referred to as partitioned resources. Reasons for this partitioning include:

- Operational fairness
- Permitting the ability to allow operations from one logical processor to bypass operations of the other logical processor that may have stalled

For example: a cache miss, a branch misprediction, or instruction dependencies may prevent a logical processor from making forward progress for some number of cycles. The partitioning prevents the stalled logical processor from blocking forward progress.

In general, the buffers for staging instructions between major pipe stages are partitioned. These buffers include μ op queues after the execution trace cache, the queues after the register rename stage, the reorder buffer which stages instructions for retirement, and the load and store buffers.

In the case of load and store buffers, partitioning also provided an easier implementation to maintain memory ordering for each logical processor and detect memory ordering violations.

2.7.1.3 Shared Resources

Most resources in a physical processor are fully shared to improve the dynamic utilization of the resource, including caches and all the execution units. Some shared resources which are linearly addressed, like the DTLB, include a logical processor ID bit to distinguish whether the entry belongs to one logical processor or the other.

The first level cache can operate in two modes depending on a context-ID bit:

- Shared mode: The L1 data cache is fully shared by two logical processors.
- Adaptive mode: In adaptive mode, memory accesses using the page directory is mapped identically across logical processors sharing the L1 data cache.

The other resources are fully shared.

2.7.2 Microarchitecture Pipeline and HT Technology

This section describes the HT Technology microarchitecture and how instructions from the two logical processors are handled between the front end and the back end of the pipeline.

Although instructions originating from two programs or two threads execute simultaneously and not necessarily in program order in the execution core and memory hierarchy, the front end and back end contain several selection points to select between instructions from the two logical processors. All selection points alternate between the two logical processors unless one logical processor cannot make use of a pipeline stage. In this case, the other logical processor has full use of every cycle of the pipeline stage. Reasons why a logical processor may not use a pipeline stage include cache misses, branch mispredictions, and instruction dependencies.

2.7.3 Front End Pipeline

The execution trace cache is shared between two logical processors. Execution trace cache access is arbitrated by the two logical processors every clock. If a cache line is fetched for one logical processor in one clock cycle, the next clock cycle a line would be fetched for the other logical processor provided that both logical processors are requesting access to the trace cache.

If one logical processor is stalled or is unable to use the execution trace cache, the other logical processor can use the full bandwidth of the trace cache until the initial logical processor's instruction fetches return from the L2 cache.

After fetching the instructions and building traces of μ ops, the μ ops are placed in a queue. This queue decouples the execution trace cache from the register rename pipeline stage. As described earlier, if both logical processors are active, the queue is partitioned so that both logical processors can make independent forward progress.

2.7.4 Execution Core

The core can dispatch up to six μ ops per cycle, provided the μ ops are ready to execute. Once the μ ops are placed in the queues waiting for execution, there is no distinction between instructions from the two logical processors. The execution core and memory hierarchy is also oblivious to which instructions belong to which logical processor.

After execution, instructions are placed in the re-order buffer. The re-order buffer decouples the execution stage from the retirement stage. The re-order buffer is partitioned such that each uses half the entries.

2.7.5 Retirement

The retirement logic tracks when instructions from the two logical processors are ready to be retired. It retires the instruction in program order for each logical processor by alternating between the two logical processors. If one logical processor is not ready to retire any instructions, then all retirement bandwidth is dedicated to the other logical processor.

Once stores have retired, the processor needs to write the store data into the level-one data cache. Selection logic alternates between the two logical processors to commit store data to the cache.

2.8 MULTICORE PROCESSORS

The Intel Pentium D processor and the Pentium Processor Extreme Edition introduce multicore features. These processors enhance hardware support for multithreading by providing two processor cores in each physical processor package. The Dual-core Intel Xeon and Intel Core Duo processors also provide two processor cores in a physical package. The multicore topology of Intel Core 2 Duo processors are similar to those of Intel Core Duo processor.

The Intel Pentium D processor provides two logical processors in a physical package, each logical processor has a separate execution core and a cache hierarchy. The Dual-core Intel Xeon processor and the Intel Pentium Processor Extreme Edition provide four logical processors in a physical package that has two execution cores. Each core provides two logical processors sharing an execution core and a cache hierarchy.

The Intel Core Duo processor provides two logical processors in a physical package. Each logical processor has a separate execution core (including first-level cache) and a smart second-level cache. The second-level cache is shared between two logical processors and optimized to reduce bus traffic when the same copy of cached data is used by two logical processors. The full capacity of the second-level cache can be used by one logical processor if the other logical processor is inactive.

The functional blocks of the dual-core processors are shown in Figure 2-9. The Quad-core Intel Xeon processors, Intel Core 2 Quad processor and Intel Core 2 Extreme quad-core processor consist of two replica of the dual-core modules. The functional blocks of the quad-core processors are also shown in Figure 2-9.

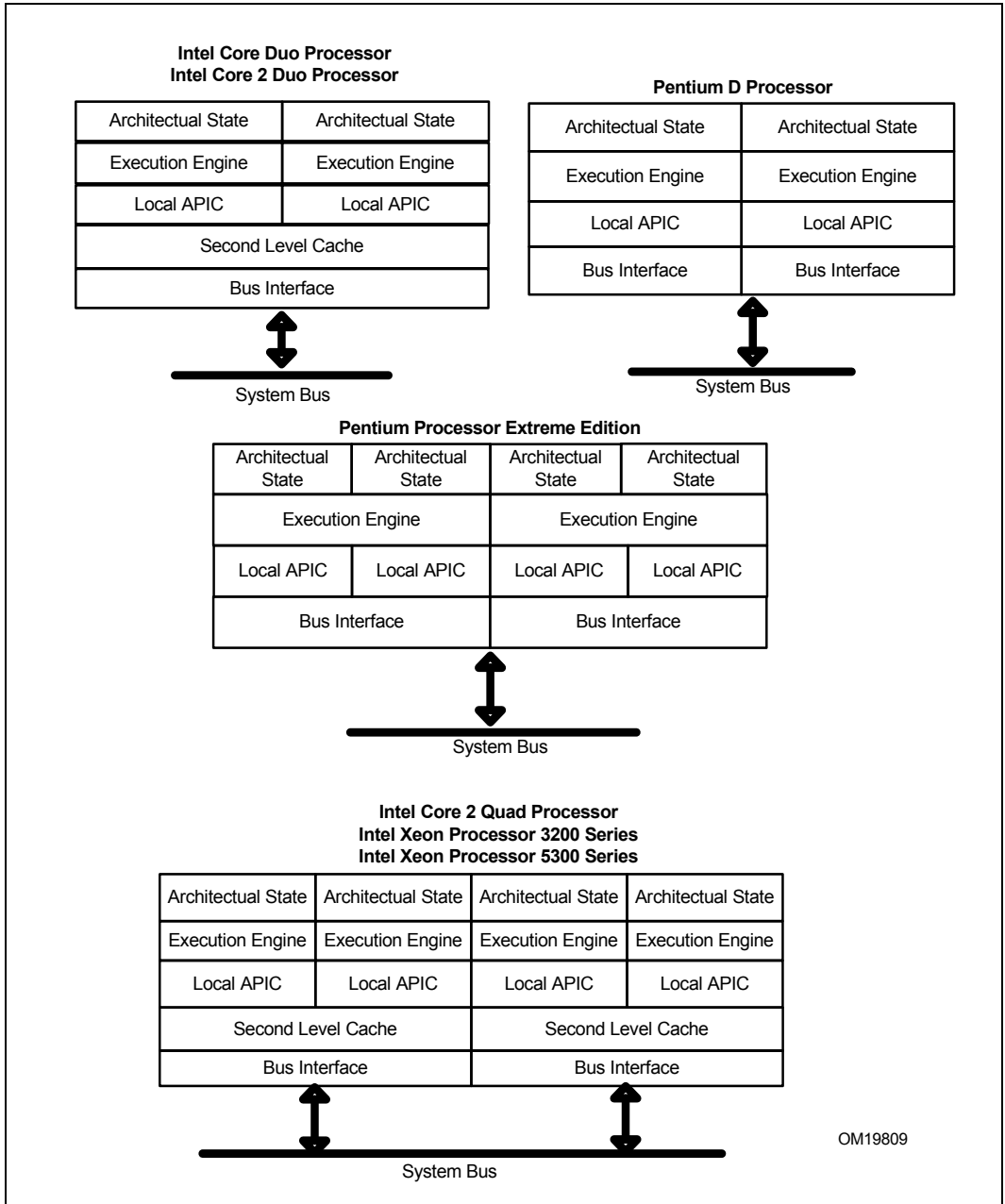


Figure 2-14. Pentium D Processor, Pentium Processor Extreme Edition, Intel Core Duo Processor, Intel Core 2 Duo Processor, and Intel Core 2 Quad Processor

2.8.1 Microarchitecture Pipeline and MultiCore Processors

In general, each core in a multicore processor resembles a single-core processor implementation of the underlying microarchitecture. The implementation of the cache hierarchy in a dual-core or multicore processor may be the same or different from the cache hierarchy implementation in a single-core processor.

CPUID should be used to determine cache-sharing topology information in a processor implementation and the underlying microarchitecture. The former is obtained by querying the deterministic cache parameter leaf (see Chapter 7, “Optimizing Cache Usage”); the latter by using the encoded values for extended family, family, extended model, and model fields. See Table 2-8.

Table 2-25. Family And Model Designations of Microarchitectures

Dual-Core Processor	Micro-architecture	Extended Family	Family	Extended Model	Model
Pentium D processor	NetBurst	0	15	0	3, 4, 6
Pentium processor Extreme Edition	NetBurst	0	15	0	3, 4, 6
Intel Core Duo processor	Improved Pentium M	0	6	0	14
Intel Core 2 Duo processor/ Intel Xeon processor 5100	Intel Core Microarchitecture	0	6	0	15
Intel Core 2 Duo processor E8000 Series/ Intel Xeon processor 5200, 5400	Enhanced Intel Core Microarchitecture	0	6	1	7

2.8.2 Shared Cache in Intel® Core™ Duo Processors

The Intel Core Duo processor has two symmetric cores that share the second-level cache and a single bus interface (see Figure 2-9). Two threads executing on two cores in an Intel Core Duo processor can take advantage of shared second-level cache, accessing a single-copy of cached data without generating bus traffic.

2.8.2.1 Load and Store Operations

When an instruction needs to read data from a memory address, the processor looks for it in caches and memory. When an instruction writes data to a memory location (write back) the processor first makes sure that the cache line that contains the memory location is owned by the first-level data cache of the initiating core (that is, the line is in exclusive or modified state). Then the processor looks for the cache line in the cache and memory sub-systems. The look-ups for the locality of load or store operation are in the following order:

1. DCU of the initiating core
2. DCU of the other core and second-level cache
3. System memory

The cache line is taken from the DCU of the other core only if it is modified, ignoring the cache line availability or state in the L2 cache. Table 2-9 lists the performance characteristics of generic load and store operations in an Intel Core Duo processor. Numeric values of Table 2-9 are in terms of processor core cycles.

Table 2-26. Characteristics of Load and Store Operations in Intel Core Duo Processors

Data Locality	Load		Store	
	Latency	Throughput	Latency	Throughput
DCU	3	1	2	1
DCU of the other core in "Modified" state	14 + bus transaction	14 + bus transaction	14 + bus transaction	~10
2nd-level cache	14	<6	14	<6
Memory	14 + bus transaction	Bus read protocol	14 + bus transaction	Bus write protocol

Throughput is expressed as the number of cycles to wait before the same operation can start again. The latency of a bus transaction is exposed in some of these operations, as indicated by entries containing "+ bus transaction". On Intel Core Duo processors, a typical bus transaction may take 5.5 bus cycles. For a 667 MHz bus and a core frequency of 2.167GHz, the total of $14 + 5.5 * 2167 / (667/4) \sim 86$ core cycles.

Sometimes a modified cache line has to be evicted to make room for a new cache line. The modified cache line is evicted in parallel to bringing in new data and does not require additional latency. However, when data is written back to memory, the eviction consumes cache bandwidth and bus bandwidth. For multiple cache misses that require the eviction of modified lines and are within a short time, there is an overall degradation in response time of these cache misses.

For store operation, reading for ownership must be completed before the data is written to the first-level data cache and the line is marked as modified. Reading for ownership and storing the data happens after instruction retirement and follows the order of retirement. The bus store latency does not affect the store instruction itself. However, several sequential stores may have cumulative latency that can effect performance.

2.9 INTEL® 64 ARCHITECTURE

Intel 64 architecture supports almost all features in the IA-32 Intel architecture and extends support to run 64-bit OS and 64-bit applications in 64-bit linear address space. Intel 64 architecture provides a new operating mode, referred to as IA-32e mode, and increases the linear address space for software to 64 bits and supports physical address space up to 40 bits.

IA-32e mode consists of two sub-modes: (1) compatibility mode enables a 64-bit operating system to run most legacy 32-bit software unmodified, (2) 64-bit mode enables a 64-bit operating system to run applications written to access 64-bit linear address space.

In the 64-bit mode of Intel 64 architecture, software may access:

- 64-bit flat linear addressing
- 8 additional general-purpose registers (GPRs)
- 8 additional registers for streaming SIMD extensions (SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2 and AESNI)
- 64-bit-wide GPRs and instruction pointers
- uniform byte-register addressing
- fast interrupt-prioritization mechanism
- a new instruction-pointer relative-addressing mode

For optimizing 64-bit applications, the features that impact software optimizations include:

- using a set of prefixes to access new registers or 64-bit register operand
- pointer size increases from 32 bits to 64 bits
- instruction-specific usages

2.10 SIMD TECHNOLOGY

SIMD computations (see Figure 2-10) were introduced to the architecture with MMX technology. MMX technology allows SIMD computations to be performed on packed byte, word, and doubleword integers. The integers are contained in a set of eight 64-bit registers called MMX registers (see Figure 2-11).

The Pentium III processor extended the SIMD computation model with the introduction of the Streaming SIMD Extensions (SSE). SSE allows SIMD computations to be performed on operands that contain four packed single-precision floating-point data elements. The operands can be in memory or in a set of eight 128-bit XMM registers (see Figure 2-11). SSE also extended SIMD computational capability by adding additional 64-bit MMX instructions.

Figure 2-10 shows a typical SIMD computation. Two sets of four packed data elements (X1, X2, X3, and X4, and Y1, Y2, Y3, and Y4) are operated on in parallel, with the same operation being performed on each corresponding pair of data elements (X1 and Y1, X2 and Y2, X3 and Y3, and X4 and Y4). The results of the four parallel computations are sorted as a set of four packed data elements.

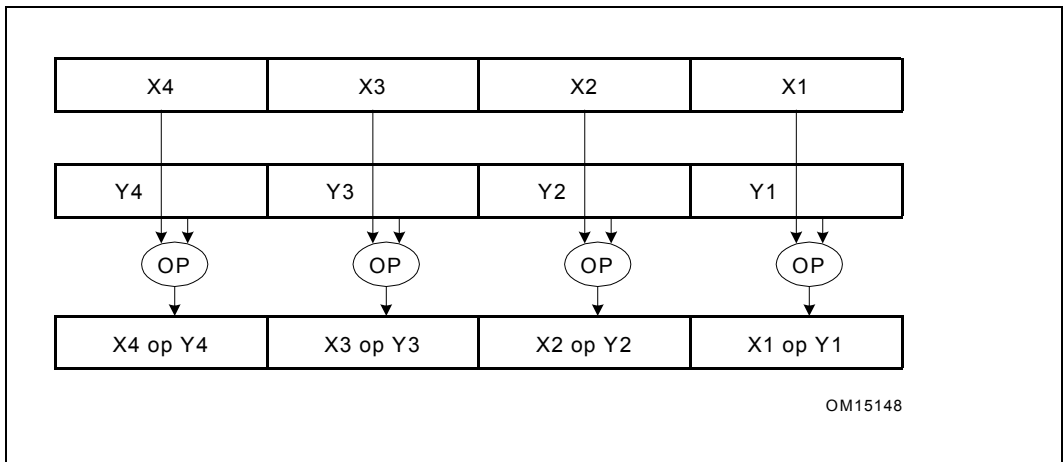


Figure 2-15. Typical SIMD Operations

The Pentium 4 processor further extended the SIMD computation model with the introduction of Streaming SIMD Extensions 2 (SSE2), Streaming SIMD Extensions 3 (SSE3), and Intel Xeon processor 5100 series introduced Supplemental Streaming SIMD Extensions 3 (SSSE3).

SSE2 works with operands in either memory or in the XMM registers. The technology extends SIMD computations to process packed double-precision floating-point data elements and 128-bit packed integers. There are 144 instructions in SSE2 that operate on two packed double-precision floating-point data elements or on 16 packed byte, 8 packed word, 4 doubleword, and 2 quadword integers.

SSE3 enhances x87, SSE and SSE2 by providing 13 instructions that can accelerate application performance in specific areas. These include video processing, complex arithmetics, and thread synchronization. SSE3 complements SSE and SSE2 with instructions that process SIMD data asymmetrically, facilitate horizontal computation, and help avoid loading cache line splits. See Figure 2-11.

SSSE3 provides additional enhancement for SIMD computation with 32 instructions on digital video and signal processing.

SSE4.1, SSE4.2 and AESNI are additional SIMD extensions that provide acceleration for applications in media processing, text/lexical processing, and block encryption/decryption.

The SIMD extensions operates the same way in Intel 64 architecture as in IA-32 architecture, with the following enhancements:

- 128-bit SIMD instructions referencing XMM register can access 16 XMM registers in 64-bit mode.
- Instructions that reference 32-bit general purpose registers can access 16 general purpose registers in 64-bit mode.

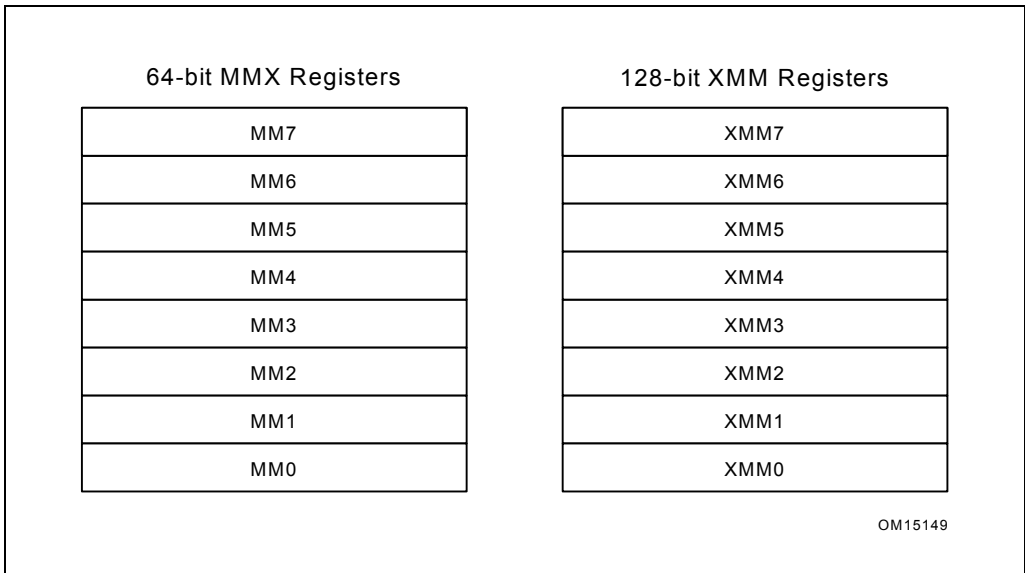


Figure 2-16. SIMD Instruction Register Usage

SIMD improves the performance of 3D graphics, speech recognition, image processing, scientific applications and applications that have the following characteristics:

- inherently parallel
- recurring memory access patterns
- localized recurring operations performed on the data
- data-independent control flow

SIMD floating-point instructions fully support the IEEE Standard 754 for Binary Floating-Point Arithmetic. They are accessible from all IA-32 execution modes: protected mode, real address mode, and Virtual 8086 mode.

SSE, SSE2, and MMX technologies are architectural extensions. Existing software will continue to run correctly, without modification on Intel microprocessors that incorporate these technologies. Existing software will also run correctly in the presence of applications that incorporate SIMD technologies.

SSE and SSE2 instructions also introduced cacheability and memory ordering instructions that can improve cache usage and application performance.

For more on SSE, SSE2, SSE3 and MMX technologies, see the following chapters in the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1*:

- Chapter 9, "Programming with Intel® MMX™ Technology"
- Chapter 10, "Programming with Streaming SIMD Extensions (SSE)"
- Chapter 11, "Programming with Streaming SIMD Extensions 2 (SSE2)"
- Chapter 12, "Programming with SSE3 and Supplemental SSE3"

2.10.1 Summary of SIMD Technologies

2.10.1.1 MMX™ Technology

MMX Technology introduced:

- 64-bit MMX registers
- Support for SIMD operations on packed byte, word, and doubleword integers

MMX instructions are useful for multimedia and communications software.

2.10.1.2 Streaming SIMD Extensions

Streaming SIMD extensions introduced:

- 128-bit XMM registers
- 128-bit data type with four packed single-precision floating-point operands
- data prefetch instructions
- non-temporal store instructions and other cacheability and memory ordering instructions
- extra 64-bit SIMD integer support

SSE instructions are useful for 3D geometry, 3D rendering, speech recognition, and video encoding and decoding.

2.10.1.3 Streaming SIMD Extensions 2

Streaming SIMD extensions 2 add the following:

- 128-bit data type with two packed double-precision floating-point operands
- 128-bit data types for SIMD integer operation on 16-byte, 8-word, 4-doubleword, or 2-quadword integers
- support for SIMD arithmetic on 64-bit integer operands
- instructions for converting between new and existing data types
- extended support for data shuffling
- Extended support for cacheability and memory ordering operations

SSE2 instructions are useful for 3D graphics, video decoding/encoding, and encryption.

2.10.1.4 Streaming SIMD Extensions 3

Streaming SIMD extensions 3 add the following:

- SIMD floating-point instructions for asymmetric and horizontal computation
- a special-purpose 128-bit load instruction to avoid cache line splits
- an x87 FPU instruction to convert to integer independent of the floating-point control word (FCW)
- instructions to support thread synchronization

SSE3 instructions are useful for scientific, video and multi-threaded applications.

2.10.1.5 Supplemental Streaming SIMD Extensions 3

The Supplemental Streaming SIMD Extensions 3 introduces 32 new instructions to accelerate eight types of computations on packed integers. These include:

- 12 instructions that perform horizontal addition or subtraction operations
- 6 instructions that evaluate the absolute values
- 2 instructions that perform multiply and add operations and speed up the evaluation of dot products
- 2 instructions that accelerate packed-integer multiply operations and produce integer values with scaling
- 2 instructions that perform a byte-wise, in-place shuffle according to the second shuffle control operand
- 6 instructions that negate packed integers in the destination operand if the signs of the corresponding element in the source operand is less than zero
- 2 instructions that align data from the composite of two operands

2.10.1.6 SSE4.1

SSE4.1 introduces 47 new instructions to accelerate video, imaging and 3D applications. SSE4.1 also improves compiler vectorization and significantly increase support for packed dword computation. These include:

- Two instructions perform packed dword multiplies.
- Two instructions perform floating-point dot products with input/output selects.
- One instruction provides a streaming hint for WC loads.
- Six instructions simplify packed blending.
- Eight instructions expand support for packed integer MIN/MAX.
- Four instructions support floating-point round with selectable rounding mode and precision exception override.
- Seven instructions improve data insertion and extractions from XMM registers
- Twelve instructions improve packed integer format conversions (sign and zero extensions).
- One instruction improves SAD (sum absolute difference) generation for small block sizes.
- One instruction aids horizontal searching operations of word integers.
- One instruction improves masked comparisons.
- One instruction adds qword packed equality comparisons.
- One instruction adds dword packing with unsigned saturation.

2.10.1.7 SSE4.2

SSE4.2 introduces 7 new instructions. These include:

- A 128-bit SIMD integer instruction for comparing 64-bit integer data elements.
- Four string/text processing instructions providing a rich set of primitives, these primitives can accelerate:
 - basic and advanced string library functions from `strlen`, `strcmp`, to `strcspn`,
 - delimiter processing, token extraction for lexing of text streams,
 - Parser, schema validation including XML processing.
- A general-purpose instruction for accelerating cyclic redundancy checksum signature calculations.
- A general-purpose instruction for calculating bit count population of integer numbers.

2.10.1.8 AESNI and PCLMULQDQ

AESNI introduces 7 new instructions, six of them are primitives for accelerating algorithms based on AES encryption/decryption standard, referred to as AESNI.

The PCLMULQDQ instruction accelerates general-purpose block encryption, which can perform carry-less multiplication for two binary numbers up to 64-bit wide.

Typically, algorithm based on AES standard involve transformation of block data over multiple iterations via several primitives. The AES standard supports cipher key of sizes 128, 192, and 256 bits. The respective cipher key sizes correspond to 10, 12, and 14 rounds of iteration.

AES encryption involves processing 128-bit input data (plaintext) through a finite number of iterative operation, referred to as "AES round", into a 128-bit encrypted block (ciphertext). Decryption follows the reverse direction of iterative operation using the "equivalent inverse cipher" instead of the "inverse cipher".

The cryptographic processing at each round involves two input data, one is the "state", the other is the "round key". Each round uses a different "round key". The round keys are derived from the cipher key using a "key schedule" algorithm. The "key schedule" algorithm is independent of the data processing of encryption/decryption, and can be carried out independently from the encryption/decryption phase.

The AES extensions provide two primitives to accelerate AES rounds on encryption, two primitives for AES rounds on decryption using the equivalent inverse cipher, and two instructions to support the AES key expansion procedure.

2.10.1.9 Intel® Advanced Vector Extensions

Intel® Advanced Vector Extensions offers comprehensive architectural enhancements over previous generations of Streaming SIMD Extensions. Intel AVX introduces the following architectural enhancements:

- Support for 256-bit wide vectors and SIMD register set.
- 256-bit floating-point instruction set enhancement with up to 2X performance gain relative to 128-bit Streaming SIMD extensions.
- Instruction syntax support for generalized three-operand syntax to improve instruction programming flexibility and efficient encoding of new instruction extensions.
- Enhancement of legacy 128-bit SIMD instruction extensions to support three-operand syntax and to simplify compiler vectorization of high-level language expressions.
- Support flexible deployment of 256-bit AVX code, 128-bit AVX code, legacy 128-bit code and scalar code.

Intel AVX instruction set and 256-bit register state management detail are described in IA-32 Intel® Architecture Software Developer's Manual, Volumes 2A, 2B and 3A.

Optimization techniques for Intel AVX is discussed in Chapter 11, “Optimization for Intel AVX”.

CHAPTER 3

GENERAL OPTIMIZATION GUIDELINES

This chapter discusses general optimization techniques that can improve the performance of applications running on Intel Core i7 processors, processors based on Intel microarchitecture code name Sandy Bridge, Intel microarchitecture code name Nehalem, Intel Core microarchitecture, Enhanced Intel Core microarchitecture, Intel NetBurst microarchitecture, Intel Core Duo, Intel Core Solo, and Pentium M processors. These techniques take advantage of microarchitectural described in Chapter 2, “Intel® 64 and IA-32 Processor Architectures.” Optimization guidelines focusing on Intel multi-core processors, Hyper-Threading Technology and 64-bit mode applications are discussed in Chapter 8, “Multicore and Hyper-Threading Technology,” and Chapter 9, “64-bit Mode Coding Guidelines.”

Practices that optimize performance focus on three areas:

- tools and techniques for code generation
- analysis of the performance characteristics of the workload and its interaction with microarchitectural sub-systems
- tuning code to the target microarchitecture (or families of microarchitecture) to improve performance

Some hints on using tools are summarized first to simplify the first two tasks. the rest of the chapter will focus on recommendations of code generation or code tuning to the target microarchitectures.

This chapter explains optimization techniques for the Intel C++ Compiler, the Intel Fortran Compiler, and other compilers.

3.1 PERFORMANCE TOOLS

Intel offers several tools to help optimize application performance, including compilers, performance analyzer and multithreading tools.

3.1.1 Intel® C++ and Fortran Compilers

Intel compilers support multiple operating systems (Windows*, Linux*, Mac OS* and embedded). The Intel compilers optimize performance and give application developers access to advanced features:

- Flexibility to target 32-bit or 64-bit Intel processors for optimization
- Compatibility with many integrated development environments or third-party compilers.

GENERAL OPTIMIZATION GUIDELINES

- Automatic optimization features to take advantage of the target processor's architecture.
- Automatic compiler optimization reduces the need to write different code for different processors.
- Common compiler features that are supported across Windows, Linux and Mac OS include:
 - General optimization settings
 - Cache-management features
 - Interprocedural optimization (IPO) methods
 - Profile-guided optimization (PGO) methods
 - Multithreading support
 - Floating-point arithmetic precision and consistency support
 - Compiler optimization and vectorization reports

3.1.2 General Compiler Recommendations

Generally speaking, a compiler that has been tuned for the target microarchitecture can be expected to match or outperform hand-coding. However, if performance problems are noted with the compiled code, some compilers (like Intel C++ and Fortran Compilers) allow the coder to insert intrinsics or inline assembly in order to exert control over what code is generated. If inline assembly is used, the user must verify that the code generated is of good quality and yields good performance.

Default compiler switches are targeted for common cases. An optimization may be made to the compiler default if it is beneficial for most programs. If the root cause of a performance problem is a poor choice on the part of the compiler, using different switches or compiling the targeted module with a different compiler may be the solution.

3.1.3 VTune™ Performance Analyzer

VTune uses performance monitoring hardware to collect statistics and coding information of your application and its interaction with the microarchitecture. This allows software engineers to measure performance characteristics of the workload for a given microarchitecture. VTune supports Intel Core i7 processors, Intel Core microarchitecture, Intel NetBurst microarchitecture, Intel Core Duo, Intel Core Solo, and Pentium M processor families.

The VTune Performance Analyzer provides two kinds of feedback:

- indication of a performance improvement gained by using a specific coding recommendation or microarchitectural feature

- information on whether a change in the program has improved or degraded performance with respect to a particular metric

The VTune Performance Analyzer also provides measures for a number of workload characteristics, including:

- retirement throughput of instruction execution as an indication of the degree of extractable instruction-level parallelism in the workload
- data traffic locality as an indication of the stress point of the cache and memory hierarchy
- data traffic parallelism as an indication of the degree of effectiveness of amortization of data access latency

NOTE

Improving performance in one part of the machine does not necessarily bring significant gains to overall performance. It is possible to degrade overall performance by improving performance for some particular metric.

Where appropriate, coding recommendations in this chapter include descriptions of the VTune Performance Analyzer events that provide measurable data on the performance gain achieved by following the recommendations. For more on using the VTune analyzer, refer to the application's online help.

3.2 PROCESSOR PERSPECTIVES

Many coding recommendations for Intel Core microarchitecture work well across Intel Core i7, Pentium M, Intel Core Solo, Intel Core Duo processors and processors based on Intel NetBurst microarchitecture. However, there are situations where a recommendation may benefit one microarchitecture more than another. Some of these are:

- Instruction decode throughput is important for processors based on Intel Core i7 processors, Intel Core microarchitecture (Pentium M, Intel Core Solo, and Intel Core Duo processors). Taking advantage of decoded ICache, Loop Stream Detector and macrofusion can further improve front end performance.
- Generating code with a 4-1-1 template (instruction with four μ ops followed by two instructions with one μ op each) helps the Pentium M processor.

Intel Core Solo and Intel Core Duo processors have an enhanced front end that is less sensitive to the 4-1-1 template. Processors based on Intel Core microarchitecture have 4 decoders and employ micro-fusion and macro-fusion so that each of three simple decoders are not restricted to handling simple instructions consisting of one μ op.

Taking advantage of micro-fusion will increase decoder throughput across Intel Core Solo, Intel Core Duo and Intel Core2 Duo processors. Taking advantage of macro-fusion can improve decoder throughput further on Intel Core 2 Duo processor family. Taking advantage of macro-fusion can improve decoder throughput in both 64-bit and 32-bit code for Intel microarchitecture code name Nehalem.

- On processors based on Intel microarchitecture code name Sandy Bridge, the code size for optimal front end performance is associated with the decode ICache.
- Dependencies for partial register writes incur large penalties when using the Pentium M processor (this applies to processors with CPUID signature family 6, model 9). On Pentium 4, Intel Xeon processors, Pentium M processor (with CPUID signature family 6, model 13), such penalties are relieved by artificial dependencies between each partial register write. Intel Core Solo, Intel Core Duo processors and processors based on Intel Core microarchitecture can experience minor delays due to partial register stalls. To avoid false dependences from partial register updates, use full register updates and extended moves.
- Use appropriate instructions that support dependence-breaking (PXOR, SUB, XOR instructions). Dependence-breaking support for XORPS is available in Intel Core Solo, Intel Core Duo processors and processors based on Intel Core microarchitecture.
- Hardware prefetching can reduce the effective memory latency for data and instruction accesses in general. But different microarchitectures may require some custom modifications to adapt to the specific hardware prefetch implementation of each microarchitecture.

3.2.1 CPUID Dispatch Strategy and Compatible Code Strategy

When optimum performance on all processor generations is desired, applications can take advantage of the CPUID instruction to identify the processor generation and integrate processor-specific instructions into the source code. The Intel C++ Compiler supports the integration of different versions of the code for different target processors. The selection of which code to execute at runtime is made based on the CPU identifiers. Binary code targeted for different processor generations can be generated under the control of the programmer or by the compiler.

For applications that target multiple generations of microarchitectures, and where minimum binary code size and single code path is important, a compatible code strategy is the best. Optimizing applications using techniques developed for the Intel Core microarchitecture and combined with some for Intel NetBurst microarchitecture are likely to improve code efficiency and scalability when running on processors based on current and future generations of Intel 64 and IA-32 processors. This compatible approach to optimization is also likely to deliver high performance on Pentium M, Intel Core Solo and Intel Core Duo processors.

3.2.2 Transparent Cache-Parameter Strategy

If the CPUID instruction supports function leaf 4, also known as deterministic cache parameter leaf, the leaf reports cache parameters for each level of the cache hierarchy in a deterministic and forward-compatible manner across Intel 64 and IA-32 processor families.

For coding techniques that rely on specific parameters of a cache level, using the deterministic cache parameter allows software to implement techniques in a way that is forward-compatible with future generations of Intel 64 and IA-32 processors, and cross-compatible with processors equipped with different cache sizes.

3.2.3 Threading Strategy and Hardware Multithreading Support

Intel 64 and IA-32 processor families offer hardware multithreading support in two forms: dual-core technology and HT Technology.

To fully harness the performance potential of hardware multithreading in current and future generations of Intel 64 and IA-32 processors, software must embrace a threaded approach in application design. At the same time, to address the widest range of installed machines, multi-threaded software should be able to run without failure on a single processor without hardware multithreading support and should achieve performance on a single logical processor that is comparable to an unthreaded implementation (if such comparison can be made). This generally requires architecting a multi-threaded application to minimize the overhead of thread synchronization. Additional guidelines on multithreading are discussed in Chapter 8, “Multicore and Hyper-Threading Technology.”

3.3 CODING RULES, SUGGESTIONS AND TUNING HINTS

This section includes rules, suggestions and hints. They are targeted for engineers who are:

- modifying source code to enhance performance (user/source rules)
- writing assemblers or compilers (assembly/compiler rules)
- doing detailed performance tuning (tuning suggestions)

Coding recommendations are ranked in importance using two measures:

- Local impact (high, medium, or low) refers to a recommendation’s affect on the performance of a given instance of code.
- Generality (high, medium, or low) measures how often such instances occur across all application domains. Generality may also be thought of as “frequency”.

These recommendations are approximate. They can vary depending on coding style, application domain, and other factors.

The purpose of the high, medium, and low (H, M, and L) priorities is to suggest the relative level of performance gain one can expect if a recommendation is implemented.

Because it is not possible to predict the frequency of a particular code instance in applications, priority hints cannot be directly correlated to application-level performance gain. In cases in which application-level performance gain has been observed, we have provided a quantitative characterization of the gain (for information only). In cases in which the impact has been deemed inapplicable, no priority is assigned.

3.4 OPTIMIZING THE FRONT END

Optimizing the front end covers two aspects:

- Maintaining steady supply of μ ops to the execution engine — Mispredicted branches can disrupt streams of μ ops, or cause the execution engine to waste execution resources on executing streams of μ ops in the non-architected code path. Much of the tuning in this respect focuses on working with the Branch Prediction Unit. Common techniques are covered in Section 3.4.1, “Branch Prediction Optimization.”
- Supplying streams of μ ops to utilize the execution bandwidth and retirement bandwidth as much as possible — For Intel Core microarchitecture and Intel Core Duo processor family, this aspect focuses maintaining high decode throughput. In Intel microarchitecture code name Sandy Bridge, this aspect focuses on keeping the *h*od code running from Decoded ICache. Techniques to maximize decode throughput for Intel Core microarchitecture are covered in Section 3.4.2, “Fetch and Decode Optimization.”

3.4.1 Branch Prediction Optimization

Branch optimizations have a significant impact on performance. By understanding the flow of branches and improving their predictability, you can increase the speed of code significantly.

Optimizations that help branch prediction are:

- Keep code and data on separate pages. This is very important; see Section 3.6, “Optimizing Memory Accesses,” for more information.
- Eliminate branches whenever possible.
- Arrange code to be consistent with the static branch prediction algorithm.
- Use the PAUSE instruction in spin-wait loops.
- Inline functions and pair up calls and returns.
- Unroll as necessary so that repeatedly-executed loops have sixteen or fewer iterations (unless this causes an excessive code size increase).

- Separate branches so that they occur no more frequently than every three μ ops where possible.

3.4.1.1 Eliminating Branches

Eliminating branches improves performance because:

- It reduces the possibility of mispredictions.
- It reduces the number of required branch target buffer (BTB) entries. Conditional branches, which are never taken, do not consume BTB resources.

There are four principal ways of eliminating branches:

- Arrange code to make basic blocks contiguous.
- Unroll loops, as discussed in Section 3.4.1.7, "Loop Unrolling."
- Use the CMOV instruction.
- Use the SETCC instruction.

The following rules apply to branch elimination:

Assembly/Compiler Coding Rule 1. (MH impact, M generality) *Arrange code to make basic blocks contiguous and eliminate unnecessary branches.*

For the Pentium M processor, every branch counts. Even correctly predicted branches have a negative effect on the amount of useful code delivered to the processor. Also, taken branches consume space in the branch prediction structures and extra branches create pressure on the capacity of the structures.

Assembly/Compiler Coding Rule 2. (M impact, ML generality) *Use the SETCC and CMOV instructions to eliminate unpredictable conditional branches where possible. Do not do this for predictable branches. Do not use these instructions to eliminate all unpredictable conditional branches (because using these instructions will incur execution overhead due to the requirement for executing both paths of a conditional branch). In addition, converting a conditional branch to SETCC or CMOV trades off control flow dependence for data dependence and restricts the capability of the out-of-order engine. When tuning, note that all Intel 64 and IA-32 processors usually have very high branch prediction rates. Consistently mispredicted branches are generally rare. Use these instructions only if the increase in computation time is less than the expected cost of a mispredicted branch.*

Consider a line of C code that has a condition dependent upon one of the constants:

```
X = (A < B) ? CONST1 : CONST2;
```

This code conditionally compares two values, A and B. If the condition is true, X is set to CONST1; otherwise it is set to CONST2. An assembly code sequence equivalent to the above C code can contain branches that are not predictable if there are no correlation in the two values.

GENERAL OPTIMIZATION GUIDELINES

Example 3-1 shows the assembly code with unpredictable branches. The unpredictable branches can be removed with the use of the SETCC instruction. Example 3-2 shows optimized code that has no branches.

Example 3-1. Assembly Code with an Unpredictable Branch

```
    cmp a, b           ; Condition
    jbe L30           ; Conditional branch
    mov ebx, const1   ; ebx holds X
    jmp L31           ; Unconditional branch
L30:
    mov ebx, const2
L31:
```

Example 3-2. Code Optimization to Eliminate Branches

```
xor  ebx, ebx        ; Clear ebx (X in the C code)
cmp  A, B
setge bl             ; When ebx = 0 or 1
                        ; OR the complement condition
sub  ebx, 1          ; ebx=11...11 or 00...00
and  ebx, CONST3; CONST3 = CONST1-CONST2
add  ebx, CONST2; ebx=CONST1 or CONST2
```

The optimized code in Example 3-2 sets EBX to zero, then compares A and B. If A is greater than or equal to B, EBX is set to one. Then EBX is decreased and AND'd with the difference of the constant values. This sets EBX to either zero or the difference of the values. By adding CONST2 back to EBX, the correct value is written to EBX. When CONST2 is equal to zero, the last instruction can be deleted.

Another way to remove branches on Pentium II and subsequent processors is to use the CMOV and FCMOV instructions. Example 3-3 shows how to change a TEST and branch instruction sequence using CMOV to eliminate a branch. If the TEST sets the equal flag, the value in EBX will be moved to EAX. This branch is data-dependent, and is representative of an unpredictable branch.

Example 3-3. Eliminating Branch with CMOV Instruction

```

test ecx, ecx
jne 1H
mov eax, ebx
1H:
; To optimize code, combine jne and mov into one cmovcc instruction that checks the equal flag
test    ecx, ecx           ; Test the flags
cmoveq  eax, ebx          ; If the equal flag is set, move
                               ; ebx to eax- the 1H: tag no longer needed

```

The CMOV and FCMOV instructions are available on the Pentium II and subsequent processors, but not on Pentium processors and earlier IA-32 processors. Be sure to check whether a processor supports these instructions with the CPUID instruction.

3.4.1.2 Spin-Wait and Idle Loops

The Pentium 4 processor introduces a new PAUSE instruction; the instruction is architecturally a NOP on Intel 64 and IA-32 processor implementations.

To the Pentium 4 and later processors, this instruction acts as a hint that the code sequence is a spin-wait loop. Without a PAUSE instruction in such loops, the Pentium 4 processor may suffer a severe penalty when exiting the loop because the processor may detect a possible memory order violation. Inserting the PAUSE instruction significantly reduces the likelihood of a memory order violation and as a result improves performance.

In Example 3-4, the code spins until memory location A matches the value stored in the register EAX. Such code sequences are common when protecting a critical section, in producer-consumer sequences, for barriers, or other synchronization.

Example 3-4. Use of PAUSE Instruction

```

lock:  cmp eax, a
       jne loop
       ; Code in critical section:
loop:  pause
       cmp eax, a
       jne loop
       jmp lock

```

3.4.1.3 Static Prediction

Branches that do not have a history in the BTB (see Section 3.4.1, “Branch Prediction Optimization”) are predicted using a static prediction algorithm. Pentium 4,

Pentium M, Intel Core Solo and Intel Core Duo processors have similar static prediction algorithms that:

- predict unconditional branches to be taken
- predict indirect branches to be NOT taken

In addition, conditional branches in processors based on the Intel NetBurst microarchitecture are predicted using the following static prediction algorithm:

- predict backward conditional branches to be taken; rule is suitable for loops
- predict forward conditional branches to be NOT taken

Pentium M, Intel Core Solo and Intel Core Duo processors do not statically predict conditional branches according to the jump direction. All conditional branches are dynamically predicted, even at first appearance.

The following rule applies to static elimination.

Assembly/Compiler Coding Rule 3. (M impact, H generality) *Arrange code to be consistent with the static branch prediction algorithm: make the fall-through code following a conditional branch be the likely target for a branch with a forward target, and make the fall-through code following a conditional branch be the unlikely target for a branch with a backward target.*

Example 3-5 illustrates the static branch prediction algorithm. The body of an IF-THEN conditional is predicted.

Example 3-5. Pentium 4 Processor Static Branch Prediction Algorithm

```

//Forward condition branches not taken (fall through)
  IF<condition> {...
  ↓
  }

IF<condition> {...
  ↓
  }

//Backward conditional branches are taken
  LOOP {...
  ↑ — }<condition>

//Unconditional branches taken
  JMP
  ----->
    
```

Examples 3-6 and Example 3-7 provide basic rules for a static prediction algorithm. In Example 3-6, the backward branch (JC BEGIN) is not in the BTB the first time

through; therefore, the BTB does not issue a prediction. The static predictor, however, will predict the branch to be taken, so a misprediction will not occur.

Example 3-6. Static Taken Prediction

```

Begin:  mov    eax, mem32
        and    eax, ebx
        imul   eax, edx
        shld   eax, 7
        jc     Begin

```

The first branch instruction (JC BEGIN) in Example 3-7 is a conditional forward branch. It is not in the BTB the first time through, but the static predictor will predict the branch to fall through. The static prediction algorithm correctly predicts that the CALL CONVERT instruction will be taken, even before the branch has any branch history in the BTB.

Example 3-7. Static Not-Taken Prediction

```

        mov    eax, mem32
        and    eax, ebx
        imul   eax, edx
        shld   eax, 7
        jc     Begin
        mov    eax, 0
Begin:  call   Convert

```

The Intel Core microarchitecture does not use the static prediction heuristic. However, to maintain consistency across Intel 64 and IA-32 processors, software should maintain the static prediction heuristic as the default.

3.4.1.4 Inlining, Calls and Returns

The return address stack mechanism augments the static and dynamic predictors to optimize specifically for calls and returns. It holds 16 entries, which is large enough to cover the call depth of most programs. If there is a chain of more than 16 nested calls and more than 16 returns in rapid succession, performance may degrade.

The trace cache in Intel NetBurst microarchitecture maintains branch prediction information for calls and returns. As long as the trace with the call or return remains in the trace cache and the call and return targets remain unchanged, the depth limit of the return address stack described above will not impede performance.

To enable the use of the return stack mechanism, calls and returns must be matched in pairs. If this is done, the likelihood of exceeding the stack depth in a manner that will impact performance is very low.

The following rules apply to inlining, calls, and returns.

Assembly/Compiler Coding Rule 4. (MH impact, MH generality) *Near calls must be matched with near returns, and far calls must be matched with far returns. Pushing the return address on the stack and jumping to the routine to be called is not recommended since it creates a mismatch in calls and returns.*

Calls and returns are expensive; use inlining for the following reasons:

- Parameter passing overhead can be eliminated.
- In a compiler, inlining a function exposes more opportunity for optimization.
- If the inlined routine contains branches, the additional context of the caller may improve branch prediction within the routine.
- A mispredicted branch can lead to performance penalties inside a small function that are larger than those that would occur if that function is inlined.

Assembly/Compiler Coding Rule 5. (MH impact, MH generality) *Selectively inline a function if doing so decreases code size or if the function is small and the call site is frequently executed.*

Assembly/Compiler Coding Rule 6. (H impact, H generality) *Do not inline a function if doing so increases the working set size beyond what will fit in the trace cache.*

Assembly/Compiler Coding Rule 7. (ML impact, ML generality) *If there are more than 16 nested calls and returns in rapid succession; consider transforming the program with inline to reduce the call depth.*

Assembly/Compiler Coding Rule 8. (ML impact, ML generality) *Favor inlining small functions that contain branches with poor prediction rates. If a branch misprediction results in a RETURN being prematurely predicted as taken, a performance penalty may be incurred.)*

Assembly/Compiler Coding Rule 9. (L impact, L generality) *If the last statement in a function is a call to another function, consider converting the call to a jump. This will save the call/return overhead as well as an entry in the return stack buffer.*

Assembly/Compiler Coding Rule 10. (M impact, L generality) *Do not put more than four branches in a 16-byte chunk.*

Assembly/Compiler Coding Rule 11. (M impact, L generality) *Do not put more than two end loop branches in a 16-byte chunk.*

3.4.1.5 Code Alignment

Careful arrangement of code can enhance cache and memory locality. Likely sequences of basic blocks should be laid out contiguously in memory. This may involve removing unlikely code, such as code to handle error conditions, from the sequence. See Section 3.7, "Prefetching," on optimizing the instruction prefetcher.

Assembly/Compiler Coding Rule 12. (M impact, H generality) *All branch targets should be 16-byte aligned.*

Assembly/Compiler Coding Rule 13. (M impact, H generality) *If the body of a conditional is not likely to be executed, it should be placed in another part of the program. If it is highly unlikely to be executed and code locality is an issue, it should be placed on a different code page.*

3.4.1.6 Branch Type Selection

The default predicted target for indirect branches and calls is the fall-through path. Fall-through prediction is overridden if and when a hardware prediction is available for that branch. The predicted branch target from branch prediction hardware for an indirect branch is the previously executed branch target.

The default prediction to the fall-through path is only a significant issue if no branch prediction is available, due to poor code locality or pathological branch conflict problems. For indirect calls, predicting the fall-through path is usually not an issue, since execution will likely return to the instruction after the associated return.

Placing data immediately following an indirect branch can cause a performance problem. If the data consists of all zeros, it looks like a long stream of ADDs to memory destinations and this can cause resource conflicts and slow down branch recovery. Also, data immediately following indirect branches may appear as branches to the branch predication hardware, which can branch off to execute other data pages. This can lead to subsequent self-modifying code problems.

Assembly/Compiler Coding Rule 14. (M impact, L generality) *When indirect branches are present, try to put the most likely target of an indirect branch immediately following the indirect branch. Alternatively, if indirect branches are common but they cannot be predicted by branch prediction hardware, then follow the indirect branch with a UD2 instruction, which will stop the processor from decoding down the fall-through path.*

Indirect branches resulting from code constructs (such as switch statements, computed GOTOs or calls through pointers) can jump to an arbitrary number of locations. If the code sequence is such that the target destination of a branch goes to the same address most of the time, then the BTB will predict accurately most of the time. Since only one taken (non-fall-through) target can be stored in the BTB, indirect branches with multiple taken targets may have lower prediction rates.

The effective number of targets stored may be increased by introducing additional conditional branches. Adding a conditional branch to a target is fruitful if:

- The branch direction is correlated with the branch history leading up to that branch; that is, not just the last target, but how it got to this branch.
- The source/target pair is common enough to warrant using the extra branch prediction capacity. This may increase the number of overall branch mispredictions, while improving the misprediction of indirect branches. The profitability is lower if the number of mispredicting branches is very large.

User/Source Coding Rule 1. (M impact, L generality) *If an indirect branch has two or more common taken targets and at least one of those targets is correlated with branch history leading up to the branch, then convert the indirect branch to a tree where one or more indirect branches are preceded by conditional branches to those targets. Apply this "peeling" procedure to the common target of an indirect branch that correlates to branch history.*

The purpose of this rule is to reduce the total number of mispredictions by enhancing the predictability of branches (even at the expense of adding more branches). The added branches must be predictable for this to be worthwhile. One reason for such predictability is a strong correlation with preceding branch history. That is, the directions taken on preceding branches are a good indicator of the direction of the branch under consideration.

Example 3-8 shows a simple example of the correlation between a target of a preceding conditional branch and a target of an indirect branch.

Example 3-8. Indirect Branch With Two Favored Targets

```
function ()
{
int n = rand();           // random integer 0 to RAND_MAX
  if ( !(n & 0x01) ) { // n will be 0 half the times
    n = 0;           // updates branch history to predict taken
  }
  // indirect branches with multiple taken targets
  // may have lower prediction rates

switch (n) {
  case 0: handle_0(); break; // common target, correlated with
                          // branch history that is forward taken
  case 1: handle_1(); break; // uncommon
  case 3: handle_3(); break; // uncommon
  default: handle_other(); // common target
}
}
```

Correlation can be difficult to determine analytically, for a compiler and for an assembly language programmer. It may be fruitful to evaluate performance with and without peeling to get the best performance from a coding effort.

An example of peeling out the most favored target of an indirect branch with correlated branch history is shown in Example 3-9.

Example 3-9. A Peeling Technique to Reduce Indirect Branch Misprediction

```

function ()
{
  int n = rand();                // Random integer 0 to RAND_MAX
  if (!(n & 0x01)) THEN
    n = 0;                       // n will be 0 half the times
  if (!n) THEN
    handle_0();                  // Peel out the most common target
                                // with correlated branch history

  {
    switch (n) {
      case 1: handle_1(); break;  // Uncommon
      case 3: handle_3(); break; // Uncommon

      default: handle_other();   // Make the favored target in
                                // the fall-through path
    }
  }
}

```

3.4.1.7 Loop Unrolling

Benefits of unrolling loops are:

- Unrolling amortizes the branch overhead, since it eliminates branches and some of the code to manage induction variables.
- Unrolling allows one to aggressively schedule (or pipeline) the loop to hide latencies. This is useful if you have enough free registers to keep variables live as you stretch out the dependence chain to expose the critical path.
- Unrolling exposes the code to various other optimizations, such as removal of redundant loads, common subexpression elimination, and so on.
- The Pentium 4 processor can correctly predict the exit branch for an inner loop that has 16 or fewer iterations (if that number of iterations is predictable and there are no conditional branches in the loop). So, if the loop body size is not excessive and the probable number of iterations is known, unroll inner loops until they have a maximum of 16 iterations. With the Pentium M processor, do not unroll loops having more than 64 iterations.

The potential costs of unrolling loops are:

- Excessive unrolling or unrolling of very large loops can lead to increased code size. This can be harmful if the unrolled loop no longer fits in the trace cache (TC).
- Unrolling loops whose bodies contain branches increases demand on BTB capacity. If the number of iterations of the unrolled loop is 16 or fewer, the branch

predictor should be able to correctly predict branches in the loop body that alternate direction.

Assembly/Compiler Coding Rule 15. (H impact, M generality) *Unroll small loops until the overhead of the branch and induction variable accounts (generally) for less than 10% of the execution time of the loop.*

Assembly/Compiler Coding Rule 16. (H impact, M generality) *Avoid unrolling loops excessively; this may thrash the trace cache or instruction cache.*

Assembly/Compiler Coding Rule 17. (M impact, M generality) *Unroll loops that are frequently executed and have a predictable number of iterations to reduce the number of iterations to 16 or fewer. Do this unless it increases code size so that the working set no longer fits in the trace or instruction cache. If the loop body contains more than one conditional branch, then unroll so that the number of iterations is 16/(# conditional branches).*

Example 3-10 shows how unrolling enables other optimizations.

Example 3-10. Loop Unrolling

```

Before unrolling:
    do i = 1, 100
        if ( i mod 2 == 0 ) then a(i) = x
            else a(i) = y
        enddo
After unrolling
    do i = 1, 100, 2
        a(i) = y
        a(i+1) = x
    enddo
    
```

In this example, the loop that executes 100 times assigns X to every even-numbered element and Y to every odd-numbered element. By unrolling the loop you can make assignments more efficiently, removing one branch in the loop body.

3.4.1.8 Compiler Support for Branch Prediction

Compilers generate code that improves the efficiency of branch prediction in the Pentium 4, Pentium M, Intel Core Duo processors and processors based on Intel Core microarchitecture. The Intel C++ Compiler accomplishes this by:

- keeping code and data on separate pages
- using conditional move instructions to eliminate branches
- generating code consistent with the static branch prediction algorithm
- inlining where appropriate
- unrolling if the number of iterations is predictable

With profile-guided optimization, the compiler can lay out basic blocks to eliminate branches for the most frequently executed paths of a function or at least improve their predictability. Branch prediction need not be a concern at the source level. For more information, see Intel C++ Compiler documentation.

3.4.2 Fetch and Decode Optimization

Intel Core microarchitecture provides several mechanisms to increase front end throughput. Techniques to take advantage of some of these features are discussed below.

3.4.2.1 Optimizing for Micro-fusion

An Instruction that operates on a register and a memory operand decodes into more μ ops than its corresponding register-register version. Replacing the equivalent work of the former instruction using the register-register version usually require a sequence of two instructions. The latter sequence is likely to result in reduced fetch bandwidth.

Assembly/Compiler Coding Rule 18. (ML impact, M generality) *For improving fetch/decode throughput, Give preference to memory flavor of an instruction over the register-only flavor of the same instruction, if such instruction can benefit from micro-fusion.*

The following examples are some of the types of micro-fusions that can be handled by all decoders:

- All stores to memory, including store immediate. Stores execute internally as two separate μ ops: store-address and store-data.
- All “read-modify” (load+op) instructions between register and memory, for example:


```
ADDPS  XMM9, QWORD PTR [RSP+40]
FADD   DOUBLE PTR [RDI+RSI*8]
XOR    RAX, QWORD PTR [RBP+32]
```
- All instructions of the form “load and jump,” for example:


```
JMP    [RDI+200]
RET
```
- CMP and TEST with immediate operand and memory

An Intel 64 instruction with RIP relative addressing is not micro-fused in the following cases:

- When an additional immediate is needed, for example:


```
CMP    [RIP+400], 27
MOV    [RIP+3000], 142
```
- When an RIP is needed for control flow purposes, for example:


```
JMP    [RIP+5000000]
```

In these cases, Intel Core Microarchitecture and Intel microarchitecture code name Sandy Bridge provides a 2 μ op flow from decoder 0, resulting in a slight loss of decode bandwidth since 2 μ op flow must be steered to decoder 0 from the decoder with which it was aligned.

RIP addressing may be common in accessing global data. Since it will not benefit from micro-fusion, compiler may consider accessing global data with other means of memory addressing.

3.4.2.2 Optimizing for Macro-fusion

Macro-fusion merges two instructions to a single μ op. Intel Core Microarchitecture performs this hardware optimization under limited circumstances.

The first instruction of the macro-fused pair must be a CMP or TEST instruction. This instruction can be REG-REG, REG-IMM, or a micro-fused REG-MEM comparison. The second instruction (adjacent in the instruction stream) should be a conditional branch.

Since these pairs are common ingredient in basic iterative programming sequences, macro-fusion improves performance even on un-recompiled binaries. All of the decoders can decode one macro-fused pair per cycle, with up to three other instructions, resulting in a peak decode bandwidth of 5 instructions per cycle.

Each macro-fused instruction executes with a single dispatch. This process reduces latency, which in this case shows up as a cycle removed from branch mispredict penalty. Software also gain all other fusion benefits: increased rename and retire bandwidth, more storage for instructions in-flight, and power savings from representing more work in fewer bits.

The following list details when you can use macro-fusion:

- CMP or TEST can be fused when comparing:
 - REG-REG. For example: `CMP EAX,ECX; JZ label`
 - REG-IMM. For example: `CMP EAX,0x80; JZ label`
 - REG-MEM. For example: `CMP EAX,[ECX]; JZ label`
 - MEM-REG. For example: `CMP [EAX],ECX; JZ label`
- TEST can fused with all conditional jumps.
- CMP can be fused with only the following conditional jumps in Intel Core microarchitecture. These conditional jumps check carry flag (CF) or zero flag (ZF). jump. The list of macro-fusion-capable conditional jumps are:

- JA or JNBE
- JAE or JNB or JNC
- JE or JZ
- JNA or JBE
- JNAE or JC or JB
- JNE or JNZ

CMP and TEST can not be fused when comparing MEM-IMM (e.g. CMP [EAX],0x80; JZ label). Macro-fusion is not supported in 64-bit mode for Intel Core microarchitecture.

- Intel microarchitecture code name Nehalem supports the following enhancements in macrofusion:
 - CMP can be fused with the following conditional jumps (that was not supported in Intel Core microarchitecture):
 - JL or JNGE
 - JGE or JNL
 - JLE or JNG
 - JG or JNLE
 - Macro-fusion is support in 64-bit mode.
- Enhanced macrofusion support in Intel microarchitecture code name Sandy Bridge is summarized in Table 3-1 with additional information in Section 2.1.2.1 and Example 3-15:

Table 3-1. Macro-Fusible Instructions in Intel microarchitecture code name Sandy Bridge

Instructions	TEST	AND	CMP	ADD	SUB	INC	DEC
JO/JNO	Y	Y	N	N	N	N	N
JC/JB/JAE/JNB	Y	Y	Y	Y	Y	N	N
JE/JZ/JNE/JNZ	Y	Y	Y	Y	Y	Y	Y
JNA/JBE/JA/JNBE	Y	Y	Y	Y	Y	N	N
JS/JNS/JP/JPE/JNP/JPO	Y	Y	N	N	N	N	N
JL/JNGE/JGE/JNL/JLE/JNG/JG/JNLE	Y	Y	Y	Y	Y	Y	Y

Assembly/Compiler Coding Rule 19. (M impact, ML generality) *Employ macro-fusion where possible using instruction pairs that support macro-fusion. Prefer TEST over CMP if possible. Use unsigned variables and unsigned jumps when*

GENERAL OPTIMIZATION GUIDELINES

possible. Try to logically verify that a variable is non-negative at the time of comparison. Avoid CMP or TEST of MEM-IMM flavor when possible. However, do not add other instructions to avoid using the MEM-IMM flavor.

Example 3-11. Macro-fusion, Unsigned Iteration Count

	Without Macro-fusion	With Macro-fusion
C code	for (int ¹ i = 0; i < 1000; i++) a++;	for (unsigned int ² i = 0; i < 1000; i++) a++;
Disassembly	for (int i = 0; i < 1000; i++) mov dword ptr [i], 0 jmp First Loop: mov eax, dword ptr [i] add eax, 1 mov dword ptr [i], eax First: cmp dword ptr [i], 3E8H ³ jge End a++; mov eax, dword ptr [a] addq eax, 1 mov dword ptr [a], eax jmp Loop End:	for (unsigned int i = 0; i < 1000; i++) mov dword ptr [i], 0 jmp First Loop: mov eax, dword ptr [i] add eax, 1 mov dword ptr [i], eax First: cmp eax, 3E8H ⁴ jae End a++; mov eax, dword ptr [a] add eax, 1 mov dword ptr [a], eax jmp Loop End:

NOTES:

1. Signed iteration count inhibits macro-fusion
2. Unsigned iteration count is compatible with macro-fusion
3. CMP MEM-IMM, JGE inhibit macro-fusion
4. CMP REG-IMM, JAE permits macro-fusion

Example 3-12. Macro-fusion, If Statement

	Without Macro-fusion	With Macro-fusion
C code	int ¹ a = 7; if (a < 77) a++; else a--;	unsigned int ² a = 7; if (a < 77) a++; else a--;

Example 3-12. Macro-fusion, If Statement (Contd.)

	Without Macro-fusion	With Macro-fusion
Disassembly	<pre>int a = 7; mov dword ptr [a], 7 if (a < 77) cmp dword ptr [a], 4DH³ jge Dec a++; mov eax, dword ptr [a] add eax, 1 mov dword ptr [a], eax else jmp End a--; Dec: mov eax, dword ptr [a] sub eax, 1 mov dword ptr [a], eax End::</pre>	<pre>unsigned int a = 7; mov dword ptr [a], 7 if (a < 77) mov eax, dword ptr [a] cmp eax, 4DH jae Dec a++; add eax, 1 mov dword ptr [a], eax else jmp End a--; Dec: sub eax, 1 mov dword ptr [a], eax End::</pre>

NOTES:

1. Signed iteration count inhibits macro-fusion
2. Unsigned iteration count is compatible with macro-fusion
3. CMP MEM-IMM, JGE inhibit macro-fusion

Assembly/Compiler Coding Rule 20. (M impact, ML generality) Software can enable macro fusion when it can be logically determined that a variable is non-negative at the time of comparison; use TEST appropriately to enable macro-fusion when comparing a variable with 0.

Example 3-13. Macro-fusion, Signed Variable

Without Macro-fusion	With Macro-fusion
<pre>test ecx, ecx jle OutSideTheIF cmp ecx, 64H jge OutSideTheIF <IF BLOCK CODE> OutSideTheIF:</pre>	<pre>test ecx, ecx jle OutSideTheIF cmp ecx, 64H jae OutSideTheIF <IF BLOCK CODE> OutSideTheIF:</pre>

For either signed or unsigned variable 'a'; "CMP a,0" and "TEST a,a" produce the same result as far as the flags are concerned. Since TEST can be macro-fused more

often, software can use “TEST a,a” to replace “CMP a,0” for the purpose of enabling macro-fusion.

Example 3-14. Macro-fusion, Signed Comparison

C Code	Without Macro-fusion	With Macro-fusion
if (a == 0)	<pre> cmp a, 0 jne lbl ... lbl: </pre>	<pre> test a, a jne lbl ... lbl: </pre>
if (a >= 0)	<pre> cmp a, 0 jl lbl; ... lbl: </pre>	<pre> test a, a jl lbl ... lbl: </pre>

Intel microarchitecture code name Sandy Bridge enables more arithmetic and logic instructions to macro-fuse with conditional branches. In loops where the ALU ports are already congested, performing one of these macro-fusions can relieve the pressure, as the macro-fused instruction consumes only port 5, instead of an ALU port plus port 5.

In Example 3-15, the “add/cmp/jnz” loop contains two ALU instructions that can be dispatched via either port 0, 1, 5. So there is higher probability of port 5 might bind to either ALU instruction causing JNZ to wait a cycle. The “sub/jnz” loop, the likelihood of ADD/SUB/JNZ can be dispatched in the same cycle is increased because only SUB is free to bind with either port 0, 1, 5.

Example 3-15. Additional Macro-fusion Benefit in Intel microarchitecture code name Sandy Bridge

Add + cmp + jnz alternative	Loop control with sub + jnz
<pre> lea rdx, buff xor rcx, rcx xor eax, eax loop: add eax, [rdx + 4 * rcx] add rcx, 1 cmp rcx, LEN jnz loop </pre>	<pre> lea rdx, buff - 4 xor rcx, LEN xor eax, eax loop: add eax, [rdx + 4 * rcx] sub rcx, 1 jnz loop </pre>

3.4.2.3 Length-Changing Prefixes (LCP)

The length of an instruction can be up to 15 bytes in length. Some prefixes can dynamically change the length of an instruction that the decoder must recognize. Typically, the pre-decode unit will estimate the length of an instruction in the byte

stream assuming the absence of LCP. When the predecoder encounters an LCP in the fetch line, it must use a slower length decoding algorithm. With the slower length decoding algorithm, the predecoder decodes the fetch in 6 cycles, instead of the usual 1 cycle. Normal queuing throughout of the machine pipeline generally cannot hide LCP penalties.

The prefixes that can dynamically change the length of an instruction include:

- operand size prefix (0x66)
- address size prefix (0x67)

The instruction `MOV DX, 01234h` is subject to LCP stalls in processors based on Intel Core microarchitecture, and in Intel Core Duo and Intel Core Solo processors. Instructions that contain `imm16` as part of their fixed encoding but do not require LCP to change the immediate size are not subject to LCP stalls. The REX prefix (4xh) in 64-bit mode can change the size of two classes of instruction, but does not cause an LCP penalty.

If the LCP stall happens in a tight loop, it can cause significant performance degradation. When decoding is not a bottleneck, as in floating-point heavy code, isolated LCP stalls usually do not cause performance degradation.

Assembly/Compiler Coding Rule 21. (MH impact, MH generality) *Favor generating code using `imm8` or `imm32` values instead of `imm16` values.*

If `imm16` is needed, load equivalent `imm32` into a register and use the word value in the register instead.

Double LCP Stalls

Instructions that are subject to LCP stalls and cross a 16-byte fetch line boundary can cause the LCP stall to trigger twice. The following alignment situations can cause LCP stalls to trigger twice:

- An instruction is encoded with a `MODR/M` and `SIB` byte, and the fetch line boundary crossing is between the `MODR/M` and the `SIB` bytes.
- An instruction starts at offset 13 of a fetch line references a memory location using register and immediate byte offset addressing mode.

The first stall is for the 1st fetch line, and the 2nd stall is for the 2nd fetch line. A double LCP stall causes a decode penalty of 11 cycles.

The following examples cause LCP stall once, regardless of their fetch-line location of the first byte of the instruction:

```
ADD DX, 01234H
ADD word ptr [EDX], 01234H
ADD word ptr 012345678H[EDX], 01234H
ADD word ptr [012345678H], 01234H
```

The following instructions cause a double LCP stall when starting at offset 13 of a fetch line:

```
ADD word ptr [EDX+ESI], 01234H
```

```
ADD word ptr 012H[EDX], 01234H
ADD word ptr 012345678H[EDX+ESI], 01234H
```

To avoid double LCP stalls, do not use instructions subject to LCP stalls that use SIB byte encoding or addressing mode with byte displacement.

False LCP Stalls

False LCP stalls have the same characteristics as LCP stalls, but occur on instructions that do not have any imm16 value.

False LCP stalls occur when (a) instructions with LCP that are encoded using the F7 opcodes, and (b) are located at offset 14 of a fetch line. These instructions are: not, neg, div, idiv, mul, and imul. False LCP experiences delay because the instruction length decoder can not determine the length of the instruction before the next fetch line, which holds the exact opcode of the instruction in its MODR/M byte.

The following techniques can help avoid false LCP stalls:

- Upcast all short operations from the F7 group of instructions to long, using the full 32 bit version.
- Ensure that the F7 opcode never starts at offset 14 of a fetch line.

Assembly/Compiler Coding Rule 22. (M impact, ML generality) *Ensure instructions using 0xF7 opcode byte does not start at offset 14 of a fetch line; and avoid using these instruction to operate on 16-bit data, upcast short data to 32 bits.*

Example 3-16. Avoiding False LCP Delays with 0xF7 Group Instructions

A Sequence Causing Delay in the Decoder	Alternate Sequence to Avoid Delay
neg word ptr a	movsx eax, word ptr a neg eax mov word ptr a, AX

3.4.2.4 Optimizing the Loop Stream Detector (LSD)

Loops that fit the following criteria are detected by the LSD and replayed from the instruction queue to feed the decoder in Intel Core microarchitecture:

- Must be less than or equal to four 16-byte fetches.
- Must be less than or equal to 18 instructions.
- Can contain no more than four taken branches and none of them can be a RET.
- Should usually have more than 64 iterations.

Loop Stream Detector in Intel microarchitecture code name Nehalem is improved by:

- Caching decoded micro-operations in the instruction decoder queue (IDQ, see Section 2.3.2) to feed the rename/alloc stage.
- The size of the LSD is increased to 28 micro-ops.

Many calculation-intensive loops, searches and software string moves match these characteristics. These loops exceed the BPU prediction capacity and always terminate in a branch misprediction.

Assembly/Compiler Coding Rule 23. (MH impact, MH generality) Break up a loop long sequence of instructions into loops of shorter instruction blocks of no more than the size of LSD.

Assembly/Compiler Coding Rule 24. (MH impact, M generality) Avoid unrolling loops containing LCP stalls, if the unrolled block exceeds the size of LSD.

3.4.2.5 Exploit LSD Micro-op Emission Bandwidth in Intel® microarchitecture code name Sandy Bridge

The LSD holds micro-ops that construct small "infinite" loops. Micro-ops from the LSD are allocated in the out-of-order engine. The loop in the LSD ends with a taken branch to the beginning of the loop. The taken branch at the end of the loop is always the last micro-op allocated in the cycle. The instruction at the beginning of the loop is always allocated at the next cycle. If the code performance is bound by front-end bandwidth, unused allocation slots result in a bubble in allocation, and can cause performance degradation.

Allocation bandwidth in Intel microarchitecture code name Sandy Bridge is four micro-ops per cycle. Performance is best, when the number of micro-ops in the LSD result in the least number of unused allocation slots. You can use loop unrolling to control the number of micro-ops that are in the LSD.

In the Example 3-17, the code sums all array elements. The original code adds one element per iteration. It has three micro-ops per iteration, all allocated in one cycle. Code throughput is one load per cycle.

When unrolling the loop once there are five micro-ops per iteration, which are allocated in two cycles. Code throughput is still one load per cycle. Therefore there is no performance gain.

When unrolling the loop twice there are seven micro-ops per iteration, still allocated in two cycles. Since two loads can be executed in each cycle this code has a potential throughput of three load operations in two cycles.

Example 3-17. Unrolling Loops in LSD to Optimize Emission Bandwidth

No Unrolling	Unroll once	Unroll Twice
lp: add eax, [rsi + 4* rcx] dec rcx jnz lp	lp: add eax, [rsi + 4* rcx] add eax, [rsi + 4* rcx +4] add rcx, -2 jnz lp	lp: add eax, [rsi + 4* rcx] add eax, [rsi + 4* rcx +4] add eax, [rsi + 4* rcx + 8] add rcx, -3 jnz lp

3.4.2.6 Optimization for Decoded ICache

The decoded ICache is a new feature in Intel microarchitecture code name Sandy Bridge. Running the code from the Decoded ICache has two advantages:

- Higher bandwidth of micro-ops feeding the out-of-order engine.
- The front end does not need to decode the code that is in the Decoded ICache. This saves power.

There is overhead in switching between the Decoded ICache and the legacy decode pipeline. If your code switches frequently between the front end and the Decoded ICache, the penalty may be higher than running only from the legacy pipeline

To ensure “hot” code is feeding from the decoded ICache,

- Make sure each hot code block is less than about 500 instructions. Specifically, do not unroll to more than 500 instructions in a loop. This should enable Decoded ICache residency even when hyper-threading is enabled.
- For applications with very large blocks of calculations inside a loop, consider loop-fission: split the loop into multiple loops that fit in the Decoded ICache, rather than a single loop that overflows.
- If an application can be sure to run with only one thread per core, it can increase hot code block size to about 1000 instructions.

Dense Read-Modify-Write Code

The Decoded ICache can hold only up to 18 micro-ops per each 32 byte aligned memory chunk. Therefore, code with a high concentration of instructions that are encoded in a small number of bytes, yet have many micro-ops, may overflow the 18 micro-op limitation and not enter the Decoded ICache. Read-modify-write (RMW) instructions are a good example of such instructions.

RMW instructions accept one memory source operand, one register source operand, and use the source memory operand as the destination. The same functionality can be achieved by two or three instructions: the first reads the memory source operand, the second performs the operation with the second register source operand, and the last writes the result back to memory. These instructions usually result in the same number of micro-ops but use more bytes to encode the same functionality.

One case where RMW instructions may be used extensively is when the compiler optimizes aggressively for code size.

Here are some possible solutions to fit the hot code in the Decoded ICache:

- Replace RMW instructions with two or three instructions that have the same functionality. For example, “`adc [rdi], rcx`” is only three bytes long; the equivalent sequence “`adc rax, [rdi]`” + “`mov [rdi], rax`” has a footprint of six bytes.
- Align the code so that the dense part is broken down among two different 32-byte chunks. This solution is useful when using a tool that aligns code automatically, and is indifferent to code changes.

- Spread the code by adding multiple byte NOPs in the loop. Note that this solution adds micro-ops for execution.

Align Unconditional Branches for Decoded ICache

For code entering the Decoded ICache, each unconditional branch is the last micro-op occupying a Decoded ICache Way. Therefore, only three unconditional branches per a 32 byte aligned chunk can enter the Decoded ICache.

Unconditional branches are frequent in jump tables and switch declarations. Below are examples for these constructs, and methods for writing them so that they fit in the Decoded ICache.

Compilers create jump tables for C++ virtual class methods or DLL dispatch tables. Each unconditional branch consumes five bytes; therefore up to seven of them can be associated with a 32-byte chunk. Thus jump tables may not fit in the Decoded ICache if the unconditional branches are too dense in each 32Byte-aligned chunk. This can cause performance degradation for code executing before and after the branch table.

The solution is to add multi-byte NOP instructions among the branches in the branch table. This may increase code size and should be used cautiously. However, these NOPs are not executed and therefore have no penalty in later pipe stages.

Switch-Case constructs represents a similar situation. Each evaluation of a case condition results in an unconditional branch. The same solution of using multi-byte NOP can apply for every three consecutive unconditional branches that fits inside an aligned 32-byte chunk.

Two Branches in a Decoded ICache Way

The Decoded ICache can hold up to two branches in a way. Dense branches in a 32 byte aligned chunk, or their ordering with other instructions may prohibit all the micro-ops of the instructions in the chunk from entering the Decoded ICache. This does not happen often. When it does happen, you can space the code with NOP instructions where appropriate. Make sure that these NOP instructions are not part of hot code

3.4.2.7 Scheduling Rules for the Pentium M Processor Decoder

The Pentium M processor has three decoders, but the decoding rules to supply μ ops at high bandwidth are less stringent than those of the Pentium III processor. This provides an opportunity to build a front-end tracker in the compiler and try to schedule instructions correctly. The decoder limitations are:

- The first decoder is capable of decoding one macroinstruction made up of four or fewer μ ops in each clock cycle. It can handle any number of bytes up to the maximum of 15. Multiple prefix instructions require additional cycles.
- The two additional decoders can each decode one macroinstruction per clock cycle (assuming the instruction is one μ op up to seven bytes in length).
- Instructions composed of more than four μ ops take multiple cycles to decode.

Assembly/Compiler Coding Rule 25. (M impact, M generality) Avoid putting explicit references to ESP in a sequence of stack operations (POP, PUSH, CALL, RET).

3.4.2.8 Other Decoding Guidelines

Assembly/Compiler Coding Rule 26. (ML impact, L generality) Use simple instructions that are less than eight bytes in length.

Assembly/Compiler Coding Rule 27. (M impact, MH generality) Avoid using prefixes to change the size of immediate and displacement.

Long instructions (more than seven bytes) limit the number of decoded instructions per cycle on the Pentium M processor. Each prefix adds one byte to the length of instruction, possibly limiting the decoder's throughput. In addition, multiple prefixes can only be decoded by the first decoder. These prefixes also incur a delay when decoded. If multiple prefixes or a prefix that changes the size of an immediate or displacement cannot be avoided, schedule them behind instructions that stall the pipe for some other reason.

3.5 OPTIMIZING THE EXECUTION CORE

The superscalar, out-of-order execution core(s) in recent generations of microarchitectures contain multiple execution hardware resources that can execute multiple μ ops in parallel. These resources generally ensure that μ ops execute efficiently and proceed with fixed latencies. General guidelines to make use of the available parallelism are:

- Follow the rules (see Section 3.4) to maximize useful decode bandwidth and front end throughput. These rules include favouring single μ op instructions and taking advantage of micro-fusion, Stack pointer tracker and macro-fusion.
- Maximize rename bandwidth. Guidelines are discussed in this section and include properly dealing with partial registers, ROB read ports and instructions which causes side-effects on flags.
- Scheduling recommendations on sequences of instructions so that multiple dependency chains are alive in the reservation station (RS) simultaneously, thus ensuring that your code utilizes maximum parallelism.
- Avoid hazards, minimize delays that may occur in the execution core, allowing the dispatched μ ops to make progress and be ready for retirement quickly.

3.5.1 Instruction Selection

Some execution units are not pipelined, this means that μ ops cannot be dispatched in consecutive cycles and the throughput is less than one per cycle.

It is generally a good starting point to select instructions by considering the number of μ ops associated with each instruction, favoring in the order of: single- μ op instructions, simple instruction with less than 4 μ ops, and last instruction requiring microsequencer ROM (μ ops which are executed out of the microsequencer involve extra overhead).

Assembly/Compiler Coding Rule 28. (M impact, H generality) *Favor single-micro-operation instructions. Also favor instruction with shorter latencies.*

A compiler may be already doing a good job on instruction selection. If so, user intervention usually is not necessary.

Assembly/Compiler Coding Rule 29. (M impact, L generality) *Avoid prefixes, especially multiple non-0F-prefixed opcodes.*

Assembly/Compiler Coding Rule 30. (M impact, L generality) *Do not use many segment registers.*

On the Pentium M processor, there is only one level of renaming of segment registers.

Assembly/Compiler Coding Rule 31. (ML impact, M generality) *Avoid using complex instructions (for example, enter, leave, or loop) that have more than four μ ops and require multiple cycles to decode. Use sequences of simple instructions instead.*

Complex instructions may save architectural registers, but incur a penalty of 4 μ ops to set up parameters for the microsequencer ROM in Intel NetBurst microarchitecture.

Theoretically, arranging instructions sequence to match the 4-1-1-1 template applies to processors based on Intel Core microarchitecture. However, with macro-fusion and micro-fusion capabilities in the front end, attempts to schedule instruction sequences using the 4-1-1-1 template will likely provide diminishing returns.

Instead, software should follow these additional decoder guidelines:

- If you need to use multiple μ op, non-microsequenced instructions, try to separate by a few single μ op instructions. The following instructions are examples of multiple- μ op instruction not requiring micro-sequencer:

ADC/SBB

CMOVcc

Read-modify-write instructions

- If a series of multiple- μ op instructions cannot be separated, try breaking the series into a different equivalent instruction sequence. For example, a series of read-modify-write instructions may go faster if sequenced as a series of read-modify + store instructions. This strategy could improve performance even if the new code sequence is larger than the original one.

3.5.1.1 Use of the INC and DEC Instructions

The INC and DEC instructions modify only a subset of the bits in the flag register. This creates a dependence on all previous writes of the flag register. This is especially

problematic when these instructions are on the critical path because they are used to change an address for a load on which many other instructions depend.

Assembly/Compiler Coding Rule 32. (M impact, H generality) *INC and DEC instructions should be replaced with ADD or SUB instructions, because ADD and SUB overwrite all flags, whereas INC and DEC do not, therefore creating false dependencies on earlier instructions that set the flags.*

3.5.1.2 Integer Divide

Typically, an integer divide is preceded by a CWD or CDQ instruction. Depending on the operand size, divide instructions use DX:AX or EDX:EAX for the dividend. The CWD or CDQ instructions sign-extend AX or EAX into DX or EDX, respectively. These instructions have denser encoding than a shift and move would be, but they generate the same number of micro-ops. If AX or EAX is known to be positive, replace these instructions with:

```
xor dx, dx
```

or

```
xor edx, edx
```

Modern compilers typically can transform high-level language expression involving integer division where the divisor is a known integer constant at compile time into a faster sequence using IMUL instruction instead. Thus programmers should minimize integer division expression with divisor whose value can not be known at compile time.

Alternately, if certain known divisor value are favored over other unknown ranges, software may consider isolating the few favored, known divisor value into constant-divisor expressions.

Section 9.2.4 describes more detail of using MUL/IMUL to replace integer divisions.

3.5.1.3 Using LEA

In Intel microarchitecture code name Sandy Bridge, there are two significant changes to the performance characteristics of LEA instruction:

- LEA can be dispatched via port 1 and 5 in most cases, doubling the throughput over prior generations. However this apply only to LEA instructions with one or two source operands.

Example 3-18. Independent Two-Operand LEA Example

```
mov    edx, N
mov    eax, X
mov    ecx, Y
```

Example 3-18. Independent Two-Operand LEA Example (Contd.)

```

loop:
  lea   ecx, [ecx = ecx *2]
  lea   eax, [eax = eax *5]
  and   ecx, 0xff
  and   eax, 0xff
  dec   edx
  jg    loop

```

- For LEA instructions with three source operands and some specific situations, instruction latency has increased to 3 cycles, and must dispatch via port 1:
 - LEA that has all three source operands: base, index, and offset
 - LEA that uses base and index registers where the base is EBP, RBP, or R13
 - LEA that uses RIP relative addressing mode
 - LEA that uses 16-bit addressing mode

Example 3-19. Alternative to Three-Operand LEA

3 operand LEA is slower	Two-operand LEA alternative	Alternative 2
<pre> #define K 1 uint32 an = 0; uint32 N= mi_N; mov ecx, N xor esi, esi; xor edx, edx; cmp ecx, 2; jb finished; dec ecx; loop1: mov edi, esi; lea esi, [K+esi+edx]; and esi, 0xFF; mov edx, edi; dec ecx; jnz loop1; finished: mov [an],esi; </pre>	<pre> #define K 1 uint32 an = 0; uint32 N= mi_N; mov ecx, N xor esi, esi; xor edx, edx; cmp ecx, 2; jb finished; dec ecx; loop1: mov edi, esi; lea esi, [K+edx]; lea esi, [esi+edx]; and esi, 0xFF; mov edx, edi; dec ecx; jnz loop1; finished: mov [an],esi; </pre>	<pre> #define K 1 uint32 an = 0; uint32 N= mi_N; mov ecx, N xor esi, esi; mov edx, K; cmp ecx, 2; jb finished; mov eax, 2 dec ecx; loop1: mov edi, esi; lea esi, [esi+edx]; and esi, 0xFF; lea edx, [edi +K]; dec ecx; jnz loop1; finished: mov [an],esi; </pre>

In some cases with processor based on Intel NetBurst microarchitecture, the LEA instruction or a sequence of LEA, ADD, SUB and SHIFT instructions can replace constant multiply instructions. The LEA instruction can also be used as a multiple operand addition instruction, for example:

```
LEA ECX, [EAX + EBX + 4 + A]
```

Using LEA in this way may avoid register usage by not tying up registers for operands of arithmetic instructions. This use may also save code space.

If the LEA instruction uses a shift by a constant amount then the latency of the sequence of μ ops is shorter if adds are used instead of a shift, and the LEA instruction may be replaced with an appropriate sequence of μ ops. This, however, increases the total number of μ ops, leading to a trade-off.

Assembly/Compiler Coding Rule 33. (ML impact, L generality) *If an LEA instruction using the scaled index is on the critical path, a sequence with ADDs may be better. If code density and bandwidth out of the trace cache are the critical factor, then use the LEA instruction.*

3.5.1.4 ADC and SBB in Intel® microarchitecture code name Sandy Bridge

The throughput of ADC and SBB in Intel microarchitecture code name Sandy Bridge is 1 cycle, compared to 1.5-2 cycles in prior generation. These two instructions are useful in numeric handling of integer data types that are wider than the maximum width of native hardware.

Example 3-20. Examples of 512-bit Additions

<pre>//Add 64-bit to 512 Number lea rsi, gLongCounter lea rdi, gStepValue mov rax, [rdi] xor rcx, rcx loop_start: mov r10, [rsi+rcx] add r10, rax mov [rsi+rcx], r10 mov r10, [rsi+rcx+8] adc r10, 0 mov [rsi+rcx+8], r10</pre>	<pre>// 512-bit Addition loop1: mov rax, [StepValue] add rax, [LongCounter] mov LongCounter, rax mov rax, [StepValue+8] adc rax, [LongCounter+8] mov LongCounter+8, rax mov rax, [StepValue+16] adc rax, [LongCounter+16] mov LongCounter+16, rax mov rax, [StepValue+24] adc rax, [LongCounter+24]</pre>
--	--

Example 3-20. Examples of 512-bit Additions (Contd.)

<pre> mov r10, [rsi+rcx+16] adc r10, 0 mov [rsi+rcx+16], r10 mov r10, [rsi+rcx+24] adc r10, 0 mov [rsi+rcx+24], r10 mov r10, [rsi+rcx+32] adc r10, 0 mov [rsi+rcx+32], r10 mov r10, [rsi+rcx+40] adc r10, 0 mov [rsi+rcx+40], r10 mov r10, [rsi+rcx+48] adc r10, 0 mov [rsi+rcx+48], r10 mov r10, [rsi+rcx+56] adc r10, 0 mov [rsi+rcx+56], r10 add rcx, 64 cmp rcx, SIZE jnz loop_start </pre>	<pre> mov LongCounter+24, rax mov rax, [StepValue+32] adc rax, [LongCounter+32] mov LongCounter+32, rax mov rax, [StepValue+40] adc rax, [LongCounter+40] mov LongCounter+40, rax mov rax, [StepValue+48] adc rax, [LongCounter+48] mov LongCounter+48, rax mov rax, [StepValue+56] adc rax, [LongCounter+56] mov LongCounter+56, rax dec rcx jnz loop1 </pre>
--	---

3.5.1.5 Bitwise Rotation

Bitwise rotation can choose between rotate with count specified in the CL register, an immediate constant and by 1 bit. Generally, The rotate by immediate and rotate by register instructions are slower than rotate by 1 bit. The rotate by 1 instruction has the same latency as a shift.

Assembly/Compiler Coding Rule 34. (ML impact, L generality) Avoid *ROTATE* by register or *ROTATE* by immediate instructions. If possible, replace with a *ROTATE* by 1 instruction.

In Intel microarchitecture code name Sandy Bridge, ROL/ROR by immediate has 1-cycle throughput, SHLD/SHRD using the same register as source and destination by an immediate constant has 1-cycle latency with 0.5 cycle throughput.

3.5.1.6 Address Calculations

For computing addresses, use the addressing modes rather than general-purpose computations. Internally, memory reference instructions can have four operands:

- Relocatable load-time constant
- Immediate constant
- Base register
- Scaled index register

Note that the latency and throughput of LEA with more than two operands are slower (see Section 3.5.1.3) in Intel microarchitecture code name Sandy Bridge. Addressing modes that uses both base and index registers will consume more read port resource in the execution engine and may experience more stalls due to availability of read port resources. Software should take care by selecting the speedy version of address calculation.

In the segmented model, a segment register may constitute an additional operand in the linear address calculation. In many cases, several integer instructions can be eliminated by fully using the operands of memory references.

3.5.1.7 Clearing Registers and Dependency Breaking Idioms

Code sequences that modifies partial register can experience some delay in its dependency chain, but can be avoided by using dependency breaking idioms.

In processors based on Intel Core microarchitecture, a number of instructions can help clear execution dependency when software uses these instruction to clear register content to zero. The instructions include

```
XOR REG, REG
SUB REG, REG
XORPS/PD XMMREG, XMMREG
PXOR XMMREG, XMMREG
SUBPS/PD XMMREG, XMMREG
PSUBB/W/D/Q XMMREG, XMMREG
```

In processors based on Intel microarchitecture code name Sandy Bridge, the instruction listed above plus equivalent AVX counter parts are also zero idioms that can be used to break dependency chains. Furthermore, they do not consume an issue port or an execution unit. So using zero idioms are preferable than moving 0's into the register. The AVX equivalent zero idioms are:

```
VXORPS/PD XMMREG, XMMREG
VXORPS/PD YMMREG, YMMREG
VPXOR XMMREG, XMMREG
VSUBPS/PD XMMREG, XMMREG
VSUBPS/PD YMMREG, YMMREG
VPSUBB/W/D/Q XMMREG, XMMREG
```


In Intel Core Solo and Intel Core Duo processors, the XOR, SUB, XORPS, or PXOR instructions can be used to clear execution dependencies on the zero evaluation of the destination register.

The Pentium 4 processor provides special support for XOR, SUB, and PXOR operations when executed within the same register. This recognizes that clearing a register does not depend on the old value of the register. The XORPS and XORPD instructions do not have this special support. They cannot be used to break dependence chains.

Assembly/Compiler Coding Rule 35. (M impact, ML generality) *Use dependency-breaking-idiom instructions to set a register to 0, or to break a false dependence chain resulting from re-use of registers. In contexts where the condition codes must be preserved, move 0 into the register instead. This requires more code space than using XOR and SUB, but avoids setting the condition codes.*

Example 3-21 of using pxor to break dependency idiom on a XMM register when performing negation on the elements of an array.

```
int a[4096], b[4096], c[4096];
For (int i = 0; i < 4096; i++)
    C[i] = - ( a[i] + b[i] );
```

Example 3-21. Clearing Register to Break Dependency While Negating Array Elements

Negation (-x = (x XOR (-1)) - (-1)) without breaking dependency	Negation (-x = 0 -x) using PXOR reg, reg breaks dependency
<pre>Lea eax, a lea ecx, b lea edi, c xor edx, edx movdqa xmm7, allone lp:</pre>	<pre>lea eax, a lea ecx, b lea edi, c xor edx, edx lp:</pre>
<pre>movdqa xmm0, [eax + edx] paddd xmm0, [ecx + edx] pxor xmm0, xmm7 psubd xmm0, xmm7 movdqa [edi + edx], xmm0 add edx, 16 cmp edx, 4096 jl lp</pre>	<pre>movdqa xmm0, [eax + edx] paddd xmm0, [ecx + edx] pxor xmm7, xmm7 psubd xmm7, xmm0 movdqa [edi + edx], xmm7 add edx, 16 cmp edx, 4096 jl lp</pre>

Assembly/Compiler Coding Rule 36. (M impact, MH generality) *Break dependences on portions of registers between instructions by operating on 32-bit registers instead of partial registers. For moves, this can be accomplished with 32-bit moves or by using MOVZX.*

On Pentium M processors, the MOVSB and MOVSD instructions both take a single μop , whether they move from a register or memory. On Pentium 4 processors, the

MOVSX takes an additional μop . This is likely to cause less delay than the partial register update problem mentioned above, but the performance gain may vary. If the additional μop is a critical problem, MOVSX can sometimes be used as alternative.

Sometimes sign-extended semantics can be maintained by zero-extending operands. For example, the C code in the following statements does not need sign extension, nor does it need prefixes for operand size overrides:

```
static short INT a, b;
IF (a == b) {
    ...
}
```

Code for comparing these 16-bit operands might be:

```
MOVZW EAX, [a]
MOVZW EBX, [b]
CMP   EAX, EBX
```

These circumstances tend to be common. However, the technique will not work if the compare is for greater than, less than, greater than or equal, and so on, or if the values in eax or ebx are to be used in another operation where sign extension is required.

Assembly/Compiler Coding Rule 37. (M impact, M generality) *Try to use zero extension or operate on 32-bit operands instead of using moves with sign extension.*

The trace cache can be packed more tightly when instructions with operands that can only be represented as 32 bits are not adjacent.

Assembly/Compiler Coding Rule 38. (ML impact, L generality) *Avoid placing instructions that use 32-bit immediates which cannot be encoded as sign-extended 16-bit immediates near each other. Try to schedule μops that have no immediate immediately before or after μops with 32-bit immediates.*

3.5.1.8 Compares

Use TEST when comparing a value in a register with zero. TEST essentially ANDs operands together without writing to a destination register. TEST is preferred over AND because AND produces an extra result register. TEST is better than CMP ..., 0 because the instruction size is smaller.

Use TEST when comparing the result of a logical AND with an immediate constant for equality or inequality if the register is EAX for cases such as:

```
IF (AVAR & 8) { }
```

The TEST instruction can also be used to detect rollover of modulo of a power of 2. For example, the C code:

```
IF ( (AVAR % 16) == 0 ) { }
```

can be implemented using:

```
TEST  EAX, 0x0F
JNZ   AfterIf
```

Using the TEST instruction between the instruction that may modify part of the flag register and the instruction that uses the flag register can also help prevent partial flag register stall.

Assembly/Compiler Coding Rule 39. (ML impact, M generality) Use the TEST instruction instead of AND when the result of the logical AND is not used. This saves μops in execution. Use a TEST of a register with itself instead of a CMP of the register to zero, this saves the need to encode the zero and saves encoding space. Avoid comparing a constant to a memory operand. It is preferable to load the memory operand and compare the constant to a register.

Often a produced value must be compared with zero, and then used in a branch. Because most Intel architecture instructions set the condition codes as part of their execution, the compare instruction may be eliminated. Thus the operation can be tested directly by a JCC instruction. The notable exceptions are MOV and LEA. In these cases, use TEST.

Assembly/Compiler Coding Rule 40. (ML impact, M generality) Eliminate unnecessary compare with zero instructions by using the appropriate conditional jump instruction when the flags are already set by a preceding arithmetic instruction. If necessary, use a TEST instruction instead of a compare. Be certain that any code transformations made do not introduce problems with overflow.

3.5.1.9 Using NOPs

Code generators generate a no-operation (NOP) to align instructions. Examples of NOPs of different lengths in 32-bit mode are shown below:

1-byte: XCHG EAX, EAX

2-byte: 66 NOP

3-byte: LEA REG, 0 (REG) (8-bit displacement)

4-byte: NOP DWORD PTR [EAX + 0] (8-bit displacement)

5-byte: NOP DWORD PTR [EAX + EAX*1 + 0] (8-bit displacement)

6-byte: LEA REG, 0 (REG) (32-bit displacement)

7-byte: NOP DWORD PTR [EAX + 0] (32-bit displacement)

8-byte: NOP DWORD PTR [EAX + EAX*1 + 0] (32-bit displacement)

9-byte: NOP WORD PTR [EAX + EAX*1 + 0] (32-bit displacement)

These are all true NOPs, having no effect on the state of the machine except to advance the EIP. Because NOPs require hardware resources to decode and execute, use the fewest number to achieve the desired padding.

The one byte NOP:[XCHG EAX,EAX] has special hardware support. Although it still consumes a μop and its accompanying resources, the dependence upon the old value of EAX is removed. This μop can be executed at the earliest possible opportunity, reducing the number of outstanding instructions and is the lowest cost NOP.

The other NOPs have no special hardware support. Their input and output registers are interpreted by the hardware. Therefore, a code generator should arrange to use the register containing the oldest value as input, so that the NOP will dispatch and release RS resources at the earliest possible opportunity.

Try to observe the following NOP generation priority:

- Select the smallest number of NOPs and pseudo-NOPs to provide the desired padding.
- Select NOPs that are least likely to execute on slower execution unit clusters.
- Select the register arguments of NOPs to reduce dependencies.

3.5.1.10 Mixing SIMD Data Types

Previous microarchitectures (before Intel Core microarchitecture) do not have explicit restrictions on mixing integer and floating-point (FP) operations on XMM registers. For Intel Core microarchitecture, mixing integer and floating-point operations on the content of an XMM register can degrade performance. Software should avoid mixed-use of integer/FP operation on XMM registers. Specifically,

- Use SIMD integer operations to feed SIMD integer operations. Use PXOR for idiom.
- Use SIMD floating point operations to feed SIMD floating point operations. Use XORPS for idiom.
- When floating point operations are bitwise equivalent, use PS data type instead of PD data type. MOVAPS and MOVAPD do the same thing, but MOVAPS takes one less byte to encode the instruction.

3.5.1.11 Spill Scheduling

The spill scheduling algorithm used by a code generator will be impacted by the memory subsystem. A spill scheduling algorithm is an algorithm that selects what values to spill to memory when there are too many live values to fit in registers. Consider the code in Example 3-22, where it is necessary to spill either A, B, or C.

Example 3-22. Spill Scheduling Code

```
LOOP
  C := ...
  B := ...
  A := A + ...
```

For modern microarchitectures, using dependence depth information in spill scheduling is even more important than in previous processors. The loop-carried dependence in A makes it especially important that A not be spilled. Not only would a store/load be placed in the dependence chain, but there would also be a data-not-ready stall of the load, costing further cycles.

Assembly/Compiler Coding Rule 41. (H impact, MH generality) *For small loops, placing loop invariants in memory is better than spilling loop-carried dependencies.*

A possibly counter-intuitive result is that in such a situation it is better to put loop invariants in memory than in registers, since loop invariants never have a load blocked by store data that is not ready.

3.5.2 Avoiding Stalls in Execution Core

Although the design of the execution core is optimized to make common cases executes quickly. A μop may encounter various hazards, delays, or stalls while making forward progress from the front end to the ROB and RS. The significant cases are:

- ROB Read Port Stalls
- Partial Register Reference Stalls
- Partial Updates to XMM Register Stalls
- Partial Flag Register Reference Stalls

3.5.2.1 ROB Read Port Stalls

As a μop is renamed, it determines whether its source operands have executed and been written to the reorder buffer (ROB), or whether they will be captured “in flight” in the RS or in the bypass network. Typically, the great majority of source operands are found to be “in flight” during renaming. Those that have been written back to the ROB are read through a set of read ports.

Since the Intel Core Microarchitecture is optimized for the common case where the operands are “in flight”, it does not provide a full set of read ports to enable all renamed μops to read all sources from the ROB in the same cycle.

When not all sources can be read, a μop can stall in the rename stage until it can get access to enough ROB read ports to complete renaming the μop . This stall is usually short-lived. Typically, a μop will complete renaming in the next cycle, but it appears to the application as a loss of rename bandwidth.

Some of the software-visible situations that can cause ROB read port stalls include:

- Registers that have become cold and require a ROB read port because execution units are doing other independent calculations.
- Constants inside registers
- Pointer and index registers

In rare cases, ROB read port stalls may lead to more significant performance degradations. There are a couple of heuristics that can help prevent over-subscribing the ROB read ports:

GENERAL OPTIMIZATION GUIDELINES

- Keep common register usage clustered together. Multiple references to the same written-back register can be “folded” inside the out of order execution core.
- Keep short dependency chains intact. This practice ensures that the registers will not have been written back when the new micro-ops are written to the RS.

These two scheduling heuristics may conflict with other more common scheduling heuristics. To reduce demand on the ROB read port, use these two heuristics only if both the following situations are met:

- short latency operations
- indications of actual ROB read port stalls can be confirmed by measurements of the performance event (the relevant event is RAT_STALLS.ROB_READ_PORT, see Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3B*)

If the code has a long dependency chain, these two heuristics should not be used because they can cause the RS to fill, causing damage that outweighs the positive effects of reducing demands on the ROB read port.

In Intel microarchitecture code name Sandy Bridge, there are more read ports available than in previous microarchitectures. Nonetheless, some read port stalls remain and can cause the delay of load and store operations with base+index addressing. When read port stalls occur, it may be due to oversubscription in the integer execution stack. A few techniques can solve the issue:

- Unrolling loops. This reduces the loop terminating operations and relieves read port pressure.
- Use instructions with less general purpose register source operands. For example, changing loads and stores from indexed addressing to base+displacement addressing or ALU ops using one register + immediate rather than 2 registers.

3.5.2.2 Writeback Bus Conflicts

The writeback bus inside the execution engine is a common resource needed to facilitate out-of-order execution of micro-ops in flight. When the writeback bus is needed at the same time by two micro-ops executing in the same stack of execution units (see Table 2-4), the younger micro-op will have to wait for the writeback bus to be available. This situation typically will be more likely for short-latency instructions experience a delay when it might have been otherwise ready for dispatching into the execution engine.

Consider a repeating sequence of independent floating-point ADDs with a single-cycle MOV bound to the same dispatch port. When the MOV finds the dispatch port available, the writeback bus can be occupied by the ADD. This delays the MOV operation.

If this problem is detected, you can sometimes change the instruction selection to use a different dispatch port and reduce the writeback contention.

3.5.2.3 Bypass between Execution Domains

Floating point (FP) loads have an extra cycle of latency. Moves between FP and SIMD stacks have another additional cycle of latency.

Example:

```
ADDPS XMM0, XMM1
PAND XMM0, XMM3
ADDPS XMM2, XMM0
```

The overall latency for the above calculation is 9 cycles:

- 3 cycles for each ADDPS instruction
- 1 cycle for the PAND instruction
- 1 cycle to bypass between the ADDPS floating point domain to the PAND integer domain
- 1 cycle to move the data from the PAND integer to the second floating point ADDPS domain

To avoid this penalty, you should organize code to minimize domain changes. Sometimes you cannot avoid bypasses.

Account for bypass cycles when counting the overall latency of your code. If your calculation is latency-bound, you can execute more instructions in parallel or break dependency chains to reduce total latency.

Code that has many bypass domains and is completely latency-bound may run slower on the Intel Core microarchitecture than it did on previous microarchitectures.

3.5.2.4 Partial Register Stalls

General purpose registers can be accessed in granularities of bytes, words, double-words; 64-bit mode also supports quadword granularity. Referencing a portion of a register is referred to as a partial register reference.

A partial register stall happens when an instruction refers to a register, portions of which were previously modified by other instructions. For example, partial register stalls occurs with a read to AX while previous instructions stored AL and AH, or a read to EAX while previous instruction modified AX.

The delay of a partial register stall is small in processors based on Intel Core and NetBurst microarchitectures, and in Pentium M processor (with CPUID signature family 6, model 13), Intel Core Solo, and Intel Core Duo processors. Pentium M processors (CPUID signature with family 6, model 9) and the P6 family incur a large penalty.

Note that in Intel 64 architecture, an update to the lower 32 bits of a 64 bit integer register is architecturally defined to zero extend the upper 32 bits. While this action may be logically viewed as a 32 bit update, it is really a 64 bit update (and therefore does not cause a partial stall).

Referencing partial registers frequently produces code sequences with either false or real dependencies. Example 3-18 demonstrates a series of false and real dependencies caused by referencing partial registers.

If instructions 4 and 6 (in Example 3-18) are changed to use a movzx instruction instead of a mov, then the dependences of instruction 4 on 2 (and transitively 1 before it), and instruction 6 on 5 are broken. This creates two independent chains of computation instead of one serial one.

Example 3-23 illustrates the use of MOVZX to avoid a partial register stall when packing three byte values into a register.

Example 3-23. Avoiding Partial Register Stalls in Integer Code

A Sequence Causing Partial Register Stall	Alternate Sequence Using MOVZX to Avoid Delay
<pre>mov al, byte ptr a[2] shl eax,16 mov ax, word ptr a movd mm0, eax ret</pre>	<pre>movzx eax, byte ptr a[2] shl eax, 16 movzx ecx, word ptr a or eax,ecx movd mm0, eax ret</pre>

In Intel microarchitecture code name Sandy Bridge, partial register access is handled in hardware by inserting a micro-op that merges the partial register with the full register in the following cases:

- After a write to one of the registers AH, BH, CH or DH and before a following read of the 2-, 4- or 8-byte form of the same register. In these cases a merge micro-op is inserted. The insertion consumes a full allocation cycle in which other micro-ops cannot be allocated.
- After a micro-op with a destination register of 1 or 2 bytes, which is not a source of the instruction (or the register's bigger form), and before a following read of a 2-,4- or 8-byte form of the same register. In these cases the merge micro-op is part of the flow. For example:
 - `MOV AX, [BX]`
 when you want to load from memory to a partial register, consider using MOVZX or MOVSX to avoid the additional merge micro-op penalty.
 - `LEA AX, [BX+CX]`

For optimal performance, use of zero idioms, before the use of the register, eliminates the need for partial register merge micro-ops

3.5.2.5 Partial XMM Register Stalls

Partial register stalls can also apply to XMM registers. The following SSE and SSE2 instructions update only part of the destination register:


```

MOVL/HPD XMM, MEM64
MOVL/HPS XMM, MEM32
MOVSS/SD between registers

```

Using these instructions creates a dependency chain between the unmodified part of the register and the modified part of the register. This dependency chain can cause performance loss.

Example 3-24 illustrates the use of MOVZX to avoid a partial register stall when packing three byte values into a register.

Follow these recommendations to avoid stalls from partial updates to XMM registers:

- Avoid using instructions which update only part of the XMM register.
- If a 64-bit load is needed, use the MOVSD or MOVQ instruction.
- If 2 64-bit loads are required to the same register from non continuous locations, use MOVSD/MOVHPD instead of MOVLPD/MOVHPD.
- When copying the XMM register, use the following instructions for full register copy, even if you only want to copy some of the source register data:

```

MOVAPS
MOVAPD
MOVDQA

```

Example 3-24. Avoiding Partial Register Stalls in SIMD Code

Using movlpd for memory transactions and movsd between register copies Causing Partial Register Stall	Using movsd for memory and movapd between register copies Avoid Delay
<pre> mov edx, x mov ecx, count movlpd xmm3, _1_ movlpd xmm2, _1pt5_ align 16 </pre>	<pre> mov edx, x mov ecx, count movsd xmm3, _1_ movsd xmm2, _1pt5_ align 16 </pre>
<pre> lp: movlpd xmm0, [edx] addsd xmm0, xmm3 movsd xmm1, xmm2 subsd xmm1, [edx] mulsd xmm0, xmm1 movsd [edx], xmm0 add edx, 8 dec ecx jnz lp </pre>	<pre> lp: movsd xmm0, [edx] addsd xmm0, xmm3 movapd xmm1, xmm2 subsd xmm1, [edx] mulsd xmm0, xmm1 movsd [edx], xmm0 add edx, 8 dec ecx jnz lp </pre>

3.5.2.6 Partial Flag Register Stalls

A “partial flag register stall” occurs when an instruction modifies a part of the flag register and the following instruction is dependent on the outcome of the flags. This happens most often with shift instructions (SAR, SAL, SHR, SHL). The flags are not modified in the case of a zero shift count, but the shift count is usually known only at execution time. The front end stalls until the instruction is retired.

Other instructions that can modify some part of the flag register include CMPXCHG8B, various rotate instructions, STC, and STD. An example of assembly with a partial flag register stall and alternative code without the stall is shown in Example 3-25.

In processors based on Intel Core microarchitecture, shift immediate by 1 is handled by special hardware such that it does not experience partial flag stall.

Example 3-25. Avoiding Partial Flag Register Stalls

Partial Flag Register Stall	Avoiding Partial Flag Register Stall
<pre>xor eax, eax mov ecx, a sar ecx, 2 setz al ;SAR can update carry causing a stall</pre>	<pre>or eax, eax mov ecx, a sar ecx, 2 test ecx, ecx ; test always updates all flags setz al ;No partial reg or flag stall,</pre>

In Intel microarchitecture code name Sandy Bridge, the cost of partial flag access is replaced by the insertion of a micro-op instead of a stall. However, it is still recommended to use less of instructions that write only to some of the flags (such as INC, DEC, SET CL) before instructions that can write flags conditionally (such as SHIFT CL).

Example 3-26 compares two techniques to implement the addition of very large integers (e.g. 1024 bits). The alternative sequence on the right side of Example 3-26 will be faster than the left side on Intel microarchitecture code name Sandy Bridge, but it will experience partial flag stalls on prior microarchitectures.

Example 3-26. Partial Flag Register Accesses in Intel microarchitecture code name Sandy Bridge

Save partial flag register to avoid stall	Simplified code sequence
<pre> lea rsi, [A] lea rdi, [B] xor rax, rax mov rcx, 16 ; 16*64 = 1024 bit lp_64bit: add rax, [rsi] adc rax, [rdi] mov [rdi], rax setc al ;save carry for next iteration movzx rax, al add rsi, 8 add rdi, 8 dec rcx jnz lp_64bit </pre>	<pre> lea rsi, [A] lea rdi, [B] xor rax, rax mov rcx, 16 lp_64bit: add rax, [rsi] adc rax, [rdi] mov [rdi], rax lea rsi, [rsi+8] lea rdi, [rdi+8] dec rcx jnz lp_64bit </pre>

3.5.2.7 Floating Point/SIMD Operands

Moves that write a portion of a register can introduce unwanted dependences. The MOVSD REG, REG instruction writes only the bottom 64 bits of a register, not all 128 bits. This introduces a dependence on the preceding instruction that produces the upper 64 bits (even if those bits are not longer wanted). The dependence inhibits register renaming, and thereby reduces parallelism.

Use MOVAPD as an alternative; it writes all 128 bits. Even though this instruction has a longer latency, the μ ops for MOVAPD use a different execution port and this port is more likely to be free. The change can impact performance. There may be exceptional cases where the latency matters more than the dependence or the execution port.

Assembly/Compiler Coding Rule 42. (M impact, ML generality) *Avoid introducing dependences with partial floating point register writes, e.g. from the MOVSD XMMREG1, XMMREG2 instruction. Use the MOVAPD XMMREG1, XMMREG2 instruction instead.*

The MOVSD XMMREG, MEM instruction writes all 128 bits and breaks a dependence.

The MOVUPD from memory instruction performs two 64-bit loads, but requires additional μ ops to adjust the address and combine the loads into a single register. This same functionality can be obtained using MOVSD XMMREG1, MEM; MOVSD XMMREG2, MEM+8; UNPCKLPD XMMREG1, XMMREG2, which uses fewer μ ops and

can be packed into the trace cache more effectively. The latter alternative has been found to provide a several percent performance improvement in some cases. Its encoding requires more instruction bytes, but this is seldom an issue for the Pentium 4 processor. The store version of MOVUPD is complex and slow, so much so that the sequence with two MOVSD and a UNPCKHPD should always be used.

Assembly/Compiler Coding Rule 43. (ML impact, L generality) *Instead of using MOVUPD XMMREG1, MEM for a unaligned 128-bit load, use MOVSD XMMREG1, MEM; MOVSD XMMREG2, MEM+8; UNPCKLPD XMMREG1, XMMREG2. If the additional register is not available, then use MOVSD XMMREG1, MEM; MOVHPD XMMREG1, MEM+8.*

Assembly/Compiler Coding Rule 44. (M impact, ML generality) *Instead of using MOVUPD MEM, XMMREG1 for a store, use MOVSD MEM, XMMREG1; UNPCKHPD XMMREG1, XMMREG1; MOVSD MEM+8, XMMREG1 instead.*

3.5.3 Vectorization

This section provides a brief summary of optimization issues related to vectorization. There is more detail in the chapters that follow.

Vectorization is a program transformation that allows special hardware to perform the same operation on multiple data elements at the same time. Successive processor generations have provided vector support through the MMX technology, Streaming SIMD Extensions (SSE), Streaming SIMD Extensions 2 (SSE2), Streaming SIMD Extensions 3 (SSE3) and Supplemental Streaming SIMD Extensions 3 (SSSE3).

Vectorization is a special case of SIMD, a term defined in Flynn's architecture taxonomy to denote a single instruction stream capable of operating on multiple data elements in parallel. The number of elements which can be operated on in parallel range from four single-precision floating point data elements in Streaming SIMD Extensions and two double-precision floating-point data elements in Streaming SIMD Extensions 2 to sixteen byte operations in a 128-bit register in Streaming SIMD Extensions 2. Thus, vector length ranges from 2 to 16, depending on the instruction extensions used and on the data type.

The Intel C++ Compiler supports vectorization in three ways:

- The compiler may be able to generate SIMD code without intervention from the user.
- The can user insert pragmas to help the compiler realize that it can vectorize the code.
- The user can write SIMD code explicitly using intrinsics and C++ classes.

To help enable the compiler to generate SIMD code, avoid global pointers and global variables. These issues may be less troublesome if all modules are compiled simultaneously, and whole-program optimization is used.

User/Source Coding Rule 2. (H impact, M generality) Use the smallest possible floating-point or SIMD data type, to enable more parallelism with the use of a (longer) SIMD vector. For example, use single precision instead of double precision where possible..

User/Source Coding Rule 3. (M impact, ML generality) Arrange the nesting of loops so that the innermost nesting level is free of inter-iteration dependencies. Especially avoid the case where the store of data in an earlier iteration happens lexically after the load of that data in a future iteration, something which is called a lexically backward dependence..

The integer part of the SIMD instruction set extensions cover 8-bit,16-bit and 32-bit operands. Not all SIMD operations are supported for 32 bits, meaning that some source code will not be able to be vectorized at all unless smaller operands are used.

User/Source Coding Rule 4. (M impact, ML generality) Avoid the use of conditional branches inside loops and consider using SSE instructions to eliminate branches.

User/Source Coding Rule 5. (M impact, ML generality) Keep induction (loop) variable expressions simple.

3.5.4 Optimization of Partially Vectorizable Code

Frequently, a program contains a mixture of vectorizable code and some routines that are non-vectorizable. A common situation of partially vectorizable code involves a loop structure which include mixtures of vectorized code and unvectorizable code. This situation is depicted in Figure 3-1.

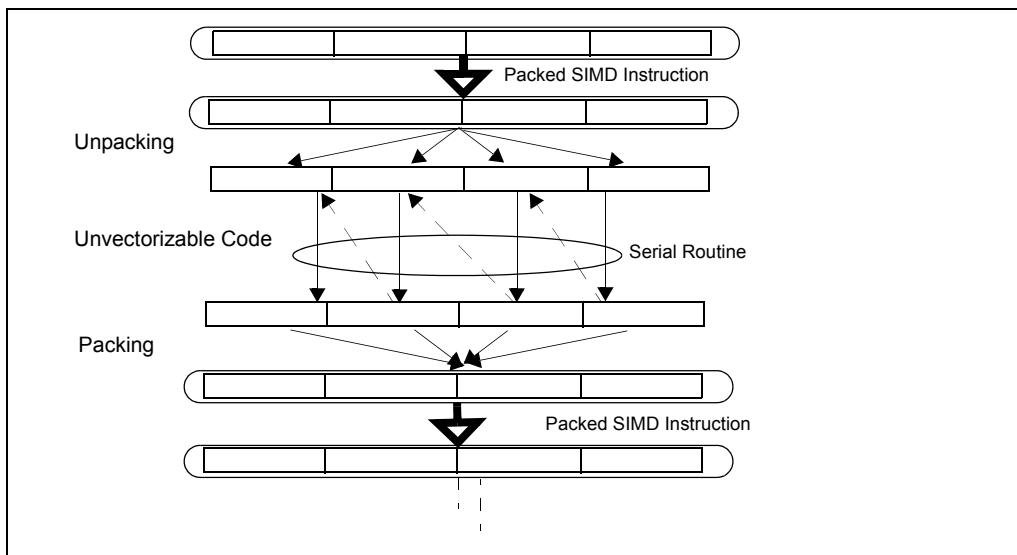


Figure 3-1. Generic Program Flow of Partially Vectorized Code

GENERAL OPTIMIZATION GUIDELINES

It generally consists of five stages within the loop:

- Prolog
- Unpacking vectorized data structure into individual elements
- Calling a non-vectorizable routine to process each element serially
- Packing individual result into vectorized data structure
- Epilog

This section discusses techniques that can reduce the cost and bottleneck associated with the packing/unpacking stages in these partially vectorize code.

Example 3-27 shows a reference code template that is representative of partially vectorizable coding situations that also experience performance issues. The unvectorizable portion of code is represented generically by a sequence of calling a serial function named "foo" multiple times. This generic example is referred to as "shuffle with store forwarding", because the problem generally involves an unpacking stage that shuffles data elements between register and memory, followed by a packing stage that can experience store forwarding issue.

There are more than one useful techniques that can reduce the store-forwarding bottleneck between the serialized portion and the packing stage. The following subsections presents alternate techniques to deal with the packing, unpacking, and parameter passing to serialized function calls.

Example 3-27. Reference Code Template for Partially Vectorizable Program

```
// Prolog //////////////////////////////////////  
push ebp  
mov ebp, esp  
  
// Unpacking //////////////////////////////////////  
sub ebp, 32  
and ebp, 0xffffffff  
movaps [ebp], xmm0
```

Example 3-27. Reference Code Template for Partially Vectorizable Program (Contd.)

```

// Serial operations on components //
sub ebp, 4

mov eax, [ebp+4]
mov [ebp], eax
call foo
mov [ebp+16+4], eax

mov eax, [ebp+8]
mov [ebp], eax
call foo
mov [ebp+16+4+4], eax

mov eax, [ebp+12]
mov [ebp], eax
call foo
mov [ebp+16+8+4], eax

mov eax, [ebp+12+4]
mov [ebp], eax
call foo
mov [ebp+16+12+4], eax

// Packing //
movaps xmm0, [ebp+16+4]

// Epilog //
pop ebp
ret

```

3.5.4.1 Alternate Packing Techniques

The packing method implemented in the reference code of Example 3-27 will experience delay as it assembles 4 doubleword result from memory into an XMM register due to store-forwarding restrictions.

Three alternate techniques for packing, using different SIMD instruction to assemble contents in XMM registers are shown in Example 3-28. All three techniques avoid

store-forwarding delay by satisfying the restrictions on data sizes between a preceding store and subsequent load operations.

Example 3-28. Three Alternate Packing Methods for Avoiding Store Forwarding Difficulty

Packing Method 1	Packing Method 2	Packing Method 3
movd xmm0, [ebp+16+4] movd xmm1, [ebp+16+8] movd xmm2, [ebp+16+12] movd xmm3, [ebp+12+16+4] punpckldq xmm0, xmm1 punpckldq xmm2, xmm3 punpckldq xmm0, xmm2	movd xmm0, [ebp+16+4] movd xmm1, [ebp+16+8] movd xmm2, [ebp+16+12] movd xmm3, [ebp+12+16+4] psllq xmm3, 32 orps xmm2, xmm3 psllq xmm1, 32 orps xmm0, xmm1 movlhps xmm0, xmm2 xmm0, xmm2	movd xmm0, [ebp+16+4] movd xmm1, [ebp+16+8] movd xmm2, [ebp+16+12] movd xmm3, [ebp+12+16+4] movlhps xmm1, xmm3 psllq xmm1, 32 movlhps xmm0, xmm2 orps xmm0, xmm1

3.5.4.2 Simplifying Result Passing

In Example 3-27, individual results were passed to the packing stage by storing to contiguous memory locations. Instead of using memory spills to pass four results, result passing may be accomplished by using either one or more registers. Using registers to simplify result passing and reduce memory spills can improve performance by varying degrees depending on the register pressure at runtime.

Example 3-29 shows the coding sequence that uses four extra XMM registers to reduce all memory spills of passing results back to the parent routine. However, software must observe the following conditions when using this technique:

- There is no register shortage.
- If the loop does not have many stores or loads but has many computations, this technique does not help performance. This technique adds work to the computational units, while the store and loads ports are idle.

Example 3-29. Using Four Registers to Reduce Memory Spills and Simplify Result Passing

```

mov eax, [ebp+4]
mov [ebp], eax
call foo
movd xmm0, eax

mov eax, [ebp+8]
mov [ebp], eax
call foo
movd xmm1, eax
    
```


Example 3-29. Using Four Registers to Reduce Memory Spills and Simplify Result Passing

```

mov eax, [ebp+12]
mov [ebp], eax
call foo
movd xmm2, eax

mov eax, [ebp+12+4]
mov [ebp], eax
call foo
movd xmm3, eax

```

3.5.4.3 Stack Optimization

In Example 3-27, an input parameter was copied in turn onto the stack and passed to the non-vectorizable routine for processing. The parameter passing from consecutive memory locations can be simplified by a technique shown in Example 3-30.

Example 3-30. Stack Optimization Technique to Simplify Parameter Passing

```

call foo
mov [ebp+16], eax

add ebp, 4
call foo
mov [ebp+16], eax

add ebp, 4
call foo
mov [ebp+16], eax

add ebp, 4
call foo

```

Stack Optimization can only be used when:

- The serial operations are function calls. The function “foo” is declared as: INT FOO(INT A). The parameter is passed on the stack.
- The order of operation on the components is from last to first.

Note the call to FOO and the advance of EBP when passing the vector elements to FOO one by one from last to first.

3.5.4.4 Tuning Considerations

Tuning considerations for situations represented by looping of Example 3-27 include

- Applying one of more of the following combinations:
 - choose an alternate packing technique
 - consider a technique to simply result-passing
 - consider the stack optimization technique to simplify parameter passing
- Minimizing the average number of cycles to execute one iteration of the loop
- Minimizing the per-iteration cost of the unpacking and packing operations

The speed improvement by using the techniques discussed in this section will vary, depending on the choice of combinations implemented and characteristics of the non-vectorizable routine. For example, if the routine "foo" is short (representative of tight, short loops), the per-iteration cost of unpacking/packing tend to be smaller than situations where the non-vectorizable code contain longer operation or many dependencies. This is because many iterations of short, tight loop can be in flight in the execution core, so the per-iteration cost of packing and unpacking is only partially exposed and appear to cause very little performance degradation.

Evaluation of the per-iteration cost of packing/unpacking should be carried out in a methodical manner over a selected number of test cases, where each case may implement some combination of the techniques discussed in this section. The per-iteration cost can be estimated by:

- evaluating the average cycles to execute one iteration of the test case
- evaluating the average cycles to execute one iteration of a base line loop sequence of non-vectorizable code

Example 3-31 shows the base line code sequence that can be used to estimate the average cost of a loop that executes non-vectorizable routines.

Example 3-31. Base Line Code Sequence to Estimate Loop Overhead

```
push ebp
mov ebp, esp
sub ebp, 4

mov [ebp], edi
call foo
```

Example 3-31. Base Line Code Sequence to Estimate Loop Overhead (Contd.)

```

mov [ebp], edi
call foo

mov [ebp], edi
call foo

mov [ebp], edi
call foo

add ebp, 4
pop ebp
ret

```

The average per-iteration cost of packing/unpacking can be derived from measuring the execution times of a large number of iterations by:

$$\frac{((\text{Cycles to run TestCase}) - (\text{Cycles to run equivalent baseline sequence}))}{(\text{Iteration count})}$$

For example, using a simple function that returns an input parameter (representative of tight, short loops), the per-iteration cost of packing/unpacking may range from slightly more than 7 cycles (the shuffle with store forwarding case, Example 3-27) to ~0.9 cycles (accomplished by several test cases). Across 27 test cases (consisting of one of the alternate packing methods, no result-simplification/simplification of either 1 or 4 results, no stack optimization or with stack optimization), the average per-iteration cost of packing/unpacking is about 1.7 cycles.

Generally speaking, packing method 2 and 3 (see Example 3-28) tend to be more robust than packing method 1; the optimal choice of simplifying 1 or 4 results will be affected by register pressure of the runtime and other relevant microarchitectural conditions.

Note that the numeric discussion of per-iteration cost of packing/packing is illustrative only. It will vary with test cases using a different base line code sequence and will generally increase if the non-vectorizable routine requires longer time to execute because the number of loop iterations that can reside in flight in the execution core decreases.

3.6 OPTIMIZING MEMORY ACCESSES

This section discusses guidelines for optimizing code and data memory accesses. The most important recommendations are:

- Execute load and store operations within available execution bandwidth.
- Enable forward progress of speculative execution.

GENERAL OPTIMIZATION GUIDELINES

- Enable store forwarding to proceed.
- Align data, paying attention to data layout and stack alignment.
- Place code and data on separate pages.
- Enhance data locality.
- Use prefetching and cacheability control instructions.
- Enhance code locality and align branch targets.
- Take advantage of write combining.

Alignment and forwarding problems are among the most common sources of large delays on processors based on Intel NetBurst microarchitecture.

3.6.1 Load and Store Execution Bandwidth

Typically, loads and stores are the most frequent operations in a workload, up to 40% of the instructions in a workload carrying load or store intent are not uncommon. Each generation of microarchitecture provides multiple buffers to support executing load and store operations while there are instructions in flight.

Software can maximize memory performance by not exceeding the issue or buffering limitations of the machine. In the Intel Core microarchitecture, only 20 stores and 32 loads may be in flight at once. In Intel microarchitecture code name Nehalem, there are 32 store buffers and 48 load buffers. Since only one load can issue per cycle, algorithms which operate on two arrays are constrained to one operation every other cycle unless you use programming tricks to reduce the amount of memory usage.

Intel Core Duo and Intel Core Solo processors have less buffers. Nevertheless the general heuristic applies to all of them.

3.6.1.1 Make Use of Load Bandwidth in Intel® microarchitecture code name Sandy Bridge

While prior microarchitecture has one load port (port 2), Intel microarchitecture code name Sandy Bridge can load from port 2 and port 3. Thus two load operations can be performed every cycle and doubling the load throughput of the code. This improves code that reads a lot of data and does not need to write out results to memory very often (Port 3 also handles store-address operation). To exploit this bandwidth, the data has to stay in the L1 data cache or it should be accessed sequentially, enabling the hardware prefetchers to bring the data to the L1 data cache in time.

Consider the following C code example of adding all the elements of an array:

```
int buff[BUFF_SIZE];  
  
int sum = 0;  
  
for (i=0;i<BUFF_SIZE;i++){
```

```

sum += buff[i];
}

```

Alternative 1 is the assembly code generated by the Intel compiler for this C code, using the optimization flag for Intel microarchitecture code name Nehalem. The compiler vectorizes execution using Intel SSE instructions. In this code, each ADD operation uses the result of the previous ADD operation. This limits the throughput to one load and ADD operation per cycle. Alternative 2 is optimized for Intel microarchitecture code name Sandy Bridge by enabling it to use the additional load bandwidth. The code removes the dependency among ADD operations, by using two registers to sum the array values. Two load and two ADD operations can be executed every cycle.

Example 3-32. Optimize for Load Port Bandwidth in Intel microarchitecture code name Sandy Bridge

Register dependency inhibits PADD execution	Reduce register dependency allow two load port to supply PADD execution
<pre> xor eax, eax pxor xmm0, xmm0 lea rsi, buff loop_start: padd xmm0, [rsi+4*rax] padd xmm0, [rsi+4*rax+16] padd xmm0, [rsi+4*rax+32] padd xmm0, [rsi+4*rax+48] padd xmm0, [rsi+4*rax+64] padd xmm0, [rsi+4*rax+80] padd xmm0, [rsi+4*rax+96] padd xmm0, [rsi+4*rax+112] add eax, 32 cmp eax, BUFF_SIZE jl loop_start </pre>	<pre> xor eax, eax pxor xmm0, xmm0 pxor xmm1, xmm1 lea rsi, buff loop_start: padd xmm0, [rsi+4*rax] padd xmm1, [rsi+4*rax+16] padd xmm0, [rsi+4*rax+32] padd xmm1, [rsi+4*rax+48] padd xmm0, [rsi+4*rax+64] padd xmm1, [rsi+4*rax+80] padd xmm0, [rsi+4*rax+96] padd xmm1, [rsi+4*rax+112] add eax, 32 cmp eax, BUFF_SIZE jl loop_start </pre>
<pre> sum_partials: movdqa xmm1, xmm0 psrldq xmm1, 8 padd xmm0, xmm1 movdqa xmm2, xmm0 psrldq xmm2, 4 padd xmm0, xmm2 mov [sum], xmm0 </pre>	<pre> sum_partials: padd xmm0, xmm1 movdqa xmm1, xmm0 psrldq xmm1, 8 padd xmm0, xmm1 movdqa xmm2, xmm0 psrldq xmm2, 4 padd xmm0, xmm2 mov [sum], xmm0 </pre>

3.6.1.2 L1D Cache Latency in Intel® microarchitecture code name Sandy Bridge

Load latency from L1D cache may vary (see Table 2-8). The best case is 4 cycles, which apply to load operations to general purpose registers using

- one register or
- a base register plus an offset that is smaller than 2048.

Consider the pointer-chasing code example in Example 3-33.

Example 3-33. Index versus Pointers in Pointer-Chasing Code

Traversing through indexes	Traversing through pointers
// C code example index = buffer.m_buff[index].next_index;	// C code example node = node->pNext;
// ASM example loop: shl rbx, 6 mov rbx, 0x20(rbx+rcx) dec rax cmp rax, -1 jne loop	// ASM example loop: mov rdx, [rdx] dec rax cmp rax, -1 jne loop

The left side implements pointer chasing via traversing an index. Compiler then generates the code shown below addressing memory using base+index with an offset. The right side shows compiler generated code from pointer de-referencing code and uses only a base register.

The code on the right side is faster than the left side across Intel microarchitecture code name Sandy Bridge and prior microarchitecture. However the code that traverses index will be slower on Intel microarchitecture code name Sandy Bridge relative to prior microarchitecture.

3.6.1.3 Handling L1D Cache Bank Conflict

In Intel microarchitecture code name Sandy Bridge, the internal organization of the L1D cache may manifest a situation when two load micro-ops whose addresses have a bank conflict. When a bank conflict is present between two load operations, the more recent one will be delayed until the conflict is resolved. A bank conflict happens when two simultaneous load operations have the same bit 2-5 of their linear address but they are not from the same set in the cache (bits 6 - 12).

Bank conflicts should be handled only if the code is bound by load bandwidth. Some bank conflicts do not cause any performance degradation since they are hidden by

other performance limiters. Eliminating such bank conflicts does not improve performance.

The following example demonstrates bank conflict and how to modify the code and avoid them. It uses two source arrays with a size that is a multiple of cache line size. When loading an element from A and the counterpart element from B the elements have the same offset in their cache lines and therefore a bank conflict may happen.

Example 3-34. Example of Bank Conflicts in L1D Cache and Remedy

<pre>int A[128]; int B[128]; int C[128]; for (i=0;i<128;i+=4){ C[i]=A[i]+B[i]; the loads from A[i] and B[i] collide C[i+1]=A[i+1]+B[i+1]; C[i+2]=A[i+2]+B[i+2]; C[i+3]=A[i+3]+B[i+3]; }</pre>	
<pre>// Code with Bank Conflicts xor rcx, rcx lea r11, A lea r12, B lea r13, C loop: lea esi, [rcx*4] movsxd rsi, esi</pre>	<pre>// Code without Bank Conflicts xor rcx, rcx lea r11, A lea r12, B lea r13, C loop: lea esi, [rcx*4] movsxd rsi, esi</pre>

Example 3-34. Example of Bank Conflicts in L1D Cache and Remedy

<pre> mov edi, [r11+rsi*4] add edi, [r12+rsi*4] mov r8d, [r11+rsi*4+4] add r8d, [r12+rsi*4+4] mov r9d, [r11+rsi*4+8] add r9d, [r12+rsi*4+8] mov r10d, [r11+rsi*4+12] add r10d, [r12+rsi*4+12] mov [r13+rsi*4], edi inc ecx mov [r13+rsi*4+4], r8d mov [r13+rsi*4+8], r9d mov [r13+rsi*4+12], r10d cmp ecx, LEN jb loop </pre>	<pre> mov edi, [r11+rsi*4] mov r8d, [r11+rsi*4+4] add edi, [r12+rsi*4] add r8d, [r12+rsi*4+4] mov r9d, [r11+rsi*4+8] mov r10d, [r11+rsi*4+12] add r9d, [r12+rsi*4+8] add r10d, [r12+rsi*4+12] inc ecx mov [r13+rsi*4], edi mov [r13+rsi*4+4], r8d mov [r13+rsi*4+8], r9d mov [r13+rsi*4+12], r10d cmp ecx, LEN jb loop </pre>
---	---

3.6.2 Minimize Register Spills

When a piece of code has more live variables than the processor can keep in general purpose registers, a common method is to hold some of the variables in memory. This method is called register spill. The effect of L1D cache latency can negatively affect the performance of this code. The effect can be more pronounced if the address of register spills uses the slower addressing modes.

One option is to spill general purpose registers to XMM registers. This method is likely to improve performance also on previous processor generations. The following example shows how to spill a register to an XMM register rather than to memory.

Example 3-35. Using XMM Register in Lieu of Memory for Register Spills

Register spills into memory	Register spills into XMM
<pre> loop: mov rdx, [rsp+0x18] movdqa xmm0, [rdx] movdqa xmm1, [rsp+0x20] pcmpeqd xmm1, xmm0 pmovmskb eax, xmm1 test eax, eax jne end_loop movzx rcx, [rbx+0x60] </pre>	<pre> movq xmm4, [rsp+0x18] mov rcx, 0x10 movq xmm5, rcx loop: movq rdx, xmm4 movdqa xmm0, [rdx] movdqa xmm1, [rsp+0x20] pcmpeqd xmm1, xmm0 pmovmskb eax, xmm1 test eax, eax jne end_loop movzx rcx, [rbx+0x60] </pre>
<pre> add qword ptr[rsp+0x18], 0x10 add rdi, 0x4 movzx rdx, di sub rcx, 0x4 add rsi, 0x1d0 cmp rdx, rcx jle loop </pre>	<pre> padd xmm4, xmm5 add rdi, 0x4 movzx rdx, di sub rcx, 0x4 add rsi, 0x1d0 cmp rdx, rcx jle loop </pre>

3.6.3 Enhance Speculative Execution and Memory Disambiguation

Prior to Intel Core microarchitecture, when code contains both stores and loads, the loads cannot be issued before the address of the store is resolved. This rule ensures correct handling of load dependencies on preceding stores.

The Intel Core microarchitecture contains a mechanism that allows some loads to be issued early speculatively. The processor later checks if the load address overlaps with a store. If the addresses do overlap, then the processor re-executes the instructions.

Example 3-36 illustrates a situation that the compiler cannot be sure that “Ptr->Array” does not change during the loop. Therefore, the compiler cannot keep “Ptr->Array” in a register as an invariant and must read it again in every iteration. Although this situation can be fixed in software by a rewriting the code to require the address of the pointer is invariant, memory disambiguation provides performance gain without rewriting the code.

Example 3-36. Loads Blocked by Stores of Unknown Address

C code	Assembly sequence
<pre> struct AA { AA ** array; }; void nullify_array (AA *Ptr, DWORD Index, AA *ThisPtr) { while (Ptr->Array[--Index] != ThisPtr) { Ptr->Array[Index] = NULL ; }; }; </pre>	<pre> nullify_loop: mov dword ptr [eax], 0 mov edx, dword ptr [edi] sub ecx, 4 cmp dword ptr [ecx+edx], esi lea eax, [ecx+edx] jne nullify_loop </pre>

3.6.4 Alignment

Alignment of data concerns all kinds of variables:

- Dynamically allocated variables
- Members of a data structure
- Global or local variables
- Parameters passed on the stack

Misaligned data access can incur significant performance penalties. This is particularly true for cache line splits. The size of a cache line is 64 bytes in the Pentium 4 and other recent Intel processors, including processors based on Intel Core microarchitecture.

An access to data unaligned on 64-byte boundary leads to two memory accesses and requires several μ ops to be executed (instead of one). Accesses that span 64-byte boundaries are likely to incur a large performance penalty, the cost of each stall generally are greater on machines with longer pipelines.

Double-precision floating-point operands that are eight-byte aligned have better performance than operands that are not eight-byte aligned, since they are less likely to incur penalties for cache and MOB splits. Floating-point operation on a memory operands require that the operand be loaded from memory. This incurs an additional μ op, which can have a minor negative impact on front end bandwidth. Additionally, memory operands may cause a data cache miss, causing a penalty.

Assembly/Compiler Coding Rule 45. (H impact, H generality) *Align data on natural operand size address boundaries. If the data will be accessed with vector instruction loads and stores, align the data on 16-byte boundaries.*

For best performance, align data as follows:

- Align 8-bit data at any address.
- Align 16-bit data to be contained within an aligned 4-byte word.
- Align 32-bit data so that its base address is a multiple of four.
- Align 64-bit data so that its base address is a multiple of eight.
- Align 80-bit data so that its base address is a multiple of sixteen.
- Align 128-bit data so that its base address is a multiple of sixteen.

A 64-byte or greater data structure or array should be aligned so that its base address is a multiple of 64. Sorting data in decreasing size order is one heuristic for assisting with natural alignment. As long as 16-byte boundaries (and cache lines) are never crossed, natural alignment is not strictly necessary (though it is an easy way to enforce this).

Example 3-37 shows the type of code that can cause a cache line split. The code loads the addresses of two DWORD arrays. 029E70FEH is not a 4-byte-aligned address, so a 4-byte access at this address will get 2 bytes from the cache line this address is contained in, and 2 bytes from the cache line that starts at 029E700H. On processors with 64-byte cache lines, a similar cache line split will occur every 8 iterations.

Example 3-37. Code That Causes Cache Line Split

```

mov     esi, 029e70feh
mov     edi, 05be5260h
Blockmove:
mov     eax, DWORD PTR [esi]
mov     ebx, DWORD PTR [esi+4]
mov     DWORD PTR [edi], eax
mov     DWORD PTR [edi+4], ebx
add     esi, 8
add     edi, 8
sub     edx, 1
jnz    Blockmove

```

Figure 3-2 illustrates the situation of accessing a data element that span across cache line boundaries.

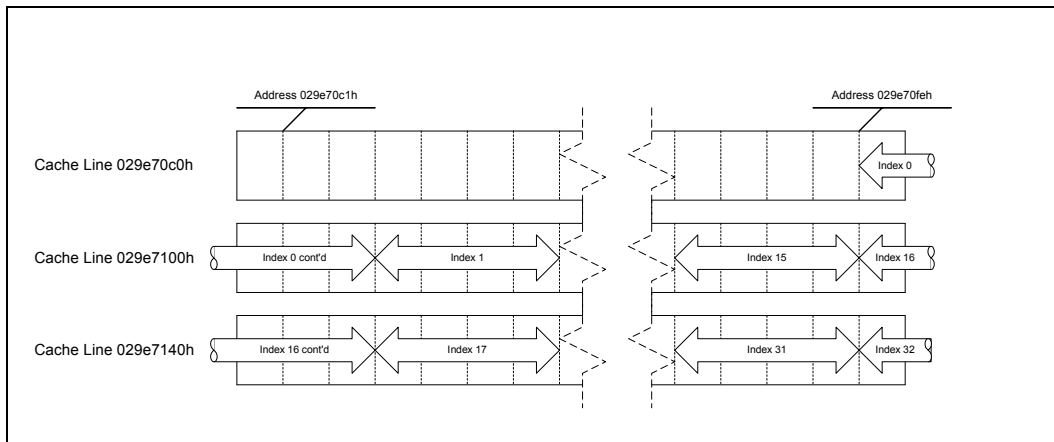


Figure 3-2. Cache Line Split in Accessing Elements in a Array

Alignment of code is less important for processors based on Intel NetBurst microarchitecture. Alignment of branch targets to maximize bandwidth of fetching cached instructions is an issue only when not executing out of the trace cache.

Alignment of code can be an issue for the Pentium M, Intel Core Duo and Intel Core 2 Duo processors. Alignment of branch targets will improve decoder throughput.

3.6.5 Store Forwarding

The processor’s memory system only sends stores to memory (including cache) after store retirement. However, store data can be forwarded from a store to a subsequent load from the same address to give a much shorter store-load latency.

There are two kinds of requirements for store forwarding. If these requirements are violated, store forwarding cannot occur and the load must get its data from the cache (so the store must write its data back to the cache first). This incurs a penalty that is largely related to pipeline depth of the underlying micro-architecture.

The first requirement pertains to the size and alignment of the store-forwarding data. This restriction is likely to have high impact on overall application performance. Typically, a performance penalty due to violating this restriction can be prevented. The store-to-load forwarding restrictions vary from one microarchitecture to another. Several examples of coding pitfalls that cause store-forwarding stalls and solutions to these pitfalls are discussed in detail in Section 3.6.5.1, “Store-to-Load-Forwarding Restriction on Size and Alignment.” The second requirement is the availability of data, discussed in Section 3.6.5.2, “Store-forwarding Restriction on Data Availability.” A good practice is to eliminate redundant load operations.

It may be possible to keep a temporary scalar variable in a register and never write it to memory. Generally, such a variable must not be accessible using indirect pointers. Moving a variable to a register eliminates all loads and stores of that variable and eliminates potential problems associated with store forwarding. However, it also increases register pressure.

Load instructions tend to start chains of computation. Since the out-of-order engine is based on data dependence, load instructions play a significant role in the engine's ability to execute at a high rate. Eliminating loads should be given a high priority.

If a variable does not change between the time when it is stored and the time when it is used again, the register that was stored can be copied or used directly. If register pressure is too high, or an unseen function is called before the store and the second load, it may not be possible to eliminate the second load.

Assembly/Compiler Coding Rule 46. (H impact, M generality) Pass parameters in registers instead of on the stack where possible. Passing arguments on the stack requires a store followed by a reload. While this sequence is optimized in hardware by providing the value to the load directly from the memory order buffer without the need to access the data cache if permitted by store-forwarding restrictions, floating point values incur a significant latency in forwarding. Passing floating point arguments in (preferably XMM) registers should save this long latency operation.

Parameter passing conventions may limit the choice of which parameters are passed in registers which are passed on the stack. However, these limitations may be overcome if the compiler has control of the compilation of the whole binary (using whole-program optimization).

3.6.5.1 Store-to-Load-Forwarding Restriction on Size and Alignment

Data size and alignment restrictions for store-forwarding apply to processors based on Intel NetBurst microarchitecture, Intel Core microarchitecture, Intel Core 2 Duo, Intel Core Solo and Pentium M processors. The performance penalty for violating store-forwarding restrictions is less for shorter-pipelined machines than for Intel NetBurst microarchitecture.

Store-forwarding restrictions vary with each microarchitecture. Intel NetBurst microarchitecture places more constraints than Intel Core microarchitecture on code generation to enable store-forwarding to make progress instead of experiencing stalls. Fixing store-forwarding problems for Intel NetBurst microarchitecture generally also avoids problems on Pentium M, Intel Core Duo and Intel Core 2 Duo processors. The size and alignment restrictions for store-forwarding in processors based on Intel NetBurst microarchitecture are illustrated in Figure 3-3.

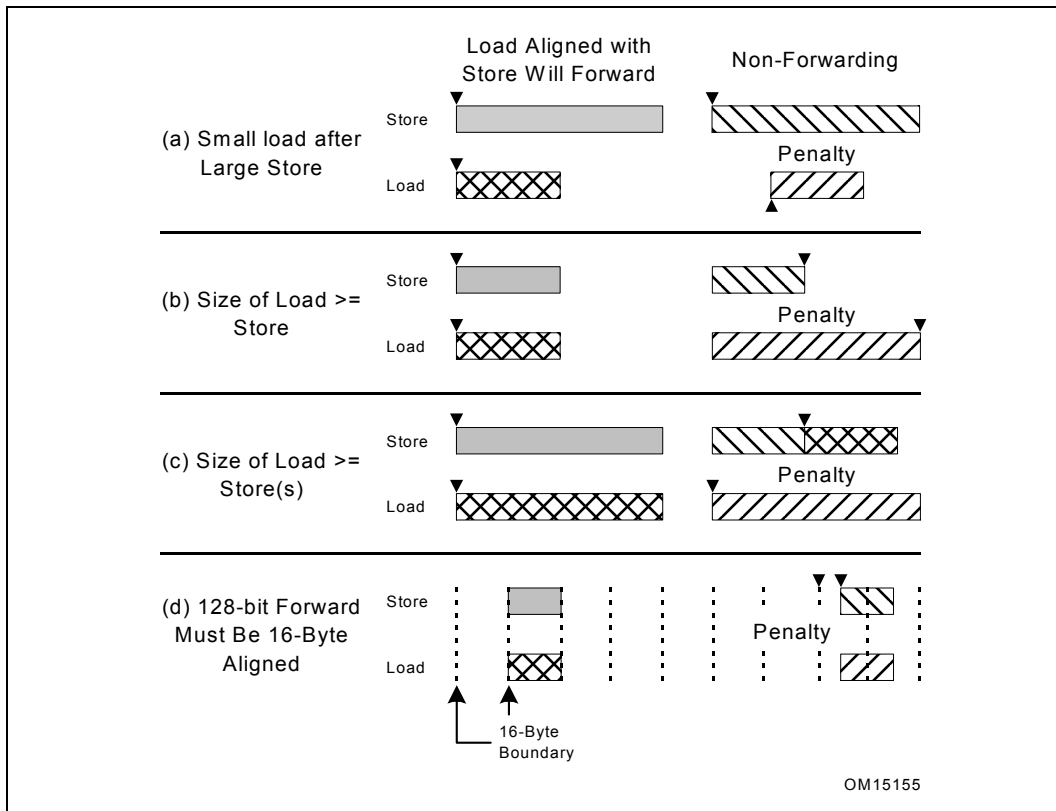


Figure 3-3. Size and Alignment Restrictions in Store Forwarding

The following rules help satisfy size and alignment restrictions for store forwarding:

Assembly/Compiler Coding Rule 47. (H impact, M generality) A load that forwards from a store must have the same address start point and therefore the same alignment as the store data.

Assembly/Compiler Coding Rule 48. (H impact, M generality) The data of a load which is forwarded from a store must be completely contained within the store data.

A load that forwards from a store must wait for the store's data to be written to the store buffer before proceeding, but other, unrelated loads need not wait.

Assembly/Compiler Coding Rule 49. (H impact, ML generality) *If it is necessary to extract a non-aligned portion of stored data, read out the smallest aligned portion that completely contains the data and shift/mask the data as necessary. This is better than incurring the penalties of a failed store-forward.*

Assembly/Compiler Coding Rule 50. (MH impact, ML generality) *Avoid several small loads after large stores to the same area of memory by using a single large read and register copies as needed.*

Example 3-38 depicts several store-forwarding situations in which small loads follow large stores. The first three load operations illustrate the situations described in Rule 50. However, the last load operation gets data from store-forwarding without problem.

Example 3-38. Situations Showing Small Loads After Large Store

```

mov [EBP], 'abcd'
mov AL, [EBP]      ; Not blocked - same alignment
mov BL, [EBP + 1]  ; Blocked
mov CL, [EBP + 2]  ; Blocked
mov DL, [EBP + 3]  ; Blocked
mov AL, [EBP]      ; Not blocked - same alignment
                  ; n.b. passes older blocked loads

```

Example 3-39 illustrates a store-forwarding situation in which a large load follows several small stores. The data needed by the load operation cannot be forwarded because all of the data that needs to be forwarded is not contained in the store buffer. Avoid large loads after small stores to the same area of memory.

Example 3-39. Non-forwarding Example of Large Load After Small Store

```

mov [EBP], 'a'
mov [EBP + 1], 'b'
mov [EBP + 2], 'c'
mov [EBP + 3], 'd'
mov EAX, [EBP] ; Blocked
                ; The first 4 small store can be consolidated into
                ; a single DWORD store to prevent this non-forwarding
                ; situation.

```

Example 3-40 illustrates a stalled store-forwarding situation that may appear in compiler generated code. Sometimes a compiler generates code similar to that

GENERAL OPTIMIZATION GUIDELINES

shown in Example 3-40 to handle a spilled byte to the stack and convert the byte to an integer value.

Example 3-40. A Non-forwarding Situation in Compiler Generated Code

```
mov DWORD PTR [esp+10h], 00000000h
mov BYTE PTR [esp+10h], bl
mov eax, DWORD PTR [esp+10h] ; Stall
and eax, 0xff                ; Converting back to byte value
```

Example 3-41 offers two alternatives to avoid the non-forwarding situation shown in Example 3-40.

Example 3-41. Two Ways to Avoid Non-forwarding Situation in Example 3-40

```
; A. Use MOVZ instruction to avoid large load after small
; store, when spills are ignored.
movzx eax, bl                ; Replaces the last three instructions
; B. Use MOVZ instruction and handle spills to the stack
mov DWORD PTR [esp+10h], 00000000h
mov BYTE PTR [esp+10h], bl
movzx eax, BYTE PTR [esp+10h] ; Not blocked
```

When moving data that is smaller than 64 bits between memory locations, 64-bit or 128-bit SIMD register moves are more efficient (if aligned) and can be used to avoid unaligned loads. Although floating-point registers allow the movement of 64 bits at a time, floating point instructions should not be used for this purpose, as data may be inadvertently modified.

As an additional example, consider the cases in Example 3-42.

Example 3-42. Large and Small Load Stalls

```
; A. Large load stall
mov     mem, eax              ; Store dword to address "MEM"
mov     mem + 4, ebx         ; Store dword to address "MEM + 4"
fld     mem                  ; Load qword at address "MEM", stalls
; B. Small Load stall
fstp   mem                  ; Store qword to address "MEM"
mov     bx, mem+2            ; Load word at address "MEM + 2", stalls
mov     cx, mem+4           ; Load word at address "MEM + 4", stalls
```

In the first case (A), there is a large load after a series of small stores to the same area of memory (beginning at memory address MEM). The large load will stall.

The FLD must wait for the stores to write to memory before it can access all the data it requires. This stall can also occur with other data types (for example, when bytes or words are stored and then words or doublewords are read from the same area of memory).

In the second case (B), there is a series of small loads after a large store to the same area of memory (beginning at memory address MEM). The small loads will stall.

The word loads must wait for the quadword store to write to memory before they can access the data they require. This stall can also occur with other data types (for example, when doublewords or words are stored and then words or bytes are read from the same area of memory). This can be avoided by moving the store as far from the loads as possible.

Store forwarding restrictions for processors based on Intel Core microarchitecture is listed in Table 3-2.

**Table 3-2. Store Forwarding Restrictions of Processors
Based on Intel Core Microarchitecture**

Store Alignment	Width of Store (bits)	Load Alignment (byte)	Width of Load (bits)	Store Forwarding Restriction
To Natural size	16	word aligned	8, 16	not stalled
To Natural size	16	not word aligned	8	stalled
To Natural size	32	dword aligned	8, 32	not stalled
To Natural size	32	not dword aligned	8	stalled
To Natural size	32	word aligned	16	not stalled
To Natural size	32	not word aligned	16	stalled
To Natural size	64	qword aligned	8, 16, 64	not stalled
To Natural size	64	not qword aligned	8, 16	stalled
To Natural size	64	dword aligned	32	not stalled
To Natural size	64	not dword aligned	32	stalled
To Natural size	128	dqword aligned	8, 16, 128	not stalled
To Natural size	128	not dqword aligned	8, 16	stalled
To Natural size	128	dword aligned	32	not stalled
To Natural size	128	not dword aligned	32	stalled
To Natural size	128	qword aligned	64	not stalled
To Natural size	128	not qword aligned	64	stalled
Unaligned, start byte 1	32	byte 0 of store	8, 16, 32	not stalled
Unaligned, start byte 1	32	not byte 0 of store	8, 16	stalled
Unaligned, start byte 1	64	byte 0 of store	8, 16, 32	not stalled

**Table 3-2. Store Forwarding Restrictions of Processors
Based on Intel Core Microarchitecture (Contd.)**

Store Alignment	Width of Store (bits)	Load Alignment (byte)	Width of Load (bits)	Store Forwarding Restriction
Unaligned, start byte 1	64	not byte 0 of store	8, 16, 32	stalled
Unaligned, start byte 1	64	byte 0 of store	64	stalled
Unaligned, start byte 7	32	byte 0 of store	8	not stalled
Unaligned, start byte 7	32	not byte 0 of store	8	not stalled
Unaligned, start byte 7	32	don't care	16, 32	stalled
Unaligned, start byte 7	64	don't care	16, 32, 64	stalled

3.6.5.2 Store-forwarding Restriction on Data Availability

The value to be stored must be available before the load operation can be completed. If this restriction is violated, the execution of the load will be delayed until the data is available. This delay causes some execution resources to be used unnecessarily, and that can lead to sizable but non-deterministic delays. However, the overall impact of this problem is much smaller than that from violating size and alignment requirements.

In processors based on Intel NetBurst microarchitecture, hardware predicts when loads are dependent on and get their data forwarded from preceding stores. These predictions can significantly improve performance. However, if a load is scheduled too soon after the store it depends on or if the generation of the data to be stored is delayed, there can be a significant penalty.

There are several cases in which data is passed through memory, and the store may need to be separated from the load:

- Spills, save and restore registers in a stack frame
- Parameter passing
- Global and volatile variables
- Type conversion between integer and floating point
- When compilers do not analyze code that is inlined, forcing variables that are involved in the interface with inlined code to be in memory, creating more memory variables and preventing the elimination of redundant loads

Assembly/Compiler Coding Rule 51. (H impact, MH generality) *Where it is possible to do so without incurring other penalties, prioritize the allocation of variables to registers, as in register allocation and for parameter passing, to minimize the likelihood and impact of store-forwarding problems. Try not to store-forward data generated from a long latency instruction - for example, MUL or DIV. Avoid store-forwarding data for variables with the shortest store-load distance.*

Avoid store-forwarding data for variables with many and/or long dependence chains, and especially avoid including a store forward on a loop-carried dependence chain.

shows an example of a loop-carried dependence chain.

Example 3-43. Loop-carried Dependence Chain

```
for ( i = 0; i < MAX; i++ ) {
    a[i] = b[i] * foo;
    foo = a[i] / 3;
} // foo is a loop-carried dependence.
```

Assembly/Compiler Coding Rule 52. (M impact, MH generality) Calculate store addresses as early as possible to avoid having stores block loads.

3.6.6 Data Layout Optimizations

User/Source Coding Rule 6. (H impact, M generality) Pad data structures defined in the source code so that every data element is aligned to a natural operand size address boundary.

If the operands are packed in a SIMD instruction, align to the packed element size (64-bit or 128-bit).

Align data by providing padding inside structures and arrays. Programmers can reorganize structures and arrays to minimize the amount of memory wasted by padding. However, compilers might not have this freedom. The C programming language, for example, specifies the order in which structure elements are allocated in memory. For more information, see Section 4.4, “Stack and Data Alignment,” and Appendix D, “Stack Alignment.”

Example 3-44 shows how a data structure could be rearranged to reduce its size.

Example 3-44. Rearranging a Data Structure

```
struct unpacked { /* Fits in 20 bytes due to padding */
    int    a;
    char   b;
    int    c;
    char   d;
    int    e;
};
```

Example 3-44. Rearranging a Data Structure (Contd.)

```

struct packed { /* Fits in 16 bytes */
    int    a;
    int    c;
    int    e;
    char   b;
    char   d;
}

```

Cache line size of 64 bytes can impact streaming applications (for example, multi-media). These reference and use data only once before discarding it. Data accesses which sparsely utilize the data within a cache line can result in less efficient utilization of system memory bandwidth. For example, arrays of structures can be decomposed into several arrays to achieve better packing, as shown in Example 3-45.

Example 3-45. Decomposing an Array

```

struct { /* 1600 bytes */
    int  a, c, e;
    char b, d;
} array_of_struct [100];

struct { /* 1400 bytes */
    int  a[100], c[100], e[100];
    char b[100], d[100];
} struct_of_array;

struct { /* 1200 bytes */
    int  a, c, e;
} hybrid_struct_of_array_ace[100];

struct { /* 200 bytes */
    char b, d;
} hybrid_struct_of_array_bd[100];

```

The efficiency of such optimizations depends on usage patterns. If the elements of the structure are all accessed together but the access pattern of the array is random, then ARRAY_OF_STRUCT avoids unnecessary prefetch even though it wastes memory.

However, if the access pattern of the array exhibits locality (for example, if the array index is being swept through) then processors with hardware prefetchers will

prefetch data from `STRUCT_OF_ARRAY`, even if the elements of the structure are accessed together.

When the elements of the structure are not accessed with equal frequency, such as when element A is accessed ten times more often than the other entries, then `STRUCT_OF_ARRAY` not only saves memory, but it also prevents fetching unnecessary data items B, C, D, and E.

Using `STRUCT_OF_ARRAY` also enables the use of the SIMD data types by the programmer and the compiler.

Note that `STRUCT_OF_ARRAY` can have the disadvantage of requiring more independent memory stream references. This can require the use of more prefetches and additional address generation calculations. It can also have an impact on DRAM page access efficiency. An alternative, `HYBRID_STRUCT_OF_ARRAY` blends the two approaches. In this case, only 2 separate address streams are generated and referenced: 1 for `HYBRID_STRUCT_OF_ARRAY_ACE` and 1 for `HYBRID_STRUCT_OF_ARRAY_BD`. The second alternative also prevents fetching unnecessary data — assuming that (1) the variables A, C and E are always used together, and (2) the variables B and D are always used together, but not at the same time as A, C and E.

The hybrid approach ensures:

- Simpler/fewer address generations than `STRUCT_OF_ARRAY`
- Fewer streams, which reduces DRAM page misses
- Fewer prefetches due to fewer streams
- Efficient cache line packing of data elements that are used concurrently

Assembly/Compiler Coding Rule 53. (H impact, M generality) *Try to arrange data structures such that they permit sequential access.*

If the data is arranged into a set of streams, the automatic hardware prefetcher can prefetch data that will be needed by the application, reducing the effective memory latency. If the data is accessed in a non-sequential manner, the automatic hardware prefetcher cannot prefetch the data. The prefetcher can recognize up to eight concurrent streams. See Chapter 7, “Optimizing Cache Usage,” for more information on the hardware prefetcher.

On Intel Core 2 Duo, Intel Core Duo, Intel Core Solo, Pentium 4, Intel Xeon and Pentium M processors, memory coherence is maintained on 64-byte cache lines (rather than 32-byte cache lines, as in earlier processors). This can increase the opportunity for false sharing.

User/Source Coding Rule 7. (M impact, L generality) *Beware of false sharing within a cache line (64 bytes) and within a sector of 128 bytes on processors based on Intel NetBurst microarchitecture.*

3.6.7 Stack Alignment

The easiest way to avoid stack alignment problems is to keep the stack aligned at all times. For example, a language that supports 8-bit, 16-bit, 32-bit, and 64-bit data quantities but never uses 80-bit data quantities can require the stack to always be aligned on a 64-bit boundary.

Assembly/Compiler Coding Rule 54. (H impact, M generality) *If 64-bit data is ever passed as a parameter or allocated on the stack, make sure that the stack is aligned to an 8-byte boundary.*

Doing this will require using a general purpose register (such as EBP) as a frame pointer. The trade-off is between causing unaligned 64-bit references (if the stack is not aligned) and causing extra general purpose register spills (if the stack is aligned). Note that a performance penalty is caused only when an unaligned access splits a cache line. This means that one out of eight spatially consecutive unaligned accesses is always penalized.

A routine that makes frequent use of 64-bit data can avoid stack misalignment by placing the code described in Example 3-46 in the function prologue and epilogue.

Example 3-46. Dynamic Stack Alignment

```

prologue:
    subl    esp, 4           ; Save frame ptr
    movl   [esp], ebp
    movl   ebp, esp        ; New frame pointer
    andl   ebp, 0xFFFFFFFF ; Aligned to 64 bits
    movl   [ebp], esp      ; Save old stack ptr
    subl   esp, FRAMESIZE  ; Allocate space
    ; ... callee saves, etc.

epilogue:
    ; ... callee restores, etc.
    movl   esp, [ebp]      ; Restore stack ptr
    movl   ebp, [esp]      ; Restore frame ptr
    addl   esp, 4
    ret

```

If for some reason it is not possible to align the stack for 64-bits, the routine should access the parameter and save it into a register or known aligned storage, thus incurring the penalty only once.

3.6.8 Capacity Limits and Aliasing in Caches

There are cases in which addresses with a given stride will compete for some resource in the memory hierarchy.

Typically, caches are implemented to have multiple ways of set associativity, with each way consisting of multiple sets of cache lines (or sectors in some cases). Multiple memory references that compete for the same set of each way in a cache can cause a capacity issue. There are aliasing conditions that apply to specific microarchitectures. Note that first-level cache lines are 64 bytes. Thus, the least significant 6 bits are not considered in alias comparisons. For processors based on Intel NetBurst microarchitecture, data is loaded into the second level cache in a sector of 128 bytes, so the least significant 7 bits are not considered in alias comparisons.

3.6.8.1 Capacity Limits in Set-Associative Caches

Capacity limits may be reached if the number of outstanding memory references that are mapped to the same set in each way of a given cache exceeds the number of ways of that cache. The conditions that apply to the first-level data cache and second level cache are listed below:

- **L1 Set Conflicts** — Multiple references map to the same first-level cache set. The conflicting condition is a stride determined by the size of the cache in bytes, divided by the number of ways. These competing memory references can cause excessive cache misses only if the number of outstanding memory references exceeds the number of ways in the working set:
 - On Pentium 4 and Intel Xeon processors with a CPUID signature of family encoding 15, model encoding of 0, 1, or 2; there will be an excess of first-level cache misses for more than 4 simultaneous competing memory references to addresses with 2-KByte modulus.
 - On Pentium 4 and Intel Xeon processors with a CPUID signature of family encoding 15, model encoding 3; there will be an excess of first-level cache misses for more than 8 simultaneous competing references to addresses that are apart by 2-KByte modulus.
 - On Intel Core 2 Duo, Intel Core Duo, Intel Core Solo, and Pentium M processors, there will be an excess of first-level cache misses for more than 8 simultaneous references to addresses that are apart by 4-KByte modulus.
- **L2 Set Conflicts** — Multiple references map to the same second-level cache set. The conflicting condition is also determined by the size of the cache or the number of ways:
 - On Pentium 4 and Intel Xeon processors, there will be an excess of second-level cache misses for more than 8 simultaneous competing references. The stride sizes that can cause capacity issues are 32 KBytes, 64 KBytes, or 128 KBytes, depending of the size of the second level cache.
 - On Pentium M processors, the stride sizes that can cause capacity issues are 128 KBytes or 256 KBytes, depending of the size of the second level cache. On Intel Core 2 Duo, Intel Core Duo, Intel Core Solo processors, stride size of 256 KBytes can cause capacity issue if the number of simultaneous accesses exceeded the way associativity of the L2 cache.

3.6.8.2 Aliasing Cases in the Pentium® M, Intel® Core™ Solo, Intel® Core™ Duo and Intel® Core™ 2 Duo Processors

Pentium M, Intel Core Solo, Intel Core Duo and Intel Core 2 Duo processors have the following aliasing case:

- **Store forwarding** — If a store to an address is followed by a load from the same address, the load will not proceed until the store data is available. If a store is followed by a load and their addresses differ by a multiple of 4 KBytes, the load stalls until the store operation completes.

Assembly/Compiler Coding Rule 55. (H impact, M generality) *Avoid having a store followed by a non-dependent load with addresses that differ by a multiple of 4 KBytes. Also, lay out data or order computation to avoid having cache lines that have linear addresses that are a multiple of 64 KBytes apart in the same working set. Avoid having more than 4 cache lines that are some multiple of 2 KBytes apart in the same first-level cache working set, and avoid having more than 8 cache lines that are some multiple of 4 KBytes apart in the same first-level cache working set.*

When declaring multiple arrays that are referenced with the same index and are each a multiple of 64 KBytes (as can happen with STRUCT_OF_ARRAY data layouts), pad them to avoid declaring them contiguously. Padding can be accomplished by either intervening declarations of other variables or by artificially increasing the dimension.

User/Source Coding Rule 8. (H impact, ML generality) *Consider using a special memory allocation library with address offset capability to avoid aliasing.*

One way to implement a memory allocator to avoid aliasing is to allocate more than enough space and pad. For example, allocate structures that are 68 KB instead of 64 KBytes to avoid the 64-KByte aliasing, or have the allocator pad and return random offsets that are a multiple of 128 Bytes (the size of a cache line).

User/Source Coding Rule 9. (M impact, M generality) *When padding variable declarations to avoid aliasing, the greatest benefit comes from avoiding aliasing on second-level cache lines, suggesting an offset of 128 bytes or more.*

4-KByte memory aliasing occurs when the code accesses two different memory locations with a 4-KByte offset between them. The 4-KByte aliasing situation can manifest in a memory copy routine where the addresses of the source buffer and destination buffer maintain a constant offset and the constant offset happens to be a multiple of the byte increment from one iteration to the next.

Example 3-47 shows a routine that copies 16 bytes of memory in each iteration of a loop. If the offsets (modular 4096) between source buffer (EAX) and destination buffer (EDX) differ by 16, 32, 48, 64, 80; loads have to wait until stores have been retired before they can continue. For example at offset 16, the load of the next iteration is 4-KByte aliased current iteration store, therefore the loop must wait until the store operation completes, making the entire loop serialized. The amount of time needed to wait decreases with larger offset until offset of 96 resolves the issue (as there is no pending stores by the time of the load with same address).

The Intel Core microarchitecture provides a performance monitoring event (see LOAD_BLOCK.OVERLAP_STORE in *Intel® 64 and IA-32 Architectures Software*

Developer's Manual, Volume 3B) that allows software tuning effort to detect the occurrence of aliasing conditions.

Example 3-47. Aliasing Between Loads and Stores Across Loop Iterations

```
LP:
    movaps xmm0, [eax+ecx]
    movaps [edx+ecx], xmm0
    add ecx, 16
    jnz lp
```

3.6.9 Mixing Code and Data

The aggressive prefetching and pre-decoding of instructions by Intel processors have two related effects:

- Self-modifying code works correctly, according to the Intel architecture processor requirements, but incurs a significant performance penalty. Avoid self-modifying code if possible.
- Placing writable data in the code segment might be impossible to distinguish from self-modifying code. Writable data in the code segment might suffer the same performance penalty as self-modifying code.

Assembly/Compiler Coding Rule 56. (M impact, L generality) *If (hopefully read-only) data must occur on the same page as code, avoid placing it immediately after an indirect jump. For example, follow an indirect jump with its mostly likely target, and place the data after an unconditional branch.*

Tuning Suggestion 1. *In rare cases, a performance problem may be caused by executing data on a code page as instructions. This is very likely to happen when execution is following an indirect branch that is not resident in the trace cache. If this is clearly causing a performance problem, try moving the data elsewhere, or inserting an illegal opcode or a PAUSE instruction immediately after the indirect branch. Note that the latter two alternatives may degrade performance in some circumstances.*

Assembly/Compiler Coding Rule 57. (H impact, L generality) *Always put code and data on separate pages. Avoid self-modifying code wherever possible. If code is to be modified, try to do it all at once and make sure the code that performs the modifications and the code being modified are on separate 4-KByte pages or on separate aligned 1-KByte subpages.*

3.6.9.1 Self-modifying Code

Self-modifying code (SMC) that ran correctly on Pentium III processors and prior implementations will run correctly on subsequent implementations. SMC and cross-

modifying code (when multiple processors in a multiprocessor system are writing to a code page) should be avoided when high performance is desired.

Software should avoid writing to a code page in the same 1-KByte subpage that is being executed or fetching code in the same 2-KByte subpage of that is being written. In addition, sharing a page containing directly or speculatively executed code with another processor as a data page can trigger an SMC condition that causes the entire pipeline of the machine and the trace cache to be cleared. This is due to the self-modifying code condition.

Dynamic code need not cause the SMC condition if the code written fills up a data page before that page is accessed as code. Dynamically-modified code (for example, from target fix-ups) is likely to suffer from the SMC condition and should be avoided where possible. Avoid the condition by introducing indirect branches and using data tables on data pages (not code pages) using register-indirect calls.

3.6.9.2 Position Independent Code

Position independent code often needs to obtain the value of the instruction pointer. Example 3-48a shows one technique to put the value of IP into the ECX register by issuing a CALL without a matching RET. Example 3-48b shows an alternative technique to put the value of IP into the ECX register using a matched pair of CALL/RET.

Example 3-48. Instruction Pointer Query Techniques

a) Using call without return to obtain IP does not corrupt the RSB call _label; return address pushed is the IP of next instruction _label: pop ecx; IP of this instruction is now put into ECX
b) Using matched call/ret pair call _blcx; ... ; ECX now contains IP of this instruction ... _blcx mov ecx, [esp]; ret

3.6.10 Write Combining

Write combining (WC) improves performance in two ways:

- On a write miss to the first-level cache, it allows multiple stores to the same cache line to occur before that cache line is read for ownership (RFO) from further out in the cache/memory hierarchy. Then the rest of line is read, and the bytes that have not been written are combined with the unmodified bytes in the returned line.
- Write combining allows multiple writes to be assembled and written further out in the cache hierarchy as a unit. This saves port and bus traffic. Saving traffic is particularly important for avoiding partial writes to uncached memory.

There are six write-combining buffers (on Pentium 4 and Intel Xeon processors with a CPUID signature of family encoding 15, model encoding 3; there are 8 write-combining buffers). Two of these buffers may be written out to higher cache levels and freed up for use on other write misses. Only four write-combining buffers are guaranteed to be available for simultaneous use. Write combining applies to memory type WC; it does not apply to memory type UC.

There are six write-combining buffers in each processor core in Intel Core Duo and Intel Core Solo processors. Processors based on Intel Core microarchitecture have eight write-combining buffers in each core.

Assembly/Compiler Coding Rule 58. (H impact, L generality) *If an inner loop writes to more than four arrays (four distinct cache lines), apply loop fission to break up the body of the loop such that only four arrays are being written to in each iteration of each of the resulting loops.*

Write combining buffers are used for stores of all memory types. They are particularly important for writes to uncached memory: writes to different parts of the same cache line can be grouped into a single, full-cache-line bus transaction instead of going across the bus (since they are not cached) as several partial writes. Avoiding partial writes can have a significant impact on bus bandwidth-bound graphics applications, where graphics buffers are in uncached memory. Separating writes to uncached memory and writes to writeback memory into separate phases can assure that the write combining buffers can fill before getting evicted by other write traffic. Eliminating partial write transactions has been found to have performance impact on the order of 20% for some applications. Because the cache lines are 64 bytes, a write to the bus for 63 bytes will result in 8 partial bus transactions.

When coding functions that execute simultaneously on two threads, reducing the number of writes that are allowed in an inner loop will help take full advantage of write-combining store buffers. For write-combining buffer recommendations for Hyper-Threading Technology, see Chapter 8, "Multicore and Hyper-Threading Technology."

Store ordering and visibility are also important issues for write combining. When a write to a write-combining buffer for a previously-unwritten cache line occurs, there will be a read-for-ownership (RFO). If a subsequent write happens to another write-combining buffer, a separate RFO may be caused for that cache line. Subsequent

writes to the first cache line and write-combining buffer will be delayed until the second RFO has been serviced to guarantee properly ordered visibility of the writes. If the memory type for the writes is write-combining, there will be no RFO since the line is not cached, and there is no such delay. For details on write-combining, see Chapter 10, "Power Optimization for Mobile Usages," of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*.

3.6.11 Locality Enhancement

Locality enhancement can reduce data traffic originating from an outer-level sub-system in the cache/memory hierarchy. This is to address the fact that the access-cost in terms of cycle-count from an outer level will be more expensive than from an inner level. Typically, the cycle-cost of accessing a given cache level (or memory system) varies across different microarchitectures, processor implementations, and platform components. It may be sufficient to recognize the relative data access cost trend by locality rather than to follow a large table of numeric values of cycle-costs, listed per locality, per processor/platform implementations, etc. The general trend is typically that access cost from an outer sub-system may be approximately 3-10X more expensive than accessing data from the immediate inner level in the cache/memory hierarchy, assuming similar degrees of data access parallelism.

Thus locality enhancement should start with characterizing the dominant data traffic locality. Section A, "Application Performance Tools," describes some techniques that can be used to determine the dominant data traffic locality for any workload.

Even if cache miss rates of the last level cache may be low relative to the number of cache references, processors typically spend a sizable portion of their execution time waiting for cache misses to be serviced. Reducing cache misses by enhancing a program's locality is a key optimization. This can take several forms:

- Blocking to iterate over a portion of an array that will fit in the cache (with the purpose that subsequent references to the data-block [or tile] will be cache hit references)
- Loop interchange to avoid crossing cache lines or page boundaries
- Loop skewing to make accesses contiguous

Locality enhancement to the last level cache can be accomplished with sequencing the data access pattern to take advantage of hardware prefetching. This can also take several forms:

- Transformation of a sparsely populated multi-dimensional array into a one-dimension array such that memory references occur in a sequential, small-stride pattern that is friendly to the hardware prefetch (see Section 2.4.4.4, "Data Prefetch")
- Optimal tile size and shape selection can further improve temporal data locality by increasing hit rates into the last level cache and reduce memory traffic resulting from the actions of hardware prefetching (see Section 7.6.11, "Hardware Prefetching and Cache Blocking Techniques")

It is important to avoid operations that work against locality-enhancing techniques. Using the lock prefix heavily can incur large delays when accessing memory, regardless of whether the data is in the cache or in system memory.

User/Source Coding Rule 10. (H impact, H generality) *Optimization techniques such as blocking, loop interchange, loop skewing, and packing are best done by the compiler. Optimize data structures either to fit in one-half of the first-level cache or in the second-level cache; turn on loop optimizations in the compiler to enhance locality for nested loops.*

Optimizing for one-half of the first-level cache will bring the greatest performance benefit in terms of cycle-cost per data access. If one-half of the first-level cache is too small to be practical, optimize for the second-level cache. Optimizing for a point in between (for example, for the entire first-level cache) will likely not bring a substantial improvement over optimizing for the second-level cache.

3.6.12 Minimizing Bus Latency

Each bus transaction includes the overhead of making requests and arbitrations. The average latency of bus read and bus write transactions will be longer if reads and writes alternate. Segmenting reads and writes into phases can reduce the average latency of bus transactions. This is because the number of incidences of successive transactions involving a read following a write, or a write following a read, are reduced.

User/Source Coding Rule 11. (M impact, ML generality) *If there is a blend of reads and writes on the bus, changing the code to separate these bus transactions into read phases and write phases can help performance.*

Note, however, that the order of read and write operations on the bus is not the same as it appears in the program.

Bus latency for fetching a cache line of data can vary as a function of the access stride of data references. In general, bus latency will increase in response to increasing values of the stride of successive cache misses. Independently, bus latency will also increase as a function of increasing bus queue depths (the number of outstanding bus requests of a given transaction type). The combination of these two trends can be highly non-linear, in that bus latency of large-stride, bandwidth-sensitive situations are such that effective throughput of the bus system for data-parallel accesses can be significantly less than the effective throughput of small-stride, bandwidth-sensitive situations.

To minimize the per-access cost of memory traffic or amortize raw memory latency effectively, software should control its cache miss pattern to favor higher concentration of smaller-stride cache misses.

User/Source Coding Rule 12. (H impact, H generality) *To achieve effective amortization of bus latency, software should favor data access patterns that result in higher concentrations of cache miss patterns, with cache miss strides that are significantly smaller than half the hardware prefetch trigger threshold.*

3.6.13 Non-Temporal Store Bus Traffic

Peak system bus bandwidth is shared by several types of bus activities, including reads (from memory), reads for ownership (of a cache line), and writes. The data transfer rate for bus write transactions is higher if 64 bytes are written out to the bus at a time.

Typically, bus writes to Writeback (WB) memory must share the system bus bandwidth with read-for-ownership (RFO) traffic. Non-temporal stores do not require RFO traffic; they do require care in managing the access patterns in order to ensure 64 bytes are evicted at once (rather than evicting several 8-byte chunks).

Although the data bandwidth of full 64-byte bus writes due to non-temporal stores is twice that of bus writes to WB memory, transferring 8-byte chunks wastes bus request bandwidth and delivers significantly lower data bandwidth. This difference is depicted in Examples 3-49 and 3-50.

Example 3-49. Using Non-temporal Stores and 64-byte Bus Write Transactions

```
#define STRIDESIZE 256
lea ecx, p64byte_Aligned
mov edx, ARRAY_LEN
xor eax, eax
sloop:
movntps XMMWORD ptr [ecx + eax], xmm0
movntps XMMWORD ptr [ecx + eax+16], xmm0
movntps XMMWORD ptr [ecx + eax+32], xmm0
movntps XMMWORD ptr [ecx + eax+48], xmm0
; 64 bytes is written in one bus transaction
add eax, STRIDESIZE
cmp eax, edx
jl sloop
```

Example 3-50. On-temporal Stores and Partial Bus Write Transactions

```

#define STRIDESIZE 256
Lea ecx, p64byte_Aligned
Mov edx, ARRAY_LEN
Xor eax, eax
sloop:
movntps XMMWORD ptr [ecx + eax], xmm0
movntps XMMWORD ptr [ecx + eax+16], xmm0
movntps XMMWORD ptr [ecx + eax+32], xmm0
; Storing 48 bytes results in 6 bus partial transactions
add eax, STRIDESIZE
cmp eax, edx

```

3.7 PREFETCHING

Recent Intel processor families employ several prefetching mechanisms to accelerate the movement of data or code and improve performance:

- Hardware instruction prefetcher
- Software prefetch for data
- Hardware prefetch for cache lines of data or instructions

3.7.1 Hardware Instruction Fetching and Software Prefetching

In processor based on Intel NetBurst microarchitecture, the hardware instruction fetcher reads instructions, 32 bytes at a time, into the 64-byte instruction streaming buffers. Instruction fetching for Intel Core microarchitecture is discussed in Section 2.2.2.

Software prefetching requires a programmer to use PREFETCH hint instructions and anticipate some suitable timing and location of cache misses.

In Intel Core microarchitecture, software PREFETCH instructions can prefetch beyond page boundaries and can perform one-to-four page walks. Software PREFETCH instructions issued on fill buffer allocations retire after the page walk completes and the DCU miss is detected. Software PREFETCH instructions can trigger all hardware prefetchers in the same manner as do regular loads.

Software PREFETCH operations work the same way as do load from memory operations, with the following exceptions:

- Software PREFETCH instructions retire after virtual to physical address translation is completed.

- If an exception, such as page fault, is required to prefetch the data, then the software prefetch instruction retires without prefetching data.

3.7.2 Software and Hardware Prefetching in Prior Microarchitectures

Pentium 4 and Intel Xeon processors based on Intel NetBurst microarchitecture introduced hardware prefetching in addition to software prefetching. The hardware prefetcher operates transparently to fetch data and instruction streams from memory without requiring programmer intervention. Subsequent microarchitectures continue to improve and add features to the hardware prefetching mechanisms. Earlier implementations of hardware prefetching mechanisms focus on prefetching data and instruction from memory to L2; more recent implementations provide additional features to prefetch data from L2 to L1.

In Intel NetBurst microarchitecture, the hardware prefetcher can track 8 independent streams.

The Pentium M processor also provides a hardware prefetcher for data. It can track 12 separate streams in the forward direction and 4 streams in the backward direction. The processor's PREFETCHNTA instruction also fetches 64-bytes into the first-level data cache without polluting the second-level cache.

Intel Core Solo and Intel Core Duo processors provide more advanced hardware prefetchers for data than Pentium M processors. Key differences are summarized in Table 2-23.

Although the hardware prefetcher operates transparently (requiring no intervention by the programmer), it operates most efficiently if the programmer specifically tailors data access patterns to suit its characteristics (it favors small-stride cache miss patterns). Optimizing data access patterns to suit the hardware prefetcher is highly recommended, and should be a higher-priority consideration than using software prefetch instructions.

The hardware prefetcher is best for small-stride data access patterns in either direction with a cache-miss stride not far from 64 bytes. This is true for data accesses to addresses that are either known or unknown at the time of issuing the load operations. Software prefetch can complement the hardware prefetcher if used carefully.

There is a trade-off to make between hardware and software prefetching. This pertains to application characteristics such as regularity and stride of accesses. Bus bandwidth, issue bandwidth (the latency of loads on the critical path) and whether access patterns are suitable for non-temporal prefetch will also have an impact.

For a detailed description of how to use prefetching, see Chapter 7, "Optimizing Cache Usage."

Chapter 5, "Optimizing for SIMD Integer Applications," contains an example that uses software prefetch to implement a memory copy algorithm.

Tuning Suggestion 2. *If a load is found to miss frequently, either insert a prefetch before it or (if issue bandwidth is a concern) move the load up to execute earlier.*

3.7.3 Hardware Prefetching for First-Level Data Cache

The hardware prefetching mechanism for L1 in Intel Core microarchitecture is discussed in Section 2.2.4.2. A similar L1 prefetch mechanism is also available to processors based on Intel NetBurst microarchitecture with CPUID signature of family 15 and model 6.

Example 3-51 depicts a technique to trigger hardware prefetch. The code demonstrates traversing a linked list and performing some computational work on 2 members of each element that reside in 2 different cache lines. Each element is of size 192 bytes. The total size of all elements is larger than can be fitted in the L2 cache.

Example 3-51. Using DCU Hardware Prefetch

Original code	Modified sequence benefit from prefetch
<pre> mov ebx, DWORD PTR [First] xor eax, eax scan_list: mov eax, [ebx+4] mov ecx, 60 do_some_work_1: add eax, eax and eax, 6 sub ecx, 1 jnz do_some_work_1 </pre>	<pre> mov ebx, DWORD PTR [First] xor eax, eax scan_list: mov eax, [ebx+4] mov eax, [ebx+4] mov eax, [ebx+4] mov ecx, 60 do_some_work_1: add eax, eax and eax, 6 sub ecx, 1 jnz do_some_work_1 </pre>
<pre> mov eax, [ebx+64] mov ecx, 30 do_some_work_2: add eax, eax and eax, 6 sub ecx, 1 jnz do_some_work_2 mov ebx, [ebx] test ebx, ebx jnz scan_list </pre>	<pre> mov eax, [ebx+64] mov ecx, 30 do_some_work_2: add eax, eax and eax, 6 sub ecx, 1 jnz do_some_work_2 mov ebx, [ebx] test ebx, ebx jnz scan_list </pre>

The additional instructions to load data from one member in the modified sequence can trigger the DCU hardware prefetch mechanisms to prefetch data in the next cache line, enabling the work on the second member to complete sooner.

Software can gain from the first-level data cache prefetchers in two cases:

- If data is not in the second-level cache, the first-level data cache prefetcher enables early trigger of the second-level cache prefetcher.
- If data is in the second-level cache and not in the first-level data cache, then the first-level data cache prefetcher triggers earlier data bring-up of sequential cache line to the first-level data cache.

There are situations that software should pay attention to a potential side effect of triggering unnecessary DCU hardware prefetches. If a large data structure with many members spanning many cache lines is accessed in ways that only a few of its members are actually referenced, but there are multiple pair accesses to the same cache line. The DCU hardware prefetcher can trigger fetching of cache lines that are not needed. In Example , references to the "Pts" array and "AltPts" will trigger DCU

prefetch to fetch additional cache lines that won't be needed. If significant negative performance impact is detected due to DCU hardware prefetch on a portion of the code, software can try to reduce the size of that contemporaneous working set to be less than half of the L2 cache.

Example 3-52. Avoid Causing DCU Hardware Prefetch to Fetch Un-needed Lines

```
while ( CurrBond != NULL )
{
  MyATOM *a1 = CurrBond->At1 ;
  MyATOM *a2 = CurrBond->At2 ;

  if ( a1->CurrStep <= a1->LastStep &&
      a2->CurrStep <= a2->LastStep
      )
  {
    a1->CurrStep++ ;
    a2->CurrStep++ ;

    double ux = a1->Pts[0].x - a2->Pts[0].x ;
    double uy = a1->Pts[0].y - a2->Pts[0].y ;
    double uz = a1->Pts[0].z - a2->Pts[0].z ;

    a1->AuxPts[0].x += ux ;
    a1->AuxPts[0].y += uy ;
    a1->AuxPts[0].z += uz ;

    a2->AuxPts[0].x += ux ;
    a2->AuxPts[0].y += uy ;
    a2->AuxPts[0].z += uz ;
  } ;
  CurrBond = CurrBond->Next ;
};
```

To fully benefit from these prefetchers, organize and access the data using one of the following methods:

Method 1:

- Organize the data so consecutive accesses can usually be found in the same 4-KByte page.
- Access the data in constant strides forward or backward IP Prefetcher.

Method 2:

- Organize the data in consecutive lines.
- Access the data in increasing addresses, in sequential cache lines.

Example demonstrates accesses to sequential cache lines that can benefit from the first-level cache prefetcher.

Example 3-53. Technique For Using L1 Hardware Prefetch

```
unsigned int *p1, j, a, b;
for (j = 0; j < num; j += 16)
{
    a = p1[j];
    b = p1[j+1];
    // Use these two values
}
```

By elevating the load operations from memory to the beginning of each iteration, it is likely that a significant part of the latency of the pair cache line transfer from memory to the second-level cache will be in parallel with the transfer of the first cache line.

The IP prefetcher uses only the lower 8 bits of the address to distinguish a specific address. If the code size of a loop is bigger than 256 bytes, two loads may appear similar in the lowest 8 bits and the IP prefetcher will be restricted. Therefore, if you have a loop bigger than 256 bytes, make sure that no two loads have the same lowest 8 bits in order to use the IP prefetcher.

3.7.4 Hardware Prefetching for Second-Level Cache

The Intel Core microarchitecture contains two second-level cache prefetchers:

- **Streamer** — Loads data or instructions from memory to the second-level cache. To use the streamer, organize the data or instructions in blocks of 128 bytes, aligned on 128 bytes. The first access to one of the two cache lines in this block while it is in memory triggers the streamer to prefetch the pair line. To software, the L2 streamer's functionality is similar to the adjacent cache line prefetch mechanism found in processors based on Intel NetBurst microarchitecture.
- **Data prefetch logic (DPL)** — DPL and L2 Streamer are triggered only by writeback memory type. They prefetch only inside page boundary (4 KBytes). Both L2 prefetchers can be triggered by software prefetch instructions and by prefetch request from DCU prefetchers. DPL can also be triggered by read for ownership (RFO) operations. The L2 Streamer can also be triggered by DPL requests for L2 cache misses.

Software can gain from organizing data both according to the instruction pointer and according to line strides. For example, for matrix calculations, columns can be

prefetched by IP-based prefetches, and rows can be prefetched by DPL and the L2 streamer.

3.7.5 Cacheability Instructions

SSE2 provides additional cacheability instructions that extend those provided in SSE. The new cacheability instructions include:

- new streaming store instructions
- new cache line flush instruction
- new memory fencing instructions

For more information, see Chapter 7, “Optimizing Cache Usage.”

3.7.6 REP Prefix and Data Movement

The REP prefix is commonly used with string move instructions for memory related library functions such as MEMCPY (using REP MOVSD) or MEMSET (using REP STOS). These STRING/MOV instructions with the REP prefixes are implemented in MS-ROM and have several implementation variants with different performance levels.

The specific variant of the implementation is chosen at execution time based on data layout, alignment and the counter (ECX) value. For example, MOVSB/STOSB with the REP prefix should be used with counter value less than or equal to three for best performance.

String MOVE/STORE instructions have multiple data granularities. For efficient data movement, larger data granularities are preferable. This means better efficiency can be achieved by decomposing an arbitrary counter value into a number of double-words plus single byte moves with a count value less than or equal to 3.

Because software can use SIMD data movement instructions to move 16 bytes at a time, the following paragraphs discuss general guidelines for designing and implementing high-performance library functions such as MEMCPY(), MEMSET(), and MEMMOVE(). Four factors are to be considered:

- **Throughput per iteration** — If two pieces of code have approximately identical path lengths, efficiency favors choosing the instruction that moves larger pieces of data per iteration. Also, smaller code size per iteration will in general reduce overhead and improve throughput. Sometimes, this may involve a comparison of the relative overhead of an iterative loop structure versus using REP prefix for iteration.
- **Address alignment** — Data movement instructions with highest throughput usually have alignment restrictions, or they operate more efficiently if the destination address is aligned to its natural data size. Specifically, 16-byte moves need to ensure the destination address is aligned to 16-byte boundaries, and 8-bytes moves perform better if the destination address is aligned to 8-byte

boundaries. Frequently, moving at doubleword granularity performs better with addresses that are 8-byte aligned.

- **REP string move vs. SIMD move** — Implementing general-purpose memory functions using SIMD extensions usually requires adding some prolog code to ensure the availability of SIMD instructions, preamble code to facilitate aligned data movement requirements at runtime. Throughput comparison must also take into consideration the overhead of the prolog when considering a REP string implementation versus a SIMD approach.
- **Cache eviction** — If the amount of data to be processed by a memory routine approaches half the size of the last level on-die cache, temporal locality of the cache may suffer. Using streaming store instructions (for example: MOVNTQ, MOVNTDQ) can minimize the effect of flushing the cache. The threshold to start using a streaming store depends on the size of the last level cache. Determine the size using the deterministic cache parameter leaf of CPUID.

Techniques for using streaming stores for implementing a MEMSET()-type library must also consider that the application can benefit from this technique only if it has no immediate need to reference the target addresses. This assumption is easily upheld when testing a streaming-store implementation on a micro-benchmark configuration, but violated in a full-scale application situation.

When applying general heuristics to the design of general-purpose, high-performance library routines, the following guidelines can be useful when optimizing an arbitrary counter value N and address alignment. Different techniques may be necessary for optimal performance, depending on the magnitude of N:

- When N is less than some small count (where the small count threshold will vary between microarchitectures -- empirically, 8 may be a good value when optimizing for Intel NetBurst microarchitecture), each case can be coded directly without the overhead of a looping structure. For example, 11 bytes can be processed using two MOVSD instructions explicitly and a MOVSB with REP counter equaling 3.
- When N is not small but still less than some threshold value (which may vary for different micro-architectures, but can be determined empirically), an SIMD implementation using run-time CPUID and alignment prolog will likely deliver less throughput due to the overhead of the prolog. A REP string implementation should favor using a REP string of doublewords. To improve address alignment, a small piece of prolog code using MOVSB/STOSB with a count less than 4 can be used to peel off the non-aligned data moves before starting to use MOVSD/STOSD.
- When N is less than half the size of last level cache, throughput consideration may favor either:
 - An approach using a REP string with the largest data granularity because a REP string has little overhead for loop iteration, and the branch misprediction overhead in the prolog/epilogue code to handle address alignment is amortized over many iterations.

- An iterative approach using the instruction with largest data granularity, where the overhead for SIMD feature detection, iteration overhead, and prolog/epilogue for alignment control can be minimized. The trade-off between these approaches may depend on the microarchitecture.

An example of MEMSET() implemented using stosd for arbitrary counter value with the destination address aligned to doubleword boundary in 32-bit mode is shown in Example 3-54.

- When N is larger than half the size of the last level cache, using 16-byte granularity streaming stores with prolog/epilog for address alignment will likely be more efficient, if the destination addresses will not be referenced immediately afterwards.

Example 3-54. REP STOSD with Arbitrary Count Size and 4-Byte-Aligned Destination

A 'C' example of Memset()	Equivalent Implementation Using REP STOSD
<pre>void memset(void *dst,int c,size_t size) { char *d = (char *)dst; size_t i; for (i=0;i<size;i++) *d++ = (char)c; }</pre>	<pre>push edi movzx eax, byte ptr [esp+12] mov ecx, eax shl ecx, 8 or ecx, eax mov ecx, eax shl ecx, 16 or eax, ecx mov edi, [esp+8] ; 4-byte aligned mov ecx, [esp+16] ; byte count shr ecx, 2 ; do dword cmp ecx, 127 jle _main test edi, 4 jz _main stosd ;peel off one dword dec ecx _main: ; 8-byte aligned rep stosd mov ecx, [esp + 16] and ecx, 3 ; do count <= 3 rep stosb ; optimal with <= 3 pop edi ret</pre>

Memory routines in the runtime library generated by Intel compilers are optimized across a wide range of address alignments, counter values, and microarchitectures. In most cases, applications should take advantage of the default memory routines provided by Intel compilers.

In some situations, the byte count of the data is known by the context (as opposed to being known by a parameter passed from a call), and one can take a simpler approach than those required for a general-purpose library routine. For example, if the byte count is also small, using REP MOVSB/STOSB with a count less than four can ensure good address alignment and loop-unrolling to finish the remaining data; using MOVSD/STOSD can reduce the overhead associated with iteration.

Using a REP prefix with string move instructions can provide high performance in the situations described above. However, using a REP prefix with string scan instructions (SCASB, SCASW, SCASD, SCASQ) or compare instructions (CMPSB, CMPSW, SMPD, SMPDQ) is not recommended for high performance. Consider using SIMD instructions instead.

3.8 FLOATING-POINT CONSIDERATIONS

When programming floating-point applications, it is best to start with a high-level programming language such as C, C++, or Fortran. Many compilers perform floating-point scheduling and optimization when it is possible. However in order to produce optimal code, the compiler may need some assistance.

3.8.1 Guidelines for Optimizing Floating-point Code

User/Source Coding Rule 13. (M impact, M generality) *Enable the compiler's use of SSE, SSE2 or SSE3 instructions with appropriate switches.*

Follow this procedure to investigate the performance of your floating-point application:

- Understand how the compiler handles floating-point code.
- Look at the assembly dump and see what transforms are already performed on the program.
- Study the loop nests in the application that dominate the execution time.
- Determine why the compiler is not creating the fastest code.
- See if there is a dependence that can be resolved.
- Determine the problem area: bus bandwidth, cache locality, trace cache bandwidth, or instruction latency. Focus on optimizing the problem area. For example, adding PREFETCH instructions will not help if the bus is already saturated. If trace cache bandwidth is the problem, added prefetch `µops` may degrade performance.

Also, in general, follow the general coding recommendations discussed in this chapter, including:

- Blocking the cache
- Using prefetch

- Enabling vectorization
- Unrolling loops

User/Source Coding Rule 14. (H impact, ML generality) *Make sure your application stays in range to avoid denormal values, underflows..*

Out-of-range numbers cause very high overhead.

User/Source Coding Rule 15. (M impact, ML generality) *Do not use double precision unless necessary. Set the precision control (PC) field in the x87 FPU control word to "Single Precision". This allows single precision (32-bit) computation to complete faster on some operations (for example, divides due to early out). However, be careful of introducing more than a total of two values for the floating point control word, or there will be a large performance penalty. See Section 3.8.3.*

User/Source Coding Rule 16. (H impact, ML generality) *Use fast float-to-int routines, FISTTP, or SSE2 instructions. If coding these routines, use the FISTTP instruction if SSE3 is available, or the CVTTSS2SI and CVTTSD2SI instructions if coding with Streaming SIMD Extensions 2.*

Many libraries generate X87 code that does more work than is necessary. The FISTTP instruction in SSE3 can convert floating-point values to 16-bit, 32-bit, or 64-bit integers using truncation without accessing the floating-point control word (FCW). The instructions CVTTSS2SI and CVTTSD2SI save many μ ops and some store-forwarding delays over some compiler implementations. This avoids changing the rounding mode.

User/Source Coding Rule 17. (M impact, ML generality) *Removing data dependence enables the out-of-order engine to extract more ILP from the code. When summing up the elements of an array, use partial sums instead of a single accumulator..*

For example, to calculate $z = a + b + c + d$, instead of:

```
X = A + B;
Y = X + C;
Z = Y + D;
```

use:

```
X = A + B;
Y = C + D;
Z = X + Y;
```

User/Source Coding Rule 18. (M impact, ML generality) *Usually, math libraries take advantage of the transcendental instructions (for example, FSIN) when evaluating elementary functions. If there is no critical need to evaluate the transcendental functions using the extended precision of 80 bits, applications should consider an alternate, software-based approach, such as a look-up-table-based algorithm using interpolation techniques. It is possible to improve*

transcendental performance with these techniques by choosing the desired numeric precision and the size of the look-up table, and by taking advantage of the parallelism of the SSE and the SSE2 instructions.

3.8.2 Floating-point Modes and Exceptions

When working with floating-point numbers, high-speed microprocessors frequently must deal with situations that need special handling in hardware or code.

3.8.2.1 Floating-point Exceptions

The most frequent cause of performance degradation is the use of masked floating-point exception conditions such as:

- arithmetic overflow
- arithmetic underflow
- denormalized operand

Refer to Chapter 4 of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1*, for definitions of overflow, underflow and denormal exceptions.

Denormalized floating-point numbers impact performance in two ways:

- directly when are used as operands
- indirectly when are produced as a result of an underflow situation

If a floating-point application never underflows, the denormals can only come from floating-point constants.

User/Source Coding Rule 19. (H impact, ML generality) *Denormalized floating-point constants should be avoided as much as possible.*

Denormal and arithmetic underflow exceptions can occur during the execution of x87 instructions or SSE/SSE2/SSE3 instructions. Processors based on Intel NetBurst microarchitecture handle these exceptions more efficiently when executing SSE/SSE2/SSE3 instructions and when speed is more important than complying with the IEEE standard. The following paragraphs give recommendations on how to optimize your code to reduce performance degradations related to floating-point exceptions.

3.8.2.2 Dealing with floating-point exceptions in x87 FPU code

Every special situation listed in Section 3.8.2.1, "Floating-point Exceptions," is costly in terms of performance. For that reason, x87 FPU code should be written to avoid these situations.

There are basically three ways to reduce the impact of overflow/underflow situations with x87 FPU code:

- Choose floating-point data types that are large enough to accommodate results without generating arithmetic overflow and underflow exceptions.
- Scale the range of operands/results to reduce as much as possible the number of arithmetic overflow/underflow situations.
- Keep intermediate results on the x87 FPU register stack until the final results have been computed and stored in memory. Overflow or underflow is less likely to happen when intermediate results are kept in the x87 FPU stack (this is because data on the stack is stored in double extended-precision format and overflow/underflow conditions are detected accordingly).
- Denormalized floating-point constants (which are read-only, and hence never change) should be avoided and replaced, if possible, with zeros of the same sign.

3.8.2.3 Floating-point Exceptions in SSE/SSE2/SSE3 Code

Most special situations that involve masked floating-point exceptions are handled efficiently in hardware. When a masked overflow exception occurs while executing SSE/SSE2/SSE3 code, processor hardware can handle it without performance penalty.

Underflow exceptions and denormalized source operands are usually treated according to the IEEE 754 specification, but this can incur significant performance delay. If a programmer is willing to trade pure IEEE 754 compliance for speed, two non-IEEE 754 compliant modes are provided to speed situations where underflows and input are frequent: FTZ mode and DAZ mode.

When the FTZ mode is enabled, an underflow result is automatically converted to a zero with the correct sign. Although this behavior is not compliant with IEEE 754, it is provided for use in applications where performance is more important than IEEE 754 compliance. Since denormal results are not produced when the FTZ mode is enabled, the only denormal floating-point numbers that can be encountered in FTZ mode are the ones specified as constants (read only).

The DAZ mode is provided to handle denormal source operands efficiently when running a SIMD floating-point application. When the DAZ mode is enabled, input denormals are treated as zeros with the same sign. Enabling the DAZ mode is the way to deal with denormal floating-point constants when performance is the objective.

If departing from the IEEE 754 specification is acceptable and performance is critical, run SSE/SSE2/SSE3 applications with FTZ and DAZ modes enabled.

NOTE

The DAZ mode is available with both the SSE and SSE2 extensions, although the speed improvement expected from this mode is fully realized only in SSE code.

3.8.3 Floating-point Modes

On the Pentium III processor, the FLDCW instruction is an expensive operation. On early generations of Pentium 4 processors, FLDCW is improved only for situations where an application alternates between two constant values of the x87 FPU control word (FCW), such as when performing conversions to integers. On Pentium M, Intel Core Solo, Intel Core Duo and Intel Core 2 Duo processors, FLDCW is improved over previous generations.

Specifically, the optimization for FLDCW in the first two generations of Pentium 4 processors allow programmers to alternate between two constant values efficiently. For the FLDCW optimization to be effective, the two constant FCW values are only allowed to differ on the following 5 bits in the FCW:

FCW[8-9] ; Precision control
FCW[10-11] ; Rounding control
FCW[12] ; Infinity control

If programmers need to modify other bits (for example: mask bits) in the FCW, the FLDCW instruction is still an expensive operation.

In situations where an application cycles between three (or more) constant values, FLDCW optimization does not apply, and the performance degradation occurs for each FLDCW instruction.

One solution to this problem is to choose two constant FCW values, take advantage of the optimization of the FLDCW instruction to alternate between only these two constant FCW values, and devise some means to accomplish the task that requires the 3rd FCW value without actually changing the FCW to a third constant value. An alternative solution is to structure the code so that, for periods of time, the application alternates between only two constant FCW values. When the application later alternates between a pair of different FCW values, the performance degradation occurs only during the transition.

It is expected that SIMD applications are unlikely to alternate between FTZ and DAZ mode values. Consequently, the SIMD control word does not have the short latencies that the floating-point control register does. A read of the MXCSR register has a fairly long latency, and a write to the register is a serializing instruction.

There is no separate control word for single and double precision; both use the same modes. Notably, this applies to both FTZ and DAZ modes.

Assembly/Compiler Coding Rule 59. (H impact, M generality) *Minimize changes to bits 8-12 of the floating point control word. Changes for more than two values (each value being a combination of the following bits: precision, rounding and infinity control, and the rest of bits in FCW) leads to delays that are on the order of the pipeline depth.*

3.8.3.1 Rounding Mode

Many libraries provide float-to-integer library routines that convert floating-point values to integer. Many of these libraries conform to ANSI C coding standards which

state that the rounding mode should be truncation. With the Pentium 4 processor, one can use the CVTTSD2SI and CVTTSS2SI instructions to convert operands with truncation without ever needing to change rounding modes. The cost savings of using these instructions over the methods below is enough to justify using SSE and SSE2 wherever possible when truncation is involved.

For x87 floating point, the FIST instruction uses the rounding mode represented in the floating-point control word (FCW). The rounding mode is generally “round to nearest”, so many compiler writers implement a change in the rounding mode in the processor in order to conform to the C and FORTRAN standards. This implementation requires changing the control word on the processor using the FLDCW instruction. For a change in the rounding, precision, and infinity bits, use the FSTCW instruction to store the floating-point control word. Then use the FLDCW instruction to change the rounding mode to truncation.

In a typical code sequence that changes the rounding mode in the FCW, a FSTCW instruction is usually followed by a load operation. The load operation from memory should be a 16-bit operand to prevent store-forwarding problem. If the load operation on the previously-stored FCW word involves either an 8-bit or a 32-bit operand, this will cause a store-forwarding problem due to mismatch of the size of the data between the store operation and the load operation.

To avoid store-forwarding problems, make sure that the write and read to the FCW are both 16-bit operations.

If there is more than one change to the rounding, precision, and infinity bits, and the rounding mode is not important to the result, use the algorithm in Example 3-55 to avoid synchronization issues, the overhead of the FLDCW instruction, and having to change the rounding mode. Note that the example suffers from a store-forwarding problem which will lead to a performance penalty. However, its performance is still better than changing the rounding, precision, and infinity bits among more than two values.

Example 3-55. Algorithm to Avoid Changing Rounding Mode

```

_fto132proc
  lea    ecx, [esp-8]
  sub    esp, 16          ; Allocate frame
  and    ecx, -8         ; Align pointer on boundary of 8
  fld    st(0)           ; Duplicate FPU stack top

  fistp  qword ptr[ecx]
  fld    qword ptr[ecx]
  mov    edx, [ecx+4]    ; High DWORD of integer
  mov    eax, [ecx]      ; Low DWIRD of integer
  test   eax, eax
  je     integer_QnaN_or_zero

```

Example 3-55. Algorithm to Avoid Changing Rounding Mode (Contd.)

```

arg_is_not_integer_QNaN:
    fsubp    st(1), st          ; TOS=d-round(d), { st(1) = st(1)-st & pop ST}
    test    edx, edx          ; What's sign of integer
    jns     positive          ; Number is negative
    fstp    dword ptr[ecx]    ; Result of subtraction
    mov     ecx, [ecx]        ; DWORD of diff(single-precision)
    add     esp, 16
    xor     ecx, 80000000h
    add     ecx, 7fffffffh    ; If diff<0 then decrement integer
    adc     eax, 0             ; INC EAX (add CARRY flag)
    ret

positive:
    positive:
    fstp    dword ptr[ecx]    ; 17-18 result of subtraction
    mov     ecx, [ecx]        ; DWORD of diff(single precision)
    add     esp, 16
    add     ecx, 7fffffffh    ; If diff<0 then decrement integer
    sbb     eax, 0            ; DEC EAX (subtract CARRY flag)
    ret

integer_QNaN_or_zero:
    test    edx, 7fffffffh
    jnz     arg_is_not_integer_QNaN
    add     esp, 16
    ret

```

Assembly/Compiler Coding Rule 60. (H impact, L generality) *Minimize the number of changes to the rounding mode. Do not use changes in the rounding mode to implement the floor and ceiling functions if this involves a total of more than two values of the set of rounding, precision, and infinity bits.*

3.8.3.2 Precision

If single precision is adequate, use it instead of double precision. This is true because:

- Single precision operations allow the use of longer SIMD vectors, since more single precision data elements can fit in a register.
- If the precision control (PC) field in the x87 FPU control word is set to single precision, the floating-point divider can complete a single-precision computation much faster than either a double-precision computation or an extended double-precision computation. If the PC field is set to double precision, this will enable those x87 FPU operations on double-precision data to complete faster than

extended double-precision computation. These characteristics affect computations including floating-point divide and square root.

Assembly/Compiler Coding Rule 61. (H impact, L generality) *Minimize the number of changes to the precision mode.*

3.8.3.3 Improving Parallelism and the Use of FXCH

The x87 instruction set relies on the floating point stack for one of its operands. If the dependence graph is a tree, which means each intermediate result is used only once and code is scheduled carefully, it is often possible to use only operands that are on the top of the stack or in memory, and to avoid using operands that are buried under the top of the stack. When operands need to be pulled from the middle of the stack, an FXCH instruction can be used to swap the operand on the top of the stack with another entry in the stack.

The FXCH instruction can also be used to enhance parallelism. Dependent chains can be overlapped to expose more independent instructions to the hardware scheduler. An FXCH instruction may be required to effectively increase the register name space so that more operands can be simultaneously live.

In processors based on Intel NetBurst microarchitecture, however, that FXCH inhibits issue bandwidth in the trace cache. It does this not only because it consumes a slot, but also because of issue slot restrictions imposed on FXCH. If the application is not bound by issue or retirement bandwidth, FXCH will have no impact.

The effective instruction window size in processors based on Intel NetBurst microarchitecture is large enough to permit instructions that are as far away as the next iteration to be overlapped. This often obviates the need to use FXCH to enhance parallelism.

The FXCH instruction should be used only when it's needed to express an algorithm or to enhance parallelism. If the size of register name space is a problem, the use of XMM registers is recommended.

Assembly/Compiler Coding Rule 62. (M impact, M generality) *Use FXCH only where necessary to increase the effective name space.*

This in turn allows instructions to be reordered and made available for execution in parallel. Out-of-order execution precludes the need for using FXCH to move instructions for very short distances.

3.8.4 x87 vs. Scalar SIMD Floating-point Trade-offs

There are a number of differences between x87 floating-point code and scalar floating-point code (using SSE and SSE2). The following differences should drive decisions about which registers and instructions to use:

- When an input operand for a SIMD floating-point instruction contains values that are less than the representable range of the data type, a denormal exception occurs. This causes a significant performance penalty. An SIMD floating-point

operation has a flush-to-zero mode in which the results will not underflow. Therefore subsequent computation will not face the performance penalty of handling denormal input operands. For example, in the case of 3D applications with low lighting levels, using flush-to-zero mode can improve performance by as much as 50% for applications with large numbers of underflows.

- Scalar floating-point SIMD instructions have lower latencies than equivalent x87 instructions. Scalar SIMD floating-point multiply instruction may be pipelined, while x87 multiply instruction is not.
- Only x87 supports transcendental instructions.
- x87 supports 80-bit precision, double extended floating point. SSE support a maximum of 32-bit precision. SSE2 supports a maximum of 64-bit precision.
- Scalar floating-point registers may be accessed directly, avoiding FXCH and top-of-stack restrictions.
- The cost of converting from floating point to integer with truncation is significantly lower with Streaming SIMD Extensions 2 and Streaming SIMD Extensions in the processors based on Intel NetBurst microarchitecture than with either changes to the rounding mode or the sequence prescribed in the Example 3-55.

Assembly/Compiler Coding Rule 63. (M impact, M generality) *Use Streaming SIMD Extensions 2 or Streaming SIMD Extensions unless you need an x87 feature. Most SSE2 arithmetic operations have shorter latency than their X87 counterpart and they eliminate the overhead associated with the management of the X87 register stack.*

3.8.4.1 Scalar SSE/SSE2 Performance on Intel® Core™ Solo and Intel® Core™ Duo Processors

On Intel Core Solo and Intel Core Duo processors, the combination of improved decoding and μ op fusion allows instructions which were formerly two, three, and four μ ops to go through all decoders. As a result, scalar SSE/SSE2 code can match the performance of x87 code executing through two floating-point units. On Pentium M processors, scalar SSE/SSE2 code can experience approximately 30% performance degradation relative to x87 code executing through two floating-point units.

In code sequences that have conversions from floating-point to integer, divide single-precision instructions, or any precision change, x87 code generation from a compiler typically writes data to memory in single-precision and reads it again in order to reduce precision. Using SSE/SSE2 scalar code instead of x87 code can generate a large performance benefit using Intel NetBurst microarchitecture and a modest benefit on Intel Core Solo and Intel Core Duo processors.

Recommendation: Use the compiler switch to generate SSE2 scalar floating-point code rather than x87 code.

When working with scalar SSE/SSE2 code, pay attention to the need for clearing the content of unused slots in an XMM register and the associated performance impact.

For example, loading data from memory with MOVSS or MOVSD causes an extra micro-op for zeroing the upper part of the XMM register.

On Pentium M, Intel Core Solo, and Intel Core Duo processors, this penalty can be avoided by using MOVLPD. However, using MOVLPD causes a performance penalty on Pentium 4 processors.

Another situation occurs when mixing single-precision and double-precision code. On processors based on Intel NetBurst microarchitecture, using CVTSS2SD has performance penalty relative to the alternative sequence:

```
XORPS    XMM1, XMM1
MOVSS    XMM1, XMM2
CVTSS2PD XMM1, XMM1
```

On Intel Core Solo and Intel Core Duo processors, using CVTSS2SD is more desirable than the alternative sequence.

3.8.4.2 x87 Floating-point Operations with Integer Operands

For processors based on Intel NetBurst microarchitecture, splitting floating-point operations (FIADD, FISUB, FIMUL, and FIDIV) that take 16-bit integer operands into two instructions (FILD and a floating-point operation) is more efficient. However, for floating-point operations with 32-bit integer operands, using FIADD, FISUB, FIMUL, and FIDIV is equally efficient compared with using separate instructions.

Assembly/Compiler Coding Rule 64. (M impact, L generality) *Try to use 32-bit operands rather than 16-bit operands for FILD. However, do not do so at the expense of introducing a store-forwarding problem by writing the two halves of the 32-bit memory operand separately.*

3.8.4.3 x87 Floating-point Comparison Instructions

The FCOMI and FCMOV instructions should be used when performing x87 floating-point comparisons. Using the FCOM, FCOMP, and FCOMPP instructions typically requires additional instruction like FSTSW. The latter alternative causes more μ ops to be decoded, and should be avoided.

3.8.4.4 Transcendental Functions

If an application needs to emulate math functions in software for performance or other reasons (see Section 3.8.1, "Guidelines for Optimizing Floating-point Code"), it may be worthwhile to inline math library calls because the CALL and the prologue/epilogue involved with such calls can significantly affect the latency of operations.

Note that transcendental functions are supported only in x87 floating point, not in Streaming SIMD Extensions or Streaming SIMD Extensions 2.

3.9 MAXIMIZING PCIE PERFORMANCE

PCIe performance can be dramatically impacted by the size and alignment of upstream reads and writes (read and write transactions issued from a PCIe agent to the host's memory). As a general rule, the best performance, in terms of both bandwidth and latency, is obtained by aligning the start addresses of upstream reads and writes on 64-byte boundaries and ensuring that the request size is a multiple of 64-bytes, with modest further increases in bandwidth when larger multiples (128, 192, 256 bytes) are employed. In particular, a partial write will cause a delay for the following request (read or write).

A second rule is to avoid multiple concurrently outstanding accesses to a single cache line. This can result in a conflict which in turn can cause serialization of accesses that would otherwise be pipelined, resulting in higher latency and/or lower bandwidth. Patterns that violate this rule include sequential accesses (reads or writes) that are not a multiple of 64-bytes, as well as explicit accesses to the same cache line address. Overlapping requests—those with different start addresses but with request lengths that result in overlap of the requests—can have the same effect. For example, a 96-byte read of address 0x00000200 followed by a 64-byte read of address 0x00000240 will cause a conflict—and a likely delay—for the second read.

Upstream writes that are a multiple of 64-byte but are non-aligned will have the performance of a series of partial and full sequential writes. For example, a write of length 128-byte to address 0x00000070 will perform similarly to 3 sequential writes of lengths 16, 64, and 48 to addresses 0x00000070, 0x00000080, and 0x00000100, respectively.

For PCIe cards implementing multi-function devices, such as dual or quad port network interface cards (NICs) or dual-GPU graphics cards, it is important to note that non-optimal behavior by one of those devices can impact the bandwidth and/or latency observed by the other devices on that card. With respect to the behavior described in this section, all traffic on a given PCIe port is treated as if it originated from a single device and function.

For the best PCIe bandwidth:

1. Align start addresses of upstream reads and writes on 64-byte boundaries.
2. Use read and write requests that are a multiple of 64-bytes.
3. Eliminate or avoid sequential and random partial line upstream writes.
4. Eliminate or avoid conflicting upstream reads, including sequential partial line reads.

Techniques for avoiding performance pitfalls include cache line aligning all descriptors and data buffers, padding descriptors that are written upstream to 64-byte alignment, buffering incoming data to achieve larger upstream write payloads, allocating data structures intended for sequential reading by the PCIe device in such a way as to enable use of (multiple of) 64-byte reads. The negative impact of unopti-

mized reads and writes depends on the specific workload and the microarchitecture on which the product is based.

CHAPTER 4

CODING FOR SIMD ARCHITECTURES

Processors based on Intel Core microarchitecture supports MMX, SSE, SSE2, SSE3, and SSSE3. Processors based on Enhanced Intel Core microarchitecture supports MMX, SSE, SSE2, SSE3, SSSE3 and SSE4.1. Processors based on Intel microarchitecture code name Nehalem supports MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1 and SSE4.2. Processors based on Intel microarchitecture code name Westmere supports MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2 and AESNI. Processors based on Intel microarchitecture code name Sandy Bridge supports MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, AESNI, PCLMULQDQ and Intel AVX.

Intel Pentium 4, Intel Xeon and Pentium M processors include support for SSE2, SSE, and MMX technology. SSE3 were introduced with the Pentium 4 processor supporting Hyper-Threading Technology at 90 nm technology. Intel Core Solo and Intel Core Duo processors support SSE3/SSE2/SSE, and MMX.

Single-instruction, multiple-data (SIMD) technologies enable the development of advanced multimedia, signal processing, and modeling applications.

Single-instruction, multiple-data techniques can be applied to text/string processing, lexing and parser applications. This is covered in Chapter 10, "SSE4.2 and SIMD Programming For Text-Processing/LexING/Parsing". Techniques for optimizing AESNI are discussed in Section 5.10.

To take advantage of the performance opportunities presented by these capabilities, do the following:

- Ensure that the processor supports MMX technology, SSE, SSE2, SSE3, SSSE3 and SSE4.1.
- Ensure that the operating system supports MMX technology and SSE (OS support for SSE2, SSE3 and SSSE3 is the same as OS support for SSE).
- Employ the optimization and scheduling strategies described in this book.
- Use stack and data alignment techniques to keep data properly aligned for efficient memory use.
- Utilize the cacheability instructions offered by SSE and SSE2, where appropriate.

4.1 CHECKING FOR PROCESSOR SUPPORT OF SIMD TECHNOLOGIES

This section shows how to check whether a processor supports MMX technology, SSE, SSE2, SSE3, SSSE3, and SSE4.1.

SIMD technology can be included in your application in three ways:

1. Check for the SIMD technology during installation. If the desired SIMD technology is available, the appropriate DLLs can be installed.
2. Check for the SIMD technology during program execution and install the proper DLLs at runtime. This is effective for programs that may be executed on different machines.
3. Create a “fat” binary that includes multiple versions of routines; versions that use SIMD technology and versions that do not. Check for SIMD technology during program execution and run the appropriate versions of the routines. This is especially effective for programs that may be executed on different machines.

4.1.1 Checking for MMX Technology Support

If MMX technology is available, then CPUID.01H:EDX[BIT 23] = 1. Use the code segment in Example 4-1 to test for MMX technology.

Example 4-1. Identification of MMX Technology with CPUID

```

...identify existence of cpuid instruction
...
...
...
mov eax, 1          ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test edx, 00800000h ; Is MMX technology bit (bit 23) in feature flags equal to 1
jnz    Found
    
```

For more information on CPUID see, *Intel® Processor Identification with CPUID Instruction*, order number 241618.

4.1.2 Checking for Streaming SIMD Extensions Support

Checking for processor support of Streaming SIMD Extensions (SSE) on your processor is similar to checking for MMX technology. However, operating system (OS) must provide support for SSE states save and restore on context switches to ensure consistent application behavior when using SSE instructions.

To check whether your system supports SSE, follow these steps:

1. Check that your processor supports the CPUID instruction.
2. Check the feature bits of CPUID for SSE existence.

Example 4-2 shows how to find the SSE feature bit (bit 25) in CPUID feature flags.

Example 4-2. Identification of SSE with CPUID

```

...Identify existence of cpuid instruction
                                ; Identify signature is genuine intel
mov eax, 1                        ; Request for feature flags
cpuid                             ; 0FH, 0A2H cpuid instruction
test EDX, 002000000h             ; Bit 25 in feature flags equal to 1
jnz      Found

```

4.1.3 Checking for Streaming SIMD Extensions 2 Support

Checking for support of SSE2 is like checking for SSE support. The OS requirements for SSE2 Support are the same as the OS requirements for SSE.

To check whether your system supports SSE2, follow these steps:

1. Check that your processor has the CPUID instruction.
2. Check the feature bits of CPUID for SSE2 technology existence.

Example 4-3 shows how to find the SSE2 feature bit (bit 26) in the CPUID feature flags.

Example 4-3. Identification of SSE2 with cpuid

```

...Identify existence of cpuid instruction
...                               ; Identify signature is genuine intel
mov eax, 1                        ; Request for feature flags
cpuid                             ; 0FH, 0A2H CPUID instruction
test EDX, 004000000h             ; Bit 26 in feature flags equal to 1
jnz      Found

```

4.1.4 Checking for Streaming SIMD Extensions 3 Support

SSE3 includes 13 instructions, 11 of those are suited for SIMD or x87 style programming. Checking for support of SSE3 instructions is similar to checking for SSE support. The OS requirements for SSE3 Support are the same as the requirements for SSE.

To check whether your system supports the x87 and SIMD instructions of SSE3, follow these steps:

1. Check that your processor has the CPUID instruction.
2. Check the ECX feature bit 0 of CPUID for SSE3 technology existence.

Example 4-4 shows how to find the SSE3 feature bit (bit 0 of ECX) in the CPUID feature flags.

Example 4-4. Identification of SSE3 with CPUID

```

...Identify existence of cpuid instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 00000001h ; Bit 0 in feature flags equal to 1
jnz    Found
    
```

Software must check for support of MONITOR and MWAIT before attempting to use MONITOR and MWAIT. Detecting the availability of MONITOR and MWAIT can be done using a code sequence similar to Example 4-4. The availability of MONITOR and MWAIT is indicated by bit 3 of the returned value in ECX.

4.1.5 Checking for Supplemental Streaming SIMD Extensions 3 Support

Checking for support of SSSE3 is similar to checking for SSE support. The OS requirements for SSSE3 support are the same as the requirements for SSE.

To check whether your system supports SSSE3, follow these steps:

1. Check that your processor has the CPUID instruction.
2. Check the feature bits of CPUID for SSSE3 technology existence.

Example 4-5 shows how to find the SSSE3 feature bit in the CPUID feature flags.

Example 4-5. Identification of SSSE3 with cpuid

```

...Identify existence of CPUID instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 000000200h ; ECX bit 9
jnz    Found
    
```

4.1.6 Checking for SSE4.1 Support

Checking for support of SSE4.1 is similar to checking for SSE support. The OS requirements for SSE4.1 support are the same as the requirements for SSE.

To check whether your system supports SSE4.1, follow these steps:

1. Check that your processor has the CPUID instruction.

2. Check the feature bit of CPUID for SSE4.1.

Example 4-6 shows how to find the SSE4.1 feature bit in the CPUID feature flags.

Example 4-6. Identification of SSE4.1 with cpuid

```

...Identify existence of CPUID instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 00008000h ; ECX bit 19
jnz    Found

```

4.1.7 Checking for SSE4.2 Support

Checking for support of SSE4.2 is similar to checking for SSE support. The OS requirements for SSE4.2 support are the same as the requirements for SSE.

To check whether your system supports SSE4.2, follow these steps:

1. Check that your processor has the CPUID instruction.
2. Check the feature bit of CPUID for SSE4.2.

Example 4-7 shows how to find the SSE4.1 feature bit in the CPUID feature flags.

Example 4-7. Identification of SSE4.2 with cpuid

```

...Identify existence of CPUID instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 00010000h ; ECX bit 20
jnz    Found

```

4.1.8 DetectiON of PCLMULQDQ and AESNI Instructions

Before an application attempts to use the following AESNI instructions: AESDEC/AESDECLAST/AESENCA/AESENCLAST/AESIMC/AESKEYGENASSIST, it must check that the processor supports the AESNI extensions. AESNI extensions is supported if CPUID.01H:ECX.AESNI[bit 25] = 1.

Prior to using PCLMULQDQ instruction, application must check if CPUID.01H:ECX.PCLMULQDQ[bit 1] = 1.

Operating systems that support handling SSE state will also support applications that use AESNI extensions and PCLMULQDQ instruction. This is the same requirement for SSE2, SSE3, SSSE3, and SSE4.

Example 4-8. Detection of AESNI Instructions

```

...Identify existence of CPUID instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 00200000h ; ECX bit 25
jnz    Found
    
```

Example 4-9. Detection of PCLMULQDQ Instruction

```

...Identify existence of CPUID instruction
...                ; Identify signature is genuine intel
mov eax, 1         ; Request for feature flags
cpuid              ; 0FH, 0A2H CPUID instruction
test ECX, 000000002h ; ECX bit 1
jnz    Found
    
```

4.1.9 Detection of AVX Instructions

Intel AVX operates on the 256-bit YMM register state. Application detection of new instruction extensions operating on the YMM state follows the general procedural flow in Figure 4-1.

Prior to using AVX, the application must identify that the operating system supports the XGETBV instruction, the YMM register state, in addition to processor’s support for YMM state management using XSAVE/XRSTOR and AVX instructions. The following simplified sequence accomplishes both and is strongly recommended.

- 1) Detect CPUID.1:ECX.OSXSAVE[bit 27] = 1 (XGETBV enabled for application use¹)
 - 2) Issue XGETBV and verify that XFEATURE_ENABLED_MASK[2:1] = ‘11b’ (XMM state and YMM state are enabled by OS).
 - 3) detect CPUID.1:ECX.AVX[bit 28] = 1 (AVX instructions supported).
- (Step 3 can be done in any order relative to 1 and 2)

1. If CPUID.01H:ECX.OSXSAVE reports 1, it also indirectly implies the processor supports XSAVE, XRSTOR, XGETBV, processor extended state bit vector XFEATURE_ENALBED_MASK register. Thus an application may streamline the checking of CPUID feature flags for XSAVE and OSXSAVE. XSETBV is a privileged instruction.

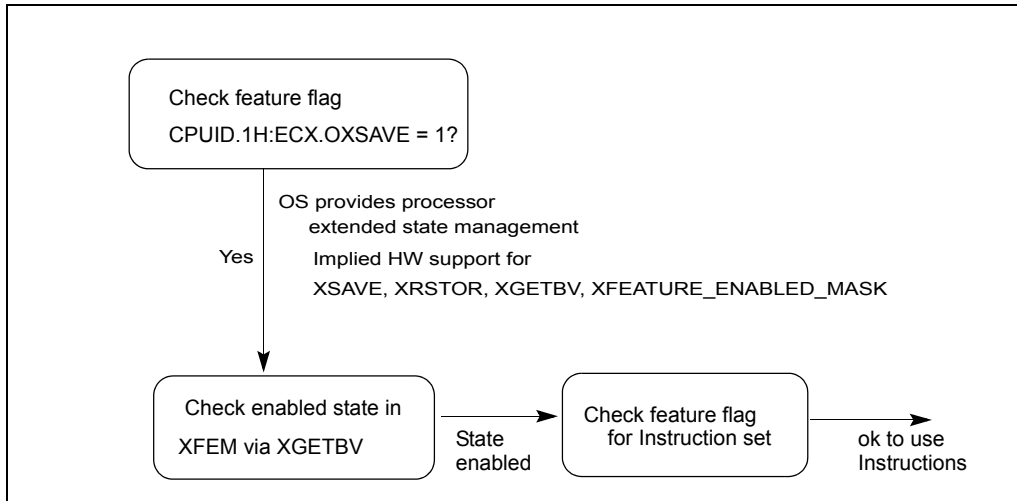


Figure 4-1. General Procedural Flow of Application Detection of AVX

The following pseudocode illustrates this recommended application AVX detection process:

Example 4-10. Detection of AVX Instruction

```

INT supports_AVX()
{
  mov    eax, 1
  cpuid
  and    ecx, 018000000H
  cmp    ecx, 018000000H; check both OSXSAVE and AVX feature flags
  jne    not_supported
; processor supports AVX instructions and XGETBV is enabled by OS
  mov    ecx, 0; specify 0 for XFEATURE_ENABLED_MASK register
  XGETBV    ; result in EDX:EAX
  and    eax, 06H
  cmp    eax, 06H; check OS has enabled both XMM and YMM state support
  jne    not_supported
  mov    eax, 1
  jmp    done
NOT_SUPPORTED:
  mov    eax, 0
done:

```

Note: It is unwise for an application to rely exclusively on CPUID.1:ECX.AVX[bit 28] or at all on CPUID.1:ECX.XSAVE[bit 26]: These indicate hardware support but not operating system support. If YMM state management is not enabled by an operating system, AVX instructions will #UD regardless of CPUID.1:ECX.AVX[bit 28]. "CPUID.1:ECX.XSAVE[bit 26] = 1" does not guarantee the OS actually uses the XSAVE process for state management.

4.1.10 Detection of VEX-Encoded AES and VPCLMULQDQ

VAESDEC/VAESDECLAST/VAESENK/VAESENCLAST/VAESIMC/VAESKEYGENASSIST instructions operate on YMM states. The detection sequence must combine checking for CPUID.1:ECX.AES[bit 25] = 1 and the sequence for detection application support for AVX.

Example 4-11. Detection of VEX-Encoded AESNI Instructions

```

INT supports_VAESNI()
{
    mov     eax, 1
    cpuid
    and     ecx, 01A000000H
    cmp     ecx, 01A000000H; check OSXSAVE AVX and AESNI feature flags
    jne     not_supported
    ; processor supports AVX and VEX-encoded AESNI and XGETBV is enabled by OS
    mov     ecx, 0; specify 0 for XFEATURE_ENABLED_MASK register
    XGETBV     ; result in EDX:EAX
    and     eax, 06H
    cmp     eax, 06H; check OS has enabled both XMM and YMM state support
    jne     not_supported
    mov     eax, 1
    jmp     done
NOT_SUPPORTED:
    mov     eax, 0
done:

```

Similarly, the detection sequence for VPCLMULQDQ must combine checking for CPUID.1:ECX.PCLMULQDQ[bit 1] = 1 and the sequence for detection application support for AVX.

This is shown in the pseudocode:

Example 4-12. Detection of VEX-Encoded AESNI Instructions

```

INT supports_VPCLMULQDQ()
{
    mov     eax, 1
    cpuid

```

Example 4-12. Detection of VEX-Encoded AESNI Instructions

```

and    ecx, 018000002H
cmp    ecx, 018000002H; check OSXSAVE AVX and PCLMULQDQ feature flags
jne    not_supported
; processor supports AVX and VEX-encoded PCLMULQDQ and XGETBV is enabled by OS
mov    ecx, 0; specify 0 for XFEATURE_ENABLED_MASK register
XGETBV    ; result in EDX:EAX
and    eax, 06H
cmp    eax, 06H; check OS has enabled both XMM and YMM state support
jne    not_supported
mov    eax, 1
jmp    done
NOT_SUPPORTED:
mov    eax, 0
done:

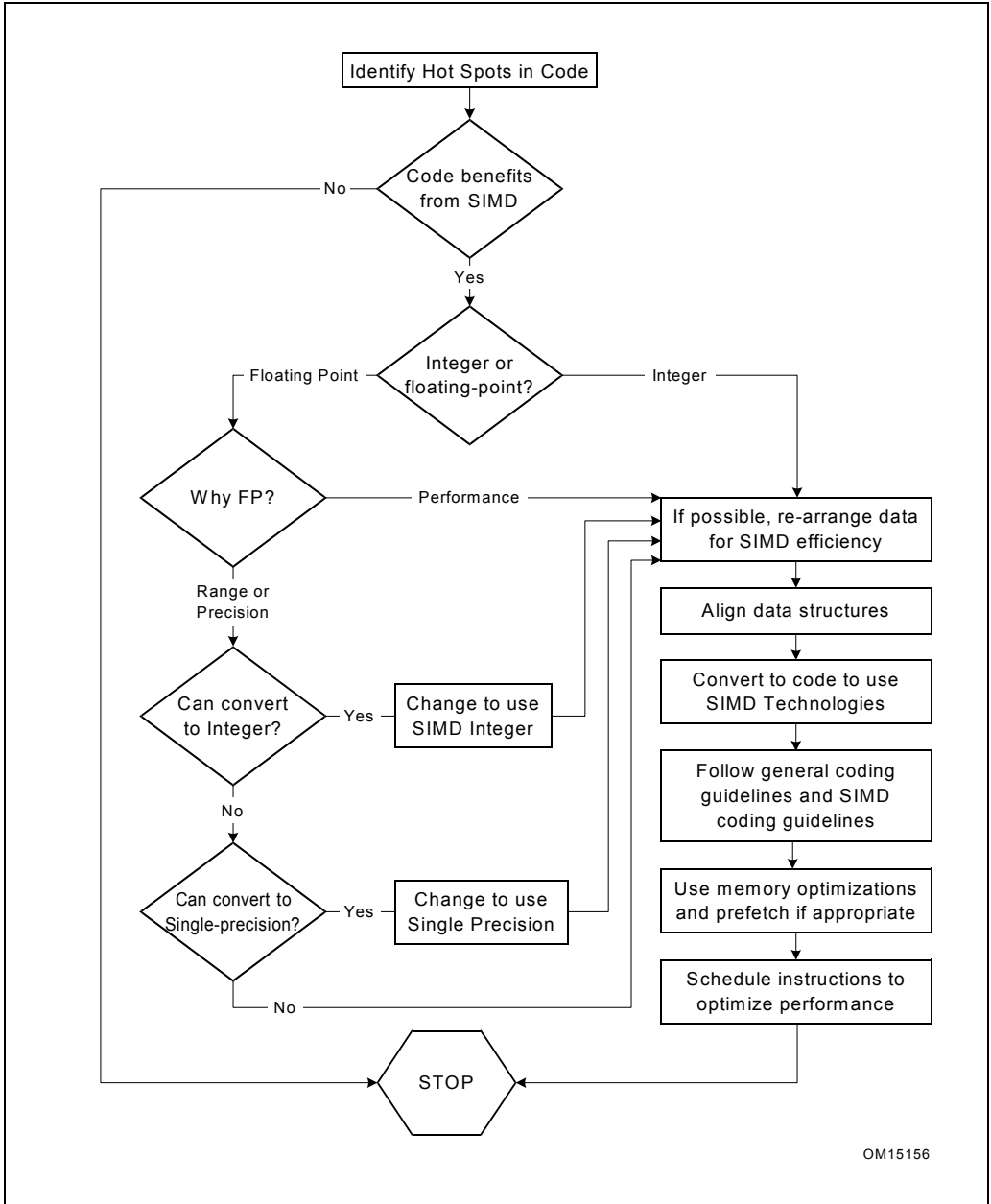
```

4.2 CONSIDERATIONS FOR CODE CONVERSION TO SIMD PROGRAMMING

The VTune Performance Enhancement Environment CD provides tools to aid in the evaluation and tuning. Before implementing them, you need answers to the following questions:

1. Will the current code benefit by using MMX technology, Streaming SIMD Extensions, Streaming SIMD Extensions 2, Streaming SIMD Extensions 3, or Supplemental Streaming SIMD Extensions 3?
2. Is this code integer or floating-point?
3. What integer word size or floating-point precision is needed?
4. What coding techniques should I use?
5. What guidelines do I need to follow?
6. How should I arrange and align the datatypes?

Figure 4-1 provides a flowchart for the process of converting code to MMX technology, SSE, SSE2, SSE3, or SSSE3.



OM15156

Figure 4-2. Converting to Streaming SIMD Extensions Chart

To use any of the SIMD technologies optimally, you must evaluate the following situations in your code:

- Fragments that are computationally intensive
- Fragments that are executed often enough to have an impact on performance
- Fragments that with little data-dependent control flow
- Fragments that require floating-point computations
- Fragments that can benefit from moving data 16 bytes at a time
- Fragments of computation that can coded using fewer instructions
- Fragments that require help in using the cache hierarchy efficiently

4.2.1 Identifying Hot Spots

To optimize performance, use the VTune Performance Analyzer to find sections of code that occupy most of the computation time. Such sections are called the hotspots. See Appendix A, “Application Performance Tools.”

The VTune analyzer provides a hotspots view of a specific module to help you identify sections in your code that take the most CPU time and that have potential performance problems. The hotspots view helps you identify sections in your code that take the most CPU time and that have potential performance problems.

The VTune analyzer enables you to change the view to show hotspots by memory location, functions, classes, or source files. You can double-click on a hotspot and open the source or assembly view for the hotspot and see more detailed information about the performance of each instruction in the hotspot.

The VTune analyzer offers focused analysis and performance data at all levels of your source code and can also provide advice at the assembly language level. The code coach analyzes and identifies opportunities for better performance of C/C++, Fortran and Java* programs, and suggests specific optimizations. Where appropriate, the coach displays pseudo-code to suggest the use of highly optimized intrinsics and functions in the Intel® Performance Library Suite. Because VTune analyzer is designed specifically for Intel architecture (IA)-based processors, including the Pentium 4 processor, it can offer detailed approaches to working with IA. See Appendix A.1.1, “Recommended Optimization Settings for Intel 64 and IA-32 Processors,” for details.

4.2.2 Determine If Code Benefits by Conversion to SIMD Execution

Identifying code that benefits by using SIMD technologies can be time-consuming and difficult. Likely candidates for conversion are applications that are highly computation intensive, such as the following:

- Speech compression algorithms and filters
- Speech recognition algorithms

- Video display and capture routines
- Rendering routines
- 3D graphics (geometry)
- Image and video processing algorithms
- Spatial (3D) audio
- Physical modeling (graphics, CAD)
- Workstation applications
- Encryption algorithms
- Complex arithmetics

Generally, good candidate code is code that contains small-sized repetitive loops that operate on sequential arrays of integers of 8, 16 or 32 bits, single-precision 32-bit floating-point data, double precision 64-bit floating-point data (integer and floating-point data items should be sequential in memory). The repetitiveness of these loops incurs costly application processing time. However, these routines have potential for increased performance when you convert them to use one of the SIMD technologies.

Once you identify your opportunities for using a SIMD technology, you must evaluate what should be done to determine whether the current algorithm or a modified one will ensure the best performance.

4.3 CODING TECHNIQUES

The SIMD features of SSE3, SSE2, SSE, and MMX technology require new methods of coding algorithms. One of them is vectorization. Vectorization is the process of transforming sequentially-executing, or scalar, code into code that can execute in parallel, taking advantage of the SIMD architecture parallelism. This section discusses the coding techniques available for an application to make use of the SIMD architecture.

To vectorize your code and thus take advantage of the SIMD architecture, do the following:

- Determine if the memory accesses have dependencies that would prevent parallel execution.
- “Strip-mine” the inner loop to reduce the iteration count by the length of the SIMD operations (for example, four for single-precision floating-point SIMD, eight for 16-bit integer SIMD on the XMM registers).
- Re-code the loop with the SIMD instructions.

Each of these actions is discussed in detail in the subsequent sections of this chapter. These sections also discuss enabling automatic vectorization using the Intel C++ Compiler.

4.3.1 Coding Methodologies

Software developers need to compare the performance improvement that can be obtained from assembly code versus the cost of those improvements. Programming directly in assembly language for a target platform may produce the required performance gain, however, assembly code is not portable between processor architectures and is expensive to write and maintain.

Performance objectives can be met by taking advantage of the different SIMD technologies using high-level languages as well as assembly. The new C/C++ language extensions designed specifically for SSSE3, SSE3, SSE2, SSE, and MMX technology help make this possible.

Figure 4-2 illustrates the trade-offs involved in the performance of hand-coded assembly versus the ease of programming and portability.

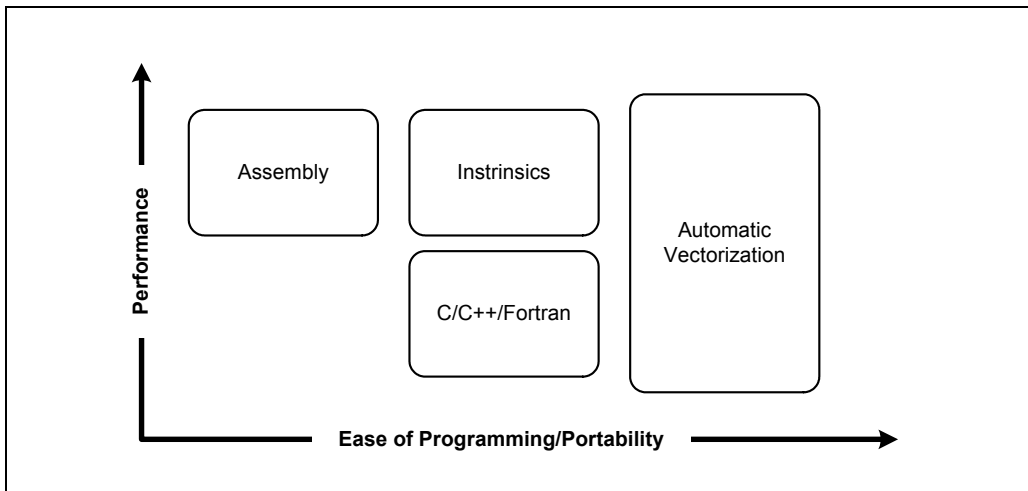


Figure 4-3. Hand-Coded Assembly and High-Level Compiler Performance Trade-offs

The examples that follow illustrate the use of coding adjustments to enable the algorithm to benefit from the SSE. The same techniques may be used for single-precision floating-point, double-precision floating-point, and integer data under SSSE3, SSE3, SSE2, SSE, and MMX technology.

As a basis for the usage model discussed in this section, consider a simple loop shown in Example 4-7.

Example 4-13. Simple Four-Iteration Loop

```
void add(float *a, float *b, float *c)
{
int i;
for (i = 0; i < 4; i++) {
    c[i] = a[i] + b[i];
}
}
```

Note that the loop runs for only four iterations. This allows a simple replacement of the code with Streaming SIMD Extensions.

For the optimal use of the Streaming SIMD Extensions that need data alignment on the 16-byte boundary, all examples in this chapter assume that the arrays passed to the routine, *A*, *B*, *C*, are aligned to 16-byte boundaries by a calling routine. For the methods to ensure this alignment, please refer to the application notes for the Pentium 4 processor.

The sections that follow provide details on the coding methodologies: inlined assembly, intrinsics, C++ vector classes, and automatic vectorization.

4.3.1.1 Assembly

Key loops can be coded directly in assembly language using an assembler or by using inlined assembly (C-asm) in C/C++ code. The Intel compiler or assembler recognize the new instructions and registers, then directly generate the corresponding code. This model offers the opportunity for attaining greatest performance, but this performance is not portable across the different processor architectures.

Example 4-8 shows the Streaming SIMD Extensions inlined assembly encoding.

Example 4-14. Streaming SIMD Extensions Using Inlined Assembly Encoding

```
void add(float *a, float *b, float *c)
{
    __asm {
        mov    eax, a
        mov    edx, b
        mov    ecx, c
        movaps xmm0, XMMWORD PTR [eax]
        addps  xmm0, XMMWORD PTR [edx]
        movaps XMMWORD PTR [ecx], xmm0
    }
}
```

4.3.1.2 Intrinsics

Intrinsics provide the access to the ISA functionality using C/C++ style coding instead of assembly language. Intel has defined three sets of intrinsic functions that are implemented in the Intel C++ Compiler to support the MMX technology, Streaming SIMD Extensions and Streaming SIMD Extensions 2. Four new C data types, representing 64-bit and 128-bit objects are used as the operands of these intrinsic functions. `__M64` is used for MMX integer SIMD, `__M128` is used for single-precision floating-point SIMD, `__M128I` is used for Streaming SIMD Extensions 2 integer SIMD, and `__M128D` is used for double precision floating-point SIMD. These types enable the programmer to choose the implementation of an algorithm directly, while allowing the compiler to perform register allocation and instruction scheduling where possible. The intrinsics are portable among all Intel architecture-based processors supported by a compiler.

The use of intrinsics allows you to obtain performance close to the levels achievable with assembly. The cost of writing and maintaining programs with intrinsics is considerably less. For a detailed description of the intrinsics and their use, refer to the Intel C++ Compiler documentation.

Example 4-9 shows the loop from Example 4-7 using intrinsics.

Example 4-15. Simple Four-Iteration Loop Coded with Intrinsics

```
#include <xmmintrin.h>
void add(float *a, float *b, float *c)
{
    __m128 t0, t1;
    t0 = _mm_load_ps(a);
    t1 = _mm_load_ps(b);
    t0 = _mm_add_ps(t0, t1);
    _mm_store_ps(c, t0);
}
```

The intrinsics map one-to-one with actual Streaming SIMD Extensions assembly code. The XMMINTRIN.H header file in which the prototypes for the intrinsics are defined is part of the Intel C++ Compiler included with the VTune Performance Enhancement Environment CD.

Intrinsics are also defined for the MMX technology ISA. These are based on the `__m64` data type to represent the contents of an mm register. You can specify values in bytes, short integers, 32-bit values, or as a 64-bit object.

The intrinsic data types, however, are not a basic ANSI C data type, and therefore you must observe the following usage restrictions:

- Use intrinsic data types only on the left-hand side of an assignment as a return value or as a parameter. You cannot use it with other arithmetic expressions (for example, "+", ">>").
- Use intrinsic data type objects in aggregates, such as unions to access the byte elements and structures; the address of an `__M64` object may be also used.
- Use intrinsic data type data only with the MMX technology intrinsics described in this guide.

For complete details of the hardware instructions, see the *Intel Architecture MMX Technology Programmer's Reference Manual*. For a description of data types, see the *Intel® 64 and IA-32 Architectures Software Developer's Manual*.

4.3.1.3 Classes

A set of C++ classes has been defined and available in Intel C++ Compiler to provide both a higher-level abstraction and more flexibility for programming with MMX technology, Streaming SIMD Extensions and Streaming SIMD Extensions 2. These classes provide an easy-to-use and flexible interface to the intrinsic functions, allowing developers to write more natural C++ code without worrying about which intrinsic or assembly language instruction to use for a given operation. Since the intrinsic functions underlie the implementation of these C++ classes, the perfor-

mance of applications using this methodology can approach that of one using the intrinsics. Further details on the use of these classes can be found in the *Intel C++ Class Libraries for SIMD Operations User's Guide*, order number 693500.

Example 4-10 shows the C++ code using a vector class library. The example assumes the arrays passed to the routine are already aligned to 16-byte boundaries.

Example 4-16. C++ Code Using the Vector Classes

```
#include <fvec.h>
void add(float *a, float *b, float *c)
{
    F32vec4 *av=(F32vec4 *) a;
    F32vec4 *bv=(F32vec4 *) b;
    F32vec4 *cv=(F32vec4 *) c;
    *cv=*av + *bv;
}
```

Here, `fvec.h` is the class definition file and `F32vec4` is the class representing an array of four floats. The “+” and “=” operators are overloaded so that the actual Streaming SIMD Extensions implementation in the previous example is abstracted out, or hidden, from the developer. Note how much more this resembles the original code, allowing for simpler and faster programming.

Again, the example is assuming the arrays, passed to the routine, are already aligned to 16-byte boundary.

4.3.1.4 Automatic Vectorization

The Intel C++ Compiler provides an optimization mechanism by which loops, such as in Example 4-7 can be automatically vectorized, or converted into Streaming SIMD Extensions code. The compiler uses similar techniques to those used by a programmer to identify whether a loop is suitable for conversion to SIMD. This involves determining whether the following might prevent vectorization:

- The layout of the loop and the data structures used
- Dependencies amongst the data accesses in each iteration and across iterations

Once the compiler has made such a determination, it can generate vectorized code for the loop, allowing the application to use the SIMD instructions.

The caveat to this is that only certain types of loops can be automatically vectorized, and in most cases user interaction with the compiler is needed to fully enable this.

Example 4-11 shows the code for automatic vectorization for the simple four-iteration loop (from Example 4-7).

Example 4-17. Automatic Vectorization for a Simple Loop

```
void add (float *restrict a,
         float *restrict b,
         float *restrict c)
{
    int i;
    for (i = 0; i < 4; i++) {
        c[i] = a[i] + b[i];
    }
}
```

Compile this code using the `-QAX` and `-QRESTRICT` switches of the Intel C++ Compiler, version 4.0 or later.

The `RESTRICT` qualifier in the argument list is necessary to let the compiler know that there are no other aliases to the memory to which the pointers point. In other words, the pointer for which it is used, provides the only means of accessing the memory in question in the scope in which the pointers live. Without the `restrict` qualifier, the compiler will still vectorize this loop using runtime data dependence testing, where the generated code dynamically selects between sequential or vector execution of the loop, based on overlap of the parameters (See documentation for the Intel C++ Compiler). The `restrict` keyword avoids the associated overhead altogether.

See Intel C++ Compiler documentation for details.

4.4 STACK AND DATA ALIGNMENT

To get the most performance out of code written for SIMD technologies data should be formatted in memory according to the guidelines described in this section. Assembly code with an unaligned accesses is a lot slower than an aligned access.

4.4.1 Alignment and Contiguity of Data Access Patterns

The 64-bit packed data types defined by MMX technology, and the 128-bit packed data types for Streaming SIMD Extensions and Streaming SIMD Extensions 2 create more potential for misaligned data accesses. The data access patterns of many algorithms are inherently misaligned when using MMX technology and Streaming SIMD Extensions. Several techniques for improving data access, such as padding, organizing data elements into arrays, etc. are described below. SSE3 provides a special-

purpose instruction LDDQU that can avoid cache line splits is discussed in Section 5.7.1.1, “Supplemental Techniques for Avoiding Cache Line Splits.”

4.4.1.1 Using Padding to Align Data

However, when accessing SIMD data using SIMD operations, access to data can be improved simply by a change in the declaration. For example, consider a declaration of a structure, which represents a point in space plus an attribute.

```
typedef struct {short x,y,z; char a} Point;
Point pt[N];
```

Assume we will be performing a number of computations on X, Y, Z in three of the four elements of a SIMD word; see Section 4.5.1, “Data Structure Layout,” for an example. Even if the first element in array PT is aligned, the second element will start 7 bytes later and not be aligned (3 shorts at two bytes each plus a single byte = 7 bytes).

By adding the padding variable PAD, the structure is now 8 bytes, and if the first element is aligned to 8 bytes (64 bits), all following elements will also be aligned. The sample declaration follows:

```
typedef struct {short x,y,z; char a; char pad;} Point;
Point pt[N];
```

4.4.1.2 Using Arrays to Make Data Contiguous

In the following code,

```
for (i=0; i<N; i++) pt[i].y *= scale;
```

the second dimension Y needs to be multiplied by a scaling value. Here, the FOR loop accesses each Y dimension in the array PT thus disallowing the access to contiguous data. This can degrade the performance of the application by increasing cache misses, by poor utilization of each cache line that is fetched, and by increasing the chance for accesses which span multiple cache lines.

The following declaration allows you to vectorize the scaling operation and further improve the alignment of the data access patterns:

```
short ptx[N], pty[N], ptz[N];
for (i=0; i<N; i++) pty[i] *= scale;
```

With the SIMD technology, choice of data organization becomes more important and should be made carefully based on the operations that will be performed on the data. In some applications, traditional data arrangements may not lead to the maximum performance.

A simple example of this is an FIR filter. An FIR filter is effectively a vector dot product in the length of the number of coefficient taps.

Consider the following code:

```
(data [ j ] *coeff [0] + data [j+1]*coeff [1]+...+data [j+num of taps-1]*coeff [num of taps-1]),
```

If in the code above the filter operation of data element I is the vector dot product that begins at data element J, then the filter operation of data element I+1 begins at data element J+1.

Assuming you have a 64-bit aligned data vector and a 64-bit aligned coefficients vector, the filter operation on the first data element will be fully aligned. For the second data element, however, access to the data vector will be misaligned. For an example of how to avoid the misalignment problem in the FIR filter, refer to Intel application notes on Streaming SIMD Extensions and filters.

Duplication and padding of data structures can be used to avoid the problem of data accesses in algorithms which are inherently misaligned. Section 4.5.1, "Data Structure Layout," discusses trade-offs for organizing data structures.

NOTE

The duplication and padding technique overcomes the misalignment problem, thus avoiding the expensive penalty for misaligned data access, at the cost of increasing the data size. When developing your code, you should consider this tradeoff and use the option which gives the best performance.

4.4.2 Stack Alignment For 128-bit SIMD Technologies

For best performance, the Streaming SIMD Extensions and Streaming SIMD Extensions 2 require their memory operands to be aligned to 16-byte boundaries. Unaligned data can cause significant performance penalties compared to aligned data. However, the existing software conventions for IA-32 (STDCALL, CDECL, FASTCALL) as implemented in most compilers, do not provide any mechanism for ensuring that certain local data and certain parameters are 16-byte aligned. Therefore, Intel has defined a new set of IA-32 software conventions for alignment to support the new `__M128*` datatypes (`__M128`, `__M128D`, and `__M218I`). These meet the following conditions:

- Functions that use Streaming SIMD Extensions or Streaming SIMD Extensions 2 data need to provide a 16-byte aligned stack frame.
- `__M128*` parameters need to be aligned to 16-byte boundaries, possibly creating "holes" (due to padding) in the argument block.

The new conventions presented in this section as implemented by the Intel C++ Compiler can be used as a guideline for an assembly language code as well. In many cases, this section assumes the use of the `__M128*` data types, as defined by the Intel C++ Compiler, which represents an array of four 32-bit floats.

For more details on the stack alignment for Streaming SIMD Extensions and SSE2, see Appendix D, "Stack Alignment."

4.4.3 Data Alignment for MMX Technology

Many compilers enable alignment of variables using controls. This aligns variable bit lengths to the appropriate boundaries. If some of the variables are not appropriately aligned as specified, you can align them using the C algorithm in Example 4-12.

Example 4-18. C Algorithm for 64-bit Data Alignment

```
/* Make newp a pointer to a 64-bit aligned array of NUM_ELEMENTS 64-bit elements. */
double *p, *newp;
p = (double*)malloc (sizeof(double)*(NUM_ELEMENTS+1));
newp = (p+7) & (~0x7);
```

The algorithm in Example 4-12 aligns an array of 64-bit elements on a 64-bit boundary. The constant of 7 is derived from one less than the number of bytes in a 64-bit element, or 8-1. Aligning data in this manner avoids the significant performance penalties that can occur when an access crosses a cache line boundary.

Another way to improve data alignment is to copy the data into locations that are aligned on 64-bit boundaries. When the data is accessed frequently, this can provide a significant performance improvement.

4.4.4 Data Alignment for 128-bit data

Data must be 16-byte aligned when loading to and storing from the 128-bit XMM registers used by SSE/SSE2/SSE3/SSSE3. This must be done to avoid severe performance penalties and, at worst, execution faults.

There are MOVE instructions (and intrinsics) that allow unaligned data to be copied to and out of XMM registers when not using aligned data, but such operations are much slower than aligned accesses. If data is not 16-byte-aligned and the programmer or the compiler does not detect this and uses the aligned instructions, a fault occurs. So keep data 16-byte-aligned. Such alignment also works for MMX technology code, even though MMX technology only requires 8-byte alignment.

The following describes alignment techniques for Pentium 4 processor as implemented with the Intel C++ Compiler.

4.4.4.1 Compiler-Supported Alignment

The Intel C++ Compiler provides the following methods to ensure that the data is aligned.

Alignment by F32vec4 or __m128 Data Types

When the compiler detects F32VEC4 or __M128 data declarations or parameters, it forces alignment of the object to a 16-byte boundary for both global and local data, as well as parameters. If the declaration is within a function, the compiler also aligns

the function's stack frame to ensure that local data and parameters are 16-byte-aligned. For details on the stack frame layout that the compiler generates for both debug and optimized ("release"-mode) compilations, refer to Intel's compiler documentation.

__declspec(align(16)) specifications

These can be placed before data declarations to force 16-byte alignment. This is useful for local or global data declarations that are assigned to 128-bit data types. The syntax for it is

```
__declspec(align(integer-constant))
```

where the INTEGER-CONSTANT is an integral power of two but no greater than 32. For example, the following increases the alignment to 16-bytes:

```
__declspec(align(16)) float buffer[400];
```

The variable BUFFER could then be used as if it contained 100 objects of type __M128 or F32VEC4. In the code below, the construction of the F32VEC4 object, X, will occur with aligned data.

```
void foo() {
    F32vec4 x = *(__m128 *) buffer;
    ...
}
```

Without the declaration of __DECLSPEC(ALIGN(16)), a fault may occur.

Alignment by Using a UNION Structure

When feasible, a UNION can be used with 128-bit data types to allow the compiler to align the data structure by default. This is preferred to forcing alignment with __DECLSPEC(ALIGN(16)) because it exposes the true program intent to the compiler in that __M128 data is being used. For example:

```
union {
    float f[400];
    __m128 m[100];
} buffer;
```

Now, 16-byte alignment is used by default due to the __M128 type in the UNION; it is not necessary to use __DECLSPEC(ALIGN(16)) to force the result.

In C++ (but not in C) it is also possible to force the alignment of a CLASS/STRUCT/UNION type, as in the code that follows:

```
struct __declspec(align(16)) my_m128
{
    float f[4];
};
```

If the data in such a CLASS is going to be used with the Streaming SIMD Extensions or Streaming SIMD Extensions 2, it is preferable to use a UNION to make this explicit. In C++, an anonymous UNION can be used to make this more convenient:

```
class my_m128 {
    union {
        __m128 m;
        float f[4];
    };
};
```

Because the UNION is anonymous, the names, M and F, can be used as immediate member names of MY__M128. Note that `__DECLSPEC(ALIGN)` has no effect when applied to a CLASS, STRUCT, or UNION member in either C or C++.

Alignment by Using `__m64` or `DOUBLE` Data

In some cases, the compiler aligns routines with `__M64` or `DOUBLE` data to 16-bytes by default. The command-line switch, `-QSFALIGN16`, limits the compiler so that it only performs this alignment on routines that contain 128-bit data. The default behavior is to use `-QSFALIGN8`. This switch instructs the compiler to align routines with 8- or 16-byte data types to 16 bytes.

For more, see the Intel C++ Compiler documentation.

4.5 IMPROVING MEMORY UTILIZATION

Memory performance can be improved by rearranging data and algorithms for SSE, SSE2, and MMX technology intrinsics. Methods for improving memory performance involve working with the following:

- Data structure layout
- Strip-mining for vectorization and memory utilization
- Loop-blocking

Using the cacheability instructions, prefetch and streaming store, also greatly enhance memory utilization. See also: Chapter 7, "Optimizing Cache Usage."

4.5.1 Data Structure Layout

For certain algorithms, like 3D transformations and lighting, there are two basic ways to arrange vertex data. The traditional method is the array of structures (AoS) arrangement, with a structure for each vertex (Example 4-13). However this method does not take full advantage of SIMD technology capabilities.

Example 4-19. AoS Data Structure

```
typedef struct{
    float x,y,z;
    int a,b,c;
    ...
} Vertex;
Vertex Vertices[NumOfVertices];
```

The best processing method for code using SIMD technology is to arrange the data in an array for each coordinate (Example 4-14). This data arrangement is called structure of arrays (SoA).

Example 4-20. SoA Data Structure

```
typedef struct{
    float x[NumOfVertices];
    float y[NumOfVertices];
    float z[NumOfVertices];
    int a[NumOfVertices];
    int b[NumOfVertices];
    int c[NumOfVertices];
    ...
} VerticesList;
VerticesList Vertices;
```

There are two options for computing data in AoS format: perform operation on the data as it stands in AoS format, or re-arrange it (swizzle it) into SoA format dynamically. See Example 4-15 for code samples of each option based on a dot-product computation.

Example 4-21. AoS and SoA Code Samples

```
; The dot product of an array of vectors (Array) and a fixed vector (Fixed) is a
; common operation in 3D lighting operations, where Array = (x0,y0,z0),(x1,y1,z1),...
; and Fixed = (xF,yF,zF)
; A dot product is defined as the scalar quantity d0 = x0*xF + y0*yF + z0*zF.
;
; AoS code
; All values marked DC are "don't-care."
```

Example 4-21. AoS and SoA Code Samples (Contd.)

```

; In the AOS model, the vertices are stored in the xyz format
movaps xmm0, Array      ; xmm0 = DC, x0, y0, z0
movaps xmm1, Fixed      ; xmm1 = DC, xF, yF, zF
mulps  xmm0, xmm1       ; xmm0 = DC, x0*xF, y0*yF, z0*zF
movhps xmm, xmm0        ; xmm = DC, DC, DC, x0*xF

addps  xmm1, xmm0       ; xmm0 = DC, DC, DC,
                        ; x0*xF+z0*zF
movaps xmm2, xmm1       ; xmm2 = DC, DC, DC, y0*yF
shufps xmm2, xmm2, 55h  ; xmm2 = DC, DC, DC, y0*yF
addps  xmm2, xmm1       ; xmm1 = DC, DC, DC,
                        ; x0*xF+y0*yF+z0*zF

; SoA code
; X = x0,x1,x2,x3
; Y = y0,y1,y2,y3
; Z = z0,z1,z2,z3
; A = xF,xF,xF,xF
; B = yF,yF,yF,yF
; C = zF,zF,zF,zF

movaps xmm0, X          ; xmm0 = x0,x1,x2,x3
movaps xmm1, Y          ; xmm1 = y0,y1,y2,y3
movaps xmm2, Z          ; xmm2 = z0,z1,z2,z3
mulps  xmm0, A          ; xmm0 = x0*xF, x1*xF, x2*xF, x3*xF
mulps  xmm1, B          ; xmm1 = y0*yF, y1*yF, y2*yF, y3*yF
mulps  xmm2, C          ; xmm2 = z0*zF, z1*zF, z2*zF, z3*zF
addps  xmm0, xmm1
addps  xmm0, xmm2       ; xmm0 = (x0*xF+y0*yF+z0*zF), ...

```

Performing SIMD operations on the original AoS format can require more calculations and some operations do not take advantage of all SIMD elements available. Therefore, this option is generally less efficient.

The recommended way for computing data in AoS format is to swizzle each set of elements to SoA format before processing it using SIMD technologies. Swizzling can either be done dynamically during program execution or statically when the data structures are generated. See Chapter 5 and Chapter 6 for examples. Performing the swizzle dynamically is usually better than using AoS, but can be somewhat inefficient because there are extra instructions during computation. Performing the swizzle statically, when data structures are being laid out, is best as there is no runtime overhead.

As mentioned earlier, the SoA arrangement allows more efficient use of the parallelism of SIMD technologies because the data is ready for computation in a more optimal vertical manner: multiplying components X0,X1,X2,X3 by XF,XF,XF,XF using

4 SIMD execution slots to produce 4 unique results. In contrast, computing directly on AoS data can lead to horizontal operations that consume SIMD execution slots but produce only a single scalar result (as shown by the many “don’t-care” (DC) slots in Example 4-15).

Use of the SoA format for data structures can lead to more efficient use of caches and bandwidth. When the elements of the structure are not accessed with equal frequency, such as when element x , y , z are accessed ten times more often than the other entries, then SoA saves memory and prevents fetching unnecessary data items a , b , and c .

Example 4-22. Hybrid SoA Data Structure

```

NumOfGroups = NumOfVertices/SIMDwidth
typedef struct{
    float x[SIMDwidth];
    float y[SIMDwidth];
    float z[SIMDwidth];
} VerticesCoordList;
typedef struct{
    int a[SIMDwidth];
    int b[SIMDwidth];
    int c[SIMDwidth];
    ...
} VerticesColorList;
VerticesCoordList VerticesCoord[NumOfGroups];
VerticesColorList VerticesColor[NumOfGroups];

```

Note that SoA can have the disadvantage of requiring more independent memory stream references. A computation that uses arrays X , Y , and Z (see Example 4-14) would require three separate data streams. This can require the use of more prefetches, additional address generation calculations, as well as having a greater impact on DRAM page access efficiency.

There is an alternative: a hybrid SoA approach blends the two alternatives (see Example 4-16). In this case, only 2 separate address streams are generated and referenced: one contains $XXXX, YYYY, ZZZZ, ZZZZ, \dots$ and the other $AAAA, BBBB, CCCC, AAAA, DDDD, \dots$. The approach prevents fetching unnecessary data, assuming the variables X, Y, Z are always used together; whereas the variables A, B, C would also be used together, but not at the same time as X, Y, Z .

The hybrid SoA approach ensures:

- Data is organized to enable more efficient vertical SIMD computation
- Simpler/less address generation than AoS
- Fewer streams, which reduces DRAM page misses

- Use of fewer prefetches, due to fewer streams
- Efficient cache line packing of data elements that are used concurrently.

With the advent of the SIMD technologies, the choice of data organization becomes more important and should be carefully based on the operations to be performed on the data. This will become increasingly important in the Pentium 4 processor and future processors. In some applications, traditional data arrangements may not lead to the maximum performance. Application developers are encouraged to explore different data arrangements and data segmentation policies for efficient computation. This may mean using a combination of AoS, SoA, and Hybrid SoA in a given application.

4.5.2 Strip-Mining

Strip-mining, also known as loop sectioning, is a loop transformation technique for enabling SIMD-encodings of loops, as well as providing a means of improving memory performance. First introduced for vectorizers, this technique consists of the generation of code when each vector operation is done for a size less than or equal to the maximum vector length on a given vector machine. By fragmenting a large loop into smaller segments or strips, this technique transforms the loop structure by:

- Increasing the temporal and spatial locality in the data cache if the data are reusable in different passes of an algorithm.
- Reducing the number of iterations of the loop by a factor of the length of each “vector,” or number of operations being performed per SIMD operation. In the case of Streaming SIMD Extensions, this vector or strip-length is reduced by 4 times: four floating-point data items per single Streaming SIMD Extensions single-precision floating-point SIMD operation are processed. Consider Example 4-17.

Example 4-23. Pseudo-code Before Strip Mining

```
typedef struct _VERTEX {
    float x, y, z, nx, ny, nz, u, v;
} Vertex_rec;

main()
{
    Vertex_rec v[Num];
    ....
    for (i=0; i<Num; i++) {
        Transform(v[i]);
    }
}
```

Example 4-23. Pseudo-code Before Strip Mining (Contd.)

```

    for (i=0; i<Num; i++) {
        Lighting(v[i]);
    }
    ....
}

```

The main loop consists of two functions: transformation and lighting. For each object, the main loop calls a transformation routine to update some data, then calls the lighting routine to further work on the data. If the size of array $V[NUM]$ is larger than the cache, then the coordinates for $V[I]$ that were cached during $TRANSFORM(V[I])$ will be evicted from the cache by the time we do $LIGHTING(V[I])$. This means that $V[I]$ will have to be fetched from main memory a second time, reducing performance.

In Example 4-18, the computation has been strip-mined to a size $STRIP_SIZE$. The value $STRIP_SIZE$ is chosen such that $STRIP_SIZE$ elements of array $V[NUM]$ fit into the cache hierarchy. By doing this, a given element $V[I]$ brought into the cache by $TRANSFORM(V[I])$ will still be in the cache when we perform $LIGHTING(V[I])$, and thus improve performance over the non-strip-mined code.

Example 4-24. Strip Mined Code

```

MAIN()
{
    Vertex_rec v[Num];
    ....
    for (i=0; i < Num; i+=strip_size) {
        FOR (J=I; J < MIN(NUM, I+STRIP_SIZE); J++) {
            TRANSFORM(V[J]);
        }
        FOR (J=I; J < MIN(NUM, I+STRIP_SIZE); J++) {
            LIGHTING(V[J]);
        }
    }
}

```

4.5.3 Loop Blocking

Loop blocking is another useful technique for memory performance optimization. The main purpose of loop blocking is also to eliminate as many cache misses as possible. This technique transforms the memory domain of a given problem into smaller chunks rather than sequentially traversing through the entire memory domain. Each chunk should be small enough to fit all the data for a given computation into the cache, thereby maximizing data reuse. In fact, one can treat loop blocking as strip mining in two or more dimensions. Consider the code in Example 4-17 and access

pattern in Figure 4-3. The two-dimensional array A is referenced in the J (column) direction and then referenced in the I (row) direction (column-major order); whereas array B is referenced in the opposite manner (row-major order). Assume the memory layout is in column-major order; therefore, the access strides of array A and B for the code in Example 4-19 would be 1 and MAX, respectively.

Example 4-25. Loop Blocking

```

A. Original Loop
float A[MAX, MAX], B[MAX, MAX]
for (i=0; i< MAX; i++) {
    for (j=0; j< MAX; j++) {
        A[i,j] = A[i,j] + B[j, i];
    }
}

B. Transformed Loop after Blocking
float A[MAX, MAX], B[MAX, MAX];
for (i=0; i< MAX; i+=block_size) {
    for (j=0; j< MAX; j+=block_size) {
        for (ii=i; ii<i+block_size; ii++) {
            for (jj=j; jj<j+block_size; jj++) {
                A[ii,jj] = A[ii,jj] + B[jj, ii];
            }
        }
    }
}

```

For the first iteration of the inner loop, each access to array B will generate a cache miss. If the size of one row of array A, that is, $A[2, 0:MAX-1]$, is large enough, by the time the second iteration starts, each access to array B will always generate a cache miss. For instance, on the first iteration, the cache line containing $B[0, 0:7]$ will be brought in when $B[0,0]$ is referenced because the float type variable is four bytes and each cache line is 32 bytes. Due to the limitation of cache capacity, this line will be evicted due to conflict misses before the inner loop reaches the end. For the next iteration of the outer loop, another cache miss will be generated while referencing $B[0, 1]$. In this manner, a cache miss occurs when each element of array B is referenced, that is, there is no data reuse in the cache at all for array B.

This situation can be avoided if the loop is blocked with respect to the cache size. In Figure 4-3, a `BLOCK_SIZE` is selected as the loop blocking factor. Suppose that `BLOCK_SIZE` is 8, then the blocked chunk of each array will be eight cache lines (32 bytes each). In the first iteration of the inner loop, $A[0, 0:7]$ and $B[0, 0:7]$ will be brought into the cache. $B[0, 0:7]$ will be completely consumed by the first iteration of the outer loop. Consequently, $B[0, 0:7]$ will only experience one cache miss after applying loop blocking optimization in lieu of eight misses for the original algorithm. As illustrated in Figure 4-3, arrays A and B are blocked into smaller rectangular

chunks so that the total size of two blocked A and B chunks is smaller than the cache size. This allows maximum data reuse.

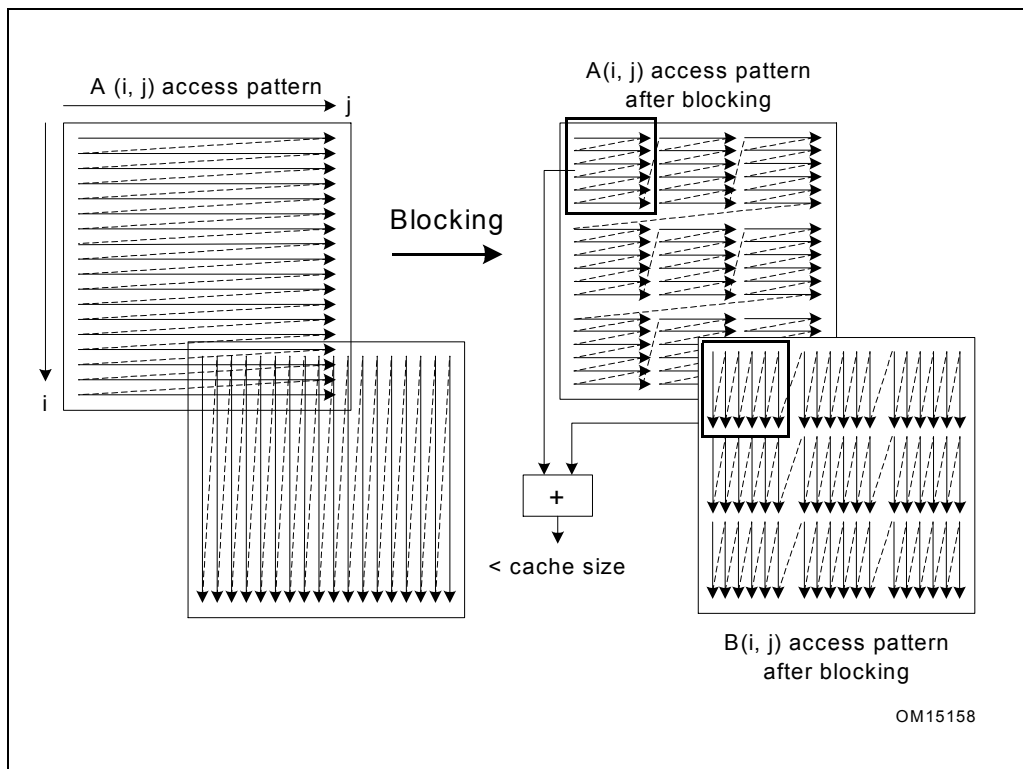


Figure 4-4. Loop Blocking Access Pattern

As one can see, all the redundant cache misses can be eliminated by applying this loop blocking technique. If MAX is huge, loop blocking can also help reduce the penalty from DTLB (data translation look-aside buffer) misses. In addition to improving the cache/memory performance, this optimization technique also saves external bus bandwidth.

4.6 INSTRUCTION SELECTION

The following section gives some guidelines for choosing instructions to complete a task.

One barrier to SIMD computation can be the existence of data-dependent branches. Conditional moves can be used to eliminate data-dependent branches. Conditional

moves can be emulated in SIMD computation by using masked compares and logicals, as shown in Example 4-20. SSE4.1 provides packed blend instruction that can vectorize data-dependent branches in a loop.

Example 4-26. Emulation of Conditional Moves

<pre> High-level code: __declspec(align(16)) short A[MAX_ELEMENT], B[MAX_ELEMENT], C[MAX_ELEMENT], D[MAX_ELEMENT], E[MAX_ELEMENT]; for (i=0; i<MAX_ELEMENT; i++) { if (A[i] > B[i]) { C[i] = D[i]; } else { C[i] = E[i]; } } </pre>
<pre> MMX assembly code processes 4 short values per iteration: xor eax, eax top_of_loop: movq mm0, [A + eax] pcmptw xmm0, [B + eax]; Create compare mask movq mm1, [D + eax] pand mm1, mm0; Drop elements where A<B pandn mm0, [E + eax]; Drop elements where A>B por mm0, mm1; Create single word movq [C + eax], mm0 add eax, 8 cmp eax, MAX_ELEMENT*2 jle top_of_loop </pre>
<pre> SSE4.1 assembly processes 8 short values per iteration: xor eax, eax </pre>
<pre> top_of_loop: movdqq xmm0, [A + eax] pcmptw xmm0, [B + eax]; Create compare mask movdqa xmm1, [E + eax] pblendv xmm1, [D + eax], xmm0; movdqa [C + eax], xmm1; add eax, 16 cmp eax, MAX_ELEMENT*2 jle top_of_loop </pre>

If there are multiple consumers of an instance of a register, group the consumers together as closely as possible. However, the consumers should not be scheduled near the producer.

4.6.1 SIMD Optimizations and Microarchitectures

Pentium M, Intel Core Solo and Intel Core Duo processors have a different microarchitecture than Intel NetBurst microarchitecture. The following sub-section discusses optimizing SIMD code targeting Intel Core Solo and Intel Core Duo processors.

The register-register variant of the following instructions has improved performance on Intel Core Solo and Intel Core Duo processor relative to Pentium M processors. This is because the instructions consist of two micro-ops instead of three. Relevant instructions are: `unpcklps`, `unpckhps`, `packsswb`, `packuswb`, `packssdw`, `pshufd`, `shuffps` and `shuffpd`.

Recommendation: When targeting code generation for Intel Core Solo and Intel Core Duo processors, favor instructions consisting of two μ ops over those with more than two μ ops.

Intel Core microarchitecture generally executes SIMD instructions more efficiently than previous microarchitectures in terms of latency and throughput, most 128-bit SIMD operations have 1 cycle throughput (except shuffle, pack, unpack operations). Many of the restrictions specific to Intel Core Duo, Intel Core Solo processors (such as 128-bit SIMD operations having 2 cycle throughput at a minimum) do not apply to Intel Core microarchitecture. The same is true of Intel Core microarchitecture relative to Intel NetBurst microarchitectures.

Enhanced Intel Core microarchitecture provides dedicated 128-bit shuffler and radix-16 divider hardware. These capabilities and SSE4.1 instructions will make vectorization using 128-bit SIMD instructions even more efficient and effective.

Recommendation: With the proliferation of 128-bit SIMD hardware in Intel Core microarchitecture and Enhanced Intel Core microarchitecture, integer SIMD code written using MMX instructions should consider more efficient implementations using 128-bit SIMD instructions.

4.7 TUNING THE FINAL APPLICATION

The best way to tune your application once it is functioning correctly is to use a profiler that measures the application while it is running on a system. VTune analyzer can help you determine where to make changes in your application to improve performance. Using the VTune analyzer can help you with various phases required for optimized performance. See Appendix A.2, "Intel® VTune™ Performance Analyzer," for details. After every effort to optimize, you should check the performance gains to see where you are making your major optimization gains.

CHAPTER 5

OPTIMIZING FOR SIMD INTEGER APPLICATIONS

SIMD integer instructions provide performance improvements in applications that are integer-intensive and can take advantage of SIMD architecture.

Guidelines in this chapter for using SIMD integer instructions (in addition to those described in Chapter 3) may be used to develop fast and efficient code that scales across processor generations.

The collection of 64-bit and 128-bit SIMD integer instructions supported by MMX technology, SSE, SSE2, SSE3, SSSE3, SSE4.1, and PCMPEQQ in SSE4.2 are referred to as SIMD integer instructions.

Code sequences in this chapter demonstrates the use of basic 64-bit SIMD integer instructions and more efficient 128-bit SIMD integer instructions.

Processors based on Intel Core microarchitecture support MMX, SSE, SSE2, SSE3, and SSSE3. Processors based on Enhanced Intel Core microarchitecture support SSE4.1 and all previous generations of SIMD integer instructions. Processors based on Intel microarchitecture code name Nehalem supports MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1 and SSE4.2.

Single-instruction, multiple-data techniques can be applied to text/string processing, lexing and parser applications. SIMD programming in string/text processing and lexing applications often require sophisticated techniques beyond those commonly used in SIMD integer programming. This is covered in Chapter 10, "SSE4.2 and SIMD Programming For Text-Processing/LexING/Parsing"

Execution of 128-bit SIMD integer instructions in Intel Core microarchitecture and Enhanced Intel Core microarchitecture are substantially more efficient than on previous microarchitectures. Thus newer SIMD capabilities introduced in SSE4.1 operate on 128-bit operands and do not introduce equivalent 64-bit SIMD capabilities. Conversion from 64-bit SIMD integer code to 128-bit SIMD integer code is highly recommended.

This chapter contains examples that will help you to get started with coding your application. The goal is to provide simple, low-level operations that are frequently used. The examples use a minimum number of instructions necessary to achieve best performance on the current generation of Intel 64 and IA-32 processors.

Each example includes a short description, sample code, and notes if necessary. These examples do not address scheduling as it is assumed the examples will be incorporated in longer code sequences.

For planning considerations of using the SIMD integer instructions, refer to Section 4.1.3.

5.1 GENERAL RULES ON SIMD INTEGER CODE

General rules and suggestions are:

- Do not intermix 64-bit SIMD integer instructions with x87 floating-point instructions. See Section 5.2, “Using SIMD Integer with x87 Floating-point.” Note that all SIMD integer instructions can be intermixed without penalty.
- Favor 128-bit SIMD integer code over 64-bit SIMD integer code. On microarchitectures prior to Intel Core microarchitecture, most 128-bit SIMD instructions have two-cycle throughput restrictions due to the underlying 64-bit data path in the execution engine. Intel Core microarchitecture executes most SIMD instructions (except shuffle, pack, unpack operations) with one-cycle throughput and provides three ports to execute multiple SIMD instructions in parallel. Enhanced Intel Core microarchitecture speeds up 128-bit shuffle, pack, unpack operations with 1 cycle throughput.
- When writing SIMD code that works for both integer and floating-point data, use the subset of SIMD convert instructions or load/store instructions to ensure that the input operands in XMM registers contain data types that are properly defined to match the instruction.

Code sequences containing cross-typed usage produce the same result across different implementations but incur a significant performance penalty. Using SSE/SSE2/SSE3/SSSE3/SSE4.1 instructions to operate on type-mismatched SIMD data in the XMM register is strongly discouraged.

- Use the optimization rules and guidelines described in Chapter 3 and Chapter 4.
- Take advantage of hardware prefetcher where possible. Use the PREFETCH instruction only when data access patterns are irregular and prefetch distance can be pre-determined. See Chapter 7, “Optimizing Cache Usage.”
- Emulate conditional moves by using blend, masked compares and logicals instead of using conditional branches.

5.2 USING SIMD INTEGER WITH X87 FLOATING-POINT

All 64-bit SIMD integer instructions use MMX registers, which share register state with the x87 floating-point stack. Because of this sharing, certain rules and considerations apply. Instructions using MMX registers cannot be freely intermixed with x87 floating-point registers. Take care when switching between 64-bit SIMD integer instructions and x87 floating-point instructions to ensure functional correctness. See Section 5.2.1.

Both Section 5.2.1 and Section 5.2.2 apply only to software that employs MMX instructions. As noted before, 128-bit SIMD integer instructions should be favored to replace MMX code and achieve higher performance. That also obviates the need to use EMMS, and the performance penalty of using EMMS when intermixing MMX and X87 instructions.

For performance considerations, there is no penalty of intermixing SIMD floating-point operations and 128-bit SIMD integer operations and x87 floating-point operations.

5.2.1 Using the EMMS Instruction

When generating 64-bit SIMD integer code, keep in mind that the eight MMX registers are aliased to x87 floating-point registers. Switching from MMX instructions to x87 floating-point instructions incurs a finite delay, so it is the best to minimize switching between these instruction types. But when switching, the EMMS instruction provides an efficient means to clear the x87 stack so that subsequent x87 code can operate properly.

As soon as an instruction makes reference to an MMX register, all valid bits in the x87 floating-point tag word are set, which implies that all x87 registers contain valid values. In order for software to operate correctly, the x87 floating-point stack should be emptied when starting a series of x87 floating-point calculations after operating on the MMX registers.

Using EMMS clears all valid bits, effectively emptying the x87 floating-point stack and making it ready for new x87 floating-point operations. The EMMS instruction ensures a clean transition between using operations on the MMX registers and using operations on the x87 floating-point stack. On the Pentium 4 processor, there is a finite overhead for using the EMMS instruction.

Failure to use the EMMS instruction (or the `_MM_EMPTY()` intrinsic) between operations on the MMX registers and x87 floating-point registers may lead to unexpected results.

NOTE

Failure to reset the tag word for FP instructions after using an MMX instruction can result in faulty execution or poor performance.

5.2.2 Guidelines for Using EMMS Instruction

When developing code with both x87 floating-point and 64-bit SIMD integer instructions, follow these steps:

1. Always call the EMMS instruction at the end of 64-bit SIMD integer code when the code transitions to x87 floating-point code.
2. Insert the EMMS instruction at the end of all 64-bit SIMD integer code segments to avoid an x87 floating-point stack overflow exception when an x87 floating-point instruction is executed.

When writing an application that uses both floating-point and 64-bit SIMD integer instructions, use the following guidelines to help you determine when to use EMMS:

- **If next instruction is x87 FP** — Use `_MM_EMPTY()` after a 64-bit SIMD integer instruction if the next instruction is an X87 FP instruction; for example, before doing calculations on floats, doubles or long doubles.
- **Don't empty when already empty** — If the next instruction uses an MMX register, `_MM_EMPTY()` incurs a cost with no benefit.
- **Group Instructions** — Try to partition regions that use X87 FP instructions from those that use 64-bit SIMD integer instructions. This eliminates the need for an EMMS instruction within the body of a critical loop.
- **Runtime initialization** — Use `_MM_EMPTY()` during runtime initialization of `__m64` and X87 FP data types. This ensures resetting the register between data type transitions. See Example 5-1 for coding usage.

Example 5-1. Resetting Register Between `__m64` and FP Data Types Code

Incorrect Usage	Correct Usage
<pre>__m64 x = _m_padd(y, z); float f = init();</pre>	<pre>__m64 x = _m_padd(y, z); float f = (_mm_empty(), init());</pre>

You must be aware that your code generates an MMX instruction, which uses MMX registers with the Intel C++ Compiler, in the following situations:

- when using a 64-bit SIMD integer intrinsic from MMX technology, SSE/SSE2/SSSE3
- when using a 64-bit SIMD integer instruction from MMX technology, SSE/SSE2/SSSE3 through inline assembly
- when referencing the `__M64` data type variable

Additional information on the x87 floating-point programming model can be found in the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1*. For more on EMMS, visit <http://developer.intel.com>.

5.3 DATA ALIGNMENT

Make sure that 64-bit SIMD integer data is 8-byte aligned and that 128-bit SIMD integer data is 16-byte aligned. Referencing unaligned 64-bit SIMD integer data can incur a performance penalty due to accesses that span 2 cache lines. Referencing unaligned 128-bit SIMD integer data results in an exception unless the `MOVDQU` (move double-quadword unaligned) instruction is used. Using the `MOVDQU` instruction on unaligned data can result in lower performance than using 16-byte aligned references. Refer to Section 4.4, "Stack and Data Alignment," for more information.

Loading 16 bytes of SIMD data efficiently requires data alignment on 16-byte boundaries. `SSSE3` provides the `PALIGNR` instruction. It reduces overhead in situations that requires software to processing data elements from non-aligned address. The

PALIGNR instruction is most valuable when loading or storing unaligned data with the address shifts by a few bytes. You can replace a set of unaligned loads with aligned loads followed by using PALIGNR instructions and simple register to register copies.

Using PALIGNRs to replace unaligned loads improves performance by eliminating cache line splits and other penalties. In routines like MEMCPY(), PALIGNR can boost the performance of misaligned cases. Example 5-2 shows a situation that benefits by using PALIGNR.

Example 5-2. FIR Processing Example in C language Code

```
void FIR(float *in, float *out, float *coeff, int count)
{int i,j;
  for ( i=0; i<count - TAP; i++ )
  {   float sum = 0;
      for ( j=0; j<TAP; j++ )
      {   sum += in[j]*coeff[j]; }
      *out++ = sum;
      in++;
  }
}
```

Example 5-3 compares an optimal SSE2 sequence of the FIR loop and an equivalent SSSE3 implementation. Both implementations unroll 4 iteration of the FIR inner loop to enable SIMD coding techniques. The SSE2 code can not avoid experiencing cache line split once every four iterations. PALIGNR allows the SSSE3 code to avoid the delays associated with cache line splits.

Example 5-3. SSE2 and SSSE3 Implementation of FIR Processing Code

Optimized for SSE2	Optimized for SSSE3
<pre>pxor xmm0, xmm0 xor ecx, ecx mov eax, dword ptr[input] mov ebx, dword ptr[coeff4]</pre>	<pre>pxor xmm0, xmm0 xor ecx, ecx mov eax, dword ptr[input] mov ebx, dword ptr[coeff4]</pre>
<pre>inner_loop: movaps xmm1, xmmword ptr[eax+ecx] mulps xmm1, xmmword ptr[ebx+4*ecx] addps xmm0, xmm1</pre>	<pre>inner_loop: movaps xmm1, xmmword ptr[eax+ecx] movaps xmm3, xmm1 mulps xmm1, xmmword ptr[ebx+4*ecx] addps xmm0, xmm1</pre>
<pre>movups xmm1, xmmword ptr[eax+ecx+4] mulps xmm1, xmmword ptr[ebx+4*ecx+16] addps xmm0, xmm1</pre>	<pre>movaps xmm2, xmmword ptr[eax+ecx+16] movaps xmm1, xmm2 palignr xmm2, xmm3, 4 mulps xmm2, xmmword ptr[ebx+4*ecx+16] addps xmm0, xmm2</pre>

Example 5-3. SSE2 and SSSE3 Implementation of FIR Processing Code (Contd.)

Optimized for SSE2	Optimized for SSSE3
<pre> movups xmm1, xmmword ptr[eax+ecx+8] mulps xmm1, xmmword ptr[ebx+4*ecx+32] addps xmm0, xmm1 </pre>	<pre> movaps xmm2, xmm1 palignr xmm2, xmm3, 8 mulps xmm2, xmmword ptr[ebx+4*ecx+32] addps xmm0, xmm2 </pre>
<pre> movups xmm1, xmmword ptr[eax+ecx+12] mulps xmm1, xmmword ptr[ebx+4*ecx+48] addps xmm0, xmm1 add ecx, 16 cmp ecx, 4*TAP jl inner_loop mov eax, dword ptr[output] movaps xmmword ptr[eax], xmm0 </pre>	<pre> movaps xmm2, xmm1 palignr xmm2, xmm3, 12 mulps xmm2, xmmword ptr[ebx+4*ecx+48] addps xmm0, xmm2 add ecx, 16 cmp ecx, 4*TAP jl inner_loop mov eax, dword ptr[output] movaps xmmword ptr[eax], xmm0 </pre>

5.4 DATA MOVEMENT CODING TECHNIQUES

In general, better performance can be achieved if data is pre-arranged for SIMD computation (see Section 4.5, “Improving Memory Utilization”). This may not always be possible.

This section covers techniques for gathering and arranging data for more efficient SIMD computation.

5.4.1 Unsigned Unpack

MMX technology provides several instructions that are used to pack and unpack data in the MMX registers. SSE2 extends these instructions so that they operate on 128-bit source and destinations.

The unpack instructions can be used to zero-extend an unsigned number. Example 5-4 assumes the source is a packed-word (16-bit) data type.

Example 5-4. Zero Extend 16-bit Values into 32 Bits Using Unsigned Unpack Instructions Code

```

; Input:
;       XMM0      8 16-bit values in source
;       XMM7 0    a local variable can be used
;                 instead of the register XMM7 if
;                 desired.
;
; Output:
;       XMM0      four zero-extended 32-bit
;                 doublewords from four low-end
;                 words
;       XMM1      four zero-extended 32-bit
;                 doublewords from four high-end
;                 words
;
movdqa   xmm1, xmm0 ; copy source
punpcklwd xmm0, xmm7 ; unpack the 4 low-end words
;                 ; into 4 32-bit doubleword
punpckhwd xmm1, xmm7 ; unpack the 4 high-end words
;                 ; into 4 32-bit doublewords

```

5.4.2 Signed Unpack

Signed numbers should be sign-extended when unpacking values. This is similar to the zero-extend shown above, except that the PSRAD instruction (packed shift right arithmetic) is used to sign extend the values.

Example 5-5 assumes the source is a packed-word (16-bit) data type.

Example 5-5. Signed Unpack Code

```

Input:
;       XMM0      source value
;
; Output:
;       XMM0      four sign-extended 32-bit doublewords
;                 from four low-end words
;       XMM1      four sign-extended 32-bit doublewords
;                 from four high-end words
;
;
;
;
;

```

Example 5-5. Signed Unpack Code (Contd.)

```

movdqa    xmm1, xmm0 ; copy source
punpcklwd xmm0, xmm0 ; unpack four low end words of the source
                    ; into the upper 16 bits of each doubleword
                    ; in the destination
punpckhwd xmm1, xmm1 ; unpack 4 high-end words of the source
                    ; into the upper 16 bits of each doubleword
                    ; in the destination

psrad     xmm0, 16   ; sign-extend the 4 low-end words of the source
                    ; into four 32-bit signed doublewords
psrad     xmm1, 16   ; sign-extend the 4 high-end words of the
                    ; source into four 32-bit signed doublewords
    
```

5.4.3 Interleaved Pack with Saturation

Pack instructions pack two values into a destination register in a predetermined order. **PACKSSDW** saturates two signed doublewords from a source operand and two signed doublewords from a destination operand into four signed words; and it packs the four signed words into a destination register. See Figure 5-1.

SSE2 extends **PACKSSDW** so that it saturates four signed doublewords from a source operand and four signed doublewords from a destination operand into eight signed words; the eight signed words are packed into the destination.

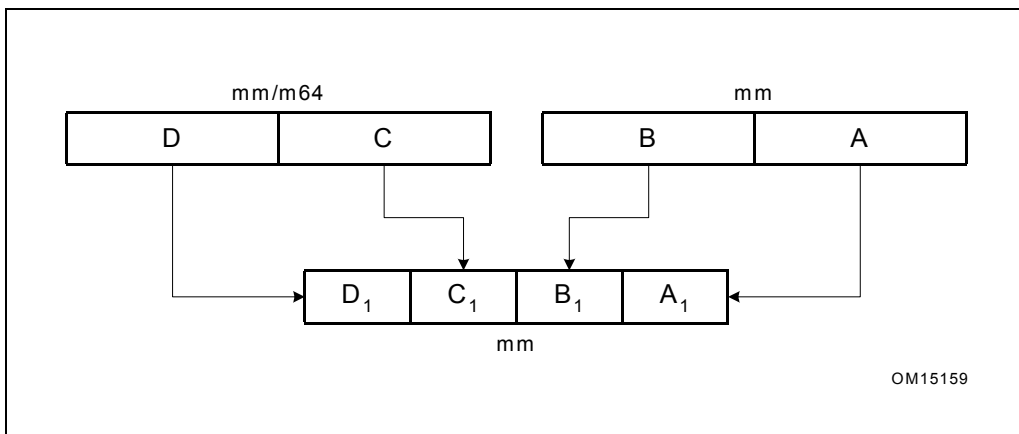


Figure 5-1. PACKSSDW mm, mm/mm64 Instruction

Figure 5-2 illustrates where two pairs of values are interleaved in a destination register; Example 5-6 shows MMX code that accomplishes the operation.

Two signed doublewords are used as source operands and the result is interleaved signed words. The sequence in Example 5-6 can be extended in SSE2 to interleave eight signed words using XMM registers.

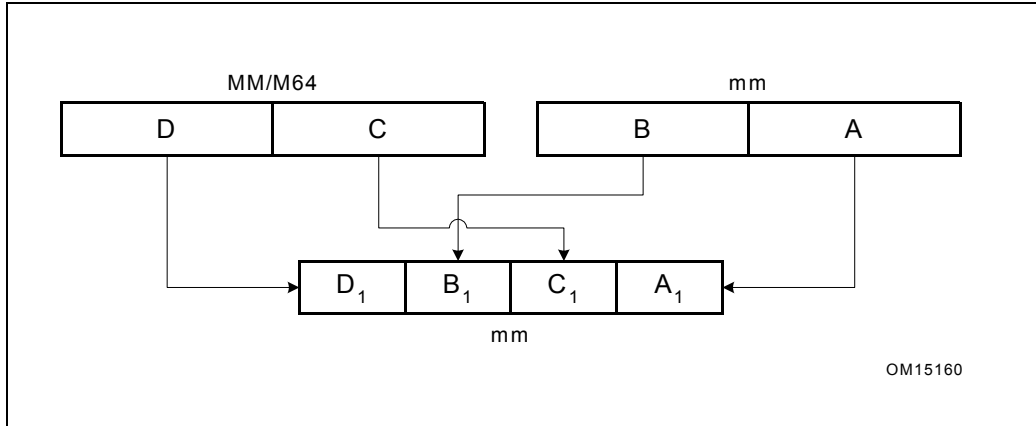


Figure 5-2. Interleaved Pack with Saturation

Example 5-6. Interleaved Pack with Saturation Code

```

; Input:
      MM0    signed source1 value
;      MM1    signed source2 value
; Output:
      MM0    the first and third words contain the
;            signed-saturated doublewords from MM0,
;            the second and fourth words contain
;            signed-saturated doublewords from MM1
;
packssdw  mm0, mm0    ; pack and sign saturate
packssdw  mm1, mm1    ; pack and sign saturate
punpcklwd mm0, mm1    ; interleave the low-end 16-bit
;                    ; values of the operands

```

Pack instructions always assume that source operands are signed numbers. The result in the destination register is always defined by the pack instruction that performs the operation. For example, PACKSSDW packs each of two signed 32-bit values of two sources into four saturated 16-bit signed values in a destination register. PACKUSWB, on the other hand, packs the four signed 16-bit values of two sources into eight saturated eight-bit unsigned values in the destination.

5.4.4 Interleaved Pack without Saturation

Example 5-7 is similar to Example 5-6 except that the resulting words are not saturated. In addition, in order to protect against overflow, only the low order 16 bits of each doubleword are used. Again, Example 5-7 can be extended in SSE2 to accomplish interleaving eight words without saturation.

Example 5-7. Interleaved Pack without Saturation Code

```

; Input:
;     MM0     signed source value
;     MM1     signed source value

; Output:
;     MM0     the first and third words contain the
;             low 16-bits of the doublewords in MM0,
;             the second and fourth words contain the
;             low 16-bits of the doublewords in MM1

pslld  mm1, 16    ; shift the 16 LSB from each of the
                  ; doubleword values to the 16 MSB
                  ; position
pand   mm0, {0,ffff,0,ffff}
                  ; mask to zero the 16 MSB
                  ; of each doubleword value
por    mm0, mm1   ; merge the two operands

```

5.4.5 Non-Interleaved Unpack

Unpack instructions perform an interleave merge of the data elements of the destination and source operands into the destination register.

The following example merges the two operands into destination registers without interleaving. For example, take two adjacent elements of a packed-word data type in SOURCE1 and place this value in the low 32 bits of the results. Then take two adjacent elements of a packed-word data type in SOURCE2 and place this value in the high 32 bits of the results. One of the destination registers will have the combination illustrated in Figure 5-3.

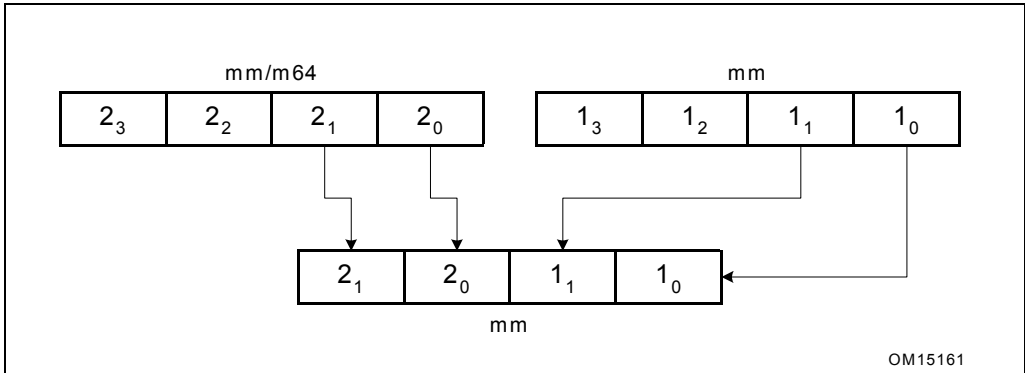


Figure 5-3. Result of Non-Interleaved Unpack Low in MM0

The other destination register will contain the opposite combination illustrated in Figure 5-4.

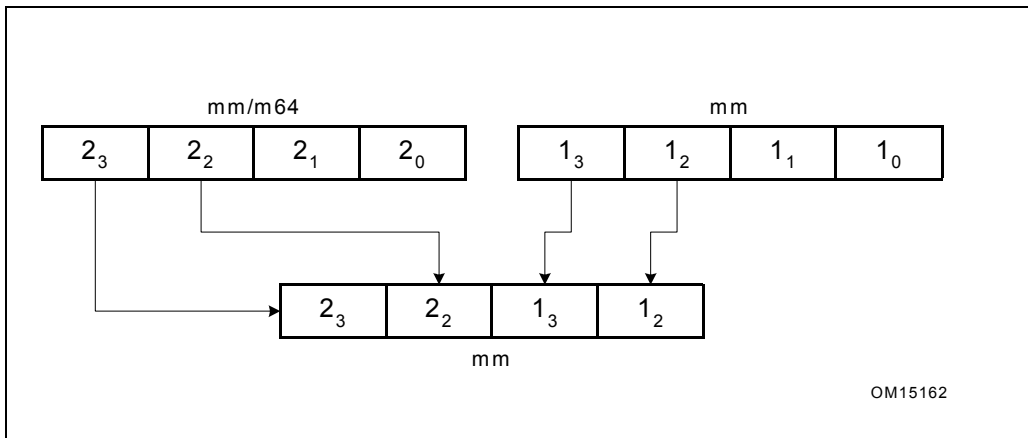


Figure 5-4. Result of Non-Interleaved Unpack High in MM1

Code in the Example 5-8 unpacks two packed-word sources in a non-interleaved way. The goal is to use the instruction which unpacks doublewords to a quadword, instead of using the instruction which unpacks words to doublewords.

Example 5-8. Unpacking Two Packed-word Sources in Non-interleaved Way Code

```

; Input:
;           MM0           packed-word source value
;           MM1           packed-word source value
; Output:
;           MM0           contains the two low-end words of the
;                           original sources, non-interleaved
;           MM2           contains the two high end words of the
;                           original sources, non-interleaved.
movq      mm2, mm0      ; copy source1
punpckldq mm0, mm1      ; replace the two high-end words of MM0 with
;                           ; two low-end words of MM1;
;                           ; leave the two low-end words of MM0 in place
punpckhdq mm2, mm1      ; move two high-end words of MM2 to the two low-end
;                           ; words of MM2; place the two high-end words of
;                           ; MM1 in two high-end words of MM2
    
```

5.4.6 Extract Data Element

The PEXTRW instruction in SSE takes the word in the designated MMX register selected by the two least significant bits of the immediate value and moves it to the lower half of a 32-bit integer register. See Figure 5-5 and Example 5-9.

With SSE2, PEXTRW can extract a word from an XMM register to the lower 16 bits of an integer register. SSE4.1 provides extraction of a byte, word, dword and qword from an XMM register into either a memory location or integer register.

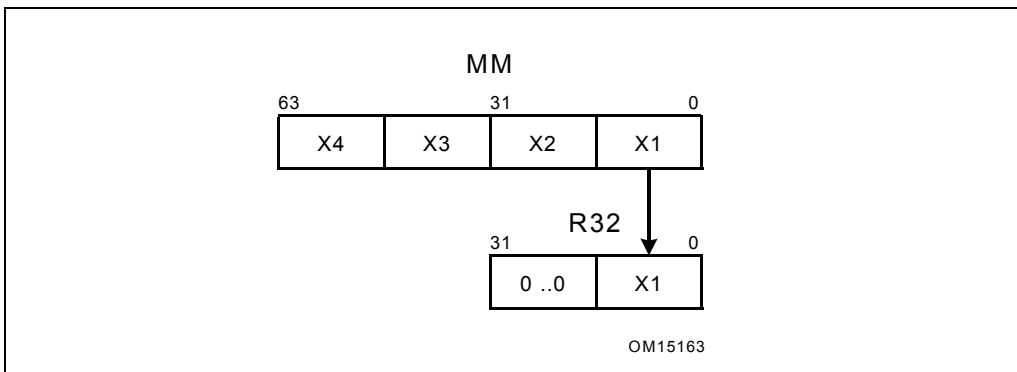


Figure 5-5. PEXTRW Instruction

Example 5-9. PEXTRW Instruction Code

```

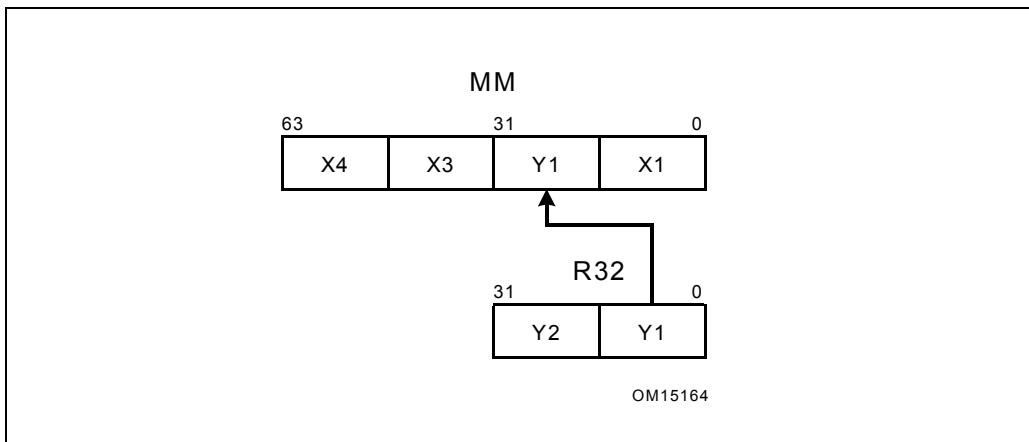
; Input:
;     eax    source value
;           immediate value: "0"
; Output:
;     edx    32-bit integer register containing the extracted word in the
;           low-order bits & the high-order bits zero-extended
movq  mm0, [eax]
pextrw edx, mm0, 0

```

5.4.7 Insert Data Element

The PINSRW instruction in SSE loads a word from the lower half of a 32-bit integer register or from memory and inserts it in an MMX technology destination register at a position defined by the two least significant bits of the immediate constant. Insertion is done in such a way that three other words from the destination register are left untouched. See Figure 5-6 and Example 5-10.

With SSE2, PINSRW can insert a word from the lower 16 bits of an integer register or memory into an XMM register. SSE4.1 provides insertion of a byte, dword and qword from either a memory location or integer register into an XMM register.

**Figure 5-6. PINSRW Instruction**

Example 5-10. PINSRW Instruction Code

```

; Input:
;     edx     pointer to source value
; Output:
;     mm0     register with new 16-bit value inserted
;
mov   eax, [edx]
pinsrw mm0, eax, 1

```

If all of the operands in a register are being replaced by a series of PINSRW instructions, it can be useful to clear the content and break the dependence chain by either using the PXOR instruction or loading the register. See Example 5-11 and Section 3.5.1.7, “Clearing Registers and Dependency Breaking Idioms.”

Example 5-11. Repeated PINSRW Instruction Code

```

; Input:
;     edx     pointer to structure containing source
;             values at offsets: of +0, +10, +13, and +24
;             immediate value: "1"
; Output:
;     MMX     register with new 16-bit value inserted
;
pxor   mm0, mm0 ; Breaks dependency on previous value of mm0
mov    eax, [edx]
pinsrw mm0, eax, 0
mov    eax, [edx+10]
pinsrw mm0, eax, 1
mov    eax, [edx+13]
pinsrw mm0, eax, 2
mov    eax, [edx+24]
pinsrw mm0, eax, 3

```

5.4.8 Non-Unit Stride Data Movement

SSE4.1 provides instructions to insert a data element from memory into an XMM register, and to extract a data element from an XMM register into memory directly. Separate instructions are provided to handle floating-point data and integer byte, word, or dword. These instructions are suited for vectorizing code that loads/stores non-unit stride data from memory, see Example 5-12.

Example 5-12. Non-Unit Stride Load/Store Using SSE4.1 Instructions

<pre>/* Goal: Non-Unit Stride Load Dwords*/ movd xmm0, [addr] pinsrd xmm0, [addr + stride], 1 pinsrd xmm0, [addr + 2*stride], 2 pinsrd xmm0, [addr + 3*stride], 3</pre>	<pre>/* Goal: Non-Unit Stride Store Dwords*/ movd [addr], xmm0 pextrd [addr + stride], xmm0, 1 pextrd [addr + 2*stride], xmm0, 2 pextrd [addr + 3*stride], xmm0, 3</pre>
---	--

Example 5-13 provides two examples: using INSERTPS and PEXTRD to perform gather operations on floating-point data; using EXTRACTPS and PEXTRD to perform scatter operations on floating-point data.

Example 5-13. Scatter and Gather Operations Using SSE4.1 Instructions

<pre>/* Goal: Gather Operation*/ movd eax, xmm0 movss xmm1, [addr + 4*eax] pextrd eax, xmm0, 1 insertps xmm1, [addr + 4*eax], 1 pextrd eax, xmm0, 2 insertps xmm1, [addr + 4*eax], 2 pextrd eax, xmm0, 3 insertps xmm1, [addr + 4*eax], 3</pre>	<pre>/* Goal: Scatter Operation*/ movd eax, xmm0 movss [addr + 4*eax], xmm1 pextrd eax, xmm0, 1 extractps [addr + 4*eax], xmm1, 1 pextrd eax, xmm0, 2 extractps [addr + 4*eax], xmm1, 2 pextrd eax, xmm0, 3 extractps [addr + 4*eax], xmm1, 3</pre>
---	---

5.4.9 Move Byte Mask to Integer

The PMOVMASKB instruction returns a bit mask formed from the most significant bits of each byte of its source operand. When used with 64-bit MMX registers, this produces an 8-bit mask, zeroing out the upper 24 bits in the destination register. When used with 128-bit XMM registers, it produces a 16-bit mask, zeroing out the upper 16 bits in the destination register.

The 64-bit version of this instruction is shown in Figure 5-7 and Example 5-14.

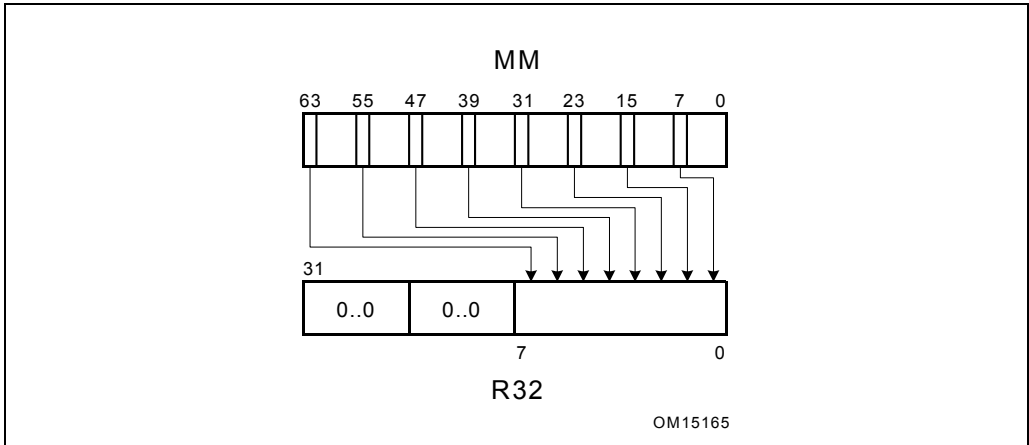


Figure 5-7. PMOVMSKB Instruction

Example 5-14. PMOVMSKB Instruction Code

```

; Input:
;     source value
; Output:
;     32-bit register containing the byte mask in the lower eight bits
;
movq  mm0, [edi]
pmovmskb eax, mm0

```

5.4.10 Packed Shuffle Word for 64-bit Registers

The PSHUFW instruction uses the immediate (IMM8) operand to select between the four words in either two MMX registers or one MMX register and a 64-bit memory location. SSE2 provides PSHUFLW to shuffle the lower four words into an XMM register. In addition to the equivalent to the PSHUFW, SSE2 also provides PSHUFW to shuffle the higher four words. Furthermore, SSE2 offers PSHUFD to shuffle four dwords into an XMM register. All of these four PSHUF instructions use an immediate byte to encode the data path of individual words within the corresponding 8 bytes from source to destination, shown in Table 5-1:

Table 5-1. PSHUF Encoding

Bits	Words
1 - 0	0
3 - 2	1
5 - 4	2
7 - 6	3

5.4.11 Packed Shuffle Word for 128-bit Registers

The PSHUFLW/PSHUFHW instruction performs a full shuffle of any source word field within the low/high 64 bits to any result word field in the low/high 64 bits, using an 8-bit immediate operand; other high/low 64 bits are passed through from the source operand.

PSHUFD performs a full shuffle of any double-word field within the 128-bit source to any double-word field in the 128-bit result, using an 8-bit immediate operand.

No more than 3 instructions, using PSHUFLW/PSHUFHW/PSHUFD, are required to implement many common data shuffling operations. Broadcast, Swap, and Reverse are illustrated in Example 5-15 and Example 5-16.

Example 5-15. Broadcast a Word Across XMM, Using 2 SSE2 Instructions

```

/* Goal: Broadcast the value from word 5 to all words */
/* Instruction      Result */
                    | 7| 6| 5| 4| 3| 2| 1| 0|
PSHUFHW (3,2,1,1) | 7| 6| 5| 5| 3| 2| 1| 0|
PSHUFD (2,2,2,2) | 5| 5| 5| 5| 5| 5| 5| 5|

```

Example 5-16. Swap/Reverse words in an XMM, Using 3 SSE2 Instructions

/* Goal: Swap the values in word 6 and word 1 */	/* Goal: Reverse the order of the words */
/* Instruction Result */	/* Instruction Result */
7 6 5 4 3 2 1 0	7 6 5 4 3 2 1 0
PSHUFD (3,0,1,2) 7 6 1 0 3 2 5 4	PSHUFLW (0,1,2,3) 7 6 5 4 0 1 2 3
PSHUFHW (3,1,2,0) 7 1 6 0 3 2 5 4	PSHUFHW (0,1,2,3) 4 5 6 7 0 1 2 3
PSHUFD (3,0,1,2) 7 1 5 4 3 2 6 0	PSHUFD (1,0,3,2) 0 1 2 3 4 5 6 7

5.4.12 Shuffle Bytes

SSSE3 provides PSHUFB; this instruction carries out byte manipulation within a 16 byte range. PSHUFB can replace up to 12 other instructions: including SHIFT, OR, AND and MOV.

Use PSHUFB if the alternative uses 5 or more instructions.

5.4.13 Conditional Data Movement

SSE4.1 provides two packed blend instructions on byte and word data elements in 128-bit operands. Packed blend instructions conditionally copies data elements from selected positions in the source to the corresponding data element using a mask specified by an immediate control byte or an implied XMM register (XMM0). The mask can be generated by a packed compare instruction for example. Thus packed blend instructions are most useful for vectorizing conditional flows within a loop and can be more efficient than inserting single element one at a time for some situations.

5.4.14 Unpacking/interleaving 64-bit Data in 128-bit Registers

The PUNPCKLQDQ/PUNPCHQDQ instructions interleave the low/high-order 64-bits of the source operand and the low/high-order 64-bits of the destination operand. It then writes the results to the destination register.

The high/low-order 64-bits of the source operands are ignored.

5.4.15 Data Movement

There are two additional instructions to enable data movement from 64-bit SIMD integer registers to 128-bit SIMD registers.

The MOVQ2DQ instruction moves the 64-bit integer data from an MMX register (source) to a 128-bit destination register. The high-order 64 bits of the destination register are zeroed-out.

The MOVDQ2Q instruction moves the low-order 64-bits of integer data from a 128-bit source register to an MMX register (destination).

5.4.16 Conversion Instructions

SSE provides Instructions to support 4-wide conversion of single-precision data to/from double-word integer data. Conversions between double-precision data to double-word integer data have been added in SSE2.

SSE4.1 provides 4 rounding instructions to convert floating-point values to integer values with rounding control specified in a more flexible manner and independent of the rounding control in MXCSR. The integer values produced by ROUNDxx instructions are maintained as floating-point data.

SSE4.1 also provides instructions to convert integer data from

- packed bytes to packed word/dword/qword format using either sign extension or zero extension,
- packed words to packed dword/qword format using either sign extension or zero extension,
- packed dword to packed qword format using either sign extension or zero extension.

5.5 GENERATING CONSTANTS

SIMD integer instruction sets do not have instructions that will load immediate constants to the SIMD registers.

The following code segments generate frequently used constants in the SIMD register. These examples can also be extended in SSE2 by substituting MMX with XMM registers. See Example 5-17.

Example 5-17. Generating Constants

```
pxor    mm0, mm0    ; generate a zero register in MM0
pcmpeq mm1, mm1    ; Generate all 1's in register MM1,
                    ; which is -1 in each of the packed
                    ; data type fields
```

Example 5-17. Generating Constants (Contd.)

```

pxor    mm0, mm0
pcmpeq mm1, mm1
psubb  mm0, mm1 [psubw mm0, mm1] (psubd mm0, mm1)
        ; three instructions above generate
        ; the constant 1 in every
        ; packed-byte [or packed-word]
        ; (or packed-dword) field

pcmpeq mm1, mm1
psrlw  mm1, 16-n (psrld mm1, 32-n)
        ; two instructions above generate
        ; the signed constant  $2^n-1$  in every
        ; packed-word (or packed-dword) field

pcmpeq mm1, mm1
psllw  mm1, n (pslld mm1, n)
        ; two instructions above generate
        ; the signed constant  $-2n$  in every
        ; packed-word (or packed-dword) field

```

NOTE

Because SIMD integer instruction sets do not support shift instructions for bytes, $2n-1$ and $-2n$ are relevant only for packed words and packed doublewords.

5.6 BUILDING BLOCKS

This section describes instructions and algorithms which implement common code building blocks.

5.6.1 Absolute Difference of Unsigned Numbers

Example 5-18 computes the absolute difference of two unsigned numbers. It assumes an unsigned packed-byte data type.

Here, we make use of the subtract instruction with unsigned saturation. This instruction receives UNSIGNED operands and subtracts them with UNSIGNED saturation.

This support exists only for packed bytes and packed words, not for packed double-words.

Example 5-18. Absolute Difference of Two Unsigned Numbers

```

; Input:
;     MM0 source operand
;     MM1 source operand
; Output:
;     MM0 absolute difference of the unsigned operands

movq   mm2, mm0 ; make a copy of mm0
psubusbmm0, mm1 ; compute difference one way
psubusbmm1, mm2 ; compute difference the other way
por    mm0, mm1 ; OR them together

```

This example will not work if the operands are signed. Note that PSADBW may also be used in some situations. See Section 5.6.9 for details.

5.6.2 Absolute Difference of Signed Numbers

Example 5-19 computes the absolute difference of two signed numbers using SSE3 instruction PABSW. This sequence is more efficient than using previous generation of SIMD instruction extensions.

Example 5-19. Absolute Difference of Signed Numbers

```

; Input:
;     XMM0 signed source operand
;     XMM1 signed source operand
; Output:
;     XMM1 absolute difference of the unsigned operands

psubw  xmm0, xmm1 ; subtract words
pabsw  xmm1, xmm0 ; results in XMM1

```

5.6.3 Absolute Value

Example 5-20 show an MMX code sequence to compute $|X|$, where X is signed. This example assumes signed words to be the operands.

With SSSE3, this sequence of three instructions can be replaced by the PABSW instruction. Additionally, SSSE3 provides a 128-bit version using XMM registers and supports byte, word and doubleword granularity.

Example 5-20. Computing Absolute Value

```

; Input:
;   MM0      signed source operand
; Output:
;   MM1      ABS(MM0)
pxor  mm1, mm1 ; set mm1 to all zeros
psubw mm1, mm0 ; make each mm1 word contain the
                ; negative of each mm0 word
pmaxswmm1, mm0 ; mm1 will contain only the positive
                ; (larger) values - the absolute value

```

NOTE

The absolute value of the most negative number (that is, 8000H for 16-bit) cannot be represented using positive numbers. This algorithm will return the original value for the absolute value (8000H).

5.6.4 Pixel Format Conversion

SSSE3 provides the PSHUFB instruction to carry out byte manipulation within a 16-byte range. PSHUFB can replace a set of up to 12 other instruction, including SHIFT, OR, AND and MOV.

Use PSHUFB if the alternative code uses 5 or more instructions. Example 5-21 shows the basic form of conversion of color pixel formats.

Example 5-21. Basic C Implementation of RGBA to BGRA Conversion

```

Standard C Code:
struct RGBA{BYTE r,g,b,a;};
struct BGRA{BYTE b,g,r,a;};

void BGRA_RGBA_Convert(BGRA *source, RGBA *dest, int num_pixels)
{
    for(int i = 0; i < num_pixels; i++){
        dest[i].r = source[i].r;
        dest[i].g = source[i].g;
        dest[i].b = source[i].b;
        dest[i].a = source[i].a;
    }
}

```

Example 5-22 and Example 5-23 show SSE2 code and SSSE3 code for pixel format conversion. In the SSSE3 example, PSHUFB replaces six SSE2 instructions.

Example 5-22. Color Pixel Format Conversion Using SSE2

```

; Optimized for SSE2
mov esi, src
  mov edi, dest
  mov ecx, iterations
  movdqa xmm0, ag_mask //{0,ff,0,ff,0,ff,0,ff,0,ff,0,ff,0,ff}
  movdqa xmm5, rb_mask //{ff,0,ff,0,ff,0,ff,0,ff,0,ff,0,ff,0}
  mov eax, remainder
convert16Pixs: // 16 pixels, 64 byte per iteration
  movdqa xmm1, [esi] // xmm1 = [r3g3b3a3,r2g2b2a2,r1g1b1a1,r0g0b0a0]
  movdqa xmm2, xmm1
  movdqa xmm7, xmm1 //xmm7 abgr
  psrld xmm2, 16 //xmm2 00ab
  pslld xmm1, 16 //xmm1 gr00

  por xmm1, xmm2 //xmm1 grab
  pand xmm7, xmm0 //xmm7 a0g0
  pand xmm1, xmm5 //xmm1 0r0b
  por xmm1, xmm7 //xmm1 argb
  movdqa [edi], xmm1

//repeats for another 3*16 bytes
...

add esi, 64
add edi, 64
sub ecx, 1
jnz convert16Pixs

```

Example 5-23. Color Pixel Format Conversion Using SSSE3

```

; Optimized for SSSE3

mov  esi, src
    mov  edi, dest
    mov  ecx, iterations
    movdqa xmm0, _shufb
// xmm0 = [15,12,13,14,11,8,9,10,7,4,5,6,3,0,1,2]
    mov  eax, remainder

convert16Pixs: // 16 pixels, 64 byte per iteration
    movdqa xmm1, [esi]
// xmm1 = [r3g3b3a3,r2g2b2a2,r1g1b1a1,r0g0b0a0]
    movdqa xmm2, [esi+16]
    pshufb xmm1, xmm0
// xmm1 = [b3g3r3a3,b2g2r2a2,b1g1r1a1,b0g0r0a0]
    movdqa [edi], xmm1

    //repeats for another 3*16 bytes

...

add  esi, 64
add  edi, 64
sub  ecx, 1
jnz  convert16Pixs

```

5.6.5 Endian Conversion

The PSHUFB instruction can also be used to reverse byte ordering within a double-word. It is more efficient than traditional techniques, such as BSWAP.

Example 5-24 (a) shows the traditional technique using four BSWAP instructions to reverse the bytes within a DWORD. Each BSWAP requires executing two μ ops. In addition, the code requires 4 loads and 4 stores for processing 4 DWORDs of data.

Example 5-24 (b) shows an SSSE3 implementation of endian conversion using PSHUFB. The reversing of four DWORDs requires one load, one store, and PSHUFB.

On Intel Core microarchitecture, reversing 4 DWORDs using PSHUFB can be approximately twice as fast as using BSWAP.

Example 5-24. Big-Endian to Little-Endian Conversion

<pre> ;(a) Using BSWAP lea eax, src lea ecx, dst mov edx, elCount start: mov edi, [eax] mov esi, [eax+4] bswap edi mov ebx, [eax+8] bswap esi mov ebp, [eax+12] mov [ecx], edi mov [ecx+4], esi bswap ebx mov [ecx+8], ebx bswap ebp mov [ecx+12], ebp add eax, 16 add ecx, 16 sub edx, 4 jnz start </pre>	<pre> ;(b) Using PSHUFB __declspec(align(16)) BYTE bswapMASK[16] = {3,2,1,0, 7,6,5,4, 11,10,9,8, 15,14,13,12}; lea eax, src lea ecx, dst mov edx, elCount movaps xmm7, bswapMASK start: movdqa xmm0, [eax] pshufb xmm0, xmm7 movdqa [ecx], xmm0 add eax, 16 add ecx, 16 sub edx, 4 jnz start </pre>
--	--

5.6.6 Clipping to an Arbitrary Range [High, Low]

This section explains how to clip a values to a range [HIGH, LOW]. Specifically, if the value is less than LOW or greater than HIGH, then clip to LOW or HIGH, respectively. This technique uses the packed-add and packed-subtract instructions with saturation (signed or unsigned), which means that this technique can only be used on packed-byte and packed-word data types.

The examples in this section use the constants PACKED_MAX and PACKED_MIN and show operations on word values. For simplicity, we use the following constants (corresponding constants are used in case the operation is done on byte values):

PACKED_MAX equals 0X7FFF7FFF7FFF7FFF

PACKED_MIN equals 0X8000800080008000

PACKED_LOW contains the value LOW in all four words of the packed-words data type

PACKED_HIGH contains the value HIGH in all four words of the packed-words data type

PACKED_USMAX all values equal 1

HIGH_US adds the HIGH value to all data elements (4 words) of PACKED_MIN

LOW_US adds the LOW value to all data elements (4 words) of PACKED_MIN

5.6.6.1 Highly Efficient Clipping

For clipping signed words to an arbitrary range, the PMAWSW and PMINSW instructions may be used. For clipping unsigned bytes to an arbitrary range, the PMAWUB and PMINUB instructions may be used.

Example 5-25 shows how to clip signed words to an arbitrary range; the code for clipping unsigned bytes is similar.

Example 5-25. Clipping to a Signed Range of Words [High, Low]

```

; Input:
;   MM0   signed source operands
; Output:
;   MM0   signed words clipped to the signed
;         range [high, low]
pminsw mm0, packed_high
pmaxswmm0, packed_low

```

With SSE4.1, Example 5-25 can be easily extended to clip signed bytes, unsigned words, signed and unsigned dwords.

Example 5-26. Clipping to an Arbitrary Signed Range [High, Low]

```

; Input:
;   MM0           signed source operands
; Output:
;   MM1           signed operands clipped to the unsigned
;               range [high, low]
paddw mm0, packed_min ; add with no saturation
;               ; 0x8000 to convert to unsigned
padduswmm0, (packed_usmax - high_us)
;               ; in effect this clips to high
psubuswmm0, (packed_usmax - high_us + low_us)
;               ; in effect this clips to low
paddw mm0, packed_low ; undo the previous two offsets

```

The code above converts values to unsigned numbers first and then clips them to an unsigned range. The last instruction converts the data back to signed data and places the data within the signed range.

Conversion to unsigned data is required for correct results when $(High - Low) < 0X8000$. If $(High - Low) \geq 0X8000$, simplify the algorithm as in Example 5-27.

Example 5-27. Simplified Clipping to an Arbitrary Signed Range

```

; Input:      MM0      signed source operands
; Output:    MM1      signed operands clipped to the unsigned
;             range [high, low]
paddssw  mm0, (packed_max - packed_high)
           ; in effect this clips to high
psubssw  mm0, (packed_usmax - packed_high + packed_low)
           ; clips to low
paddw    mm0, low     ; undo the previous two offsets

```

This algorithm saves a cycle when it is known that $(\text{High} - \text{Low}) \geq 0x8000$. The three-instruction algorithm does not work when $(\text{High} - \text{Low}) < 0x8000$ because $0xffff$ minus any number $< 0x8000$ will yield a number greater in magnitude than $0x8000$ (which is a negative number).

When the second instruction, `psubssw MM0, (0xffff - High + Low)` in the three-step algorithm (Example 5-27) is executed, a negative number is subtracted. The result of this subtraction causes the values in MM0 to be increased instead of decreased, as should be the case, and an incorrect answer is generated.

5.6.6.2 Clipping to an Arbitrary Unsigned Range [High, Low]

Example 5-28 clips an unsigned value to the unsigned range [High, Low]. If the value is less than low or greater than high, then clip to low or high, respectively. This technique uses the packed-add and packed-subtract instructions with unsigned saturation, thus the technique can only be used on packed-bytes and packed-words data types.

Figure 5-28 illustrates operation on word values.

Example 5-28. Clipping to an Arbitrary Unsigned Range [High, Low]

```

; Input:
;
; Output:    MM0      unsigned source operands
;             MM1      unsigned operands clipped to the unsigned
;             range [HIGH, LOW]
paddusw  mm0, 0xffff - high
           ; in effect this clips to high
psubusw  mm0, (0xffff - high + low)
           ; in effect this clips to low
paddw    mm0, low
           ; undo the previous two offsets

```

5.6.7 Packed Max/Min of Byte, Word and Dword

The P`MAXSW` instruction returns the maximum between four signed words in either of two SIMD registers, or one SIMD register and a memory location.

The P`MINSW` instruction returns the minimum between the four signed words in either of two SIMD registers, or one SIMD register and a memory location.

The P`MAXUB` instruction returns the maximum between the eight unsigned bytes in either of two SIMD registers, or one SIMD register and a memory location.

The P`MINUB` instruction returns the minimum between the eight unsigned bytes in either of two SIMD registers, or one SIMD register and a memory location.

SSE2 extended P`MAXSW`/P`MAXUB`/P`MINSW`/P`MINUB` to 128-bit operations. SSE4.1 adds 128-bit operations for signed bytes, unsigned word, signed and unsigned dword.

5.6.8 Packed Multiply Integers

The P`MULHUW`/P`MULHW` instruction multiplies the unsigned/signed words in the destination operand with the unsigned/signed words in the source operand. The high-order 16 bits of the 32-bit intermediate results are written to the destination operand. The P`MULLW` instruction multiplies the signed words in the destination operand with the signed words in the source operand. The low-order 16 bits of the 32-bit intermediate results are written to the destination operand.

SSE2 extended P`MULHUW`/P`MULHW`/P`MULLW` to 128-bit operations and adds P`MULUDQ`.

The P`MULUDQ` instruction performs an unsigned multiply on the lower pair of double-word operands within 64-bit chunks from the two sources; the full 64-bit result from each multiplication is returned to the destination register.

This instruction is added in both a 64-bit and 128-bit version; the latter performs 2 independent operations, on the low and high halves of a 128-bit register.

SSE4.1 adds 128-bit operations of P`MULDQ` and P`MULLD`. The P`MULLD` instruction multiplies the signed dwords in the destination operand with the signed dwords in the source operand. The low-order 32 bits of the 64-bit intermediate results are written to the destination operand. The P`MULDQ` instruction multiplies the two low-order, signed dwords in the destination operand with the two low-order, signed dwords in the source operand and stores two 64-bit results in the destination operand.

5.6.9 Packed Sum of Absolute Differences

The P`SADBW` instruction computes the absolute value of the difference of unsigned bytes for either two SIMD registers, or one SIMD register and a memory location. The differences of 8 pairs of unsigned bytes are then summed to produce a word

result in the lower 16-bit field, and the upper three words are set to zero. With SSE2, PSADBW is extended to compute two word results.

The subtraction operation presented above is an absolute difference. That is, $T = \text{ABS}(X-Y)$. Byte values are stored in temporary space, all values are summed together, and the result is written to the lower word of the destination register.

Motion estimation involves searching reference frames for best matches. Sum absolute difference (SAD) on two blocks of pixels is a common ingredient in video processing algorithms to locate matching blocks of pixels. PSADBW can be used as building blocks for finding best matches by way of calculating SAD results on 4x4, 8x4, 8x8 blocks of pixels.

5.6.10 MPSADBw and PHMINPOSUw

The MPSADBW instruction in SSE4.1 performs eight SAD operations. Each SAD operation produces a word result from 4 pairs of unsigned bytes. With 8 SAD result in an XMM register, PHMINPOSUM can help search for the best match between eight 4x4 pixel blocks.

For motion estimation algorithms, MPSADBW is likely to improve over PSADBW in several ways:

- Simplified data movement to construct packed data format for SAD computation on pixel blocks.
- Higher throughput in terms of SAD results per iteration (less iteration required per frame).
- MPSADBW results are amenable to efficient search using PHMINPOSUW.

Examples of MPSADBW vs. PSADBW for 4x4 and 8x8 block search can be found in the white paper listed in the reference section of Chapter 1.

5.6.11 Packed Average (Byte/Word)

The PAVGB and PAVGW instructions add the unsigned data elements of the source operand to the unsigned data elements of the destination register, along with a carry-in. The results of the addition are then independently shifted to the right by one bit position. The high order bits of each element are filled with the carry bits of the corresponding sum.

The destination operand is an SIMD register. The source operand can either be an SIMD register or a memory operand.

The PAVGB instruction operates on packed unsigned bytes and the PAVGW instruction operates on packed unsigned words.

5.6.12 Complex Multiply by a Constant

Complex multiplication is an operation which requires four multiplications and two additions. This is exactly how the PMADDWD instruction operates. In order to use this instruction, you need to format the data into multiple 16-bit values. The real and imaginary components should be 16-bits each. Consider Example 5-29, which assumes that the 64-bit MMX registers are being used:

- Let the input data be DR and DI, where DR is real component of the data and DI is imaginary component of the data.
- Format the constant complex coefficients in memory as four 16-bit values [CR - CI CI CR]. Remember to load the values into the MMX register using MOVQ.
- The real component of the complex product is $PR = DR*CR - DI*CI$ and the imaginary component of the complex product is $PI = DR*CI + DI*CR$.
- The output is a packed doubleword. If needed, a pack instruction can be used to convert the result to 16-bit (thereby matching the format of the input).

Example 5-29. Complex Multiply by a Constant

; Input:		
;	MM0	complex value, Dr, Di
;	MM1	constant complex coefficient in the form
;		[Cr -Ci Ci Cr]
; Output:		
;	MM0	two 32-bit dwords containing [Pr Pi]
;		
punpckldq	mm0, mm0	; makes [dr di dr di]
pmaddwd	mm0, mm1	; done, the result is
		; [(Dr*Cr-Di*Ci)(Dr*Ci+Di*Cr)]

5.6.13 Packed 64-bit Add/Subtract

The PADDQ/PSUBQ instructions add/subtract quad-word operands within each 64-bit chunk from the two sources; the 64-bit result from each computation is written to the destination register. Like the integer ADD/SUB instruction, PADDQ/PSUBQ can operate on either unsigned or signed (two's complement notation) integer operands.

When an individual result is too large to be represented in 64-bits, the lower 64-bits of the result are written to the destination operand and therefore the result wraps around. These instructions are added in both a 64-bit and 128-bit version; the latter performs 2 independent operations, on the low and high halves of a 128-bit register.

5.6.14 128-bit Shifts

The PSLLDQ/PSRLDQ instructions shift the first operand to the left/right by the number of bytes specified by the immediate operand. The empty low/high-order bytes are cleared (set to zero).

If the value specified by the immediate operand is greater than 15, then the destination is set to all zeros.

5.6.15 PTEST and Conditional Branch

SSE4.1 offers PTEST instruction that can be used in vectorizing loops with conditional branches. PTEST is an 128-bit version of the general-purpose instruction TEST. The ZF or CF field of the EFLAGS register are modified as a result of PTEST.

Example 5-30(a) depicts a loop that requires a conditional branch to handle the special case of divide-by-zero. In order to vectorize such loop, any iteration that may encounter divide-by-zero must be treated outside the vectorizable iterations.

Example 5-30. Using PTEST to Separate Vectorizable and non-Vectorizable Loop Iterations

<pre>(a) /* Loops requiring infrequent exception handling*/ float a[CNT]; unsigned int i; for (i=0;i<CNT;i++) { if (a[i] != 0.0) { a[i] = 1.0f/a[i]; } else { call DivException(); } } </pre>	<pre>(b) /* PTEST enables early out to handle infrequent, non- vectorizable portion*/ xor eax,eax movaps xmm7, [all_ones] xorps xmm6, xmm6 lp: movaps xmm0, a[eax] cmpeqps xmm6, xmm0 ; convert each non-zero to ones ptest xmm6, xmm7 jnc zero_present; carry will be set if all 4 were non- zero movaps xmm1,[_1_of_] divps xmm1, xmm0 movaps a[eax], xmm1 add eax, 16 cmp eax, CNT jnz lp jmp end zero_present: // execute one by one, call // exception when value is zero </pre>
---	--

Example 5-30(b) shows an assembly sequence that uses PTEST to cause an early-out branch whenever any one of the four floating-point values in xmm0 is zero. The fall-through path enables the rest of the floating-point calculations to be vectorized because none of the four values are zero.

5.6.16 Vectorization of Heterogeneous Computations across Loop Iterations

Vectorization techniques on un-rolled loops generally rely on repetitive, homogeneous operations between each loop iteration. Using SSE4.1's PTEST and variable blend instructions, vectorization of heterogeneous operations across loop iterations may be possible.

Example 5-31(a) depicts a simple heterogeneous loop. The heterogeneous operation and conditional branch makes simple loop-unrolling technique infeasible for vectorization.

Example 5-31. Using PTEST and Variable BLEND to Vectorize Heterogeneous Loops

<pre>(a) /* Loops with heterogeneous operation across iterations*/ float a[CNT]; unsigned int i; for (i=0;i<CNT;i++) { if (a[i] > b[i]) { a[i] += b[i]; } else { a[i] -= b[i]; } }</pre>	<pre>(b) /* Vectorize Condition Flow with PTEST, BLENDVPS*/ xor eax,eax lp: movaps xmm0, a[eax] movaps xmm1, b[eax] movaps xmm2, xmm0 // compare a and b values cmpgtps xmm0, xmm1 // xmm3 - will hold -b movaps xmm3, [SIGN_BIT_MASK] xorps xmm3, xmm1 // select values for the add operation, // true condition produce a+b, false will become a+(-b) // blend mask is xmm0 blendvps xmm1,xmm3, xmm0 addps xmm2, xmm1 movaps a[eax], xmm2 add eax, 16 cmp eax, CNT jnz lp</pre>
---	--

Example 5-31(b) depicts an assembly sequence that uses BLENDVPS and PTEST to vectorize the handling of heterogeneous computations occurring across four consecutive loop iterations.

5.6.17 Vectorization of Control Flows in Nested Loops

The PTEST and BLENDVPx instructions can be used as building blocks to vectorize more complex control-flow statements, where each control flow statement is creating a “working” mask used as a predicate of which the conditional code under the mask will operate.

The Mandelbrot-set map evaluation is useful to illustrate a situation with more complex control flows in nested loops. The Mandelbrot-set is a set of height values mapped to a 2-D grid. The height value is the number of Mandelbrot iterations (defined over the complex number space as $I_n = I_{n-1}^2 + I_0$) needed to get $|I_n| > 2$. It is common to limit the map generation by setting some maximum threshold value of the height, all other points are assigned with a height equal to the threshold. Example 5-32 shows an example of Mandelbrot map evaluation implemented in C.

Example 5-32. Baseline C Code for Mandelbrot Set Map Evaluation

```
#define DIMX (64)
#define DIMY (64)
#define X_STEP (0.5f/DIMX)
#define Y_STEP (0.4f/(DIMY/2))
int map[DIMX][DIMY];

void mandelbrot_C()
{
    int ij;
    float x,y;
    for (i=0,x=-1.8f;i<DIMX;i++,x+=X_STEP)
    {
        for (j=0,y=-0.2f;j<DIMY/2;j++,y+=Y_STEP)
        {float sx,sy;
            int iter = 0;
            sx = x;
            sy = y;
            while (iter < 256)
            {
                if (sx*sx + sy*sy >= 4.0f)    break;
                float old_sx = sx;
                sx = x + sx*sx - sy*sy;
                sy = y + 2*old_sx*sy;
                iter++;
            }
            map[i][j] = iter;
        }
    }
}
```

Example 5-33 shows a vectorized implementation of Mandelbrot map evaluation. Vectorization is not done on the inner most loop, because the presence of the break statement implies the iteration count will vary from one pixel to the next. The vectorized version take into account the parallel nature of 2-D, vectorize over four iterations of Y values of 4 consecutive pixels, and conditionally handles three scenarios:

- In the inner most iteration, when all 4 pixels do not reach break condition, vectorize 4 pixels.
- When one or more pixels reached break condition, use blend intrinsics to accumulate the complex height vector for the remaining pixels not reaching the break condition and continue the inner iteration of the complex height vector;
- When all four pixels reached break condition, exit the inner loop.

Example 5-33. Vectorized Mandelbrot Set Map Evaluation Using SSE4.1 Intrinsics

```

__declspec(align(16)) float _INIT_Y_4[4] = {0,Y_STEP,2*Y_STEP,3*Y_STEP};
F32vec4 _F_STEP_Y(4*Y_STEP);
I32vec4 _I_ONE_ = _mm_set1_epi32(1);
F32vec4 _F_FOUR_(4.0f);
F32vec4 _F_TWO_(2.0f);

void mandelbrot_C()
{
    int ij;
    F32vec4 x,y;

    for (i = 0, x = F32vec4(-1.8f); i < DIMX; i ++, x += F32vec4(X_STEP))
    {
        for (j = DIMY/2, y = F32vec4(-0.2f) +
            *(F32vec4*)_INIT_Y_4; j < DIMY; j += 4, y += _F_STEP_Y)
        {
            F32vec4 sx,sy;
            I32vec4 iter = _mm_setzero_si128();
            int scalar_iter = 0;
            sx = x;
            sy = y;
            while (scalar_iter < 256)
            {
                int mask = 0;
                F32vec4 old_sx = sx;
                __m128 vmask = _mm_cmpnlt_ps(sx*sx + sy*sy,_F_FOUR_);
                // if all data points in our vector are hitting the "exit" condition,
                // the vectorized loop can exit
                if (_mm_test_all_ones(_mm_castps_si128(vmask)))
                    break;
                (continue)
            }
        }
    }
}

```

Example 5-33. Vectorized Mandelbrot Set Map Evaluation Using SSE4.1 Intrinsics

```

// if non of the data points are out, we don't need the extra code which blends the results
    if (_mm_test_all_zeros(_mm_castps_si128(vmask),
        _mm_castps_si128(vmask)))
    {
        sx = x + sx*sx - sy*sy;
        sy = y + _F_TWO_*old_sx*sy;
        iter += _ONE_;
    }
    else
    {
// Blended flavour of the code, this code blends values from previous iteration with the values
// from current iteration. Only values which did not hit the "exit" condition are being stored;
// values which are already "out" are maintaining their value
        sx = _mm_blendv_ps(x + sx*sx - sy*sy,sx,vmask);
        sy = _mm_blendv_ps(y + _F_TWO_*old_sx*sy,sy,vmask);
        iter = l32vec4(_mm_blendv_epi8(iter + _ONE_,
            iter,_mm_castps_si128(vmask)));
    }
    scalar_iter++;
}
_mm_storeu_si128((__m128i*)&map[i][j],iter);
}
}
}

```

5.7 MEMORY OPTIMIZATIONS

You can improve memory access using the following techniques:

- Avoiding partial memory accesses
- Increasing the bandwidth of memory fills and video fills
- Prefetching data with Streaming SIMD Extensions. See Chapter 7, "Optimizing Cache Usage."

MMX registers and XMM registers allow you to move large quantities of data without stalling the processor. Instead of loading single array values that are 8, 16, or 32 bits long, consider loading the values in a single quadword or double quadword and then incrementing the structure or array pointer accordingly.

Any data that will be manipulated by SIMD integer instructions should be loaded using either:

- An SIMD integer instruction that loads a 64-bit or 128-bit operand (for example: `MOVQ MM0, M64`)

- The register-memory form of any SIMD integer instruction that operates on a quadword or double quadword memory operand (for example, PMADDW MM0, M64).

All SIMD data should be stored using an SIMD integer instruction that stores a 64-bit or 128-bit operand (for example: MOVQ M64, MM0)

The goal of the above recommendations is twofold. First, the loading and storing of SIMD data is more efficient using the larger block sizes. Second, following the above recommendations helps to avoid mixing of 8-, 16-, or 32-bit load and store operations with SIMD integer technology load and store operations to the same SIMD data.

This prevents situations in which small loads follow large stores to the same area of memory, or large loads follow small stores to the same area of memory. The Pentium II, Pentium III, and Pentium 4 processors may stall in such situations. See Chapter 3 for details.

5.7.1 Partial Memory Accesses

Consider a case with a large load after a series of small stores to the same area of memory (beginning at memory address MEM). The large load stalls in the case shown in Example 5-34.

Example 5-34. A Large Load after a Series of Small Stores (Penalty)

mov	mem, eax	; store dword to address "mem"
mov	mem + 4, ebx	; store dword to address "mem + 4"
:		
:		
movq	mm0, mem	; load qword at address "mem", stalls

MOVQ must wait for the stores to write memory before it can access all data it requires. This stall can also occur with other data types (for example, when bytes or words are stored and then words or doublewords are read from the same area of memory). When you change the code sequence as shown in Example 5-35, the processor can access the data without delay.

Example 5-35. Accessing Data Without Delay

```

movd  mm1, ebx      ; build data into a qword first
                        ; before storing it to memory
movd  mm2, eax
psllq mm1, 32
por   mm1, mm2
movq  mem, mm1      ; store SIMD variable to "mem" as
                        ; a qword
:
:
movq  mm0, mem      ; load qword SIMD "mem", no stall

```

Consider a case with a series of small loads after a large store to the same area of memory (beginning at memory address MEM), as shown in Example 5-36. Most of the small loads stall because they are not aligned with the store. See Section 3.6.5, “Store Forwarding,” for details.

Example 5-36. A Series of Small Loads After a Large Store

```

movq  mem, mm0      ; store qword to address "mem"
:
:
mov   bx, mem + 2   ; load word at "mem + 2" stalls
mov   cx, mem + 4   ; load word at "mem + 4" stalls

```

The word loads must wait for the quadword store to write to memory before they can access the data they require. This stall can also occur with other data types (for example: when doublewords or words are stored and then words or bytes are read from the same area of memory).

When you change the code sequence as shown in Example 5-37, the processor can access the data without delay.

Example 5-37. Eliminating Delay for a Series of Small Loads after a Large Store

```

movq  mem, mm0      ; store qword to address "mem"
:
:
movq  mm1, mem      ; load qword at address "mem"
movd  eax, mm1      ; transfer "mem + 2" to eax from
                        ; MMX register, not memory

```

Example 5-37. Eliminating Delay for a Series of Small Loads after a Large Store

```

psrlq  mm1, 32
shr    eax, 16
movd   ebx, mm1    ; transfer "mem + 4" to bx from
                   ; MMX register, not memory
and    ebx, 0ffffh

```

These transformations, in general, increase the number of instructions required to perform the desired operation. For Pentium II, Pentium III, and Pentium 4 processors, the benefit of avoiding forwarding problems outweighs the performance penalty due to the increased number of instructions.

5.7.1.1 Supplemental Techniques for Avoiding Cache Line Splits

Video processing applications sometimes cannot avoid loading data from memory addresses that are not aligned to 16-byte boundaries. An example of this situation is when each line in a video frame is averaged by shifting horizontally half a pixel.

Example shows a common operation in video processing that loads data from memory address not aligned to a 16-byte boundary. As video processing traverses each line in the video frame, it experiences a cache line split for each 64 byte chunk loaded from memory.

Example 5-38. An Example of Video Processing with Cache Line Splits

```

// Average half-pels horizontally (on // the "x" axis),
// from one reference frame only.

nextLinesLoop:
movdqu xmm0, XMMWORD PTR [edx] // may not be 16B aligned
movdqu xmm0, XMMWORD PTR [edx+1]
movdqu xmm1, XMMWORD PTR [edx+eax]
movdqu xmm1, XMMWORD PTR [edx+eax+1]

pavgb xmm0, xmm1
pavgb xmm2, xmm3
movdqa XMMWORD PTR [ecx], xmm0
movdqa XMMWORD PTR [ecx+eax], xmm2
// (repeat ...)

```

SSE3 provides an instruction LDDQU for loading from memory address that are not 16-byte aligned. LDDQU is a special 128-bit unaligned load designed to avoid cache line splits. If the address of the load is aligned on a 16-byte boundary, LDDQU loads the 16 bytes requested. If the address of the load is not aligned on a 16-byte

boundary, LDDQU loads a 32-byte block starting at the 16-byte aligned address immediately below the address of the load request. It then provides the requested 16 bytes. If the address is aligned on a 16-byte boundary, the effective number of memory requests is implementation dependent (one, or more).

LDDQU is designed for programming usage of loading data from memory without storing modified data back to the same address. Thus, the usage of LDDQU should be restricted to situations where no store-to-load forwarding is expected. For situations where store-to-load forwarding is expected, use regular store/load pairs (either aligned or unaligned based on the alignment of the data accessed).

Example 5-39. Video Processing Using LDDQU to Avoid Cache Line Splits

```
// Average half-pels horizontally (on // the "x" axis),
// from one reference frame only.
nextLinesLoop:
lddqu xmm0, XMMWORD PTR [edx] // may not be 16B aligned
lddqu xmm0, XMMWORD PTR [edx+1]
lddqu xmm1, XMMWORD PTR [edx+eax]
lddqu xmm1, XMMWORD PTR [edx+eax+1]
pavgbxmm0, xmm1
pavgbxmm2, xmm3
movdqaXMMWORD PTR [ecx], xmm0 //results stored elsewhere
movdqaXMMWORD PTR [ecx+eax], xmm2
// (repeat ...)
```

5.7.2 Increasing Bandwidth of Memory Fills and Video Fills

It is beneficial to understand how memory is accessed and filled. A memory-to-memory fill (for example a memory-to-video fill) is defined as a 64-byte (cache line) load from memory which is immediately stored back to memory (such as a video frame buffer).

The following are guidelines for obtaining higher bandwidth and shorter latencies for sequential memory fills (video fills). These recommendations are relevant for all Intel architecture processors with MMX technology and refer to cases in which the loads and stores do not hit in the first- or second-level cache.

5.7.2.1 Increasing Memory Bandwidth Using the MOVDQ Instruction

Loading any size data operand will cause an entire cache line to be loaded into the cache hierarchy. Thus, any size load looks more or less the same from a memory bandwidth perspective. However, using many smaller loads consumes more microarchitectural resources than fewer larger stores. Consuming too many resources can

cause the processor to stall and reduce the bandwidth that the processor can request of the memory subsystem.

Using MOVDQ to store the data back to UC memory (or WC memory in some cases) instead of using 32-bit stores (for example, MOVD) will reduce by three-quarters the number of stores per memory fill cycle. As a result, using the MOVDQ in memory fill cycles can achieve significantly higher effective bandwidth than using MOVD.

5.7.2.2 Increasing Memory Bandwidth by Loading and Storing to and from the Same DRAM Page

DRAM is divided into pages, which are not the same as operating system (OS) pages. The size of a DRAM page is a function of the total size of the DRAM and the organization of the DRAM. Page sizes of several Kilobytes are common. Like OS pages, DRAM pages are constructed of sequential addresses. Sequential memory accesses to the same DRAM page have shorter latencies than sequential accesses to different DRAM pages.

In many systems the latency for a page miss (that is, an access to a different page instead of the page previously accessed) can be twice as large as the latency of a memory page hit (access to the same page as the previous access). Therefore, if the loads and stores of the memory fill cycle are to the same DRAM page, a significant increase in the bandwidth of the memory fill cycles can be achieved.

5.7.2.3 Increasing UC and WC Store Bandwidth by Using Aligned Stores

Using aligned stores to fill UC or WC memory will yield higher bandwidth than using unaligned stores. If a UC store or some WC stores cross a cache line boundary, a single store will result in two transactions on the bus, reducing the efficiency of the bus transactions. By aligning the stores to the size of the stores, you eliminate the possibility of crossing a cache line boundary, and the stores will not be split into separate transactions.

5.7.3 Reverse Memory Copy

Copying blocks of memory from a source location to a destination location in reverse order presents a challenge for software to make the most out of the machine's capabilities while avoiding microarchitectural hazards. The basic, un-optimized C code is shown in Example 5-40.

The simple C code in Example 5-40 is sub-optimal, because it loads and stores one byte at a time (even in situations that hardware prefetcher might have brought data in from system memory to cache).

Example 5-40. Un-optimized Reverse Memory Copy in C

```

unsigned char* src;
unsigned char* dst;
while (len > 0)
{
  *dst-- = *src++;
  --len;
}

```

Using MOVDQA or MOVDQU, software can load and store up to 16 bytes at a time but must either ensure 16 byte alignment requirement (if using MOVDQA) or minimize the delays MOVDQU may encounter if data span across cache line boundary.

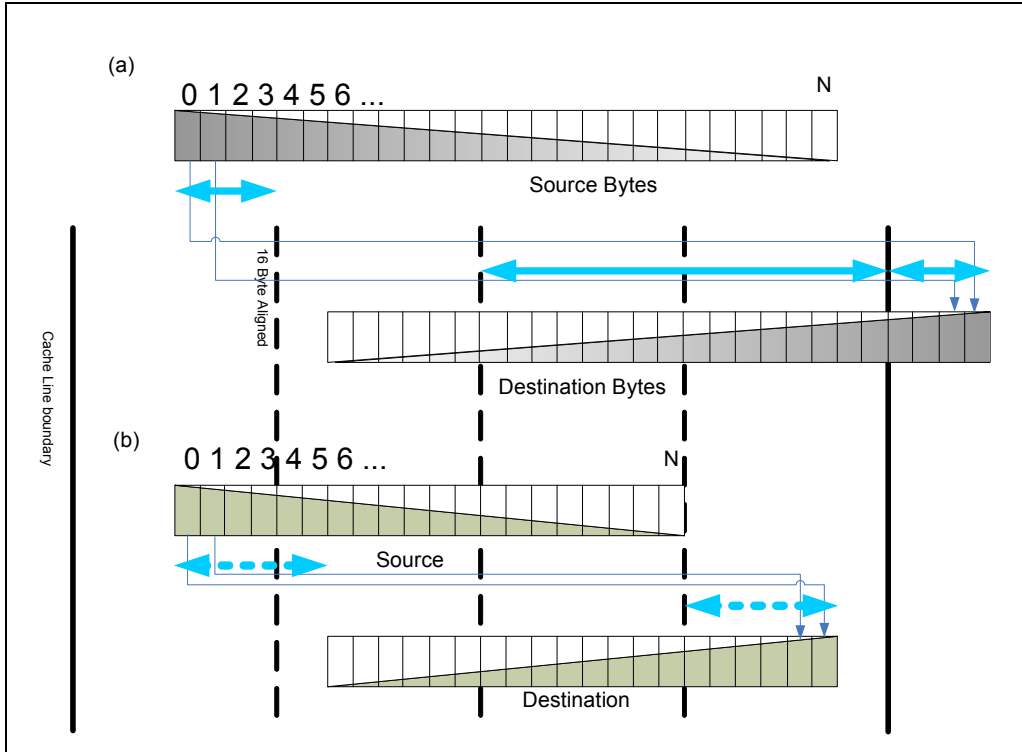


Figure 5-8. Data Alignment of Loads and Stores in Reverse Memory Copy

Given the general problem of arbitrary byte count to copy, arbitrary offsets of leading source byte and destination bytes, address alignment relative to 16 byte and cache

line boundaries, these alignment situations can be a bit complicated. Figure 5-8 (a) and (b) depict the alignment situations of reverse memory copy of N bytes.

The general guidelines for dealing with unaligned loads and stores are (in order of importance):

- Avoid stores that span cache line boundaries,
- Minimize the number of loads that span cacheline boundaries,
- Favor 16-byte aligned loads and stores over unaligned versions.

In Figure 5-8 (a), the guidelines above can be applied to the reverse memory copy problem as follows:

1. Peel off several leading destination bytes until it aligns on 16 Byte boundary, then the ensuing destination bytes can be written to using MOVAPS until the remaining byte count falls below 16 bytes.
2. After the leading source bytes have been peeled (corresponding to step 1 above), the source alignment in Figure 5-8 (a) allows loading 16 bytes at a time using MOVAPS until the remaining byte count falls below 16 bytes.

Switching the byte ordering of each 16 bytes of data can be accomplished by a 16-byte mask with PSHUFB. The pertinent code sequence is shown in Example 5-41.

Example 5-41. Using PSHUFB to Reverse Byte Ordering 16 Bytes at a Time

```

__declspec(align(16)) static const unsigned char BswapMask[16] =
{15,14,13,12,11,10,9,8,7,6,5,4,3,2,1,0};
    mov esi, src
    mov edi, dst
    mov ecx, len
    movaps xmm7, BswapMask
start:
    movdqa xmm0, [esi]
    pshufb xmm0, xmm7
    movdqa [edi-16], xmm0
sub edi, 16
    add esi, 16
    sub ecx, 16
    cmp ecx, 32
    jae start
    //handle left-overs

```

In Figure 5-8 (b), we also start with peeling the destination bytes:

1. Peel off several leading destination bytes until it aligns on 16 Byte boundary, then the ensuing destination bytes can be written to using MOVAPS until the remaining

byte count falls below 16 bytes. However, the remaining source bytes are not aligned on 16 byte boundaries, replacing MOVDQA with MOVDQU for loads will inevitably run into cache line splits.

- To achieve higher data throughput than loading unaligned bytes with MOVDQU, the 16 bytes of data targeted to each of 16 bytes of aligned destination addresses can be assembled using two aligned loads. This technique is illustrated in Figure 5-9.

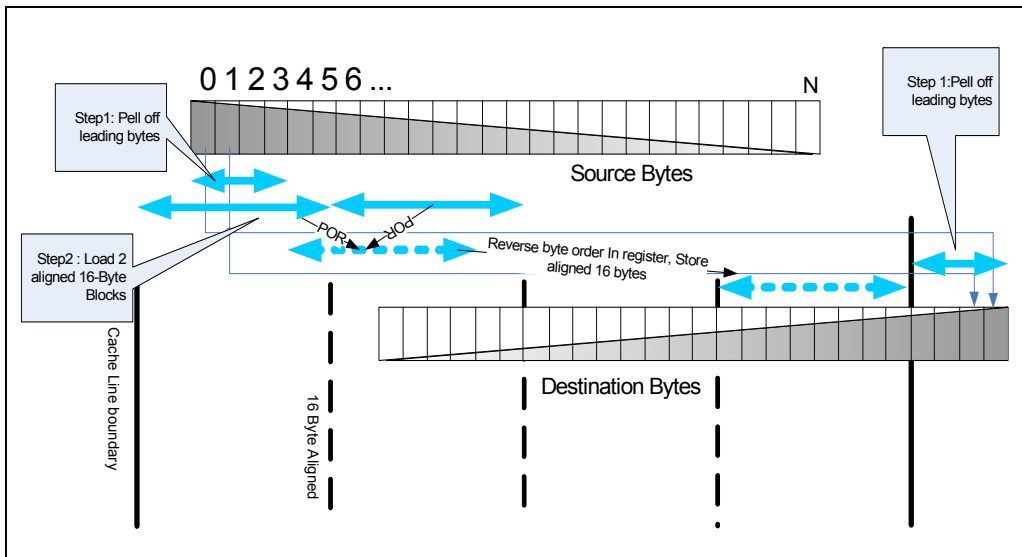


Figure 5-9. A Technique to Avoid Cacheline Split Loads in Reverse Memory Copy Using Two Aligned Loads

5.8 CONVERTING FROM 64-BIT TO 128-BIT SIMD INTEGERS

SSE2 defines a superset of 128-bit integer instructions currently available in MMX technology; the operation of the extended instructions remains. The superset simply operates on data that is twice as wide. This simplifies porting of 64-bit integer applications. However, there are few considerations:

- Computation instructions which use a memory operand that may not be aligned to a 16-byte boundary must be replaced with an unaligned 128-bit load (MOVDQU) followed by the same computation operation that uses instead register operands.

Use of 128-bit integer computation instructions with memory operands that are not 16-byte aligned will result in a #GP. Unaligned 128-bit loads and stores are

not as efficient as corresponding aligned versions; this fact can reduce the performance gains when using the 128-bit SIMD integer extensions.

- General guidelines on the alignment of memory operands are:
 - The greatest performance gains can be achieved when all memory streams are 16-byte aligned.
 - Reasonable performance gains are possible if roughly half of all memory streams are 16-byte aligned and the other half are not.
 - Little or no performance gain may result if all memory streams are not aligned to 16-bytes. In this case, use of the 64-bit SIMD integer instructions may be preferable.
- Loop counters need to be updated because each 128-bit integer instruction operates on twice the amount of data as its 64-bit integer counterpart.
- Extension of the PSHUFW instruction (shuffle word across 64-bit integer operand) across a full 128-bit operand is emulated by a combination of the following instructions: PSHUFW, PSHUFLW, and PSHUFD.
- Use of the 64-bit shift by bit instructions (PSRLQ, PSLQ) are extended to 128 bits by:
 - Use of PSRLQ and PSLQ, along with masking logic operations
 - A Code sequence rewritten to use the PSRLDQ and PSLLDQ instructions (shift double quad-word operand by bytes)

5.8.1 SIMD Optimizations and Microarchitectures

Pentium M, Intel Core Solo and Intel Core Duo processors have a different microarchitecture than Intel NetBurst microarchitecture. The following sections discuss optimizing SIMD code that targets Intel Core Solo and Intel Core Duo processors.

On Intel Core Solo and Intel Core Duo processors, `lddqu` behaves identically to `movdqu` by loading 16 bytes of data irrespective of address alignment.

5.8.1.1 Packed SSE2 Integer versus MMX Instructions

In general, 128-bit SIMD integer instructions should be favored over 64-bit MMX instructions on Intel Core Solo and Intel Core Duo processors. This is because:

- Improved decoder bandwidth and more efficient μop flows relative to the Pentium M processor
- Wider width of the XMM registers can benefit code that is limited by either decoder bandwidth or execution latency. XMM registers can provide twice the space to store data for in-flight execution. Wider XMM registers can facilitate loop-unrolling or in reducing loop overhead by halving the number of loop iterations.

In microarchitectures prior to Intel Core microarchitecture, execution throughput of 128-bit SIMD integration operations is basically the same as 64-bit MMX operations. Some shuffle/unpack/shift operations do not benefit from the front end improvements. The net impact of using 128-bit SIMD integer instruction on Intel Core Solo and Intel Core Duo processors is likely to be slightly positive overall, but there may be a few situations where their use will generate an unfavorable performance impact.

Intel Core microarchitecture generally executes 128-bit SIMD instructions more efficiently than previous microarchitectures in terms of latency and throughput, many of the limitations specific to Intel Core Duo, Intel Core Solo processors do not apply. The same is true of Intel Core microarchitecture relative to Intel NetBurst microarchitectures.

Enhanced Intel Core microarchitecture provides even more powerful 128-bit SIMD execution capabilities and more comprehensive sets of SIMD instruction extensions than Intel Core microarchitecture. The integer SIMD instructions offered by SSE4.1 operates on 128-bit XMM register only. All of these highly encourages software to favor 128-bit vectorizable code to take advantage of processors based on Enhanced Intel Core microarchitecture and Intel Core microarchitecture.

5.8.1.2 Work-around for False Dependency Issue

In processor based on Intel microarchitecture code name Nehalem, using PMOVSX and PMOVZX instructions to combine data type conversion and data movement in the same instruction will create a false-dependency due to hardware causes. A simple work-around to avoid the false dependency issue is to use PMOVSX, PMOVZX instruction solely for data type conversion and issue separate instruction to move data to destination or from origin.

Example 5-42. PMOVSX/PMOVZX Work-around to Avoid False Dependency

```
#issuing the instruction below will create a false dependency on xmm0
    pmovzxbd xmm0, dword ptr [eax]
// the above instruction may be blocked if xmm0 are updated by other instructions in flight
.....

#Alternate solution to avoid false dependency
    movd xmm0, dword ptr [eax] ; 000 hardware can hoist loads to hide latency
    pmovsxbd xmm0, xmm0
```

5.9 TUNING PARTIALLY VECTORIZABLE CODE

Some loop structured code are more difficult to vectorize than others. Example 5-43 depicts a loop carrying out table look-up operation and some arithmetic computation.

Example 5-43. Table Look-up Operations in C Code

```
// pIn1      integer input arrays.
// pOut      integer output array.
// count     size of array.
// LookUpTable integer values.
TABLE_SIZE  size of the look-up table.
for (unsigned i=0; i < count; i++)
{
    pOut[i] =
        ( ( LookUpTable[pIn1[i] % TABLE_SIZE] + pIn1[i] + 17 ) | 17
          ) % 256;
}
```

Although some of the arithmetic computations and input/output to data array in each iteration can be easily vectorizable, but the table look-up via an index array is not. This creates different approaches to tuning. A compiler can take a scalar approach to execute each iteration sequentially. Hand-tuning of such loops may use a couple of different techniques to handle the non-vectorizable table look-up operation. One vectorization technique is to load the input data for four iteration at once, then use SSE2 instruction to shift out individual index out of an XMM register to carry out table look-up sequentially. The shift technique is depicted by Example 5-44. Another technique is to use PEXTRD in SSE4.1 to extract the index from an XMM directly and then carry out table look-up sequentially. The PEXTRD technique is depicted by Example 5-45.

Example 5-44. Shift Techniques on Non-Vectorizable Table Look-up

```

int modulo[4] = {256-1, 256-1, 256-1, 256-1};
int c[4] = {17, 17, 17, 17};
    mov     esi, pln1
    mov     ebx, pOut
    mov     ecx, count
    mov     edx, pLookUpTablePTR
    movaps  xmm6, modulo
    movaps  xmm5, c
loop:
// vectorizable multiple consecutive data accesses
    movaps  xmm4, [esi]      // read 4 indices from pln1
    movaps  xmm7, xmm4
    pand    xmm7, tableSize
//Table look-up is not vectorizable, shift out one data element to look up table one by one
    movd    eax, xmm7        // get first index
    movd    xmm0, word ptr[edx + eax*4]
    psrldq  xmm7, 4
    movd    eax, xmm7        // get 2nd index
    movd    xmm1, word ptr[edx + eax*4]
    psrldq  xmm7, 4
    movd    eax, xmm7        // get 3rdindex
    movd    xmm2, word ptr[edx + eax*4]
    psrldq  xmm7, 4
    movd    eax, xmm7        // get fourth index
    movd    xmm3, word ptr[edx + eax*4]
//end of scalar part
//packing
    movlhps xmm1,xmm3
    psllq   xmm1,32
    movlhps xmm0,xmm2
    orps    xmm0,xmm1
//end of packing

```

(continue)

Example 5-44. Shift Techniques on Non-Vectorizable Table Look-up (Contd.)

```

//Vectorizable computation operations
padd    xmm0, xmm4 //+pln1
padd    xmm0, xmm5 // +17
por     xmm0, xmm5
andps   xmm0, xmm6 //mod
movaps  [ebx], xmm0
//end of vectorizable operation

add     ebx, 16
add     esi, 16
add     edi, 16
sub     ecx, 1
test    ecx, ecx
jne     lloop

```

Example 5-45. PEXTRD Techniques on Non-Vectorizable Table Look-up

```

int modulo[4] = {256-1, 256-1, 256-1, 256-1};
int c[4] = {17, 17, 17, 17};
mov     esi, pln1
mov     ebx, pOut
mov     ecx, count
mov     edx, pLookUpTablePTR
movaps  xmm6, modulo
movaps  xmm5, c
lloop:
// vectorizable multiple consecutive data accesses
movaps  xmm4, [esi]      // read 4 indices from pln1
movaps  xmm7, xmm4
pand    xmm7, tableSize
//Table look-up is not vectorizable, extract one data element to look up table one by one
movd    eax, xmm7      // get first index
mov     eax, [edx + eax*4]
movd    xmm0, eax

```

(continue)

Example 5-45. PEXTRD Techniques on Non-Vectorizable Table Look-up (Contd.)

```

pextrd    eax, xmm7, 1    // extract 2nd index
mov       eax, [edx + eax*4]
pinsrd    xmm0, eax, 1
pextrd    eax, xmm7, 2    // extract 2nd index
mov       eax, [edx + eax*4]
pinsrd    xmm0, eax, 2
pextrd    eax, xmm7, 3    // extract 2nd index
mov       eax, [edx + eax*4]
pinsrd    xmm0, eax, 2
//end of scalar part
//packing not needed
//Vectorizable operations
padd      xmm0, xmm4 //+pln1
padd      xmm0, xmm5 // +17
por       xmm0, xmm5
andps     xmm0, xmm6 //mod
movaps    [ebx], xmm0

add       ebx, 16
add       esi, 16
add       edi, 16
sub       ecx, 1
test      ecx, ecx
jne       lloop

```

The effectiveness of these two hand-tuning techniques on partially vectorizable code depends on the relative cost of transforming data layout format using various forms of pack and unpack instructions.

The shift technique requires additional instructions to pack scalar table values into an XMM to transition into vectorized arithmetic computations. The net performance gain or loss of this technique will vary with the characteristics of different microarchitectures. The alternate PEXTRD technique uses less instruction to extract each index, does not require extraneous packing of scalar data into packed SIMD data format to begin vectorized arithmetic computation.

5.10 PARALLEL MODE AES ENCRYPTION AND DECRYPTION

To deliver optimal encryption and decryption throughput using AESNI, software can optimize by re-ordering the computations and working on multiple blocks in parallel. This can speed up encryption (and decryption) in parallel modes of operation such as

ECB, CTR, and CBC-Decrypt (comparing to CBC-Encrypt which is serial mode of operation). See details in Recommendation for Block Cipher Modes of Operation?. The Related Documentation section provides a pointer to this document.

In Intel microarchitecture code name Sandy Bridge, the AES round instructions (AESENC / AESECLAST / AESDEC / AESDECLAST) have a throughput of one cycle and latency of eight cycles. This allows independent AES instructions for multiple blocks to be dispatched every cycle, if data can be provided sufficiently fast. Compared to the prior Intel microarchitecture code name Westmere, where these instructions have throughput of two cycles and a latency of six cycles, the AES encryption/decryption throughput can be significantly increased, for parallel modes of operation.

To achieve optimal parallel operation with multiple blocks, write the AES software sequences in a way that it computes one AES round on multiple blocks, using one Round Key, and then it continues to compute the subsequent round for multiple blocks, using another Round Key.

For such software optimization, you need to define the number of blocks that are processed in parallel. In Intel microarchitecture code name Sandy Bridge, the optimal parallelization parameter is eight blocks, compared to four blocks on prior microarchitecture.

5.10.1 AES Counter Mode of Operation

Example 5-46 is an example of a function that implements the Counter Mode (CTR mode) of operations while operating on eight blocks in parallel. The following pseudo-code encrypts n data blocks of 16 byte each (PT[i]):

Example 5-46. Pseudo-Code Flow of AES Counter Mode Operation

```

CTRBLK := NONCE || IV || ONE
FOR i := 1 to n-1 DO
    CT[i] := PT[i] XOR AES(CTRBLK)
    CTRBLK := CTRBLK + 1) % 256;
END
CT[n] := PT[n] XOR TRUNC(AES(CTRBLK)) CTRBLK := NONCE || IV || ONE
FOR i := 1 to n-1 DO
    CT[i] := PT[i] XOR AES(CTRBLK)// CT [i] is the i-th ciphertext block
    CTRBLK := CTRBLK + 1
END
CT[n] := PT[n] XOR TRUNC(AES(CTRBLK))

```

Example 5-47 in the following pages show the assembly implementation of the above code, optimized for Intel microarchitecture code name Sandy Bridge.

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel

```

/*****
/* This function encrypts an input buffer using AES in CTR mode      */
/* The parameters:                                                  */
/*  const unsigned char *in - pointer to the plaintext for encryption or */
/*  ciphertext for decryption                                       */
/*  unsigned char *out - pointer to the buffer where the encrypted*/
/*                        data will be stored                       */
/*  const unsigned char ivec[8] - 8 bytes of the initialization vector */
/*  const unsigned char nonce[4] - 4 bytes of the nonce             */
/*  const unsigned long length - the length of the input in bytes   */
/*  int number_of_rounds - number of AES round. 10 = AES128, 12 = AES192, 14 = AES256 */
/*  unsigned char *key_schedule - pointer to the AES key schedule   */
*****/
//void AES_128_CTR_encrypt_parallelize_8_blocks_unrolled (
//      const unsigned char *in,
//      unsigned char *out,
//      const unsigned char ivec[8],
//      const unsigned char nonce[4],
//      const unsigned long length,
//      unsigned char *key_schedule)
.align 16,0x90
.align 16
ONE:          .quad 0x00000000,0x00000001
.align 16
FOUR:        .quad 0x00000004,0x00000004
.align 16
EIGHT:       .quad 0x00000008,0x00000008
.align 16
TWO_N_ONE:   .quad 0x00000002,0x00000001
.align 16
TWO_N_TWO:   .quad 0x00000002,0x00000002
.align 16
LOAD_HIGH_BROADCAST_AND_BSWAP:  .byte 15,14,13,12,11,10,9,8
                                   .byte 15,14,13,12,11,10,9,8

align 16
BSWAP_EPI_64:  .byte 7,6,5,4,3,2,1,0
                .byte 15,14,13,12,11,10,9,8

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

.globl AES_CTR_encrypt

AES_CTR_encrypt:
# parameter 1: %rdi      # parameter 2: %rsi
# parameter 3: %rdx      # parameter 4: %rcx
# parameter 5: %r8       # parameter 6: %r9
# parameter 7: 8 + %rsp
movq  %r8, %r10
    movl  8(%rsp), %r12d
    shrq  $4, %r8
    shlq  $60, %r10
    je    NO_PARTS
    addq  $1, %r8
NO_PARTS:
    movq  %r8, %r10
    shlq  $61, %r10
    shrq  $61, %r10

    pinsrq $1, (%rdx), %xmm0
    pinsrd $1, (%rcx), %xmm0
    psrldq $4, %xmm0
    movdqa %xmm0, %xmm4
    pshufb (LOAD_HIGH_BROADCAST_AND_BSWAP), %xmm4
    paddq (TWO_N_ONE), %xmm4
    movdqa %xmm4, %xmm1
    paddq (TWO_N_TWO), %xmm4
    movdqa %xmm4, %xmm2
    paddq (TWO_N_TWO), %xmm4
    movdqa %xmm4, %xmm3
    paddq (TWO_N_TWO), %xmm4
    pshufb (BSWAP_EPI_64), %xmm1
    pshufb (BSWAP_EPI_64), %xmm2
    pshufb (BSWAP_EPI_64), %xmm3
    pshufb (BSWAP_EPI_64), %xmm4

    shrq  $3, %r8
    je    REMAINDER
    subq  $128, %rsi
    subq  $128, %rdi

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

LOOP:
  addq  $128,%rsi
  addq  $128,%rdi

  movdqa %xmm0,%xmm7
  movdqa %xmm0,%xmm8
  movdqa %xmm0,%xmm9
  movdqa %xmm0,%xmm10
  movdqa %xmm0,%xmm11
  movdqa %xmm0,%xmm12
  movdqa %xmm0,%xmm13
  movdqa %xmm0,%xmm14

  shufpd $2,%xmm1,%xmm7
  shufpd $0,%xmm1,%xmm8
  shufpd $2,%xmm2,%xmm9
  shufpd $0,%xmm2,%xmm10
  shufpd $2,%xmm3,%xmm11
  shufpd $0,%xmm3,%xmm12
  shufpd $2,%xmm4,%xmm13
  shufpd $0,%xmm4,%xmm14

  pshufb (BSWAP_EPI_64),%xmm1
  pshufb (BSWAP_EPI_64),%xmm2
  pshufb (BSWAP_EPI_64),%xmm3
  pshufb (BSWAP_EPI_64),%xmm4

  movdqa (%r9),%xmm5
  movdqa 16(%r9),%xmm6

  paddq (EIGHT),%xmm1
  paddq (EIGHT),%xmm2
  paddq (EIGHT),%xmm3
  paddq (EIGHT),%xmm4

  pxor  %xmm5,%xmm7
  pxor  %xmm5,%xmm8
  pxor  %xmm5,%xmm9
  pxor  %xmm5,%xmm10

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

pxor  %xmm5, %xmm11
pxor  %xmm5, %xmm12
pxor  %xmm5, %xmm13
pxor  %xmm5, %xmm14

pshufb (BSWAP_EPI_64), %xmm1
pshufb (BSWAP_EPI_64), %xmm2
pshufb (BSWAP_EPI_64), %xmm3
pshufb (BSWAP_EPI_64), %xmm4

aesenc %xmm6, %xmm7
aesenc %xmm6, %xmm8
aesenc %xmm6, %xmm9
aesenc %xmm6, %xmm10
aesenc %xmm6, %xmm11
aesenc %xmm6, %xmm12
aesenc %xmm6, %xmm13
aesenc %xmm6, %xmm14

movdqa 32(%r9), %xmm5
movdqa 48(%r9), %xmm6

aesenc %xmm5, %xmm7
aesenc %xmm5, %xmm8
aesenc %xmm5, %xmm9
aesenc %xmm5, %xmm10
aesenc %xmm5, %xmm11
aesenc %xmm5, %xmm12
aesenc %xmm5, %xmm13
aesenc %xmm5, %xmm14

aesenc %xmm6, %xmm7
aesenc %xmm6, %xmm8
aesenc %xmm6, %xmm9
aesenc %xmm6, %xmm10
aesenc %xmm6, %xmm11
aesenc %xmm6, %xmm12
aesenc %xmm6, %xmm13
aesenc %xmm6, %xmm14

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

movdqa 64(%r9), %xmm5
movdqa 80(%r9), %xmm6

aesenc %xmm5, %xmm7
aesenc %xmm5, %xmm8
aesenc %xmm5, %xmm9
aesenc %xmm5, %xmm10
aesenc %xmm5, %xmm11
aesenc %xmm5, %xmm12
aesenc %xmm5, %xmm13
aesenc %xmm5, %xmm14

aesenc %xmm6, %xmm7
aesenc %xmm6, %xmm8
aesenc %xmm6, %xmm9
aesenc %xmm6, %xmm10
aesenc %xmm6, %xmm11
aesenc %xmm6, %xmm12
aesenc %xmm6, %xmm13
aesenc %xmm6, %xmm14
movdqa 96(%r9), %xmm5
movdqa 112(%r9), %xmm6

aesenc %xmm5, %xmm7
aesenc %xmm5, %xmm8
aesenc %xmm5, %xmm9
aesenc %xmm5, %xmm10
aesenc %xmm5, %xmm11
aesenc %xmm5, %xmm12
aesenc %xmm5, %xmm13
aesenc %xmm5, %xmm14

aesenc %xmm6, %xmm7
aesenc %xmm6, %xmm8
aesenc %xmm6, %xmm9
aesenc %xmm6, %xmm10
aesenc %xmm6, %xmm11
aesenc %xmm6, %xmm12
aesenc %xmm6, %xmm13
aesenc %xmm6, %xmm14

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

movdqa 128(%r9), %xmm5
movdqa 144(%r9), %xmm6
movdqa 160(%r9), %xmm15
cmp    $12, %r12d

```

```

aesenc %xmm5, %xmm7
aesenc %xmm5, %xmm8
aesenc %xmm5, %xmm9
aesenc %xmm5, %xmm10
aesenc %xmm5, %xmm11
aesenc %xmm5, %xmm12
aesenc %xmm5, %xmm13
aesenc %xmm5, %xmm14

```

```

aesenc %xmm6, %xmm7
aesenc %xmm6, %xmm8
aesenc %xmm6, %xmm9
aesenc %xmm6, %xmm10
aesenc %xmm6, %xmm11
aesenc %xmm6, %xmm12
aesenc %xmm6, %xmm13
aesenc %xmm6, %xmm14

```

```

jb    LAST

```

```

movdqa 160(%r9), %xmm5
movdqa 176(%r9), %xmm6
movdqa 192(%r9), %xmm15
cmp    $14, %r12d

```

```

aesenc %xmm5, %xmm7
aesenc %xmm5, %xmm8
aesenc %xmm5, %xmm9
aesenc %xmm5, %xmm10
aesenc %xmm5, %xmm11
aesenc %xmm5, %xmm12
aesenc %xmm5, %xmm13
aesenc %xmm5, %xmm14

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

aesenc %xmm6,%xmm7
aesenc %xmm6,%xmm8
aesenc %xmm6,%xmm9
aesenc %xmm6,%xmm10
aesenc %xmm6,%xmm11
aesenc %xmm6,%xmm12
aesenc %xmm6,%xmm13
aesenc %xmm6,%xmm14

```

```

jb     LAST

```

```

movdqa 192(%r9), %xmm5
movdqa 208(%r9), %xmm6
movdqa 224(%r9), %xmm15

```

```

aesenc %xmm5,%xmm7
aesenc %xmm5,%xmm8
aesenc %xmm5,%xmm9
aesenc %xmm5,%xmm10
aesenc %xmm5,%xmm11
aesenc %xmm5,%xmm12
aesenc %xmm5,%xmm13
aesenc %xmm5,%xmm14

```

```

aesenc %xmm6,%xmm7
aesenc %xmm6,%xmm8
aesenc %xmm6,%xmm9
aesenc %xmm6,%xmm10
aesenc %xmm6,%xmm11
aesenc %xmm6,%xmm12
aesenc %xmm6,%xmm13
aesenc %xmm6,%xmm14

```

```

LAST:

```

```

aesenclast %xmm15,%xmm7
aesenclast %xmm15,%xmm8
aesenclast %xmm15,%xmm9
aesenclast %xmm15,%xmm10
aesenclast %xmm15,%xmm11
aesenclast %xmm15,%xmm12

```

```

        (continue)

```

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

aesencast %xmm15, %xmm13
aesencast %xmm15, %xmm14

```

```

pxor (%rdi), %xmm7
pxor 16(%rdi), %xmm8
pxor 32(%rdi), %xmm9
pxor 48(%rdi), %xmm10
pxor 64(%rdi), %xmm11
pxor 80(%rdi), %xmm12
pxor 96(%rdi), %xmm13
pxor 112(%rdi), %xmm14

```

```

dec %r8

```

```

movdqu %xmm7, (%rsi)
movdqu %xmm8, 16(%rsi)
movdqu %xmm9, 32(%rsi)
movdqu %xmm10, 48(%rsi)
movdqu %xmm11, 64(%rsi)
movdqu %xmm12, 80(%rsi)
movdqu %xmm13, 96(%rsi)
movdqu %xmm14, 112(%rsi)
jne LOOP

```

```

addq $128,%rsi
addq $128,%rdi

```

REMAINDER:

```

cmp $0, %r10
je END
shufpd $2, %xmm1, %xmm0

```

IN_LOOP:

```

movdqa %xmm0, %xmm11
pshufb (BSWAP_EPI_64), %xmm0
pxor (%r9), %xmm11
paddq (ONE), %xmm0
aesenc 16(%r9), %xmm11
aesenc 32(%r9), %xmm11
pshufb (BSWAP_EPI_64), %xmm0

```

(continue)

Example 5-47. AES128-CTR Implementation with Eight Block in Parallel (Contd.)

```

aesenc 48(%r9), %xmm11
aesenc 64(%r9), %xmm11
aesenc 80(%r9), %xmm11
aesenc 96(%r9), %xmm11
aesenc 112(%r9), %xmm11
aesenc 128(%r9), %xmm11
aesenc 144(%r9), %xmm11
movdqa 160(%r9), %xmm2
cmp    $12, %r12d
jnb   IN_LAST
aesenc 160(%r9), %xmm11
aesenc 176(%r9), %xmm11
movdqa 192(%r9), %xmm2
cmp    $14, %r12d
jnb   IN_LAST
aesenc 192(%r9), %xmm11
aesenc 208(%r9), %xmm11
movdqa 224(%r9), %xmm2
IN_LAST:
aesenclast %xmm2, %xmm11
pxor    (%rdi), %xmm11
movdqu  %xmm11, (%rsi)
addq   $16, %rdi
addq   $16, %rsi
dec    %r10
jne    IN_LOOP
END:
ret

```

5.10.2 AES Key Expansion Alternative

In Intel microarchitecture code name Sandy Bridge, the throughput of AESKEYGEN-ASSIST is two cycles with higher latency than the AESENC/AESDEC instructions. Software may consider implementing the AES key expansion by using the AESENCLAST instruction with the second operand (i.e., the round key) being the RCON value, duplicated four times in the register. The AESENCLAST instruction performs the SubBytes step and the xor-with-RCON step, while the ROTWORD step can be done using a PSHUFB instruction. Following are code examples of AES128 key expansion using either method.

Example 5-48. AES128 Key Expansion

<pre>// Use AESKEYGENASSIST .align 16,0x90 .globl AES_128_Key_Expansion AES_128_Key_Expansion: # parameter 1: %rdi # parameter 2: %rsi movl \$10, 240(%rsi) movdqu (%rdi), %xmm1 movdqa %xmm1, (%rsi) ASSISTS: aeskeygenassist \$1, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 16(%rsi) aeskeygenassist \$2, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 32(%rsi) aeskeygenassist \$4, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 48(%rsi) aeskeygenassist \$8, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 64(%rsi) aeskeygenassist \$16, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 80(%rsi) aeskeygenassist \$32, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 96(%rsi) aeskeygenassist \$64, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 112(%rsi) aeskeygenassist \$0x80, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 128(%rsi) aeskeygenassist \$0x1b, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128</pre> <p style="text-align: right;">(continue)</p>	<pre>// Use AESENCLAST mask: .long 0x0c0f0e0d,0x0c0f0e0d,0x0c0f0e0d,0x0c0f0e0d con1: .long 1,1,1,1 con2: .long 0x1b,0x1b,0x1b,0x1b .align 16,0x90 .globl AES_128_Key_Expansion AES_128_Key_Expansion: # parameter 1: %rdi # parameter 2: %rsi movdqu (%rdi), %xmm1 movdqa %xmm1, (%rsi) movdqa %xmm1, %xmm2 movdqa (con1), %xmm0 movdqa (mask), %xmm15 mov \$8, %ax LOOP1: add \$16, %rsi dec %ax pshufb %xmm15,%xmm2 aesenclast %xmm0, %xmm2 pslld \$1, %xmm0 movdqa %xmm1, %xmm3 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pxor %xmm2, %xmm1</pre> <p style="text-align: right;">(continue)</p>
--	--

Example 5-48. AES128 Key Expansion (Contd.)

<pre> movdqa %xmm1, 144(%rsi) aeskeygenassist \$0x36, %xmm1, %xmm2 call PREPARE_ROUNDKEY_128 movdqa %xmm1, 160(%rsi) ret PREPARE_ROUNDKEY_128: pshufd \$255, %xmm2, %xmm2 movdqa %xmm1, %xmm3 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pxor %xmm2, %xmm1 ret </pre>	<pre> movdqa %xmm1, (%rsi) movdqa %xmm1, %xmm2 jne LOOP1 movdqa (con2), %xmm0 pshufb %xmm15,%xmm2 aesencast %xmm0, %xmm2 pslld \$1, %xmm0 movdqa %xmm1, %xmm3 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pxor %xmm2, %xmm1 movdqa %xmm1, 16(%rsi) movdqa %xmm1, %xmm2 pshufb %xmm15,%xmm2 aesencast %xmm0, %xmm2 movdqa %xmm1, %xmm3 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pslldq \$4, %xmm3 pxor %xmm3, %xmm1 pxor %xmm2, %xmm1 movdqa %xmm1, 32(%rsi) movdqa %xmm1, %xmm2 ret </pre>
---	---

CHAPTER 6

OPTIMIZING FOR SIMD FLOATING-POINT APPLICATIONS

This chapter discusses rules for optimizing for the single-instruction, multiple-data (SIMD) floating-point instructions available in SSE, SSE2, SSE3, and SSE4.1. The chapter also provides examples that illustrate the optimization techniques for single-precision and double-precision SIMD floating-point applications.

6.1 GENERAL RULES FOR SIMD FLOATING-POINT CODE

The rules and suggestions in this section help optimize floating-point code containing SIMD floating-point instructions. Generally, it is important to understand and balance port utilization to create efficient SIMD floating-point code. Basic rules and suggestions include the following:

- Follow all guidelines in Chapter 3 and Chapter 4.
- Mask exceptions to achieve higher performance. When exceptions are unmasked, software performance is slower.
- Utilize the flush-to-zero and denormals-are-zero modes for higher performance to avoid the penalty of dealing with denormals and underflows.
- Use the reciprocal instructions followed by iteration for increased accuracy. These instructions yield reduced accuracy but execute much faster. Note the following:
 - If reduced accuracy is acceptable, use them with no iteration.
 - If near full accuracy is needed, use a Newton-Raphson iteration.
 - If full accuracy is needed, then use divide and square root which provide more accuracy, but slow down performance.

6.2 PLANNING CONSIDERATIONS

Whether adapting an existing application or creating a new one, using SIMD floating-point instructions to achieve optimum performance gain requires programmers to consider several issues. In general, when choosing candidates for optimization, look for code segments that are computationally intensive and floating-point intensive. Also consider efficient use of the cache architecture.

The sections that follow answer the questions that should be raised before implementation:

- Can data layout be arranged to increase parallelism or cache utilization?

- Which part of the code benefits from SIMD floating-point instructions?
- Is the current algorithm the most appropriate for SIMD floating-point instructions?
- Is the code floating-point intensive?
- Do either single-precision floating-point or double-precision floating-point computations provide enough range and precision?
- Does the result of computation affected by enabling flush-to-zero or denormals-to-zero modes?
- Is the data arranged for efficient utilization of the SIMD floating-point registers?
- Is this application targeted for processors without SIMD floating-point instructions?

See also: Section 4.2, "Considerations for Code Conversion to SIMD Programming."

6.3 USING SIMD FLOATING-POINT WITH X87 FLOATING-POINT

Because the XMM registers used for SIMD floating-point computations are separate registers and are not mapped to the existing x87 floating-point stack, SIMD floating-point code can be mixed with x87 floating-point or 64-bit SIMD integer code.

With Intel Core microarchitecture, 128-bit SIMD integer instructions provides substantially higher efficiency than 64-bit SIMD integer instructions. Software should favor using SIMD floating-point and integer SIMD instructions with XMM registers where possible.

6.4 SCALAR FLOATING-POINT CODE

There are SIMD floating-point instructions that operate only on the lowest order element in the SIMD register. These instructions are known as scalar instructions. They allow the XMM registers to be used for general-purpose floating-point computations.

In terms of performance, scalar floating-point code can be equivalent to or exceed x87 floating-point code and has the following advantages:

- SIMD floating-point code uses a flat register model, whereas x87 floating-point code uses a stack model. Using scalar floating-point code eliminates the need to use FXCH instructions. These have performance limits on the Intel Pentium 4 processor.
- Mixing with MMX technology code without penalty.
- Flush-to-zero mode.
- Shorter latencies than x87 floating-point.

When using scalar floating-point instructions, it is not necessary to ensure that the data appears in vector form. However, the optimizations regarding alignment, scheduling, instruction selection, and other optimizations covered in Chapter 3 and Chapter 4 should be observed.

6.5 DATA ALIGNMENT

SIMD floating-point data is 16-byte aligned. Referencing unaligned 128-bit SIMD floating-point data will result in an exception unless `MOVUPS` or `MOVUPD` (move unaligned packed single or unaligned packed double) is used. The unaligned instructions used on aligned or unaligned data will also suffer a performance penalty relative to aligned accesses.

See also: Section 4.4, “Stack and Data Alignment.”

6.5.1 Data Arrangement

Because SSE and SSE2 incorporate SIMD architecture, arranging data to fully use the SIMD registers produces optimum performance. This implies contiguous data for processing, which leads to fewer cache misses. Correct data arrangement can potentially quadruple data throughput when using SSE or double throughput when using SSE2. Performance gains can occur because four data elements can be loaded with 128-bit load instructions into XMM registers using SSE (`MOVAPS`). Similarly, two data elements can be loaded with 128-bit load instructions into XMM registers using SSE2 (`MOVAPD`).

Refer to the Section 4.4, “Stack and Data Alignment,” for data arrangement recommendations. Duplicating and padding techniques overcome misalignment problems that occur in some data structures and arrangements. This increases the data space but avoids penalties for misaligned data access.

For some applications (for example: 3D geometry), traditional data arrangement requires some changes to fully utilize the SIMD registers and parallel techniques. Traditionally, the data layout has been an array of structures (AoS). To fully utilize the SIMD registers in such applications, a new data layout has been proposed — a structure of arrays (SoA) resulting in more optimized performance.

6.5.1.1 Vertical versus Horizontal Computation

The majority of the floating-point arithmetic instructions in SSE/SSE2 provide greater performance gain on vertical data processing for parallel data elements. This means each element of the destination is the result of an arithmetic operation performed from the source elements in the same vertical position (Figure 6-1).

To supplement these homogeneous arithmetic operations on parallel data elements, SSE and SSE2 provides data movement instructions (e.g., `SHUFPS`, `UNPCKLPS`,

UNPCKHPS, MOVLHPS, MOVHLPS, etc.) that facilitate moving data elements horizontally.

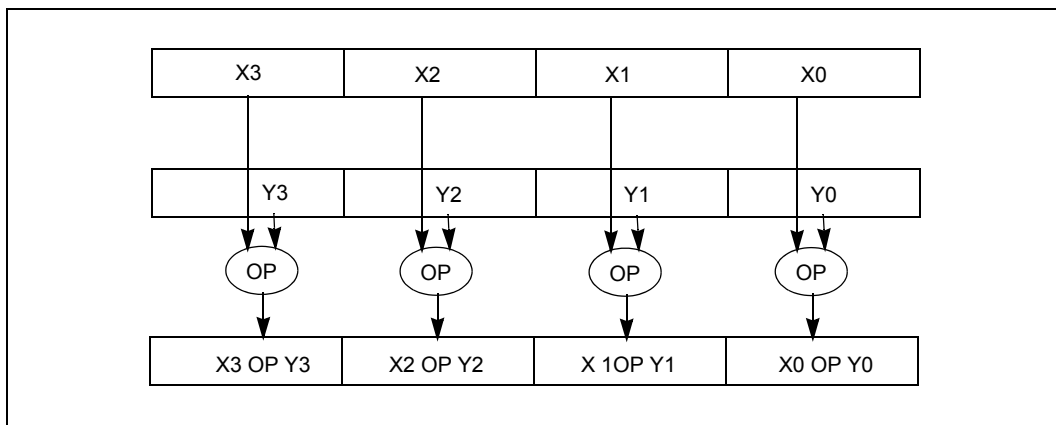


Figure 6-1. Homogeneous Operation on Parallel Data Elements

The organization of structured data have a significant impact on SIMD programming efficiency and performance. This can be illustrated using two common type of data structure organizations:

- **Array of Structure:** This refers to the arrangement of an array of data structures. Within the data structure, each member is a scalar. This is shown in Figure 6-2. Typically, a repetitive sequence of computation is applied to each element of an array, i.e. a data structure. Computational sequence for the scalar members of the structure is likely to be non-homogeneous within each iteration. AoS is generally associated with a horizontal computation model.

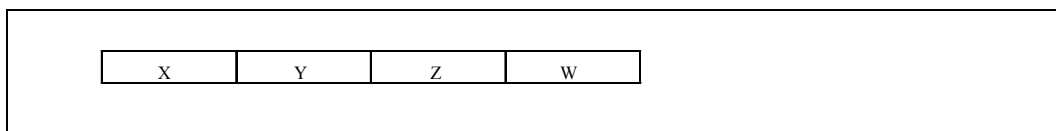


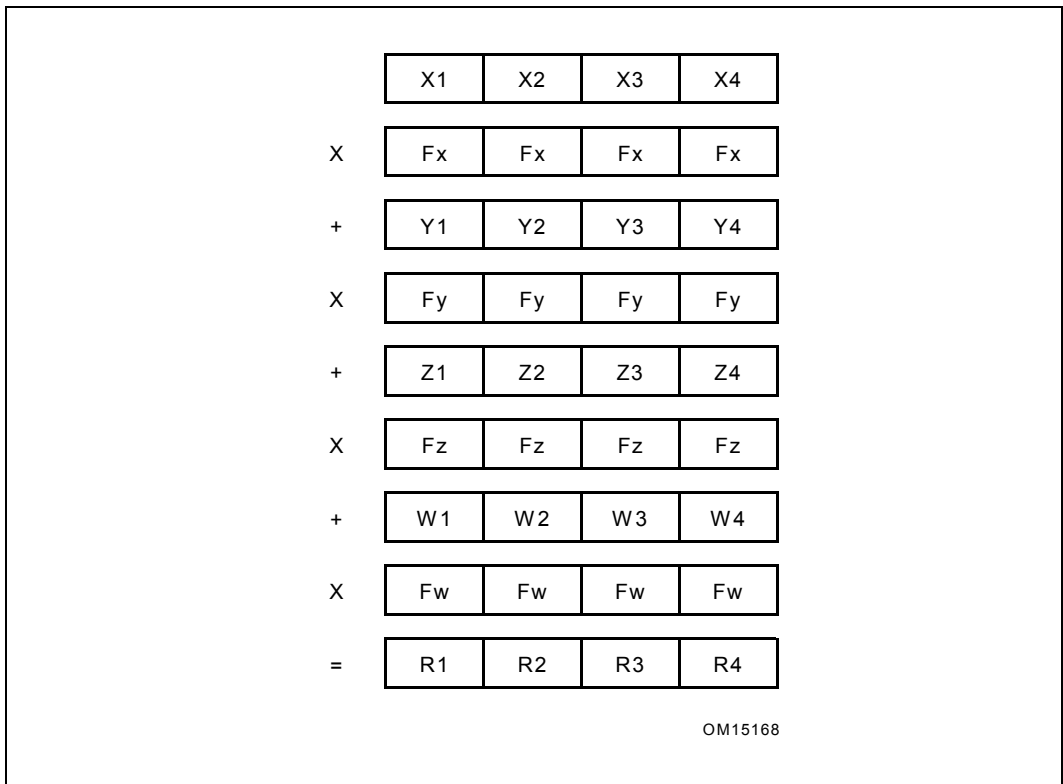
Figure 6-2. Horizontal Computation Model

- **Structure of Array:** Here, each member of the data structure is an array. Each element of the array is a scalar. This is shown Table 6-1. Repetitive computational sequence is applied to scalar elements and homogeneous operation can be easily achieved across consecutive iterations within the same structural member. Consequently, SoA is generally amenable to the vertical computation model.

Table 6-1. SoA Form of Representing Vertices Data

Vx array	X1	X2	X3	X4	Xn
Vy array	Y1	Y2	Y3	Y4	Yn
Vz array	Z1	Z2	Z3	Y4	Zn
Vw array	W1	W2	W3	W4	Wn

Using SIMD instructions with vertical computation on SOA arrangement can achieve higher efficiency and performance than AOS and horizontal computation. This can be seen with dot-product operation on vectors. The dot product operation on SoA arrangement is shown in Figure 6-3.

**Figure 6-3. Dot Product Operation**

Example 6-1 shows how one result would be computed for seven instructions if the data were organized as AoS and using SSE alone: four results would require 28 instructions.

Example 6-1. Pseudocode for Horizontal (xyz, AoS) Computation

```

mulps    ; x*x', y*y', z*z'
movaps   ; reg->reg move, since next steps overwrite
shufps  ; get b,a,d,c from a,b,c,d
addps   ; get a+b,a+b,c+d,c+d
movaps   ; reg->reg move
shufps  ; get c+d,c+d,a+b,a+b from prior addps
addps   ; get a+b+c+d,a+b+c+d,a+b+c+d,a+b+c+d

```

Now consider the case when the data is organized as SoA. Example 6-2 demonstrates how four results are computed for five instructions.

Example 6-2. Pseudocode for Vertical (xxxx, yyyy, zzzz, SoA) Computation

```

mulps   ; x*x' for all 4 x-components of 4 vertices
mulps   ; y*y' for all 4 y-components of 4 vertices
mulps   ; z*z' for all 4 z-components of 4 vertices
addps   ; x*x' + y*y'
addps   ; x*x'+y*y'+z*z'

```

For the most efficient use of the four component-wide registers, reorganizing the data into the SoA format yields increased throughput and hence much better performance for the instructions used.

As seen from this simple example, vertical computation can yield 100% use of the available SIMD registers to produce four results. (The results may vary for other situations.) If the data structures are represented in a format that is not “friendly” to vertical computation, it can be rearranged “on the fly” to facilitate better utilization of the SIMD registers. This operation is referred to as “swizzling” operation and the reverse operation is referred to as “deswizzling.”

6.5.1.2 Data Swizzling

Swizzling data from SoA to AoS format can apply to a number of application domains, including 3D geometry, video and imaging. Two different swizzling techniques can be adapted to handle floating-point and integer data. Example 6-3 illustrates a swizzle function that uses SHUFPS, MOVLHPS, MOVHLPS instructions.

Example 6-3. Swizzling Data Using SHUFPS, MOVLHPS, MOVHLP

```

typedef struct _VERTEX_AOS {
    float x, y, z, color;
} Vertex_aos;           // AoS structure declaration
typedef struct _VERTEX_SOA {
    float x[4], float y[4], float z[4];
    float color[4];
} Vertex_soa;          // SoA structure declaration
void swizzle_asm (Vertex_aos *in, Vertex_soa *out)
{
    // in mem: x1y1z1w1-x2y2z2w2-x3y3z3w3-x4y4z4w4-
    // SWIZZLE XYZW --> XXXX
    asm {
        mov ebx, in           // get structure addresses
        mov edx, out

        movaps xmm1, [ebx]    // x4 x3 x2 x1
        movaps xmm2, [ebx + 16] // y4 y3 y2 y1
        movaps xmm3, [ebx + 32] // z4 z3 z2 z1
        movaps xmm4, [ebx + 48] // w4 w3 w2 w1
        movaps xmm7, xmm4 // xmm7= w4 z4 y4 x4
        movhlps xmm7, xmm3 // xmm7= w4 z4 w3 z3
        movaps xmm6, xmm2 // xmm6= w2 z2 y2 x2
        movhlps xmm3, xmm4 // xmm3= y4 x4 y3 x3
        movhlps xmm2, xmm1 // xmm2= w2 z2 w1 z1
        movhlps xmm1, xmm6 // xmm1= y2 x2 y1 x1

        movaps xmm6, xmm2 // xmm6= w2 z2 w1 z1
        movaps xmm5, xmm1 // xmm5= y2 x2 y1 x1
        shufps xmm2, xmm7, 0xDD // xmm2= w4 w3 w2 w1 => v4
        shufps xmm1, xmm3, 0x88 // xmm1= x4 x3 x2 x1 => v1
        shufps xmm5, xmm3, 0xDD // xmm5= y4 y3 y2 y1 => v2
        shufps xmm6, xmm7, 0x88 // xmm6= z4 z3 z2 z1 => v3

        movaps [edx], xmm1    // store X
        movaps [edx+16], xmm5 // store Y
        movaps [edx+32], xmm6 // store Z
        movaps [edx+48], xmm2 // store W
    }
}

```

Example 6-4 shows a similar data-swizzling algorithm using SIMD instructions in the integer domain.

Example 6-4. Swizzling Data Using UNPCKxxx Instructions

```

void swizzle_asm (Vertex_aos *in, Vertex_soa *out)
{
// in mem: x1y1z1w1-x2y2z2w2-x3y3z3w3-x4y4z4w4-
// SWIZZLE XYZW --> XXXX
asm {
    mov ebx, in                // get structure addresses
    mov edx, out

    movdqa    xmm1, [ebx + 0*16] //w0 z0 y0 x0
    movdqa    xmm2, [ebx + 1*16] //w1 z1 y1 x1
    movdqa    xmm3, [ebx + 2*16] //w2 z2 y2 x2
    movdqa    xmm4, [ebx + 3*16] //w3 z3 y3 x3
    movdqa    xmm5, xmm1
    punpckldq xmm1, xmm2        // y1 y0 x1 x0
    punpckhdq xmm5, xmm2        // w1 w0 z1 z0
    movdqa    xmm2, xmm3
    punpckldq xmm3, xmm4        // y3 y2 x3 x2
    punpckldq xmm2, xmm4        // w3 w2 z3 z2
    movdqa    xmm4, xmm1
    punpcklqdq xmm1, xmm3       // x3 x2 x1 x0
    punpckhqdq xmm4, xmm3       // y3 y2 y1 y0
    movdqa    xmm3, xmm5
    punpcklqdq xmm5, xmm2       // z3 z2 z1 z0
    punpckhqdq xmm3, xmm2       // w3 w2 w1 w0

    movdqa    [edx+0*16], xmm1 //x3 x2 x1 x0
    movdqa    [edx+1*16], xmm4 //y3 y2 y1 y0
    movdqa    [edx+2*16], xmm5 //z3 z2 z1 z0
    movdqa    [edx+3*16], xmm3 //w3 w2 w1 w0
}

```

The technique in Example 6-3 (loading 16 bytes, using SHUFPS and copying halves of XMM registers) is preferable over an alternate approach of loading halves of each vector using MOVLPS/MOVHPS on newer microarchitectures. This is because loading 8 bytes using MOVLPS/MOVHPS can create code dependency and reduce the throughput of the execution engine.

The performance considerations of Example 6-3 and Example 6-4 often depends on the characteristics of each microarchitecture. For example, in Intel Core microarchitecture, executing a SHUFPS tend to be slower than a PUNPCKxxx instruction. In Enhanced Intel Core microarchitecture, SHUFPS and PUNPCKxxx instruction all executes with 1 cycle throughput due to the 128-bit shuffle execution unit. Then the next important consideration is that there is only one port that can execute PUNPCKxxx vs. MOVLHPS/MOVHLPS can execute on multiple ports. The performance of both techniques improves on Intel Core microarchitecture over previous microar-

chitectures due to 3 ports for executing SIMD instructions. Both techniques improves further on Enhanced Intel Core microarchitecture due to the 128-bit shuffle unit.

6.5.1.3 Data Deswizzling

In the deswizzle operation, we want to arrange the SoA format back into AoS format so the XXXX, YYYY, ZZZZ are rearranged and stored in memory as XYZ. Example 6-5 illustrates one deswizzle function for floating-point data:

Example 6-5. Deswizzling Single-Precision SIMD Data

```
void deswizzle_asm(Vertex_soa *in, Vertex_aos *out)
{
  __asm {
    mov     ecx, in                // load structure addresses
    mov     edx, out
    movaps  xmm0, [ecx]           //x3 x2 x1 x0
    movaps  xmm1, [ecx + 16]      //y3 y2 y1 y0
    movaps  xmm2, [ecx + 32]      //z3 z2 z1 z0
    movaps  xmm3, [ecx + 48]      //w3 w2 w1 w0

    movaps  xmm5, xmm0
    movaps  xmm7, xmm2
    unpcklps xmm0, xmm1          // y1 x1 y0 x0
    unpcklps xmm2, xmm3          // w1 z1 w0 z0
    movdqa  xmm4, xmm0
    movlhps xmm0, xmm2           // w0 z0 y0 x0
    movlhps xmm4, xmm2           // w1 z1 y1 x1

    unpckhps xmm5, xmm1          // y3 x3 y2 x2
    unpckhps xmm7, xmm3          // w3 z3 w2 z2
    movdqa  xmm6, xmm5
    movlhps xmm5, xmm7           // w2 z2 y2 x2
    movlhps xmm6, xmm7           // w3 z3 y3 x3
    movaps  [edx+0*16], xmm0      //w0 z0 y0 x0
    movaps  [edx+1*16], xmm4      //w1 z1 y1 x1
    movaps  [edx+2*16], xmm5      //w2 z2 y2 x2
    movaps  [edx+3*16], xmm6      //w3 z3 y3 x3
  }
}
```

Example 6-6 shows a similar deswizzle function using SIMD integer instructions. Both of these techniques demonstrate loading 16 bytes and performing horizontal data movement in registers. This approach is likely to be more efficient than alterna-

tive techniques of storing 8-byte halves of XMM registers using MOVLPS and MOVHPS.

Example 6-6. Deswizzling Data Using SIMD Integer Instructions

```

void deswizzle_rgb(Vertex_soa *in, Vertex_aos *out)
{
  //---deswizzle rgb---
  // assume: xmm1=rrrr, xmm2=gggg, xmm3=bbbb, xmm4=aaaa
  __asm {
    mov     ecx, in                // load structure addresses
    mov     edx, out
    movdqa  xmm0, [ecx]           // load r4 r3 r2 r1 => xmm1
    movdqa  xmm1, [ecx+16]       // load g4 g3 g2 g1 => xmm2
    movdqa  xmm2, [ecx+32]       // load b4 b3 b2 b1 => xmm3
    movdqa  xmm3, [ecx+48]       // load a4 a3 a2 a1 => xmm4
  // Start deswizzling here
    movdqa  xmm5, xmm0
    movdqa  xmm7, xmm2
    punpckldq  xmm0, xmm1        // g2 r2 g1 r1
    punpckldq  xmm2, xmm3        // a2 b2 a1 b1
    movdqa  xmm4, xmm0
    punpcklqdq  xmm0, xmm2       // a1 b1 g1 r1 => v1
    punpckhqdq  xmm4, xmm2       // a2 b2 g2 r2 => v2
    punpckhdq  xmm5, xmm1        // g4 r4 g3 r3
    punpckhdq  xmm7, xmm3        // a4 b4 a3 b3
    movdqa  xmm6, xmm5
    punpcklqdq  xmm5, xmm7       // a3 b3 g3 r3 => v3
    punpckhqdq  xmm6, xmm7       // a4 b4 g4 r4 => v4
  }
  movdqa  [edx], xmm0           // v1
  movdqa  [edx+16], xmm4        // v2
  movdqa  [edx+32], xmm5        // v3
  movdqa  [edx+48], xmm6        // v4
  // DESWIZZLING ENDS HERE
}
}

```

6.5.1.4 Horizontal ADD Using SSE

Although vertical computations generally make use of SIMD performance better than horizontal computations, in some cases, code must use a horizontal operation.

MOVLHPS/MOVLHPS and shuffle can be used to sum data horizontally. For example, starting with four 128-bit registers, to sum up each register horizontally while having the final results in one register, use the MOVLHPS/MOVLHPS to align the upper and lower parts of each register. This allows you to use a vertical add. With the resulting partial horizontal summation, full summation follows easily.

Figure 6-4 presents a horizontal add using MOVLHPS/MOVLHPS. Example 6-7 and Example 6-8 provide the code for this operation.

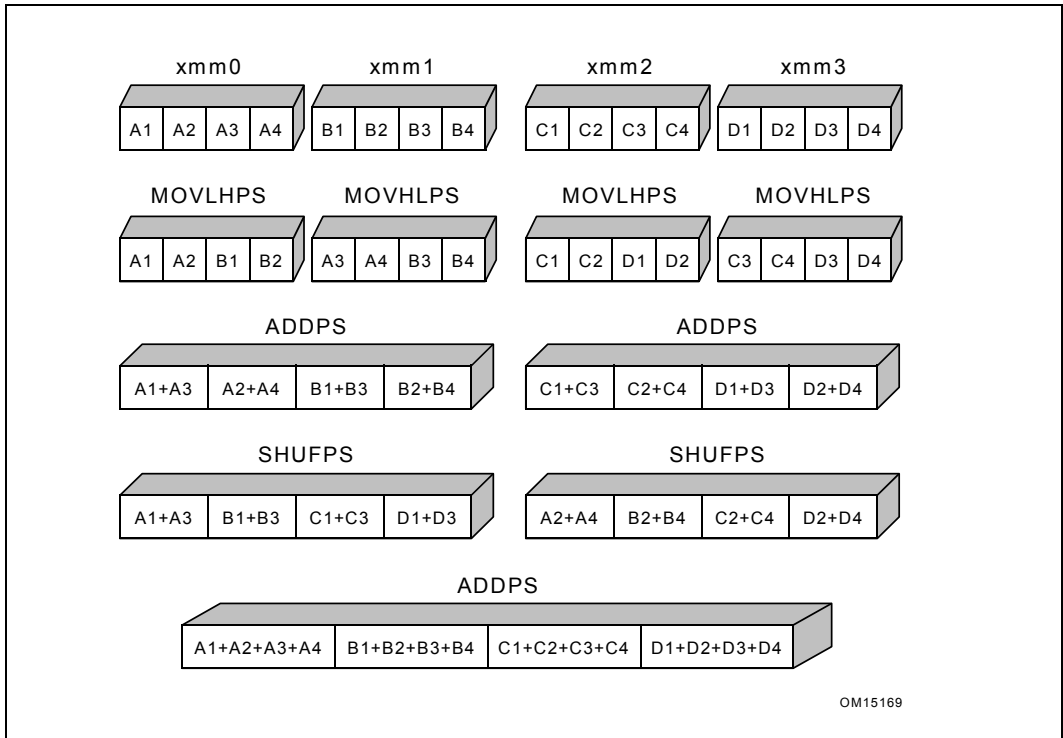


Figure 6-4. Horizontal Add Using MOVLHPS/MOVLHPS

Example 6-7. Horizontal Add Using MOVHLPS/MOVLHPS

```

void horiz_add(Vertex_soa *in, float *out) {
  __asm {
    mov    ecx, in           // load structure addresses
    mov    edx, out
    movaps xmm0, [ecx]      // load A1 A2 A3 A4 => xmm0
    movaps xmm1, [ecx+16]   // load B1 B2 B3 B4 => xmm1
    movaps xmm2, [ecx+32]   // load C1 C2 C3 C4 => xmm2
    movaps xmm3, [ecx+48]   // load D1 D2 D3 D4 => xmm3

// START HORIZONTAL ADD
    movaps xmm5, xmm0       // xmm5= A1,A2,A3,A4
    movlhps xmm5, xmm1      // xmm5= A1,A2,B1,B2
    movlhps xmm1, xmm0      // xmm1= A3,A4,B3,B4
    addps  xmm5, xmm1       // xmm5= A1+A3,A2+A4,B1+B3,B2+B4

    movaps xmm4, xmm2
    movlhps xmm2, xmm3      // xmm2= C1,C2,D1,D2
    movlhps xmm3, xmm4      // xmm3= C3,C4,D3,D4
    addps  xmm3, xmm2       // xmm3= C1+C3,C2+C4,D1+D3,D2+D4
    movaps xmm6, xmm3       // xmm6= C1+C3,C2+C4,D1+D3,D2+D4
    shufps xmm3, xmm5, 0xDD //xmm6=A1+A3,B1+B3,C1+C3,D1+D3

    shufps xmm5, xmm6, 0x88

    addps  xmm6, xmm5       // xmm5= A2+A4,B2+B4,C2+C4,D2+D4
    addps  xmm6, xmm5       // xmm6= D,C,B,A

// END HORIZONTAL ADD
    movaps [edx], xmm6
  }
}

```

Example 6-8. Horizontal Add Using Intrinsics with MOVHLPS/MOVLHPS

```

void horiz_add_intrin(Vertex_soa *in, float *out)
{
  __m128 v, v2, v3, v4;
  __m128 tmm0, tmm1, tmm2, tmm3, tmm4, tmm5, tmm6;
                                     // Temporary variables
  tmm0 = _mm_load_ps(in->x);         // tmm0 = A1 A2 A3 A4
}

```

Example 6-8. Horizontal Add Using Intrinsics with MOVHLPS/MOVLHPS (Contd.)

```

tmm1 = _mm_load_ps(in->y);           // tmm1 = B1 B2 B3 B4
tmm2 = _mm_load_ps(in->z);           // tmm2 = C1 C2 C3 C4
tmm3 = _mm_load_ps(in->w);           // tmm3 = D1 D2 D3 D4
tmm5 = tmm0;                         // tmm0 = A1 A2 A3 A4
tmm5 = _mm_movelh_ps(tmm5, tmm1);    // tmm5 = A1 A2 B1 B2
tmm1 = _mm_movehl_ps(tmm1, tmm0);    // tmm1 = A3 A4 B3 B4
tmm5 = _mm_add_ps(tmm5, tmm1);       // tmm5 = A1+A3 A2+A4 B1+B3 B2+B4
tmm4 = tmm2;

tmm2 = _mm_movelh_ps(tmm2, tmm3);    // tmm2 = C1 C2 D1 D2
tmm3 = _mm_movehl_ps(tmm3, tmm4);    // tmm3 = C3 C4 D3 D4
tmm3 = _mm_add_ps(tmm3, tmm2);       // tmm3 = C1+C3 C2+C4 D1+D3 D2+D4
tmm6 = tmm3;                         // tmm6 = C1+C3 C2+C4 D1+D3 D2+D4
tmm6 = _mm_shuffle_ps(tmm3, tmm5, 0xDD);

// tmm6 = A1+A3 B1+B3 C1+C3 D1+D3
tmm5 = _mm_shuffle_ps(tmm5, tmm6, 0x88);
// tmm5 = A2+A4 B2+B4 C2+C4 D2+D4

tmm6 = _mm_add_ps(tmm6, tmm5);

// tmm6 = A1+A2+A3+A4 B1+B2+B3+B4
// C1+C2+C3+C4 D1+D2+D3+D4

_mm_store_ps(out, tmm6);
}

```

6.5.2 Use of CVTTPS2PI/CVTSS2SI Instructions

The CVTTPS2PI and CVTSS2SI instructions encode the truncate/chop rounding mode implicitly in the instruction. They take precedence over the rounding mode specified in the MXCSR register. This behavior can eliminate the need to change the rounding mode from round-nearest, to truncate/chop, and then back to round-nearest to resume computation.

Avoid frequent changes to the MXCSR register since there is a penalty associated with writing this register. Typically, when using CVTTPS2P/CVTSS2SI, rounding control in MXCSR can always be set to round-nearest.

6.5.3 Flush-to-Zero and Denormals-are-Zero Modes

The flush-to-zero (FTZ) and denormals-are-zero (DAZ) modes are not compatible with the IEEE Standard 754. They are provided to improve performance for applications where underflow is common and where the generation of a denormalized result is not necessary.

See also: Section 3.8.2, "Floating-point Modes and Exceptions."

6.6 SIMD OPTIMIZATIONS AND MICROARCHITECTURES

Pentium M, Intel Core Solo and Intel Core Duo processors have a different microarchitecture than Intel NetBurst microarchitecture. Intel Core microarchitecture offers significantly more efficient SIMD floating-point capability than previous microarchitectures. In addition, instruction latency and throughput of SSE3 instructions are significantly improved in Intel Core microarchitecture over previous microarchitectures.

6.6.1 SIMD Floating-point Programming Using SSE3

SSE3 enhances SSE and SSE2 with nine instructions targeted for SIMD floating-point programming. In contrast to many SSE/SSE2 instructions offering homogeneous arithmetic operations on parallel data elements and favoring the vertical computation model, SSE3 offers instructions that performs asymmetric arithmetic operation and arithmetic operation on horizontal data elements.

ADDSUBPS and ADDSUBPD are two instructions with asymmetric arithmetic processing capability (see Figure 6-5). HADDPS, HADDPD, HSUBPS and HSUBPD offers horizontal arithmetic processing capability (see Figure 6-6). In addition: MOVSLDUP, MOVSHDUP and MOVDDUP load data from memory (or XMM register) and replicate data elements at once.

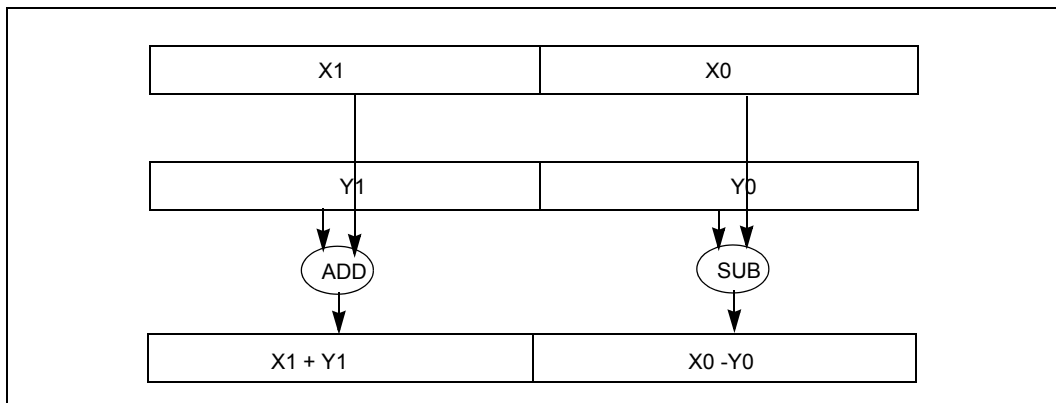


Figure 6-5. Asymmetric Arithmetic Operation of the SSE3 Instruction

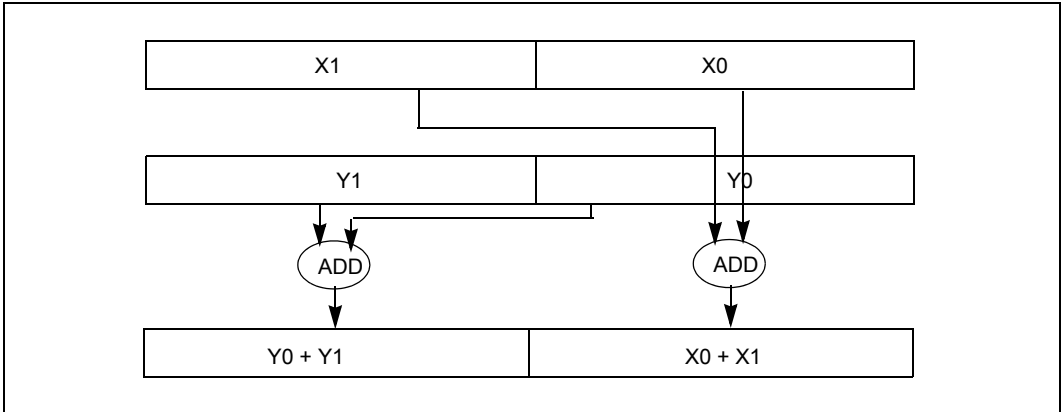


Figure 6-6. Horizontal Arithmetic Operation of the SSE3 Instruction HADDPD

6.6.1.1 SSE3 and Complex Arithmetics

The flexibility of SSE3 in dealing with AOS-type of data structure can be demonstrated by the example of multiplication and division of complex numbers. For example, a complex number can be stored in a structure consisting of its real and imaginary part. This naturally leads to the use of an array of structure. Example 6-9 demonstrates using SSE3 instructions to perform multiplications of single-precision complex numbers. Example 6-10 demonstrates using SSE3 instructions to perform division of complex numbers.

Example 6-9. Multiplication of Two Pair of Single-precision Complex Number

```
// Multiplication of (ak + i bk) * (ck + i dk)
// a + i b can be stored as a data structure
movsldup xmm0, Src1; load real parts into the destination,
                ; a1, a1, a0, a0

movaps xmm1, src2; load the 2nd pair of complex values,
                ; i.e. d1, c1, d0, c0
mulps xmm0, xmm1; temporary results, a1d1, a1c1, a0d0,
                ; a0c0

shufps xmm1, xmm1, b1; reorder the real and imaginary
                ; parts, c1, d1, c0, d0
movshdup xmm2, Src1; load the imaginary parts into the
                ; destination, b1, b1, b0, b0
```

Example 6-9. Multiplication of Two Pair of Single-precision Complex Number

```

mulps  xmm2, xmm1; temporary results, b1c1, b1d1, b0c0,
        ; b0d0
addsubps xmm0, xmm2; b1c1+a1d1, a1c1 -b1d1, b0c0+a0d0,
        ; a0c0-b0d0

```

Example 6-10. Division of Two Pair of Single-precision Complex Numbers

```

// Division of (ak + i bk) / (ck + i dk)
movshdup xmm0, Src1; load imaginary parts into the
        ; destination, b1, b1, b0, b0
movaps  xmm1, src2; load the 2nd pair of complex values,
        ; i.e. d1, c1, d0, c0
mulps  xmm0, xmm1; temporary results, b1d1, b1c1, b0d0,
        ; b0c0

shufps  xmm1, xmm1, b1; reorder the real and imaginary
        ; parts, c1, d1, c0, d0
movsldup xmm2, Src1; load the real parts into the
        ; destination, a1, a1, a0, a0

mulps  xmm2, xmm1; temp results, a1c1, a1d1, a0c0, a0d0
addsubps xmm0, xmm2; a1c1+b1d1, b1c1-a1d1, a0c0+b0d0,
        ; b0c0-a0d0

mulps  xmm1, xmm1 ; c1c1, d1d1, c0c0, d0d0
movps  xmm2, xmm1; c1c1, d1d1, c0c0, d0d0
shufps  xmm2, xmm2, b1; d1d1, c1c1, d0d0, c0c0
addps  xmm2, xmm1; c1c1+d1d1, c1c1+d1d1, c0c0+d0d0,
        ; c0c0+d0d0

divps  xmm0, xmm2
shufps  xmm0, xmm0, b1 ;(b1c1-a1d1)/(c1c1+d1d1),
        ;(a1c1+b1d1)/(c1c1+d1d1),
        ;(b0c0-a0d0)/( c0c0+d0d0),
        ;(a0c0+b0d0)/( c0c0+d0d0)

```

In both examples, the complex numbers are store in arrays of structures. MOVSLDUP, MOVSHDUP and the asymmetric ADDSUBPS allow performing complex arithmetics on two pair of single-precision complex number simultaneously and without any unnecessary swizzling between data elements.

Due to microarchitectural differences, software should implement multiplication of complex double-precision numbers using SSE3 instructions on processors based on

Intel Core microarchitecture. In Intel Core Duo and Intel Core Solo processors, software should use scalar SSE2 instructions to implement double-precision complex multiplication. This is because the data path between SIMD execution units is 128 bits in Intel Core microarchitecture, and only 64 bits in previous microarchitectures. Processors based on the Enhanced Intel Core microarchitecture generally executes SSE3 instruction more efficiently than previous microarchitectures, they also have a 128-bit shuffle unit that will benefit complex arithmetic operations further than Intel Core microarchitecture did.

Example 6-11 shows two equivalent implementations of double-precision complex multiply of two pair of complex numbers using vector SSE2 versus SSE3 instructions. Example 6-12 shows the equivalent scalar SSE2 implementation.

Example 6-11. Double-Precision Complex Multiplication of Two Pairs

SSE2 Vector Implementation	SSE3 Vector Implementation
<pre>movapd xmm0, [eax] ;y x movapd xmm1, [eax+16] ;w z unpcklpd xmm1, xmm1 ;z z movapd xmm2, [eax+16] ;w z unpckhpd xmm2, xmm2 ;w w mulpd xmm1, xmm0 ;z*y z*x mulpd xmm2, xmm0 ;w*y w*x xorpd xmm2, xmm7 ;-w*y +w*x shufpd xmm2, xmm2, 1 ;w*x -w*y addpd xmm2, xmm1 ;z*y+w*x z*x-w*y movapd [ecx], xmm2</pre>	<pre>movapd xmm0, [eax] ;y x movapd xmm1, [eax+16] ;z z movapd xmm2, xmm1 unpcklpd xmm1, xmm1 unpckhpd xmm2, xmm2 mulpd xmm1, xmm0 ;z*y z*x mulpd xmm2, xmm0 ;w*y w*x shufpd xmm2, xmm2, 1 ;w*x w*y addsubpd xmm1, xmm2 ;w*x+z*y z*x-w*y movapd [ecx], xmm1</pre>

Example 6-12. Double-Precision Complex Multiplication Using Scalar SSE2

<pre>movsd xmm0, [eax] ;x movsd xmm5, [eax+8] ;y movsd xmm1, [eax+16] ;z movsd xmm2, [eax+24] ;w movsd xmm3, xmm1 ;z movsd xmm4, xmm2 ;w mulsd xmm1, xmm0 ;z*x mulsd xmm2, xmm0 ;w*x mulsd xmm3, xmm5 ;z*y</pre>

Example 6-12. Double-Precision Complex Multiplication Using Scalar SSE2 (Contd.)

```

mulsd  xmm4, xmm5    ;w*y
subsd  xmm1, xmm4    ;z*x - w*y
addsd  xmm3, xmm2    ;z*y + w*x
movsd  [ecx], xmm1
movsd  [ecx+8], xmm3

```

6.6.1.2 Packed Floating-Point Performance in Intel Core Duo Processor

Most packed SIMD floating-point code will speed up on Intel Core Solo processors relative to Pentium M processors. This is due to improvement in decoding packed SIMD instructions.

The improvement of packed floating-point performance on the Intel Core Solo processor over Pentium M processor depends on several factors. Generally, code that is decoder-bound and/or has a mixture of integer and packed floating-point instructions can expect significant gain. Code that is limited by execution latency and has a “cycles per instructions” ratio greater than one will not benefit from decoder improvement.

When targeting complex arithmetics on Intel Core Solo and Intel Core Duo processors, using single-precision SSE3 instructions can deliver higher performance than alternatives. On the other hand, tasks requiring double-precision complex arithmetics may perform better using scalar SSE2 instructions on Intel Core Solo and Intel Core Duo processors. This is because scalar SSE2 instructions can be dispatched through two ports and executed using two separate floating-point units.

Packed horizontal SSE3 instructions (HADDPS and HSUBPS) can simplify the code sequence for some tasks. However, these instruction consist of more than five micro-ops on Intel Core Solo and Intel Core Duo processors. Care must be taken to ensure the latency and decoding penalty of the horizontal instruction does not offset any algorithmic benefits.

6.6.2 Dot Product and Horizontal SIMD Instructions

Sometimes the AOS type of data organization are more natural in many algebraic formula, one common example is the dot product operation. Dot product operation can be implemented using SSE/SSE2 instruction sets. SSE3 added a few horizontal add/subtract instructions for applications that rely on the horizontal computation model. SSE4.1 provides additional enhancement with instructions that are capable of directly evaluating dot product operations of vectors of 2, 3 or 4 components.

Example 6-13. Dot Product of Vector Length 4 Using SSE/SSE2**Using SSE/SSE2 to compute one dot product**

```

movaps xmm0, [eax] // a4, a3, a2, a1
mulps  xmm0, [eax+16] // a4*b4, a3*b3, a2*b2, a1*b1
movhyps xmm1, xmm0 // X, X, a4*b4, a3*b3, upper half not needed
addps  xmm0, xmm1 // X, X, a2*b2+a4*b4, a1*b1+a3*b3,
pshufd xmm1, xmm0, 1 // X, X, X, a2*b2+a4*b4
addss  xmm0, xmm1 // a1*b1+a3*b3+a2*b2+a4*b4
movss  [ecx], xmm0

```

Example 6-14. Dot Product of Vector Length 4 Using SSE3**Using SSE3 to compute one dot product**

```

movaps xmm0, [eax]
mulps  xmm0, [eax+16] // a4*b4, a3*b3, a2*b2, a1*b1
haddps xmm0, xmm0 // a4*b4+a3*b3, a2*b2+a1*b1, a4*b4+a3*b3, a2*b2+a1*b1
movaps xmm1, xmm0 // a4*b4+a3*b3, a2*b2+a1*b1, a4*b4+a3*b3, a2*b2+a1*b1
psrlq  xmm0, 32 // 0, a4*b4+a3*b3, 0, a4*b4+a3*b3
addss  xmm0, xmm1 // -, -, a1*b1+a3*b3+a2*b2+a4*b4
movss  [eax], xmm0

```

Example 6-15. Dot Product of Vector Length 4 Using SSE4.1**Using SSE4.1 to compute one dot product**

```

movaps xmm0, [eax]
dpps  xmm0, [eax+16], 0xf1 // 0, 0, 0, a1*b1+a3*b3+a2*b2+a4*b4
movss [eax], xmm0

```

Example 6-13, Example 6-14, and Example 6-15 compare the basic code sequence to compute one dot-product result for a pair of vectors.

The selection of an optimal sequence in conjunction with an application's memory access patterns may favor different approaches. For example, if each dot product result is immediately consumed by additional computational sequences, it may be more optimal to compare the relative speed of these different approaches. If dot products can be computed for an array of vectors and kept in the cache for subsequent computations, then more optimal choice may depend on the relative throughput of the sequence of instructions.

In Intel Core microarchitecture, Example 6-14 has higher throughput than Example 6-13. Due to the relatively longer latency of HADDPS, the speed of Example 6-14 is slightly slower than Example 6-13.

In Enhanced Intel Core microarchitecture, Example 6-15 is faster in both speed and throughput than Example 6-13 and Example 6-14. Although the latency of DPPS is also relatively long, it is compensated by the reduction of number of instructions in Example 6-15 to do the same amount of work.

Unrolling can further improve the throughput of each of three dot product implementations. Example 6-16 shows two unrolled versions using the basic SSE2 and SSE3 sequences. The SSE4.1 version can also be unrolled and using INSERTPS to pack 4 dot-product results.

Example 6-16. Unrolled Implementation of Four Dot Products

SSE2 Implementation	SSE3 Implementation
<pre> movaps xmm0, [eax] mulps xmm0, [eax+16] ;w0*w1 z0*z1 y0*y1 x0*x1 movaps xmm2, [eax+32] mulps xmm2, [eax+16+32] ;w2*w3 z2*z3 y2*y3 x2*x3 movaps xmm3, [eax+64] mulps xmm3, [eax+16+64] ;w4*w5 z4*z5 y4*y5 x4*x5 movaps xmm4, [eax+96] mulps xmm4, [eax+16+96] ;w6*w7 z6*z7 y6*y7 x6*x7 </pre>	<pre> movaps xmm0, [eax] mulps xmm0, [eax+16] movaps xmm1, [eax+32] mulps xmm1, [eax+16+32] movaps xmm2, [eax+64] mulps xmm2, [eax+16+64] movaps xmm3, [eax+96] mulps xmm3, [eax+16+96] haddps xmm0, xmm1 haddps xmm2, xmm3 haddps xmm0, xmm2 movaps [ecx], xmm0 </pre>

Example 6-16. Unrolled Implementation of Four Dot Products (Contd.)

SSE2 Implementation	SSE3 Implementation
<pre> movaps xmm1, xmm0 unpcklps xmm0, xmm2 ; y2*y3 y0*y1 x2*x3 x0*x1 unpckhps xmm1, xmm2 ; w2*w3 w0*w1 z2*z3 z0*z1 movaps xmm5, xmm3 unpcklps xmm3, xmm4 ; y6*y7 y4*y5 x6*x7 x4*x5 unpckhps xmm5, xmm4 ; w6*w7 w4*w5 z6*z7 z4*z5 addps xmm0, xmm1 addps xmm5, xmm3 movaps xmm1, xmm5 movhlps xmm1, xmm0 movlhps xmm0, xmm5 addps xmm0, xmm1 movaps [ecx], xmm0 </pre>	

6.6.3 Vector Normalization

Normalizing vectors is a common operation in many floating-point applications. Example 6-17 shows an example in C of normalizing an array of (x, y, z) vectors.

Example 6-17. Normalization of an Array of Vectors

```

for (i=0; i<CNT; i++)
{ float size = nodes[i].vec.dot();
  if (size != 0.0)
  { size = 1.0f/sqrtf(size); }
  else
  { size = 0.0; }
  nodes[i].vec.x *= size;
  nodes[i].vec.y *= size;
  nodes[i].vec.z *= size;
}

```

Example 6-18 shows an assembly sequence that normalizes the x, y, z components of a vector.

Example 6-18. Normalize (x, y, z) Components of an Array of Vectors Using SSE2

```

Vec3 *p = &nodes[i].vec;
__asm
{
    mov     eax, p
    xorps  xmm2, xmm2
    movups  xmm1, [eax] // loads the (x, y, z) of input vector plus x of next vector
    movaps  xmm7, xmm1 // save a copy of data from memory (to restore the unnormalized
value)
    movaps  xmm5, _mask // mask to select (x, y, z) values from an xmm register to normalize
    andps  xmm1, xmm5 // mask 1st 3 elements
    movaps  xmm6, xmm1 // save a copy of (x, y, z) to compute normalized vector later
    mulps  xmm1, xmm1 // 0, z*z, y*y, x*x
    pshufd  xmm3, xmm1, 0x1b // x*x, y*y, z*z, 0
    addps  xmm1, xmm3 // x*x, z*z+y*y, z*z+y*y, x*x
    pshufd  xmm3, xmm1, 0x41 // z*z+y*y, x*x, x*x, z*z+y*y
    addps  xmm1, xmm3 // x*x+y*y+z*z, x*x+y*y+z*z, x*x+y*y+z*z, x*x+y*y+z*z
    comisd  xmm1, xmm2 // compare size to 0
    jz zero
    movaps  xmm3, xmm4 // preloaded unitary vector (1.0, 1.0, 1.0, 1.0)
    sqrtps  xmm1, xmm1
    divps  xmm3, xmm1
    jmp     store
zero:
    movaps  xmm3, xmm2
store:
    mulps  xmm3, xmm6 //normalize the vector in the lower 3 elements
    andnps  xmm5, xmm7 // mask off the lower 3 elements to keep the un-normalized value
    orps   xmm3, xmm5 // order the un-normalized component after the normalized vector
    movaps  [eax ], xmm3 // writes normalized x, y, z; followed by unmodified value

```

Example 6-19 shows an assembly sequence using SSE4.1 to normalize the x, y, z components of a vector.

Example 6-19. Normalize (x, y, z) Components of an Array of Vectors Using SSE4.1

```

Vec3 *p = &nodes[i].vec;
__asm
{
    mov     eax, p
    xorps  xmm2, xmm2
    movups xmm1, [eax] // loads the (x, y, z) of input vector plus x of next vector
    movaps xmm7, xmm1 // save a copy of data from memory
    dpps   xmm1, xmm1, 0x7f // x*x+y*y+z*z, x*x+y*y+z*z, x*x+y*y+z*z, x*x+y*y+z*z
    comisd xmm1, xmm2 // compare size to 0
    jz     zero
    movaps xmm3, xmm4 // preloaded unitary vector (1.0, 1.0, 1.0, 1.0)
    sqrtps xmm1, xmm1
    divps  xmm3, xmm1
    jmp    store
zero:
    movaps xmm3, xmm2
store:
    mulps  xmm3, xmm6 //normalize the vector in the lower 3 elements
    blendps xmm3, xmm7, 0x8 // copy the un-normalized component next to the normalized
vector
    movaps [eax ], xmm3

```

In Example 6-18 and Example 6-19, the throughput of these instruction sequences are basically limited by the long-latency instructions of DIVPS and SQRTPS. In Example 6-19, the use of DPPS replaces eight SSE2 instructions to evaluate and broadcast the dot-product result to four elements of an XMM register. This could result in improvement of the relative speed of Example 6-19 over Example 6-18.

6.6.4 Using Horizontal SIMD Instruction Sets and Data Layout

SSE and SSE2 provide packed add/subtract, multiply/divide instructions that are ideal for situations that can take advantage of vertical computation model, such as SOA data layout. SSE3 and SSE4.1 added horizontal SIMD instructions including horizontal add/subtract, dot-product operations. These more recent SIMD extensions provide tools to solve problems involving data layouts or operations that do not conform to the vertical SIMD computation model.

In this section, we consider a vector-matrix multiplication problem and discuss the relevant factors for choosing various horizontal SIMD instructions.

Example 6-20 shows the vector-matrix data layout in AOS, where the input and out vectors are stored as an array of structure.

Example 6-20. Data Organization in Memory for AOS Vector-Matrix Multiplication

```

Matrix M4x4 (pMat):  M00 M01 M02 M03
                    M10 M11 M12 M13
                    M20 M21 M22 M23
                    M30 M31 M32 M33
4 input vertices V4x1 (pVert):  V0x V0y V0z V0w
                                V1x V1y V1z V1w
                                V2x V2y V2z V2w
                                V3x V3y V3z V3w
Output vertices O4x1 (pOutVert): O0x O0y O0z O0w
                                O1x O1y O1z O1w
                                O2x O2y O2z O2w
                                O3x O3y O3z O3w

```

Example 6-21 shows an example using HADDPS and MULPS to perform vector-matrix multiplication with data layout in AOS. After three HADDPS completing the summations of each output vector component, the output components are arranged in AOS.

Example 6-21. AOS Vector-Matrix Multiplication with HADDPS

```

mov     eax, pMat
mov     ebx, pVert
mov     ecx, pOutVert
xor     edx, edx
movaps  xmm5,[eax+16] // load row M1?
movaps  xmm6,[eax+2*16] // load row M2?
movaps  xmm7,[eax+3*16] // load row M3?
loop:
movaps  xmm4, [ebx + edx] // load input vector
movaps  xmm0, xmm4
mulps   xmm0, [eax] // m03*vw, m02*vz, m01*vy, m00*vx,
movaps  xmm1, xmm4
mulps   xmm1, xmm5 // m13*vw, m12*vz, m11*vy, m10*vx,

```

Example 6-21. AOS Vector-Matrix Multiplication with HADDPS

```

movaps  xmm2, xmm4
mulps   xmm2, xmm6 // m23*vw, m22*vz, m21*vy, m20*vx
movaps  xmm3, xmm4
mulps   xmm3, xmm7 // m33*vw, m32*vz, m31*vy, m30*vx,
haddps  xmm0, xmm1
haddps  xmm2, xmm3
haddps  xmm0, xmm2
movaps  [ecx + edx], xmm0 // store a vector of length 4
add     edx, 16
cmp     edx, top
jb     lloop

```

Example 6-22 shows an example using DPPS to perform vector-matrix multiplication in AOS.

Example 6-22. AOS Vector-Matrix Multiplication with DPPS

```

mov     eax, pMat
mov     ebx, pVert
mov     ecx, pOutVert
xor     edx, edx
movaps  xmm5, [eax+16] // load row M1?
movaps  xmm6, [eax+2*16] // load row M2?
movaps  xmm7, [eax+3*16] // load row M3?
lloop:
movaps  xmm4, [ebx + edx] // load input vector
movaps  xmm0, xmm4
dpps   xmm0, [eax], 0xf1 // calculate dot product of length 4, store to lowest dword
movaps  xmm1, xmm4
dpps   xmm1, xmm5, 0xf1
movaps  xmm2, xmm4
dpps   xmm2, xmm6, 0xf1
movaps  xmm3, xmm4
dpps   xmm3, xmm7, 0xf1
movss  [ecx + edx + 0*4], xmm0 // store one element of vector length 4
movss  [ecx + edx + 1*4], xmm1
movss  [ecx + edx + 2*4], xmm2
movss  [ecx + edx + 3*4], xmm3
add     edx, 16
cmp     edx, top
jb     lloop

```

Example 6-21 and Example 6-22 both work with AOS data layout using different horizontal processing techniques provided by SSE3 and SSE4.1. The effectiveness of either techniques will vary, depending on the degree of exposures of long-latency instruction in the inner loop, the overhead/efficiency of data movement, and the latency of HADDPS vs. DPPS.

On processors that support both HADDPS and DPPS, the choice between either technique may depend on application-specific considerations. If the output vectors are written back to memory directly in a batch situation, Example 6-21 may be preferable over Example 6-22, because the latency of DPPS is long and storing each output vector component individually is less than ideal for storing an array of vectors.

There may be partially-vectorizable situations that the individual output vector component is consumed immediately by other non-vectorizable computations. Then, using DPPS producing individual component may be more suitable than dispersing the packed output vector produced by three HADDPS as in Example 6-21.

6.6.4.1 SOA and Vector Matrix Multiplication

If the native data layout of a problem conforms to SOA, then vector-matrix multiply can be coded using MULPS, ADDPS without using the longer-latency horizontal arithmetic instructions, or packing scalar components into packed format (Example 6-22). To achieve higher throughput with SOA data layout, there are either pre-requisite data preparation or swizzling/deswizzling on-the-fly that must be comprehended. For example, an SOA data layout for vector-matrix multiplication is shown in Example 6-23. Each matrix element is replicated four times to minimize data movement overhead for producing packed results.

Example 6-23. Data Organization in Memory for SOA Vector-Matrix Multiplication

```

Matrix M16x4 (pMat):
    M00 M00 M00 M00 M01 M01 M01 M01 M02 M02 M02 M02 M03 M03 M03 M03
    M10 M10 M10 M10 M11 M11 M11 M11 M12 M12 M12 M12 M13 M13 M13 M13
    M20 M20 M20 M20 M21 M21 M21 M21 M22 M22 M22 M22 M23 M23 M23 M23
    M30 M30 M30 M30 M31 M31 M31 M31 M32 M32 M32 M32 M33 M33 M33 M33
4 input vertices V4x1 (pVert):  V0x V1x V2x V3x
                                V0y V1y V2y V3y
                                V0z V1z V2z V3z
                                V0w V1w V2w V3w
Output vertices O4x1 (pOutVert): 00x 01x 02x 03x
                                00y 01y 02y 03y
                                00z 01z 02z 03z
                                00w 01w 02w 03w
    
```

The corresponding vector-matrix multiply example in SOA (unrolled for four iteration of vectors) is shown in Example 6-24.

Example 6-24. Vector-Matrix Multiplication with Native SOA Data Layout

```

mov     ebx, pVert
mov     ecx, pOutVert
xor     edx, edx
movaps  xmm5,[eax+16] // load row M1?
movaps  xmm6,[eax+2*16] // load row M2?
movaps  xmm7,[eax+3*16] // load row M3?
loop_vert:
mov     eax, pMat
xor     edi, edi
movaps  xmm0,[ebx ] // load V3x, V2x, V1x, V0x
movaps  xmm1,[ebx ] // load V3y, V2y, V1y, V0y
movaps  xmm2,[ebx ] // load V3z, V2z, V1z, V0z
movaps  xmm3,[ebx ] // load V3w, V2w, V1w, V0w
loop_mat:
movaps  xmm4,[eax] // m00, m00, m00, m00,
mulps   xmm4, xmm0 // m00*V3x, m00*V2x, m00*V1x, m00*V0x,
movaps  xmm4,[eax + 16] // m01, m01, m01, m01,
mulps   xmm5, xmm1 // m01*V3y, m01*V2y, m01*V1y, m01*V0y,
addps   xmm4, xmm5
movaps  xmm5,[eax + 32] // m02, m02, m02, m02,
mulps   xmm5, xmm2 // m02*V3z, m02*V2z, m02*V1z, m02*V0z,
addps   xmm4, xmm5
movaps  xmm5,[eax + 48] // m03, m03, m03, m03,
mulps   xmm5, xmm3 // m03*V3w, m03*V2w, m03*V1w, m03*V0w,
addps   xmm4, xmm5
movaps  [ecx + edx], xmm4
add     eax, 64
add     edx, 16
add     edi, 1
cmp     edi, 4
jb     loop_mat
add     ebx, 64
cmp     edx, top
jb     loop_vert

```


CHAPTER 7

OPTIMIZING CACHE USAGE

Over the past decade, processor speed has increased. Memory access speed has increased at a slower pace. The resulting disparity has made it important to tune applications in one of two ways: either (a) a majority of data accesses are fulfilled from processor caches, or (b) effectively masking memory latency to utilize peak memory bandwidth as much as possible.

Hardware prefetching mechanisms are enhancements in microarchitecture to facilitate the latter aspect, and will be most effective when combined with software tuning. The performance of most applications can be considerably improved if the data required can be fetched from the processor caches or if memory traffic can take advantage of hardware prefetching effectively.

Standard techniques to bring data into the processor before it is needed involve additional programming which can be difficult to implement and may require special steps to prevent performance degradation. Streaming SIMD Extensions addressed this issue by providing various prefetch instructions.

Streaming SIMD Extensions introduced the various non-temporal store instructions. SSE2 extends this support to new data types and also introduce non-temporal store support for the 32-bit integer registers.

This chapter focuses on:

- Hardware Prefetch Mechanism, Software Prefetch and Cacheability Instructions — Discusses microarchitectural feature and instructions that allow you to affect data caching in an application.
- Memory Optimization Using Hardware Prefetching, Software Prefetch and Cacheability Instructions — Discusses techniques for implementing memory optimizations using the above instructions.

NOTE

In a number of cases presented, the prefetching and cache utilization described are specific to current implementations of Intel NetBurst microarchitecture but are largely applicable for the future processors.

- Using deterministic cache parameters to manage cache hierarchy.

7.1 GENERAL PREFETCH CODING GUIDELINES

The following guidelines will help you to reduce memory traffic and utilize peak memory system bandwidth more effectively when large amounts of data movement must originate from the memory system:

OPTIMIZING CACHE USAGE

- Take advantage of the hardware prefetcher's ability to prefetch data that are accessed in linear patterns, in either a forward or backward direction.
- Take advantage of the hardware prefetcher's ability to prefetch data that are accessed in a regular pattern with access strides that are substantially smaller than half of the trigger distance of the hardware prefetch (see Table 2-23).
- Use a current-generation compiler, such as the Intel C++ Compiler that supports C++ language-level features for Streaming SIMD Extensions. Streaming SIMD Extensions and MMX technology instructions provide intrinsics that allow you to optimize cache utilization. Examples of Intel compiler intrinsics include: `_mm_prefetch`, `_mm_stream` and `_mm_load`, `_mm_sfence`. For details, refer to Intel C++ Compiler User's Guide documentation.
- Facilitate compiler optimization by:
 - Minimize use of global variables and pointers.
 - Minimize use of complex control flow.
 - Use the `const` modifier, avoid register modifier.
 - Choose data types carefully (see below) and avoid type casting.
- Use cache blocking techniques (for example, strip mining) as follows:
 - Improve cache hit rate by using cache blocking techniques such as strip-mining (one dimensional arrays) or loop blocking (two dimensional arrays)
 - Explore using hardware prefetching mechanism if your data access pattern has sufficient regularity to allow alternate sequencing of data accesses (for example: tiling) for improved spatial locality. Otherwise use `PREFETCHNTA`.
- Balance single-pass versus multi-pass execution:
 - Single-pass, or unlayered execution passes a single data element through an entire computation pipeline.
 - Multi-pass, or layered execution performs a single stage of the pipeline on a batch of data elements before passing the entire batch on to the next stage.
 - If your algorithm is single-pass use `PREFETCHNTA`. If your algorithm is multi-pass use `PREFETCHT0`.
- Resolve memory bank conflict issues. Minimize memory bank conflicts by applying array grouping to group contiguously used data together or by allocating data within 4-KByte memory pages.
- Resolve cache management issues. Minimize the disturbance of temporal data held within processor's caches by using streaming store instructions.
- Optimize software prefetch scheduling distance:
 - Far ahead enough to allow interim computations to overlap memory access time.
 - Near enough that prefetched data is not replaced from the data cache.

- Use software prefetch concatenation. Arrange prefetches to avoid unnecessary prefetches at the end of an inner loop and to prefetch the first few iterations of the inner loop inside the next outer loop.
- Minimize the number of software prefetches. Prefetch instructions are not completely free in terms of bus cycles, machine cycles and resources; excessive usage of prefetches can adversely impact application performance.
- Interleave prefetches with computation instructions. For best performance, software prefetch instructions must be interspersed with computational instructions in the instruction sequence (rather than clustered together).

7.2 HARDWARE PREFETCHING OF DATA

Pentium M, Intel Core Solo, and Intel Core Duo processors and processors based on Intel Core microarchitecture and Intel NetBurst microarchitecture provide hardware data prefetch mechanisms which monitor application data access patterns and prefetches data automatically. This behavior is automatic and does not require programmer intervention.

For processors based on Intel NetBurst microarchitecture, characteristics of the hardware data prefetcher are:

1. It requires two successive cache misses in the last level cache to trigger the mechanism; these two cache misses must satisfy the condition that strides of the cache misses are less than the trigger distance of the hardware prefetch mechanism (see Table 2-23).
2. Attempts to stay 256 bytes ahead of current data access locations.
3. Follows only one stream per 4-KByte page (load or store).
4. Can prefetch up to 8 simultaneous, independent streams from eight different 4-KByte regions
5. Does not prefetch across 4-KByte boundary. This is independent of paging modes.
6. Fetches data into second/third-level cache.
7. Does not prefetch UC or WC memory types.
8. Follows load and store streams. Issues Read For Ownership (RFO) transactions for store streams and Data Reads for load streams.

Other than items 2 and 4 discussed above, most other characteristics also apply to Pentium M, Intel Core Solo and Intel Core Duo processors. The hardware prefetcher implemented in the Pentium M processor fetches data to a second level cache. It can track 12 independent streams in the forward direction and 4 independent streams in the backward direction. The hardware prefetcher of Intel Core Solo processor can track 16 forward streams and 4 backward streams. On the Intel Core Duo processor, the hardware prefetcher in each core fetches data independently.

Hardware prefetch mechanisms of processors based on Intel Core microarchitecture are discussed in Section 3.7.3 and Section 3.7.4. Despite differences in hardware implementation technique, the overall benefit of hardware prefetching to software are similar between Intel Core microarchitecture and prior microarchitectures.

7.3 PREFETCH AND CACHEABILITY INSTRUCTIONS

The PREFETCH instruction, inserted by the programmers or compilers, accesses a minimum of two cache lines of data on the Pentium 4 processor prior to the data actually being needed (one cache line of data on the Pentium M processor). This hides the latency for data access in the time required to process data already resident in the cache.

Many algorithms can provide information in advance about the data that is to be required. In cases where memory accesses are in long, regular data patterns; the automatic hardware prefetcher should be favored over software prefetches.

The cacheability control instructions allow you to control data caching strategy in order to increase cache efficiency and minimize cache pollution.

Data reference patterns can be classified as follows:

- **Temporal** — Data will be used again soon
- **Spatial** — Data will be used in adjacent locations (for example, on the same cache line).
- **Non-temporal** — Data which is referenced once and not reused in the immediate future (for example, for some multimedia data types, as the vertex buffer in a 3D graphics application).

These data characteristics are used in the discussions that follow.

7.4 PREFETCH

This section discusses the mechanics of the software PREFETCH instructions. In general, software prefetch instructions should be used to supplement the practice of tuning an access pattern to suit the automatic hardware prefetch mechanism.

7.4.1 Software Data Prefetch

The PREFETCH instruction can hide the latency of data access in performance-critical sections of application code by allowing data to be fetched in advance of actual usage. PREFETCH instructions do not change the user-visible semantics of a program, although they may impact program performance. PREFETCH merely provides a hint to the hardware and generally does not generate exceptions or faults.

PREFETCH loads either non-temporal data or temporal data in the specified cache level. This data access type and the cache level are specified as a hint. Depending on the implementation, the instruction fetches 32 or more aligned bytes (including the specified address byte) into the instruction-specified cache levels.

PREFETCH is implementation-specific; applications need to be tuned to each implementation to maximize performance.

NOTE

Using the PREFETCH instruction is recommended only if data does not fit in cache.

PREFETCH provides a hint to the hardware; it does not generate exceptions or faults except for a few special cases (see Section 7.4.3, "Prefetch and Load Instructions"). However, excessive use of PREFETCH instructions may waste memory bandwidth and result in a performance penalty due to resource constraints.

Nevertheless, PREFETCH can lessen the overhead of memory transactions by preventing cache pollution and by using caches and memory efficiently. This is particularly important for applications that share critical system resources, such as the memory bus. See an example in Section 7.7.2.1, "Video Encoder."

PREFETCH is mainly designed to improve application performance by hiding memory latency in the background. If segments of an application access data in a predictable manner (for example, using arrays with known strides), they are good candidates for using PREFETCH to improve performance.

Use the PREFETCH instructions in:

- Predictable memory access patterns
- Time-consuming innermost loops
- Locations where the execution pipeline may stall if data is not available

7.4.2 Prefetch Instructions - Pentium® 4 Processor Implementation

Streaming SIMD Extensions include four PREFETCH instructions variants, one non-temporal and three temporal. They correspond to two types of operations, temporal and non-temporal.

NOTE

At the time of PREFETCH, if data is already found in a cache level that is closer to the processor than the cache level specified by the instruction, no data movement occurs.

The non-temporal instruction is:

- **PREFETCHNTA**— Fetch the data into the second-level cache, minimizing cache pollution.

Temporal instructions are:

- **PREFETCHNT0** — Fetch the data into all cache levels; that is, to the second-level cache for the Pentium 4 processor.
- **PREFETCHNT1** — This instruction is identical to PREFETCHT0.
- **PREFETCHNT2** — This instruction is identical to PREFETCHT0.

7.4.3 Prefetch and Load Instructions

The Pentium 4 processor has a decoupled execution and memory architecture that allows instructions to be executed independently with memory accesses (if there are no data and resource dependencies). Programs or compilers can use dummy load instructions to imitate PREFETCH functionality; but preloading is not completely equivalent to using PREFETCH instructions. PREFETCH provides greater performance than preloading.

Currently, PREFETCH provides greater performance than preloading because:

- Has no destination register, it only updates cache lines.
- Does not stall the normal instruction retirement.
- Does not affect the functional behavior of the program.
- Has no cache line split accesses.
- Does not cause exceptions except when the LOCK prefix is used. The LOCK prefix is not a valid prefix for use with PREFETCH.
- Does not complete its own execution if that would cause a fault.

Currently, the advantage of PREFETCH over preloading instructions are processor-specific. This may change in the future.

There are cases where a PREFETCH will not perform the data prefetch. These include:

- PREFETCH causes a DTLB (Data Translation Lookaside Buffer) miss. This applies to Pentium 4 processors with CPUID signature corresponding to family 15, model 0, 1, or 2. PREFETCH resolves DTLB misses and fetches data on Pentium 4 processors with CPUID signature corresponding to family 15, model 3.
- An access to the specified address that causes a fault/exception.
- If the memory subsystem runs out of request buffers between the first-level cache and the second-level cache.
- PREFETCH targets an uncacheable memory region (for example, USWC and UC).
- The LOCK prefix is used. This causes an invalid opcode exception.

7.5 CACHEABILITY CONTROL

This section covers the mechanics of cacheability control instructions.

7.5.1 The Non-temporal Store Instructions

This section describes the behavior of streaming stores and reiterates some of the information presented in the previous section.

In Streaming SIMD Extensions, the MOVNTPS, MOVNTPD, MOVNTQ, MOVNTDQ, MOVNTI, MASKMOVQ and MASKMOVDQU instructions are streaming, non-temporal stores. With regard to memory characteristics and ordering, they are similar to the Write-Combining (WC) memory type:

- **Write combining** — Successive writes to the same cache line are combined.
- **Write collapsing** — Successive writes to the same byte(s) result in only the last write being visible.
- **Weakly ordered** — No ordering is preserved between WC stores or between WC stores and other loads or stores.
- **Uncacheable and not write-allocating** — Stored data is written around the cache and will not generate a read-for-ownership bus request for the corresponding cache line.

7.5.1.1 Fencing

Because streaming stores are weakly ordered, a fencing operation is required to ensure that the stored data is flushed from the processor to memory. Failure to use an appropriate fence may result in data being “trapped” within the processor and will prevent visibility of this data by other processors or system agents.

WC stores require software to ensure coherence of data by performing the fencing operation. See Section 7.5.5, “FENCE Instructions.”

7.5.1.2 Streaming Non-temporal Stores

Streaming stores can improve performance by:

- Increasing store bandwidth if the 64 bytes that fit within a cache line are written consecutively (since they do not require read-for-ownership bus requests and 64 bytes are combined into a single bus write transaction).
- Reducing disturbance of frequently used cached (temporal) data (since they write around the processor caches).

Streaming stores allow cross-aliasing of memory types for a given memory region. For instance, a region may be mapped as write-back (WB) using page attribute tables (PAT) or memory type range registers (MTRRs) and yet is written using a streaming store.

7.5.1.3 Memory Type and Non-temporal Stores

Memory type can take precedence over a non-temporal hint, leading to the following considerations:

- If the programmer specifies a non-temporal store to strongly-ordered uncacheable memory (for example, Uncacheable (UC) or Write-Protect (WP) memory types), then the store behaves like an uncacheable store. The non-temporal hint is ignored and the memory type for the region is retained.
- If the programmer specifies the weakly-ordered uncacheable memory type of Write-Combining (WC), then the non-temporal store and the region have the same semantics and there is no conflict.
- If the programmer specifies a non-temporal store to cacheable memory (for example, Write-Back (WB) or Write-Through (WT) memory types), two cases may result:
 - **CASE 1** — If the data is present in the cache hierarchy, the instruction will ensure consistency. A particular processor may choose different ways to implement this. The following approaches are probable: (a) updating data in-place in the cache hierarchy while preserving the memory type semantics assigned to that region or (b) evicting the data from the caches and writing the new non-temporal data to memory (with WC semantics).

The approaches (separate or combined) can be different for future processors. Pentium 4, Intel Core Solo and Intel Core Duo processors implement the latter policy (of evicting data from all processor caches). The Pentium M processor implements a combination of both approaches.

If the streaming store hits a line that is present in the first-level cache, the store data is combined in place within the first-level cache. If the streaming store hits a line present in the second-level, the line and stored data is flushed from the second-level to system memory.

- **CASE 2** — If the data is not present in the cache hierarchy and the destination region is mapped as WB or WT; the transaction will be weakly ordered and is subject to all WC memory semantics. This non-temporal store will not write-allocate. Different implementations may choose to collapse and combine such stores.

7.5.1.4 Write-Combining

Generally, WC semantics require software to ensure coherence with respect to other processors and other system agents (such as graphics cards). Appropriate use of synchronization and a fencing operation must be performed for producer-consumer usage models (see Section 7.5.5, "FENCE Instructions"). Fencing ensures that all system agents have global visibility of the stored data. For instance, failure to fence may result in a written cache line staying within a processor, and the line would not be visible to other agents.

For processors which implement non-temporal stores by updating data in-place that already resides in the cache hierarchy, the destination region should also be mapped as WC. Otherwise, if mapped as WB or WT, there is a potential for speculative processor reads to bring the data into the caches. In such a case, non-temporal stores would then update in place and data would not be flushed from the processor by a subsequent fencing operation.

The memory type visible on the bus in the presence of memory type aliasing is implementation-specific. As one example, the memory type written to the bus may reflect the memory type for the first store to the line, as seen in program order. Other alternatives are possible. This behavior should be considered reserved and dependence on the behavior of any particular implementation risks future incompatibility.

7.5.2 Streaming Store Usage Models

The two primary usage domains for streaming store are coherent requests and non-coherent requests.

7.5.2.1 Coherent Requests

Coherent requests are normal loads and stores to system memory, which may also hit cache lines present in another processor in a multiprocessor environment. With coherent requests, a streaming store can be used in the same way as a regular store that has been mapped with a WC memory type (PAT or MTRR). An SFENCE instruction must be used within a producer-consumer usage model in order to ensure coherency and visibility of data between processors.

Within a single-processor system, the CPU can also re-read the same memory location and be assured of coherence (that is, a single, consistent view of this memory location). The same is true for a multiprocessor (MP) system, assuming an accepted MP software producer-consumer synchronization policy is employed.

7.5.2.2 Non-coherent requests

Non-coherent requests arise from an I/O device, such as an AGP graphics card, that reads or writes system memory using non-coherent requests, which are not reflected on the processor bus and thus will not query the processor's caches. An SFENCE instruction must be used within a producer-consumer usage model in order to ensure coherency and visibility of data between processors. In this case, if the processor is writing data to the I/O device, a streaming store can be used with a processor with any behavior of Case 1 (Section 7.5.1.3) only if the region has also been mapped with a WC memory type (PAT, MTRR).

NOTE

Failure to map the region as WC may allow the line to be speculatively read into the processor caches (via the wrong path of a mispredicted branch).

In case the region is not mapped as WC, the streaming might update in-place in the cache and a subsequent SFENCE would not result in the data being written to system memory. Explicitly mapping the region as WC in this case ensures that any data read from this region will not be placed in the processor's caches. A read of this memory location by a non-coherent I/O device would return incorrect/out-of-date results.

For a processor which solely implements Case 2 (Section 7.5.1.3), a streaming store can be used in this non-coherent domain without requiring the memory region to also be mapped as WB, since any cached data will be flushed to memory by the streaming store.

7.5.3 Streaming Store Instruction Descriptions

MOVNTQ/MOVNTDQ (non-temporal store of packed integer in an MMX technology or Streaming SIMD Extensions register) store data from a register to memory. They are implicitly weakly-ordered, do no write-allocate, and so minimize cache pollution.

MOVNTPS (non-temporal store of packed single precision floating point) is similar to MOVNTQ. It stores data from a Streaming SIMD Extensions register to memory in 16-byte granularity. Unlike MOVNTQ, the memory address must be aligned to a 16-byte boundary or a general protection exception will occur. The instruction is implicitly weakly-ordered, does not write-allocate, and thus minimizes cache pollution.

MASKMOVQ/MASKMOVDQU (non-temporal byte mask store of packed integer in an MMX technology or Streaming SIMD Extensions register) store data from a register to the location specified by the EDI register. The most significant bit in each byte of the second mask register is used to selectively write the data of the first register on a per-byte basis. The instructions are implicitly weakly-ordered (that is, successive stores may not write memory in original program-order), do not write-allocate, and thus minimize cache pollution.

7.5.4 The Streaming Load Instruction

SSE4.1 introduces the MOVNTDQA instruction. MOVNTDQA loads 16 bytes from memory using a non-temporal hint if the memory source is WC type. For WC memory type, the non-temporal hint may be implemented by loading into a temporary internal buffer with the equivalent of an aligned cache line without filling this data to the cache. Subsequent MOVNTDQA reads to unread portions of the buffered WC data will cause 16 bytes of data transferred from the temporary internal buffer to an XMM register if data is available.

If used appropriately, MOVNTDQA can help software achieve significantly higher throughput when loading data in WC memory region into the processor than other means.

Chapter 1 provides a reference to an application note on using MOVNTDQA. Additional information and requirements to use MOVNTDQA appropriately can be found in Chapter 12, "Programming with SSE3, SSSE3 and SSE4" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1*, and the instruction reference pages of MOVNTDQA in *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A*.

7.5.5 FENCE Instructions

The following fence instructions are available: SFENCE, IFENCE, and MFENCE.

7.5.5.1 SFENCE Instruction

The SFENCE (STORE FENCE) instruction makes it possible for every STORE instruction that precedes an SFENCE in program order to be globally visible before any STORE that follows the SFENCE. SFENCE provides an efficient way of ensuring ordering between routines that produce weakly-ordered results.

The use of weakly-ordered memory types can be important under certain data sharing relationships (such as a producer-consumer relationship). Using weakly-ordered memory can make assembling the data more efficient, but care must be taken to ensure that the consumer obtains the data that the producer intended to see.

Some common usage models may be affected by weakly-ordered stores. Examples are:

- Library functions, which use weakly-ordered memory to write results
- Compiler-generated code, which also benefits from writing weakly-ordered results
- Hand-crafted code

The degree to which a consumer of data knows that the data is weakly-ordered can vary for different cases. As a result, SFENCE should be used to ensure ordering between routines that produce weakly-ordered data and routines that consume this data.

7.5.5.2 LFENCE Instruction

The LFENCE (LOAD FENCE) instruction makes it possible for every LOAD instruction that precedes the LFENCE instruction in program order to be globally visible before any LOAD instruction that follows the LFENCE.

The LFENCE instruction provides a means of segregating LOAD instructions from other LOADs.

7.5.5.3 MFENCE Instruction

The MFENCE (MEMORY FENCE) instruction makes it possible for every LOAD/STORE instruction preceding MFENCE in program order to be globally visible before any LOAD/STORE following MFENCE. MFENCE provides a means of segregating certain memory instructions from other memory references.

The use of a LFENCE and SFENCE is not equivalent to the use of a MFENCE since the load and store fences are not ordered with respect to each other. In other words, the load fence can be executed before prior stores and the store fence can be executed before prior loads.

MFENCE should be used whenever the cache line flush instruction (CLFLUSH) is used to ensure that speculative memory references generated by the processor do not interfere with the flush. See Section 7.5.6, "CLFLUSH Instruction."

7.5.6 CLFLUSH Instruction

The CLFLUSH instruction invalidates the cache line associated with the linear address that contain the byte address of the memory location, in all levels of the processor cache hierarchy (data and instruction). This invalidation is broadcast throughout the coherence domain. If, at any level of the cache hierarchy, a line is inconsistent with memory (dirty), it is written to memory before invalidation. Other characteristics include:

- The data size affected is the cache coherency size, which is 64 bytes on Pentium 4 processor.
- The memory attribute of the page containing the affected line has no effect on the behavior of this instruction.
- The CLFLUSH instruction can be used at all privilege levels and is subject to all permission checking and faults associated with a byte load.

CLFLUSH is an unordered operation with respect to other memory traffic, including other CLFLUSH instructions. Software should use a memory fence for cases where ordering is a concern.

As an example, consider a video usage model where a video capture device is using non-coherent AGP accesses to write a capture stream directly to system memory. Since these non-coherent writes are not broadcast on the processor bus, they will not flush copies of the same locations that reside in the processor caches. As a result, before the processor re-reads the capture buffer, it should use CLFLUSH to ensure that stale copies of the capture buffer are flushed from the processor caches. Due to speculative reads that may be generated by the processor, it is important to observe appropriate fencing (using MFENCE).

Example 7-1 provides pseudo-code for CLFLUSH usage.

Example 7-1. Pseudo-code Using CLFLUSH

```
while (!buffer_ready) {}
mfence
  for(i=0;i<num_cachelines;i+=cacheline_size) {
    clflush (char *)((unsigned int)buffer + i)
  }
mfence
  prefnta buffer[0];
  VAR = buffer[0];
```

7.6 MEMORY OPTIMIZATION USING PREFETCH

The Pentium 4 processor has two mechanisms for data prefetch: software-controlled prefetch and an automatic hardware prefetch.

7.6.1 Software-Controlled Prefetch

The software-controlled prefetch is enabled using the four PREFETCH instructions introduced with Streaming SIMD Extensions instructions. These instructions are hints to bring a cache line of data in to various levels and modes in the cache hierarchy. The software-controlled prefetch is not intended for prefetching code. Using it can incur significant penalties on a multiprocessor system when code is shared.

Software prefetching has the following characteristics:

- Can handle irregular access patterns which do not trigger the hardware prefetcher.
- Can use less bus bandwidth than hardware prefetching; see below.
- Software prefetches must be added to new code, and do not benefit existing applications.

7.6.2 Hardware Prefetch

Automatic hardware prefetch can bring cache lines into the unified last-level cache based on prior data misses. It will attempt to prefetch two cache lines ahead of the prefetch stream. Characteristics of the hardware prefetcher are:

- It requires some regularity in the data access patterns.
 - If a data access pattern has constant stride, hardware prefetching is effective if the access stride is less than half of the trigger distance of hardware prefetcher (see Table 2-23).

- If the access stride is not constant, the automatic hardware prefetcher can mask memory latency if the strides of two successive cache misses are less than the trigger threshold distance (small-stride memory traffic).
- The automatic hardware prefetcher is most effective if the strides of two successive cache misses remain less than the trigger threshold distance and close to 64 bytes.
- There is a start-up penalty before the prefetcher triggers and there may be fetches an array finishes. For short arrays, overhead can reduce effectiveness.
 - The hardware prefetcher requires a couple misses before it starts operating.
 - Hardware prefetching generates a request for data beyond the end of an array, which is not be utilized. This behavior wastes bus bandwidth. In addition this behavior results in a start-up penalty when fetching the beginning of the next array. Software prefetching may recognize and handle these cases.
- It will not prefetch across a 4-KByte page boundary. A program has to initiate demand loads for the new page before the hardware prefetcher starts prefetching from the new page.
- The hardware prefetcher may consume extra system bandwidth if the application's memory traffic has significant portions with strides of cache misses greater than the trigger distance threshold of hardware prefetch (large-stride memory traffic).
- The effectiveness with existing applications depends on the proportions of small-stride versus large-stride accesses in the application's memory traffic. An application with a preponderance of small-stride memory traffic with good temporal locality will benefit greatly from the automatic hardware prefetcher.
- In some situations, memory traffic consisting of a preponderance of large-stride cache misses can be transformed by re-arrangement of data access sequences to alter the concentration of small-stride cache misses at the expense of large-stride cache misses to take advantage of the automatic hardware prefetcher.

7.6.3 Example of Effective Latency Reduction with Hardware Prefetch

Consider the situation that an array is populated with data corresponding to a constant-access-stride, circular pointer chasing sequence (see Example 7-2). The potential of employing the automatic hardware prefetching mechanism to reduce the effective latency of fetching a cache line from memory can be illustrated by varying the access stride between 64 bytes and the trigger threshold distance of hardware prefetch when populating the array for circular pointer chasing.

Example 7-2. Populating an Array for Circular Pointer Chasing with Constant Stride

```

register char ** p;
char *next;    // Populating pArray for circular pointer
               // chasing with constant access stride
               // p = (char **) *p; loads a value pointing to next load
p = (char **) &pArray;

for (i = 0; i < aperture; i += stride) {
    p = (char **) &pArray[i];
    if (i + stride >= g_array_aperture) {
        next = &pArray[0];
    }

    else {
        next = &pArray[i + stride];
    }

    *p = next; // populate the address of the next node
}

```

The effective latency reduction for several microarchitecture implementations is shown in Figure 7-1. For a constant-stride access pattern, the benefit of the automatic hardware prefetcher begins at half the trigger threshold distance and reaches maximum benefit when the cache-miss stride is 64 bytes.

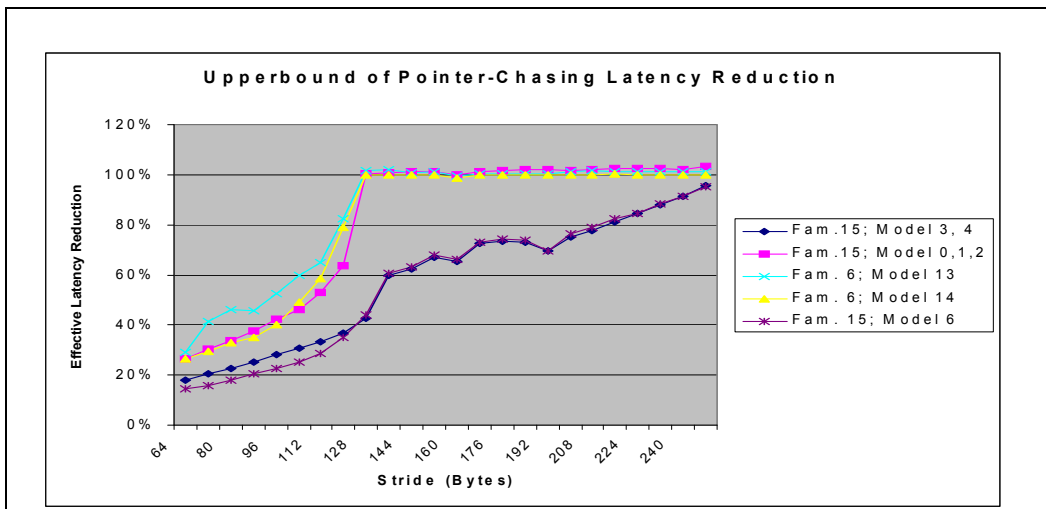


Figure 7-1. Effective Latency Reduction as a Function of Access Stride

7.6.4 Example of Latency Hiding with S/W Prefetch Instruction

Achieving the highest level of memory optimization using PREFETCH instructions requires an understanding of the architecture of a given machine. This section translates the key architectural implications into several simple guidelines for programmers to use.

Figure 7-2 and Figure 7-3 show two scenarios of a simplified 3D geometry pipeline as an example. A 3D-geometry pipeline typically fetches one vertex record at a time and then performs transformation and lighting functions on it. Both figures show two separate pipelines, an execution pipeline, and a memory pipeline (front-side bus).

Since the Pentium 4 processor (similar to the Pentium II and Pentium III processors) completely decouples the functionality of execution and memory access, the two pipelines can function concurrently. Figure 7-2 shows "bubbles" in both the execution and memory pipelines. When loads are issued for accessing vertex data, the execution units sit idle and wait until data is returned. On the other hand, the memory bus sits idle while the execution units are processing vertices. This scenario severely decreases the advantage of having a decoupled architecture.

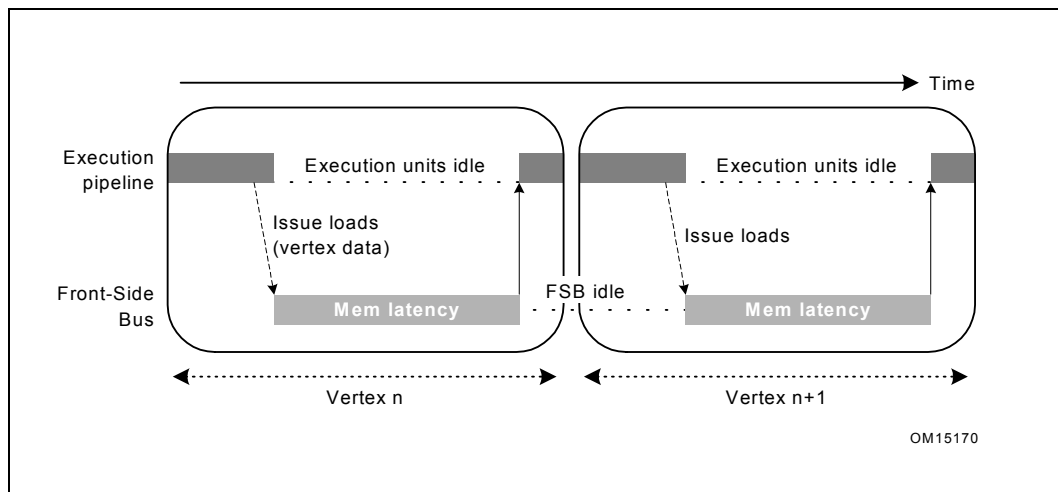


Figure 7-2. Memory Access Latency and Execution Without Prefetch

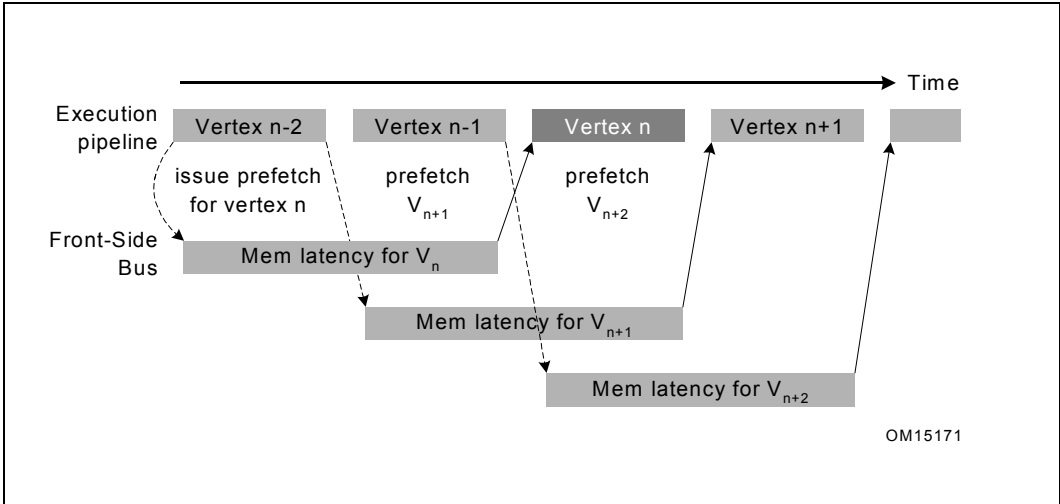


Figure 7-3. Memory Access Latency and Execution With Prefetch

The performance loss caused by poor utilization of resources can be completely eliminated by correctly scheduling the PREFETCH instructions. As shown in Figure 7-3, prefetch instructions are issued two vertex iterations ahead. This assumes that only one vertex gets processed in one iteration and a new data cache line is needed for each iteration. As a result, when iteration n , vertex V_n , is being processed; the requested data is already brought into cache. In the meantime, the front-side bus is transferring the data needed for iteration $n+1$, vertex V_{n+1} . Because there is no dependence between V_{n+1} data and the execution of V_n , the latency for data access of V_{n+1} can be entirely hidden behind the execution of V_n . Under such circumstances, no “bubbles” are present in the pipelines and thus the best possible performance can be achieved.

Prefetching is useful for inner loops that have heavy computations, or are close to the boundary between being compute-bound and memory-bandwidth-bound. It is probably not very useful for loops which are predominately memory bandwidth-bound.

When data is already located in the first level cache, prefetching can be useless and could even slow down the performance because the extra μ ops either back up waiting for outstanding memory accesses or may be dropped altogether. This behavior is platform-specific and may change in the future.

7.6.5 Software Prefetching Usage Checklist

The following checklist covers issues that need to be addressed and/or resolved to use the software PREFETCH instruction properly:

- Determine software prefetch scheduling distance.
- Use software prefetch concatenation.
- Minimize the number of software prefetches.
- Mix software prefetch with computation instructions.
- Use cache blocking techniques (for example, strip mining).
- Balance single-pass versus multi-pass execution.
- Resolve memory bank conflict issues.
- Resolve cache management issues.

Subsequent sections discuss the above items.

7.6.6 Software Prefetch Scheduling Distance

Determining the ideal prefetch placement in the code depends on many architectural parameters, including: the amount of memory to be prefetched, cache lookup latency, system memory latency, and estimate of computation cycle. The ideal distance for prefetching data is processor- and platform-dependent. If the distance is too short, the prefetch will not hide the latency of the fetch behind computation. If the prefetch is too far ahead, prefetched data may be flushed out of the cache by the time it is required.

Since prefetch distance is not a well-defined metric, for this discussion, we define a new term, prefetch scheduling distance (PSD), which is represented by the number of iterations. For large loops, prefetch scheduling distance can be set to 1 (that is, schedule prefetch instructions one iteration ahead). For small loop bodies (that is, loop iterations with little computation), the prefetch scheduling distance must be more than one iteration.

A simplified equation to compute PSD is deduced from the mathematical model. For a simplified equation, complete mathematical model, and methodology of prefetch distance determination, see Appendix E, “Summary of Rules and Suggestions.”

Example 7-3 illustrates the use of a prefetch within the loop body. The prefetch scheduling distance is set to 3, ESI is effectively the pointer to a line, EDX is the address of the data being referenced and XMM1-XMM4 are the data used in computation. Example 7-4 uses two independent cache lines of data per iteration. The PSD would need to be increased/decreased if more/less than two cache lines are used per iteration.

Example 7-3. Prefetch Scheduling Distance

```
top_loop:
    prefetchnta [edx + esi + 128*3]
    prefetchnta [edx*4 + esi + 128*3]
    .....
```


Example 7-3. Prefetch Scheduling Distance (Contd.)

```

movaps xmm1, [edx + esi]
movaps xmm2, [edx*4 + esi]
movaps xmm3, [edx + esi + 16]
movaps xmm4, [edx*4 + esi + 16]
.....
.....

add     esi, 128
cmp     esi, ecx
jl      top_loop

```

7.6.7 Software Prefetch Concatenation

Maximum performance can be achieved when the execution pipeline is at maximum throughput, without incurring any memory latency penalties. This can be achieved by prefetching data to be used in successive iterations in a loop. De-pipelining memory generates bubbles in the execution pipeline.

To explain this performance issue, a 3D geometry pipeline that processes 3D vertices in strip format is used as an example. A strip contains a list of vertices whose predefined vertex order forms contiguous triangles. It can be easily observed that the memory pipe is de-pipelined on the strip boundary due to ineffective prefetch arrangement. The execution pipeline is stalled for the first two iterations for each strip. As a result, the average latency for completing an iteration will be 165 (FIX) clocks. See Appendix E, "Summary of Rules and Suggestions", for a detailed description.

This memory de-pipelining creates inefficiency in both the memory pipeline and execution pipeline. This de-pipelining effect can be removed by applying a technique called prefetch concatenation. With this technique, the memory access and execution can be fully pipelined and fully utilized.

For nested loops, memory de-pipelining could occur during the interval between the last iteration of an inner loop and the next iteration of its associated outer loop. Without paying special attention to prefetch insertion, loads from the first iteration of an inner loop can miss the cache and stall the execution pipeline waiting for data returned, thus degrading the performance.

In Example 7-4, the cache line containing `A[II][0]` is not prefetched at all and always misses the cache. This assumes that no array `A[][]` footprint resides in the cache. The penalty of memory de-pipelining stalls can be amortized across the inner loop iterations. However, it may become very harmful when the inner loop is short. In addition, the last prefetch in the last PSD iterations are wasted and consume

machine resources. Prefetch concatenation is introduced here in order to eliminate the performance issue of memory de-pipelining.

Example 7-4. Using Prefetch Concatenation

```
for (ii = 0; ii < 100; ii++) {
  for (jj = 0; jj < 32; jj+=8) {
    prefetch a[ii][jj+8]
    computation a[ii][jj]
  }
}
```

Prefetch concatenation can bridge the execution pipeline bubbles between the boundary of an inner loop and its associated outer loop. Simply by unrolling the last iteration out of the inner loop and specifying the effective prefetch address for data used in the following iteration, the performance loss of memory de-pipelining can be completely removed. Example 7-5 gives the rewritten code.

Example 7-5. Concatenation and Unrolling the Last Iteration of Inner Loop

```
for (ii = 0; ii < 100; ii++) {
  for (jj = 0; jj < 24; jj+=8) { /* N-1 iterations */
    prefetch a[ii][jj+8]
    computation a[ii][jj]
  }
  prefetch a[ii+1][0]
  computation a[ii][jj] /* Last iteration */
}
```

This code segment for data prefetching is improved and only the first iteration of the outer loop suffers any memory access latency penalty, assuming the computation time is larger than the memory latency. Inserting a prefetch of the first data element needed prior to entering the nested loop computation would eliminate or reduce the start-up penalty for the very first iteration of the outer loop. This uncomplicated high-level code optimization can improve memory performance significantly.

7.6.8 Minimize Number of Software Prefetches

Prefetch instructions are not completely free in terms of bus cycles, machine cycles and resources, even though they require minimal clock and memory bandwidth.

Excessive prefetching may lead to performance penalties because of issue penalties in the front end of the machine and/or resource contention in the memory subsystem. This effect may be severe in cases where the target loops are small and/or cases where the target loop is issue-bound.

One approach to solve the excessive prefetching issue is to unroll and/or software-pipeline loops to reduce the number of prefetches required. Figure 7-4 presents a code example which implements prefetch and unrolls the loop to remove the redundant prefetch instructions whose prefetch addresses hit the previously issued prefetch instructions. In this particular example, unrolling the original loop once saves six prefetch instructions and nine instructions for conditional jumps in every other iteration.

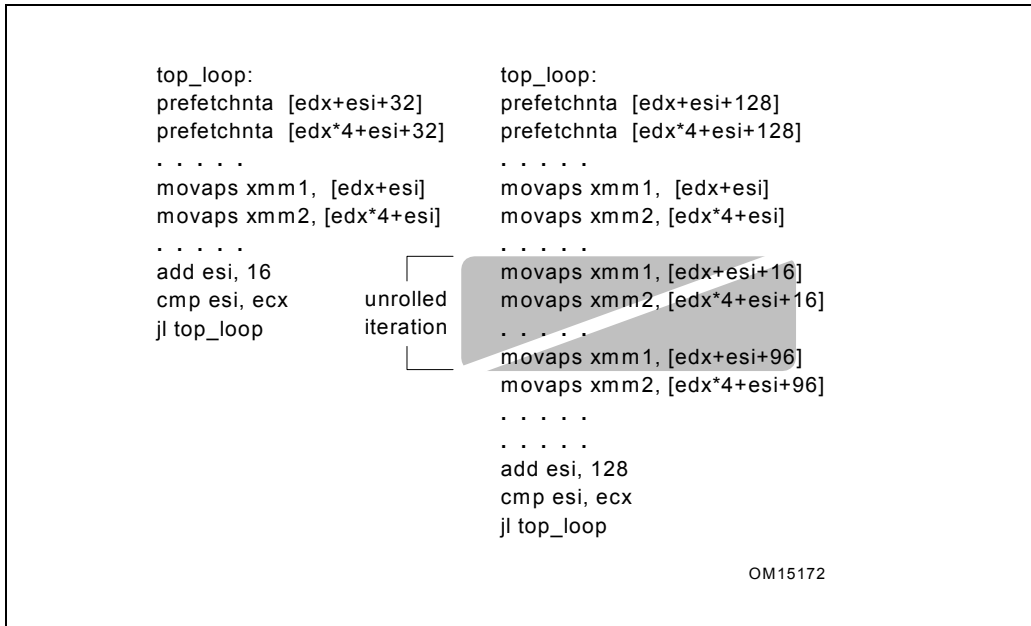


Figure 7-4. Prefetch and Loop Unrolling

Figure 7-5 demonstrates the effectiveness of software prefetches in latency hiding.

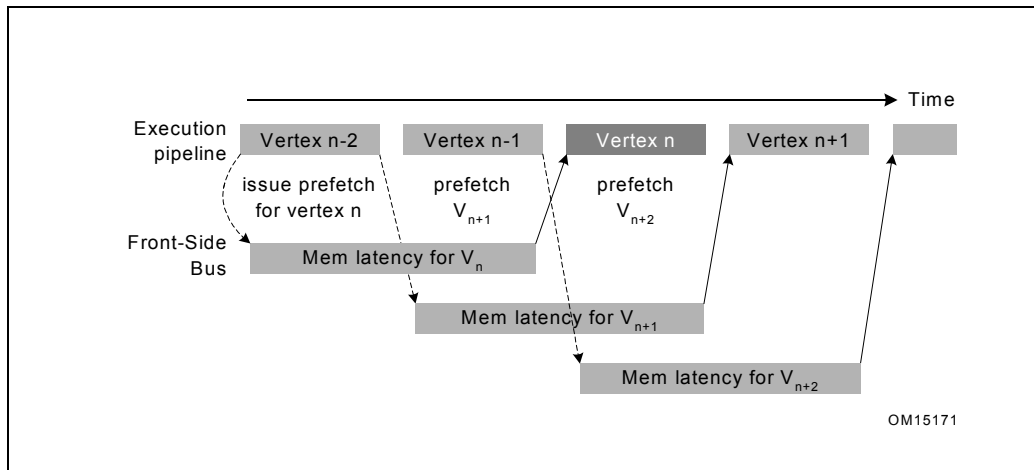


Figure 7-5. Memory Access Latency and Execution With Prefetch

The X axis in Figure 7-5 indicates the number of computation clocks per loop (each iteration is independent). The Y axis indicates the execution time measured in clocks per loop. The secondary Y axis indicates the percentage of bus bandwidth utilization. The tests vary by the following parameters:

- **Number of load/store streams** — Each load and store stream accesses one 128-byte cache line each per iteration.
- **Amount of computation per loop** — This is varied by increasing the number of dependent arithmetic operations executed.
- **Number of the software prefetches per loop** — For example, one every 16 bytes, 32 bytes, 64 bytes, 128 bytes.

As expected, the leftmost portion of each of the graphs in Figure 7-5 shows that when there is not enough computation to overlap the latency of memory access, prefetch does not help and that the execution is essentially memory-bound. The graphs also illustrate that redundant prefetches do not increase performance.

7.6.9 Mix Software Prefetch with Computation Instructions

It may seem convenient to cluster all of PREFETCH instructions at the beginning of a loop body or before a loop, but this can lead to severe performance degradation. In order to achieve the best possible performance, PREFETCH instructions must be interspersed with other computational instructions in the instruction sequence rather than clustered together. If possible, they should also be placed apart from loads. This improves the instruction level parallelism and reduces the potential instruction resource stalls. In addition, this mixing reduces the pressure on the memory access

resources and in turn reduces the possibility of the prefetch retiring without fetching data.

Figure 7-6 illustrates distributing PREFETCH instructions. A simple and useful heuristic of prefetch spreading for a Pentium 4 processor is to insert a PREFETCH instruction every 20 to 25 clocks. Rearranging PREFETCH instructions could yield a noticeable speedup for the code which stresses the cache resource.

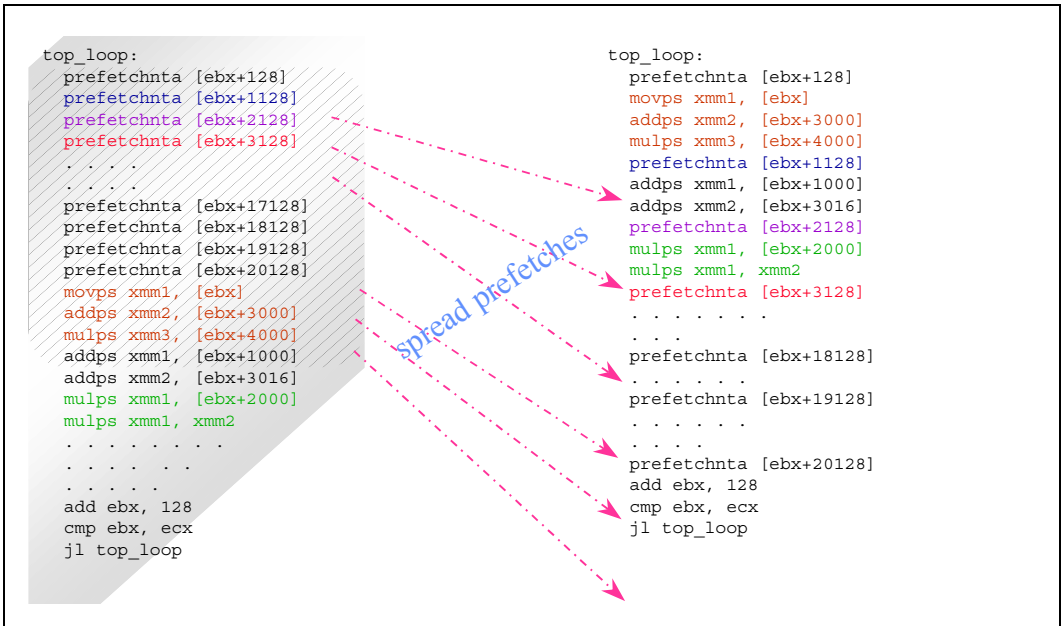


Figure 7-6. Spread Prefetch Instructions

NOTE

To avoid instruction execution stalls due to the over-utilization of the resource, PREFETCH instructions must be interspersed with computational instructions

7.6.10 Software Prefetch and Cache Blocking Techniques

Cache blocking techniques (such as strip-mining) are used to improve temporal locality and the cache hit rate. Strip-mining is one-dimensional temporal locality optimization for memory. When two-dimensional arrays are used in programs, loop blocking technique (similar to strip-mining but in two dimensions) can be applied for a better memory performance.

If an application uses a large data set that can be reused across multiple passes of a loop, it will benefit from strip mining. Data sets larger than the cache will be processed in groups small enough to fit into cache. This allows temporal data to reside in the cache longer, reducing bus traffic.

Data set size and temporal locality (data characteristics) fundamentally affect how PREFETCH instructions are applied to strip-mined code. Figure 7-7 shows two simplified scenarios for temporally-adjacent data and temporally-non-adjacent data.

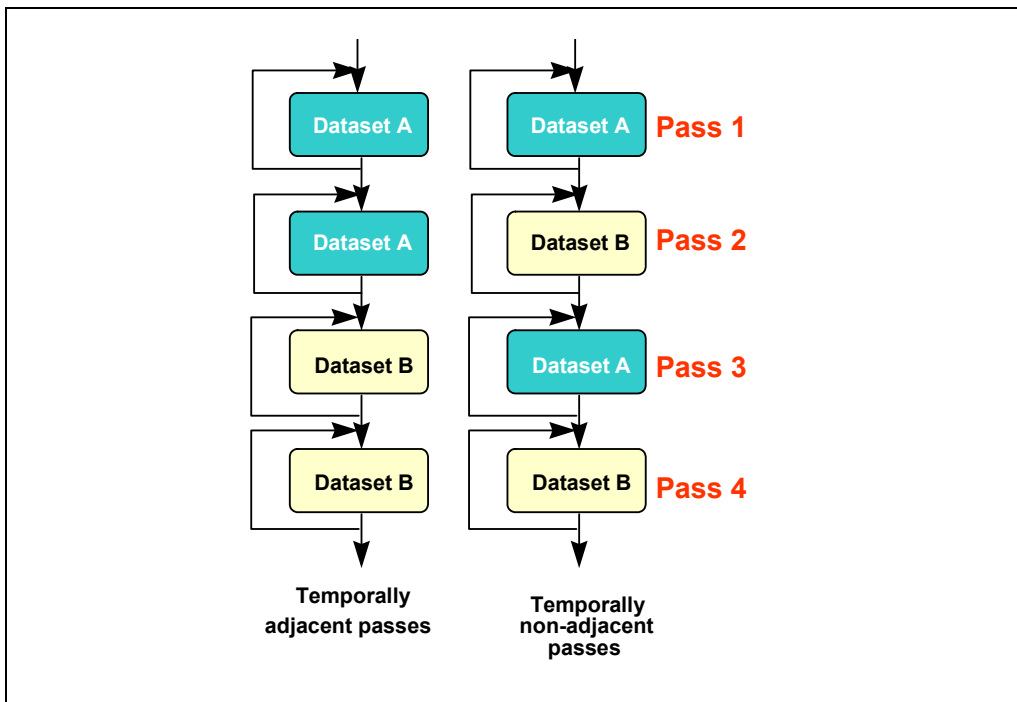


Figure 7-7. Cache Blocking - Temporally Adjacent and Non-adjacent Passes

In the temporally-adjacent scenario, subsequent passes use the same data and find it already in second-level cache. Prefetch issues aside, this is the preferred situation. In the temporally non-adjacent scenario, data used in pass m is displaced by pass $(m+1)$, requiring data *re-fetch* into the first level cache and perhaps the second level cache if a later pass reuses the data. If both data sets fit into the second-level cache, load operations in passes 3 and 4 become less expensive.

Figure 7-8 shows how prefetch instructions and strip-mining can be applied to increase performance in both of these scenarios.

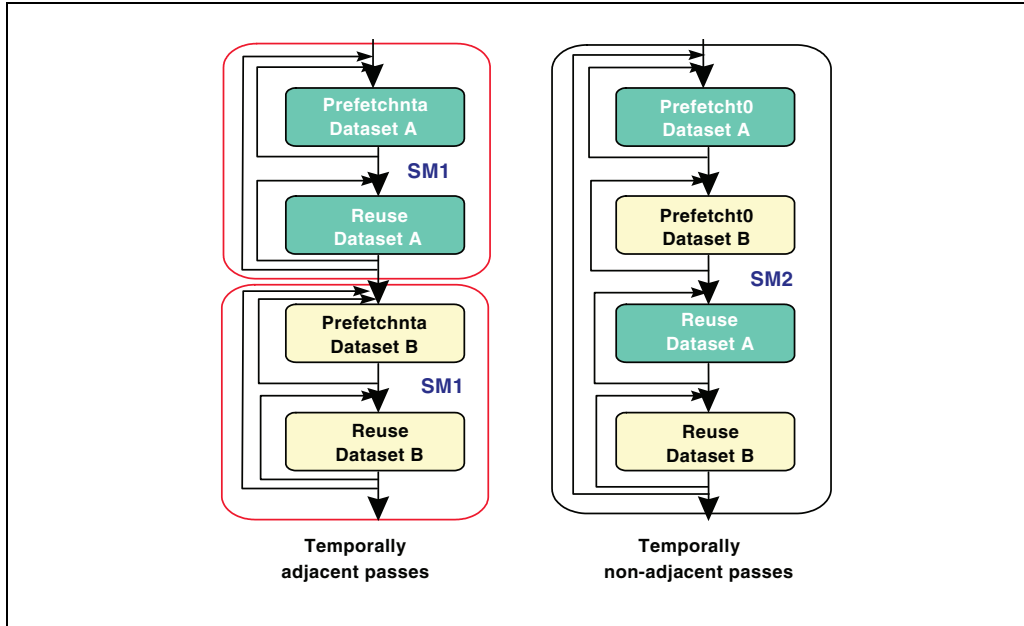


Figure 7-8. Examples of Prefetch and Strip-mining for Temporally Adjacent and Non-Adjacent Passes Loops

For Pentium 4 processors, the left scenario shows a graphical implementation of using PREFETCHNTA to prefetch data into selected ways of the second-level cache only (SM1 denotes strip mine one way of second-level), minimizing second-level cache pollution. Use PREFETCHNTA if the data is only touched once during the entire execution pass in order to minimize cache pollution in the higher level caches. This provides instant availability, assuming the prefetch was issued far ahead enough, when the read access is issued.

In scenario to the right (see Figure 7-8), keeping the data in one way of the second-level cache does not improve cache locality. Therefore, use PREFETCHT0 to prefetch the data. This amortizes the latency of the memory references in passes 1 and 2, and keeps a copy of the data in second-level cache, which reduces memory traffic and latencies for passes 3 and 4. To further reduce the latency, it might be worth considering extra PREFETCHNTA instructions prior to the memory references in passes 3 and 4.

In Example 7-6, consider the data access patterns of a 3D geometry engine first without strip-mining and then incorporating strip-mining. Note that 4-wide SIMD instructions of Pentium III processor can process 4 vertices per every iteration.

Without strip-mining, all the x,y,z coordinates for the four vertices must be re-fetched from memory in the second pass, that is, the lighting loop. This causes

OPTIMIZING CACHE USAGE

under-utilization of cache lines fetched during transformation loop as well as bandwidth wasted in the lighting loop.

Example 7-6. Data Access of a 3D Geometry Engine without Strip-mining

```
while (nvtx < MAX_NUM_VTX) {
    prefetchnta vertexi data           // v =[x,y,z,nx,ny,nz,tu,tv]
    prefetchnta vertexi+1 data
    prefetchnta vertexi+2 data
    prefetchnta vertexi+3 data
    TRANSFORMATION code                // use only x,y,z,tu,tv of a vertex
    nvtx+=4
}
while (nvtx < MAX_NUM_VTX) {
    prefetchnta vertexi data           // v =[x,y,z,nx,ny,nz,tu,tv]
                                        // x,y,z fetched again

    prefetchnta vertexi+1 data
    prefetchnta vertexi+2 data
    prefetchnta vertexi+3 data
    compute the light vectors          // use only x,y,z
    LOCAL LIGHTING code                // use only nx,ny,nz
    nvtx+=4
}
```

Now consider the code in Example 7-7 where strip-mining has been incorporated into the loops.

Example 7-7. Data Access of a 3D Geometry Engine with Strip-mining

```
while (nstrip < NUM_STRIP) {
/* Strip-mine the loop to fit data into one way of the second-level
  cache */
  while (nvtx < MAX_NUM_VTX_PER_STRIP) {
    prefetchnta vertexi data           // v =[x,y,z,nx,ny,nz,tu,tv]
    prefetchnta vertexi+1 data
    prefetchnta vertexi+2 data
    prefetchnta vertexi+3 data
    TRANSFORMATION code
    nvtx+=4
  }
  while (nvtx < MAX_NUM_VTX_PER_STRIP) {
    /* x y z coordinates are in the second-level cache, no prefetch is
       required */
```


Example 7-7. Data Access of a 3D Geometry Engine with Strip-mining

```

compute the light vectors
POINT LIGHTING code
  nvtx+=4
}
}

```

With strip-mining, all vertex data can be kept in the cache (for example, one way of second-level cache) during the strip-mined transformation loop and reused in the lighting loop. Keeping data in the cache reduces both bus traffic and the number of prefetches used.

Table 7-1 summarizes the steps of the basic usage model that incorporates only software prefetch with strip-mining. The steps are:

- Do strip-mining: partition loops so that the dataset fits into second-level cache.
- Use PREFETCHNTA if the data is only used once or the dataset fits into 32 KBytes (one way of second-level cache). Use PREFETCHT0 if the dataset exceeds 32 KBytes.

The above steps are platform-specific and provide an implementation example. The variables NUM_STRIP and MAX_NUM_VX_PER_STRIP can be heuristically determined for peak performance for specific application on a specific platform.

Table 7-1. Software Prefetching Considerations into Strip-mining Code

Read-Once Array References	Read-Multiple-Times Array References	
	Adjacent Passes	Non-Adjacent Passes
Prefetchnta	Prefetch0, SM1	Prefetch0, SM1 (2nd Level Pollution)
Evict one way; Minimize pollution	Pay memory access cost for the first pass of each array; Amortize the first pass with subsequent passes	Pay memory access cost for the first pass of every strip; Amortize the first pass with subsequent passes

7.6.11 Hardware Prefetching and Cache Blocking Techniques

Tuning data access patterns for the automatic hardware prefetch mechanism can minimize the memory access costs of the first-pass of the read-multiple-times and some of the read-once memory references. An example of the situations of read-once memory references can be illustrated with a matrix or image transpose, reading from a column-first orientation and writing to a row-first orientation, or vice versa.

Example 7-8 shows a nested loop of data movement that represents a typical matrix/image transpose problem. If the dimension of the array are large, not only the footprint of the dataset will exceed the last level cache but cache misses will

occur at large strides. If the dimensions happen to be powers of 2, aliasing condition due to finite number of way-associativity (see “Capacity Limits and Aliasing in Caches” in Chapter) will exacerbate the likelihood of cache evictions.

Example 7-8. Using HW Prefetch to Improve Read-Once Memory Traffic

```

a) Un-optimized image transpose
// dest and src represent two-dimensional arrays
for( i = 0; i < NUMCOLS; i ++){
    // inner loop reads single column
    for( j = 0; j < NUMROWS ; j ++){
        // Each read reference causes large-stride cache miss
        dest[i*NUMROWS +j] = src[j*NUMROWS + i];
    }
}
}

b)
// tilewidth = L2SizeInBytes/2/TileHeight/Sizeof(element)
for( i = 0; i < NUMCOLS; i += tilewidth) {
    for( j = 0; j < NUMROWS ; j ++){
        // access multiple elements in the same row in the inner loop
        // access pattern friendly to hw prefetch and improves hit rate
        for( k = 0; k < tilewidth; k ++){
            dest[j+ (i+k)* NUMROWS] = src[i+k+ j* NUMROWS];
        }
    }
}

```

Example 7-8 (b) shows applying the techniques of tiling with optimal selection of tile size and tile width to take advantage of hardware prefetch. With tiling, one can choose the size of two tiles to fit in the last level cache. Maximizing the width of each tile for memory read references enables the hardware prefetcher to initiate bus requests to read some cache lines before the code actually reference the linear addresses.

7.6.12 Single-pass versus Multi-pass Execution

An algorithm can use single- or multi-pass execution defined as follows:

- Single-pass, or unlayered execution passes a single data element through an entire computation pipeline.
- Multi-pass, or layered execution performs a single stage of the pipeline on a batch of data elements, before passing the batch on to the next stage.

A characteristic feature of both single-pass and multi-pass execution is that a specific trade-off exists depending on an algorithm's implementation and use of a single-pass or multiple-pass execution. See Figure 7-9.

Multi-pass execution is often easier to use when implementing a general purpose API, where the choice of code paths that can be taken depends on the specific combination of features selected by the application (for example, for 3D graphics, this might include the type of vertex primitives used and the number and type of light sources).

With such a broad range of permutations possible, a single-pass approach would be complicated, in terms of code size and validation. In such cases, each possible permutation would require a separate code sequence. For example, an object with features A, B, C, D can have a subset of features enabled, say, A, B, D. This stage would use one code path; another combination of enabled features would have a different code path. It makes more sense to perform each pipeline stage as a separate pass, with conditional clauses to select different features that are implemented within each stage. By using strip-mining, the number of vertices processed by each stage (for example, the batch size) can be selected to ensure that the batch stays within the processor caches through all passes. An intermediate cached buffer is used to pass the batch of vertices from one stage or pass to the next one.

Single-pass execution can be better suited to applications which limit the number of features that may be used at a given time. A single-pass approach can reduce the amount of data copying that can occur with a multi-pass engine. See Figure 7-9.

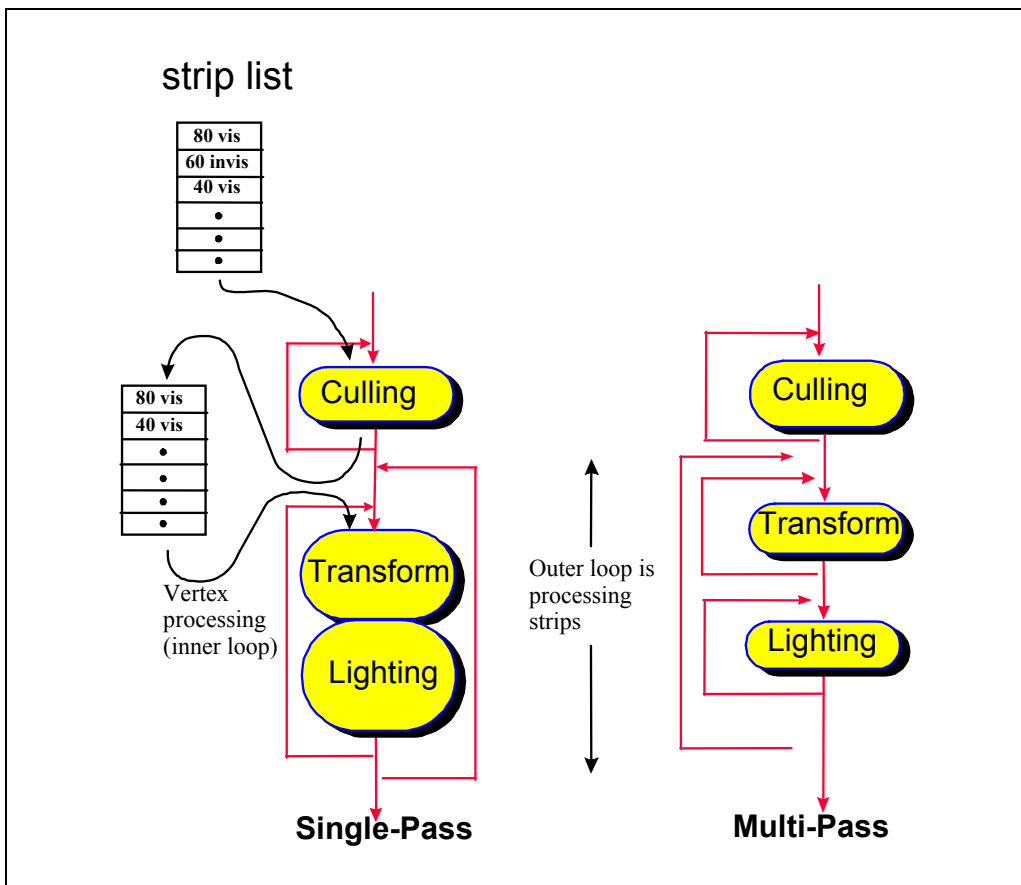


Figure 7-9. Single-Pass Vs. Multi-Pass 3D Geometry Engines

The choice of single-pass or multi-pass can have a number of performance implications. For instance, in a multi-pass pipeline, stages that are limited by bandwidth (either input or output) will reflect more of this performance limitation in overall execution time. In contrast, for a single-pass approach, bandwidth-limitations can be distributed/amortized across other computation-intensive stages. Also, the choice of which prefetch hints to use are also impacted by whether a single-pass or multi-pass approach is used.

7.7 MEMORY OPTIMIZATION USING NON-TEMPORAL STORES

Non-temporal stores can also be used to manage data retention in the cache. Uses for non-temporal stores include:

- To combine many writes without disturbing the cache hierarchy
- To manage which data structures remain in the cache and which are transient

Detailed implementations of these usage models are covered in the following sections.

7.7.1 Non-temporal Stores and Software Write-Combining

Use non-temporal stores in the cases when the data to be stored is:

- Write-once (non-temporal)
- Too large and thus cause cache thrashing

Non-temporal stores do not invoke a cache line allocation, which means they are not write-allocate. As a result, caches are not polluted and no dirty writeback is generated to compete with useful data bandwidth. Without using non-temporal stores, bus bandwidth will suffer when caches start to be thrashed because of dirty writebacks.

In Streaming SIMD Extensions implementation, when non-temporal stores are written into writeback or write-combining memory regions, these stores are weakly-ordered and will be combined internally inside the processor's write-combining buffer and be written out to memory as a line burst transaction. To achieve the best possible performance, it is recommended to align data along the cache line boundary and write them consecutively in a cache line size while using non-temporal stores. If the consecutive writes are prohibitive due to programming constraints, then software write-combining (SWWC) buffers can be used to enable line burst transaction.

You can declare small SWWC buffers (a cache line for each buffer) in your application to enable explicit write-combining operations. Instead of writing to non-temporal memory space immediately, the program writes data into SWWC buffers and combines them inside these buffers. The program only writes a SWWC buffer out using non-temporal stores when the buffer is filled up, that is, a cache line (128 bytes for the Pentium 4 processor). Although the SWWC method requires explicit instructions for performing temporary writes and reads, this ensures that the transaction on the front-side bus causes line transaction rather than several partial transactions. Application performance gains considerably from implementing this technique. These SWWC buffers can be maintained in the second-level and re-used throughout the program.

7.7.2 Cache Management

Streaming instructions (PREFETCH and STORE) can be used to manage data and minimize disturbance of temporal data held within the processor's caches.

In addition, the Pentium 4 processor takes advantage of Intel C++ Compiler support for C++ language-level features for the Streaming SIMD Extensions. Streaming SIMD Extensions and MMX technology instructions provide intrinsics that allow you to optimize cache utilization. Examples of such Intel compiler intrinsics are `_MM_PREFETCH`, `_MM_STREAM`, `_MM_LOAD`, `_MM_SFENCE`. For detail, refer to the Intel C++ Compiler User's Guide documentation.

The following examples of using prefetching instructions in the operation of video encoder and decoder as well as in simple 8-byte memory copy, illustrate performance gain from using the prefetching instructions for efficient cache management.

7.7.2.1 Video Encoder

In a video encoder, some of the data used during the encoding process is kept in the processor's second-level cache. This is done to minimize the number of reference streams that must be re-read from system memory. To ensure that other writes do not disturb the data in the second-level cache, streaming stores (MOVNTQ) are used to write around all processor caches.

The prefetching cache management implemented for the video encoder reduces the memory traffic. The second-level cache pollution reduction is ensured by preventing single-use video frame data from entering the second-level cache. Using a non-temporal PREFETCH (PREFETCHNTA) instruction brings data into only one way of the second-level cache, thus reducing pollution of the second-level cache.

If the data brought directly to second-level cache is not re-used, then there is a performance gain from the non-temporal prefetch over a temporal prefetch. The encoder uses non-temporal prefetches to avoid pollution of the second-level cache, increasing the number of second-level cache hits and decreasing the number of polluting write-backs to memory. The performance gain results from the more efficient use of the second-level cache, not only from the prefetch itself.

7.7.2.2 Video Decoder

In the video decoder example, completed frame data is written to local memory of the graphics card, which is mapped to WC (Write-combining) memory type. A copy of reference data is stored to the WB memory at a later time by the processor in order to generate future data. The assumption is that the size of the reference data is too large to fit in the processor's caches. A streaming store is used to write the data around the cache, to avoid displaying other temporal data held in the caches. Later, the processor re-reads the data using PREFETCHNTA, which ensures maximum bandwidth, yet minimizes disturbance of other cached temporal data by using the non-temporal (NTA) version of prefetch.

7.7.2.3 Conclusions from Video Encoder and Decoder Implementation

These two examples indicate that by using an appropriate combination of non-temporal prefetches and non-temporal stores, an application can be designed to lessen the overhead of memory transactions by preventing second-level cache pollution, keeping useful data in the second-level cache and reducing costly write-back transactions. Even if an application does not gain performance significantly from having data ready from prefetches, it can improve from more efficient use of the second-level cache and memory. Such design reduces the encoder's demand for such critical resource as the memory bus. This makes the system more balanced, resulting in higher performance.

7.7.2.4 Optimizing Memory Copy Routines

Creating memory copy routines for large amounts of data is a common task in software optimization. Example 7-9 presents a basic algorithm for a the simple memory copy.

Example 7-9. Basic Algorithm of a Simple Memory Copy

```
#define N 512000
double a[N], b[N];
for (i = 0; i < N; i++) {
    b[i] = a[i];
}
```

This task can be optimized using various coding techniques. One technique uses software prefetch and streaming store instructions. It is discussed in the following paragraph and a code example shown in Example 7-10.

The memory copy algorithm can be optimized using the Streaming SIMD Extensions with these considerations:

- Alignment of data
- Proper layout of pages in memory
- Cache size
- Interaction of the transaction lookaside buffer (TLB) with memory accesses
- Combining prefetch and streaming-store instructions.

The guidelines discussed in this chapter come into play in this simple example. TLB priming is required for the Pentium 4 processor just as it is for the Pentium III processor, since software prefetch instructions will not initiate page table walks on either processor.

Example 7-10. A Memory Copy Routine Using Software Prefetch

```

#define PAGESIZE 4096;
#define NUMPERPAGE 512          // # of elements to fit a page

double a[N], b[N], temp;
for (kk=0; kk<N; kk+=NUMPERPAGE) {
    temp = a[kk+NUMPERPAGE];    // TLB priming
    // use block size = page size,
    // prefetch entire block, one cache line per loop
    for (j=kk+16; j<kk+NUMPERPAGE; j+=16) {
        _mm_prefetch((char*)&a[j], _MM_HINT_NTA);
    }
    // copy 128 byte per loop
    for (j=kk; j<kk+NUMPERPAGE; j+=16) {
        _mm_stream_ps((float*)&b[j],
            _mm_load_ps((float*)&a[j]));
        _mm_stream_ps((float*)&b[j+2],
            _mm_load_ps((float*)&a[j+2]));
        _mm_stream_ps((float*)&b[j+4],
            _mm_load_ps((float*)&a[j+4]));
        _mm_stream_ps((float*)&b[j+6],
            _mm_load_ps((float*)&a[j+6]));
        _mm_stream_ps((float*)&b[j+8],
            _mm_load_ps((float*)&a[j+8]));
        _mm_stream_ps((float*)&b[j+10],
            _mm_load_ps((float*)&a[j+10]));
        _mm_stream_ps((float*)&b[j+12],
            _mm_load_ps((float*)&a[j+12]));
        _mm_stream_ps((float*)&b[j+14],
            _mm_load_ps((float*)&a[j+14]));
    } // finished copying one block
} // finished copying N elements
_mm_sfence();

```

7.7.2.5 TLB Priming

The TLB is a fast memory buffer that is used to improve performance of the translation of a virtual memory address to a physical memory address by providing fast access to page table entries. If memory pages are accessed and the page table entry

is not resident in the TLB, a TLB miss results and the page table must be read from memory.

The TLB miss results in a performance degradation since another memory access must be performed (assuming that the translation is not already present in the processor caches) to update the TLB. The TLB can be preloaded with the page table entry for the next desired page by accessing (or touching) an address in that page. This is similar to prefetch, but instead of a data cache line the page table entry is being loaded in advance of its use. This helps to ensure that the page table entry is resident in the TLB and that the prefetch happens as requested subsequently.

7.7.2.6 Using the 8-byte Streaming Stores and Software Prefetch

Example 7-10 presents the copy algorithm that uses second level cache. The algorithm performs the following steps:

1. Uses blocking technique to transfer 8-byte data from memory into second-level cache using the `_MM_PREFETCH` intrinsic, 128 bytes at a time to fill a block. The size of a block should be less than one half of the size of the second-level cache, but large enough to amortize the cost of the loop.
2. Loads the data into an XMM register using the `_MM_LOAD_PS` intrinsic.
3. Transfers the 8-byte data to a different memory location via the `_MM_STREAM` intrinsics, bypassing the cache. For this operation, it is important to ensure that the page table entry prefetched for the memory is preloaded in the TLB.

In Example 7-10, eight `_MM_LOAD_PS` and `_MM_STREAM_PS` intrinsics are used so that all of the data prefetched (a 128-byte cache line) is written back. The prefetch and streaming-stores are executed in separate loops to minimize the number of transitions between reading and writing data. This significantly improves the bandwidth of the memory accesses.

The `TEMP = A[KK+CACHESIZE]` instruction is used to ensure the page table entry for array, and `A` is entered in the TLB prior to prefetching. This is essentially a prefetch itself, as a cache line is filled from that memory location with this instruction. Hence, the prefetching starts from `KK+4` in this loop.

This example assumes that the destination of the copy is not temporally adjacent to the code. If the copied data is destined to be reused in the near future, then the streaming store instructions should be replaced with regular 128 bit stores (`_MM_STORE_PS`). This is required because the implementation of streaming stores on Pentium 4 processor writes data directly to memory, maintaining cache coherency.

7.7.2.7 Using 16-byte Streaming Stores and Hardware Prefetch

An alternate technique for optimizing a large region memory copy is to take advantage of hardware prefetcher, 16-byte streaming stores, and apply a segmented

approach to separate bus read and write transactions. See Section 3.6.12, “Minimizing Bus Latency.”

The technique employs two stages. In the first stage, a block of data is read from memory to the cache sub-system. In the second stage, cached data are written to their destination using streaming stores.

Example 7-11. Memory Copy Using Hardware Prefetch and Bus Segmentation

```
void block_prefetch(void *dst,void *src)
{ _asm {
    mov edi,dst
    mov esi,src
    mov edx,SIZE
    align 16
main_loop:
    xor ecx,ecx
    align 16
}

prefetch_loop:
    movaps xmm0, [esi+ecx]
    movaps xmm0, [esi+ecx+64]
    add ecx,128
    cmp ecx,BLOCK_SIZE
    jne prefetch_loop
    xor ecx,ecx
    align 16
cpy_loop:

    movdqa xmm0,[esi+ecx]
    movdqa xmm1,[esi+ecx+16]
    movdqa xmm2,[esi+ecx+32]
    movdqa xmm3,[esi+ecx+48]
    movdqa xmm4,[esi+ecx+64]
    movdqa xmm5,[esi+ecx+16+64]
    movdqa xmm6,[esi+ecx+32+64]
    movdqa xmm7,[esi+ecx+48+64]
    movntdq [edi+ecx],xmm0
    movntdq [edi+ecx+16],xmm1
    movntdq [edi+ecx+32],xmm2
```

Example 7-11. Memory Copy Using Hardware Prefetch and Bus Segmentation (Contd.)

```

movntdq [edi+ecx+48],xmm3
movntdq [edi+ecx+64],xmm4
movntdq [edi+ecx+80],xmm5
movntdq [edi+ecx+96],xmm6
movntdq [edi+ecx+112],xmm7
add ecx,128
cmp ecx,BLOCK_SIZE
jne cpy_loop

add esi,ecx
add edi,ecx
sub edx,ecx
jnz main_loop
sfence
}
}

```

7.7.2.8 Performance Comparisons of Memory Copy Routines

The throughput of a large-region, memory copy routine depends on several factors:

- Coding techniques that implements the memory copy task
- Characteristics of the system bus (speed, peak bandwidth, overhead in read/write transaction protocols)
- Microarchitecture of the processor

A comparison of the two coding techniques discussed above and two un-optimized techniques is shown in Table 7-2.

Table 7-2. Relative Performance of Memory Copy Routines

Processor, CPUID Signature and FSB Speed	Byte Sequential	DWORD Sequential	SW prefetch + 8 byte streaming store	4KB-Block HW prefetch + 16 byte streaming stores
Pentium M processor, 0x6Dn, 400	1.3X	1.2X	1.6X	2.5X
Intel Core Solo and Intel Core Duo processors, 0x6En, 667	3.3X	3.5X	2.1X	4.7X

Table 7-2. Relative Performance of Memory Copy Routines (Contd.)

Processor, CPUID Signature and FSB Speed	Byte Sequential	DWORD Sequential	SW prefetch + 8 byte streaming store	4KB-Block HW prefetch + 16 byte streaming stores
Pentium D processor, 0xF4n, 800	3.4X	3.3X	4.9X	5.7X

The baseline for performance comparison is the throughput (bytes/sec) of 8-MByte region memory copy on a first-generation Pentium M processor (CPUID signature 0x69n) with a 400-MHz system bus using byte-sequential technique similar to that shown in Example 7-9. The degree of improvement relative to the performance baseline for some recent processors and platforms with higher system bus speed using different coding techniques are compared.

The second coding technique moves data at 4-Byte granularity using REP string instruction. The third column compares the performance of the coding technique listed in Example 7-10. The fourth column of performance compares the throughput of fetching 4-KBytes of data at a time (using hardware prefetch to aggregate bus read transactions) and writing to memory via 16-Byte streaming stores.

Increases in bus speed is the primary contributor to throughput improvements. The technique shown in Example 7-11 will likely take advantage of the faster bus speed in the platform more efficiently. Additionally, increasing the block size to multiples of 4-KBytes while keeping the total working set within the second-level cache can improve the throughput slightly.

The relative performance figure shown in Table 7-2 is representative of clean microarchitectural conditions within a processor (e.g. looping a simple sequence of code many times). The net benefit of integrating a specific memory copy routine into an application (full-featured applications tend to create many complicated micro-architectural conditions) will vary for each application.

7.7.3 Deterministic Cache Parameters

If CPUID supports the deterministic parameter leaf, software can use the leaf to query each level of the cache hierarchy. Enumeration of each cache level is by specifying an index value (starting from 0) in the ECX register (see "CPUID-CPU Identification" in Chapter 3 of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A*).

The list of parameters is shown in Table 7-3.

Table 7-3. Deterministic Cache Parameters Leaf

Bit Location	Name	Meaning
EAX[4:0]	Cache Type	0 = Null - No more caches 1 = Data Cache 2 = Instruction Cache 3 = Unified Cache 4-31 = Reserved
EAX[7:5]	Cache Level	Starts at 1
EAX[8]	Self Initializing cache level	1: does not need SW initialization
EAX[9]	Fully Associative cache	1: Yes
EAX[13:10]	Reserved	
EAX[25:14]	Maximum number of logical processors sharing this cache	Plus encoding
EAX[31:26]	Maximum number of cores in a package	Plus 1 encoding
EBX[11:0]	System Coherency Line Size (L)	Plus 1 encoding (Bytes)
EBX[21:12]	Physical Line partitions (P)	Plus 1 encoding
EBX[31:22]	Ways of associativity (W)	Plus 1 encoding
ECX[31:0]	Number of Sets (S)	Plus 1 encoding
EDX	Reserved	
CPUID leaves > 3 < 80000000 are only visible when IA32_CR_MISC_ENABLES.BOOT_NT4 (bit 22) is clear (Default).		

The deterministic cache parameter leaf provides a means to implement software with a degree of forward compatibility with respect to enumerating cache parameters. Deterministic cache parameters can be used in several situations, including:

- Determine the size of a cache level.
- Adapt cache blocking parameters to different sharing topology of a cache-level across Hyper-Threading Technology, multicore and single-core processors.
- Determine multithreading resource topology in an MP system (See Chapter 7, "Multiple-Processor Management," of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*).
- Determine cache hierarchy topology in a platform using multicore processors (See topology enumeration white paper and reference code listed at the end of CHAPTER 1).
- Manage threads and processor affinities.
- Determine prefetch stride.

The size of a given level of cache is given by:

$$(\text{\# of Ways}) * (\text{Partitions}) * (\text{Line_size}) * (\text{Sets}) = (\text{EBX}[31:22] + 1) * (\text{EBX}[21:12] + 1) * (\text{EBX}[11:0] + 1) * (\text{ECX} + 1)$$

7.7.3.1 Cache Sharing Using Deterministic Cache Parameters

Improving cache locality is an important part of software optimization. For example, a cache blocking algorithm can be designed to optimize block size at runtime for single-processor implementations and a variety of multiprocessor execution environments (including processors supporting HT Technology, or multicore processors).

The basic technique is to place an upper limit of the blocksize to be less than the size of the target cache level divided by the number of logical processors serviced by the target level of cache. This technique is applicable to multithreaded application programming. The technique can also benefit single-threaded applications that are part of a multi-tasking workloads.

7.7.3.2 Cache Sharing in Single-Core or Multicore

Deterministic cache parameters are useful for managing shared cache hierarchy in multithreaded applications for more sophisticated situations. A given cache level may be shared by logical processors in a processor or it may be implemented to be shared by logical processors in a physical processor package.

Using the deterministic cache parameter leaf and initial APIC_ID associated with each logical processor in the platform, software can extract information on the number and the topological relationship of logical processors sharing a cache level.

7.7.3.3 Determine Prefetch Stride

The prefetch stride (see description of CPUID.01H.EBX) provides the length of the region that the processor will prefetch with the PREFETCHh instructions (PREFETCHT0, PREFETCHT1, PREFETCHT2 and PREFETCHNTA). Software will use the length as the stride when prefetching into a particular level of the cache hierarchy as identified by the instruction used. The prefetch size is relevant for cache types of Data Cache (1) and Unified Cache (3); it should be ignored for other cache types. Software should not assume that the coherency line size is the prefetch stride.

If the prefetch stride field is zero, then software should assume a default size of 64 bytes is the prefetch stride. Software should use the following algorithm to determine what prefetch size to use depending on whether the deterministic cache parameter mechanism is supported or the legacy mechanism:

- If a processor supports the deterministic cache parameters and provides a non-zero prefetch size, then that prefetch size is used.
- If a processor supports the deterministic cache parameters and does not provides a prefetch size then default size for each level of the cache hierarchy is 64 bytes.

- If a processor does not support the deterministic cache parameters but provides a legacy prefetch size descriptor (0xF0 - 64 byte, 0xF1 - 128 byte) will be the prefetch size for all levels of the cache hierarchy.
- If a processor does not support the deterministic cache parameters and does not provide a legacy prefetch size descriptor, then 32-bytes is the default size for all levels of the cache hierarchy.

CHAPTER 8

MULTICORE AND HYPER-THREADING TECHNOLOGY

This chapter describes software optimization techniques for multithreaded applications running in an environment using either multiprocessor (MP) systems or processors with hardware-based multithreading support. Multiprocessor systems are systems with two or more sockets, each mated with a physical processor package. Intel 64 and IA-32 processors that provide hardware multithreading support include dual-core processors, quad-core processors and processors supporting HT Technology¹.

Computational throughput in a multithreading environment can increase as more hardware resources are added to take advantage of thread-level or task-level parallelism. Hardware resources can be added in the form of more than one physical-processor, processor-core-per-package, and/or logical-processor-per-core. Therefore, there are some aspects of multithreading optimization that apply across MP, multicore, and HT Technology. There are also some specific microarchitectural resources that may be implemented differently in different hardware multithreading configurations (for example: execution resources are not shared across different cores but shared by two logical processors in the same core if HT Technology is enabled). This chapter covers guidelines that apply to these situations.

This chapter covers

- Performance characteristics and usage models
- Programming models for multithreaded applications
- Software optimization techniques in five specific areas

8.1 PERFORMANCE AND USAGE MODELS

The performance gains of using multiple processors, multicore processors or HT Technology are greatly affected by the usage model and the amount of parallelism in the control flow of the workload. Two common usage models are:

- Multithreaded applications
- Multitasking using single-threaded applications

1. The presence of hardware multithreading support in Intel 64 and IA-32 processors can be detected by checking the feature flag CPUID .01H:EDX[28]. A return value of in bit 28 indicates that at least one form of hardware multithreading is present in the physical processor package. The number of logical processors present in each package can also be obtained from CPUID. The application must check how many logical processors are enabled and made available to application at runtime by making the appropriate operating system calls. See the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A* for information.

8.1.1 Multithreading

When an application employs multithreading to exploit task-level parallelism in a workload, the control flow of the multi-threaded software can be divided into two parts: parallel tasks and sequential tasks.

Amdahl’s law describes an application’s performance gain as it relates to the degree of parallelism in the control flow. It is a useful guide for selecting the code modules, functions, or instruction sequences that are most likely to realize the most gains from transforming sequential tasks and control flows into parallel code to take advantage multithreading hardware support.

Figure 8-1 illustrates how performance gains can be realized for any workload according to Amdahl’s law. The bar in Figure 8-1 represents an individual task unit or the collective workload of an entire application.

In general, the speed-up of running multiple threads on an MP systems with N physical processors, over single-threaded execution, can be expressed as:

$$\text{RelativeResponse} = \frac{T_{\text{sequential}}}{T_{\text{parallel}}} = \left(1 - P + \frac{P}{N} + O\right)$$

where P is the fraction of workload that can be parallelized, and O represents the overhead of multithreading and may vary between different operating systems. In this case, performance gain is the inverse of the relative response.

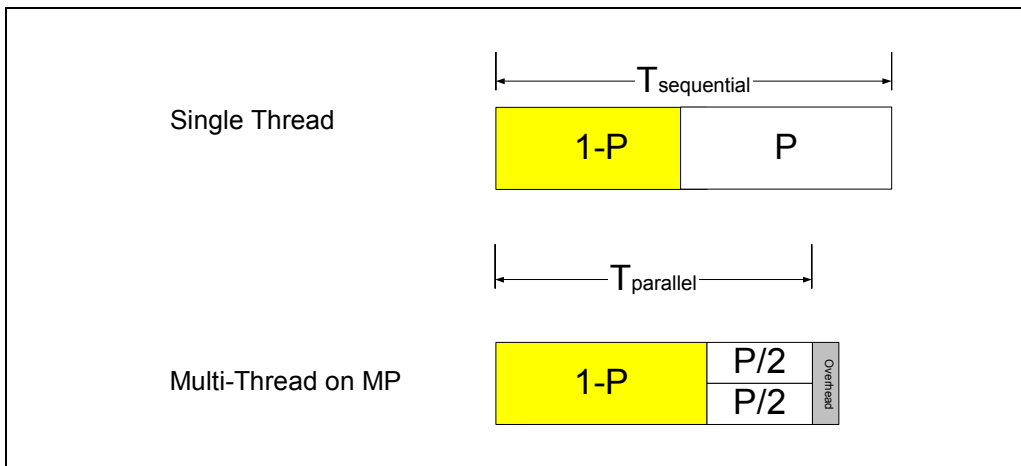


Figure 8-1. Amdahl’s Law and MP Speed-up

When optimizing application performance in a multithreaded environment, control flow parallelism is likely to have the largest impact on performance scaling with respect to the number of physical processors and to the number of logical processors per physical processor.

If the control flow of a multi-threaded application contains a workload in which only 50% can be executed in parallel, the maximum performance gain using two physical processors is only 33%, compared to using a single processor. Using four processors can deliver no more than a 60% speed-up over a single processor. Thus, it is critical to maximize the portion of control flow that can take advantage of parallelism. Improper implementation of thread synchronization can significantly increase the proportion of serial control flow and further reduce the application's performance scaling.

In addition to maximizing the parallelism of control flows, interaction between threads in the form of thread synchronization and imbalance of task scheduling can also impact overall processor scaling significantly.

Excessive cache misses are one cause of poor performance scaling. In a multi-threaded execution environment, they can occur from:

- Aliased stack accesses by different threads in the same process
- Thread contentions resulting in cache line evictions
- False-sharing of cache lines between different processors

Techniques that address each of these situations (and many other areas) are described in sections in this chapter.

8.1.2 Multitasking Environment

Hardware multithreading capabilities in Intel 64 and IA-32 processors can exploit task-level parallelism when a workload consists of several single-threaded applications and these applications are scheduled to run concurrently under an MP-aware operating system. In this environment, hardware multithreading capabilities can deliver higher throughput for the workload, although the relative performance of a single task (in terms of time of completion relative to the same task when in a single-threaded environment) will vary, depending on how much shared execution resources and memory are utilized.

For development purposes, several popular operating systems (for example Microsoft Windows* XP Professional and Home, Linux* distributions using kernel 2.4.19 or later²) include OS kernel code that can manage the task scheduling and the balancing of shared execution resources within each physical processor to maximize the throughput.

Because applications run independently under a multitasking environment, thread synchronization issues are less likely to limit the scaling of throughput. This is because the control flow of the workload is likely to be 100% parallel³ (if no inter-processor communication is taking place and if there are no system bus constraints).

-
2. This code is included in Red Hat* Linux Enterprise AS 2.1.
 3. A software tool that attempts to measure the throughput of a multitasking workload is likely to introduce control flows that are not parallel. Thread synchronization issues must be considered as an integral part of its performance measuring methodology.

With a multitasking workload, however, bus activities and cache access patterns are likely to affect the scaling of the throughput. Running two copies of the same application or same suite of applications in a lock-step can expose an artifact in performance measuring methodology. This is because an access pattern to the first level data cache can lead to excessive cache misses and produce skewed performance results. Fix this problem by:

- Including a per-instance offset at the start-up of an application
- Introducing heterogeneity in the workload by using different datasets with each instance of the application
- Randomizing the sequence of start-up of applications when running multiple copies of the same suite

When two applications are employed as part of a multitasking workload, there is little synchronization overhead between these two processes. It is also important to ensure each application has minimal synchronization overhead within itself.

An application that uses lengthy spin loops for intra-process synchronization is less likely to benefit from HT Technology in a multitasking workload. This is because critical resources will be consumed by the long spin loops.

8.2 PROGRAMMING MODELS AND MULTITHREADING

Parallelism is the most important concept in designing a multithreaded application and realizing optimal performance scaling with multiple processors. An optimized multithreaded application is characterized by large degrees of parallelism or minimal dependencies in the following areas:

- Workload
- Thread interaction
- Hardware utilization

The key to maximizing workload parallelism is to identify multiple tasks that have minimal inter-dependencies within an application and to create separate threads for parallel execution of those tasks.

Concurrent execution of independent threads is the essence of deploying a multithreaded application on a multiprocessing system. Managing the interaction between threads to minimize the cost of thread synchronization is also critical to achieving optimal performance scaling with multiple processors.

Efficient use of hardware resources between concurrent threads requires optimization techniques in specific areas to prevent contentions of hardware resources. Coding techniques for optimizing thread synchronization and managing other hardware resources are discussed in subsequent sections.

Parallel programming models are discussed next.

8.2.1 Parallel Programming Models

Two common programming models for transforming independent task requirements into application threads are:

- Domain decomposition
- Functional decomposition

8.2.1.1 Domain Decomposition

Usually large compute-intensive tasks use data sets that can be divided into a number of small subsets, each having a large degree of computational independence. Examples include:

- Computation of a discrete cosine transformation (DCT) on two-dimensional data by dividing the two-dimensional data into several subsets and creating threads to compute the transform on each subset
- Matrix multiplication; here, threads can be created to handle the multiplication of half of matrix with the multiplier matrix

Domain Decomposition is a programming model based on creating identical or similar threads to process smaller pieces of data independently. This model can take advantage of duplicated execution resources present in a traditional multiprocessor system. It can also take advantage of shared execution resources between two logical processors in HT Technology. This is because a data domain thread typically consumes only a fraction of the available on-chip execution resources.

Section 8.3.5, “Key Practices of Execution Resource Optimization,” discusses additional guidelines that can help data domain threads use shared execution resources cooperatively and avoid the pitfalls creating contentions of hardware resources between two threads.

8.2.2 Functional Decomposition

Applications usually process a wide variety of tasks with diverse functions and many unrelated data sets. For example, a video codec needs several different processing functions. These include DCT, motion estimation and color conversion. Using a functional threading model, applications can program separate threads to do motion estimation, color conversion, and other functional tasks.

Functional decomposition will achieve more flexible thread-level parallelism if it is less dependent on the duplication of hardware resources. For example, a thread executing a sorting algorithm and a thread executing a matrix multiplication routine are not likely to require the same execution unit at the same time. A design recognizing this could advantage of traditional multiprocessor systems as well as multiprocessor systems using processors supporting HT Technology.

8.2.3 Specialized Programming Models

Intel Core Duo processor and processors based on Intel Core microarchitecture offer a second-level cache shared by two processor cores in the same physical package. This provides opportunities for two application threads to access some application data while minimizing the overhead of bus traffic.

Multi-threaded applications may need to employ specialized programming models to take advantage of this type of hardware feature. One such scenario is referred to as producer-consumer. In this scenario, one thread writes data into some destination (hopefully in the second-level cache) and another thread executing on the other core in the same physical package subsequently reads data produced by the first thread.

The basic approach for implementing a producer-consumer model is to create two threads; one thread is the producer and the other is the consumer. Typically, the producer and consumer take turns to work on a buffer and inform each other when they are ready to exchange buffers. In a producer-consumer model, there is some thread synchronization overhead when buffers are exchanged between the producer and consumer. To achieve optimal scaling with the number of cores, the synchronization overhead must be kept low. This can be done by ensuring the producer and consumer threads have comparable time constants for completing each incremental task prior to exchanging buffers.

Example 8-1 illustrates the coding structure of single-threaded execution of a sequence of task units, where each task unit (either the producer or consumer) executes serially (shown in Figure 8-2). In the equivalent scenario under multi-threaded execution, each producer-consumer pair is wrapped as a thread function and two threads can be scheduled on available processor resources simultaneously.

Example 8-1. Serial Execution of Producer and Consumer Work Items

```
for (i = 0; i < number_of_iterations; i++) {
    producer (i, buff); // pass buffer index and buffer address
    consumer (i, buff);
}
```

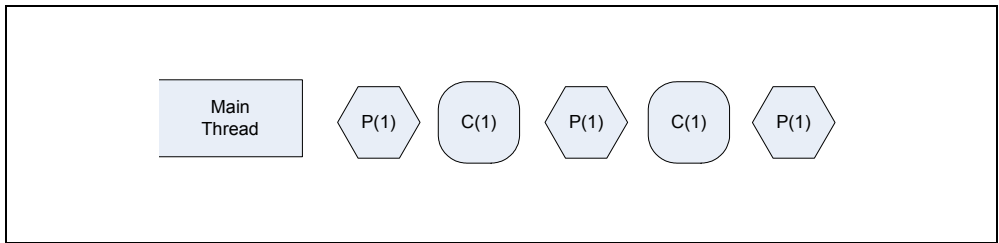


Figure 8-2. Single-threaded Execution of Producer-consumer Threading Model

8.2.3.1 Producer-Consumer Threading Models

Figure 8-3 illustrates the basic scheme of interaction between a pair of producer and consumer threads. The horizontal direction represents time. Each block represents a task unit, processing the buffer assigned to a thread.

The gap between each task represents synchronization overhead. The decimal number in the parenthesis represents a buffer index. On an Intel Core Duo processor, the producer thread can store data in the second-level cache to allow the consumer thread to continue work requiring minimal bus traffic.

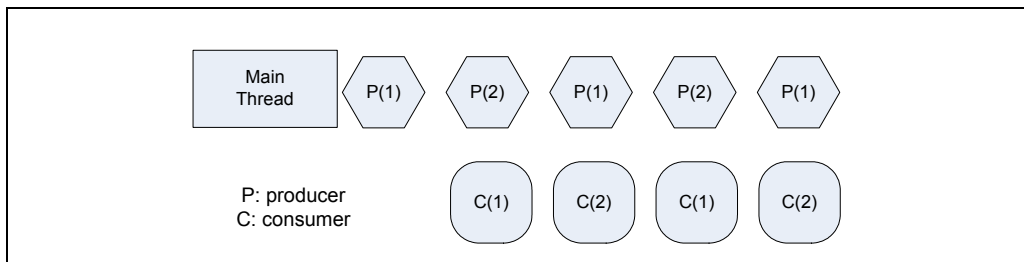


Figure 8-3. Execution of Producer-consumer Threading Model on a Multicore Processor

The basic structure to implement the producer and consumer thread functions with synchronization to communicate buffer index is shown in Example 8-2.

Example 8-2. Basic Structure of Implementing Producer Consumer Threads

```
(a) Basic structure of a producer thread function
void producer_thread()
{
    int iter_num = workamount - 1; // make local copy
    int mode1 = 1; // track usage of two buffers via 0 and 1
    produce(bufs[0],count); // placeholder function
    while (iter_num--) {

        Signal(&signal1,1); // tell the other thread to commence
        produce(bufs[mode1],count); // placeholder function
        WaitForSignal(&end1);
        mode1 = 1 - mode1; // switch to the other buffer
    }
}
```

Example 8-2. Basic Structure of Implementing Producer Consumer Threads (Contd.)

```

}
b) Basic structure of a consumer thread
void consumer_thread()
{
    int mode2 = 0; // first iteration start with buffer 0, than alternate
    int iter_num = workamount - 1;
    while (iter_num--> 0) {

        WaitForSignal(&signal1);
        consume(bufs[mode2],count); // placeholder function
        Signal(&end1,1);
        mode2 = 1 - mode2;
    }
    consume(bufs[mode2],count);
}
    
```

It is possible to structure the producer-consumer model in an interlaced manner such that it can minimize bus traffic and be effective on multicore processors without shared second-level cache.

In this interlaced variation of the producer-consumer model, each scheduling quanta of an application thread comprises of a producer task and a consumer task. Two identical threads are created to execute in parallel. During each scheduling quanta of a thread, the producer task starts first and the consumer task follows after the completion of the producer task; both tasks work on the same buffer. As each task completes, one thread signals to the other thread notifying its corresponding task to use its designated buffer. Thus, the producer and consumer tasks execute in parallel in two threads. As long as the data generated by the producer reside in either the first or second level cache of the same core, the consumer can access them without incurring bus traffic. The scheduling of the interlaced producer-consumer model is shown in Figure 8-4.

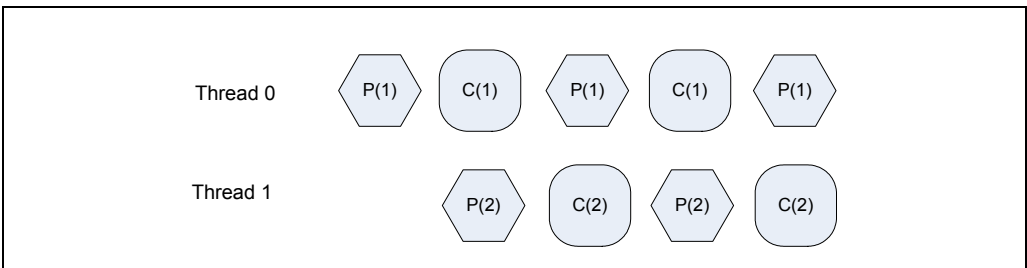


Figure 8-4. Interlaced Variation of the Producer Consumer Model

Example 8-3 shows the basic structure of a thread function that can be used in this interlaced producer-consumer model.

Example 8-3. Thread Function for an Interlaced Producer Consumer Model

```
// master thread starts first iteration, other thread must wait
// one iteration
void producer_consumer_thread(int master)
{
    int mode = 1 - master; // track which thread and its designated
                          // buffer index
    unsigned int iter_num = workamount >> 1;
    unsigned int i=0;

    iter_num += master & workamount & 1;

    if (master) // master thread starts the first iteration
    {
        produce(bufs[mode],count);
        Signal(sigp[1-mode],1); // notify producer task in follower
                               // thread that it can proceed

        consume(bufs[mode],count);
        Signal(sigc[1-mode],1);
        i = 1;
    }

    for (; i < iter_num; i++)
    {
        WaitForSignal(sigp[mode]);
        produce(bufs[mode],count); // notify the producer task in
                                   // other thread

        Signal(sigp[1-mode],1);

        WaitForSignal(sigc[mode]);
        consume(bufs[mode],count);
        Signal(sigc[1-mode],1);
    }
}
```

8.2.4 Tools for Creating Multithreaded Applications

Programming directly to a multithreading application programming interface (API) is not the only method for creating multithreaded applications. New tools (such as the Intel compiler) have become available with capabilities that make the challenge of creating multithreaded application easier.

Features available in the latest Intel compilers are:

- Generating multithreaded code using OpenMP* directives⁴
- Generating multithreaded code automatically from unmodified high-level code⁵

8.2.4.1 Programming with OpenMP Directives

OpenMP provides a standardized, non-proprietary, portable set of Fortran and C++ compiler directives supporting shared memory parallelism in applications. OpenMP supports directive-based processing. This uses special preprocessors or modified compilers to interpret parallelism expressed in Fortran comments or C/C++ pragmas. Benefits of directive-based processing include:

- The original source can be compiled unmodified.
- It is possible to make incremental code changes. This preserves algorithms in the original code and enables rapid debugging.
- Incremental code changes help programmers maintain serial consistency. When the code is run on one processor, it gives the same result as the unmodified source code.
- Offering directives to fine tune thread scheduling imbalance.
- Intel's implementation of OpenMP runtime can add minimal threading overhead relative to hand-coded multithreading.

8.2.4.2 Automatic Parallelization of Code

While OpenMP directives allow programmers to quickly transform serial applications into parallel applications, programmers must identify specific portions of the application code that contain parallelism and add compiler directives. Intel Compiler 6.0 supports a new (-QPARALLEL) option, which can identify loop structures that contain parallelism. During program compilation, the compiler automatically attempts to decompose the parallelism into threads for parallel processing. No other intervention or programmer is needed.

4. Intel Compiler 5.0 and later supports OpenMP directives. Visit <http://developer.intel.com/software/products> for details.

5. Intel Compiler 6.0 supports auto-parallelization.

8.2.4.3 Supporting Development Tools

Intel® Threading Analysis Tools include Intel® Thread Checker and Intel® Thread Profiler.

8.2.4.4 Intel® Thread Checker

Use Intel Thread Checker to find threading errors (which include data races, stalls and deadlocks) and reduce the amount of time spent debugging threaded applications.

Intel Thread Checker product is an Intel VTune Performance Analyzer plug-in data collector that executes a program and automatically locates threading errors. As the program runs, Intel Thread Checker monitors memory accesses and other events and automatically detects situations which could cause unpredictable threading-related results.

8.2.4.5 Intel® Thread Profiler

Intel Thread Profiler is a plug-in data collector for the Intel VTune Performance Analyzer. Use it to analyze threading performance and identify parallel performance bottlenecks. It graphically illustrates what each thread is doing at various levels of detail using a hierarchical summary. It can identify inactive threads, critical paths and imbalances in thread execution. Data is collapsed into relevant summaries, sorted to identify parallel regions or loops that require attention.

8.2.4.6 Intel® Threading Building Block

Intel Threading Building Block (Intel TBB) is a C++ template library that abstracts threads to tasks to create reliable, portable and scalable parallel applications. Use Intel TBB to implement task-based parallel applications and enhance developer productivity for scalable software on multi-core platforms. Intel TBB is the most efficient way to implement parallel applications and unleash multi-core platform performance compared with other threading methods like native threads and thread wrappers.

8.3 OPTIMIZATION GUIDELINES

This section summarizes optimization guidelines for tuning multithreaded applications. Five areas are listed (in order of importance):

- Thread synchronization
- Bus utilization
- Memory optimization
- Front end optimization

- Execution resource optimization

Practices associated with each area are listed in this section. Guidelines for each area are discussed in greater depth in sections that follow.

Most of the coding recommendations improve performance scaling with processor cores; and scaling due to HT Technology. Techniques that apply to only one environment are noted.

8.3.1 Key Practices of Thread Synchronization

Key practices for minimizing the cost of thread synchronization are summarized below:

- Insert the PAUSE instruction in fast spin loops and keep the number of loop repetitions to a minimum to improve overall system performance.
- Replace a spin-lock that may be acquired by multiple threads with pipelined locks such that no more than two threads have write accesses to one lock. If only one thread needs to write to a variable shared by two threads, there is no need to acquire a lock.
- Use a thread-blocking API in a long idle loop to free up the processor.
- Prevent “false-sharing” of per-thread-data between two threads.
- Place each synchronization variable alone, separated by 128 bytes or in a separate cache line.

See Section 8.4, “Thread Synchronization,” for details.

8.3.2 Key Practices of System Bus Optimization

Managing bus traffic can significantly impact the overall performance of multi-threaded software and MP systems. Key practices of system bus optimization for achieving high data throughput and quick response are:

- Improve data and code locality to conserve bus command bandwidth.
- Avoid excessive use of software prefetch instructions and allow the automatic hardware prefetcher to work. Excessive use of software prefetches can significantly and unnecessarily increase bus utilization if used inappropriately.
- Consider using overlapping multiple back-to-back memory reads to improve effective cache miss latencies.
- Use full write transactions to achieve higher data throughput.

See Section 8.5, “System Bus Optimization,” for details.

8.3.3 Key Practices of Memory Optimization

Key practices for optimizing memory operations are summarized below:

- Use cache blocking to improve locality of data access. Target one quarter to one half of cache size when targeting processors supporting HT Technology.
- Minimize the sharing of data between threads that execute on different physical processors sharing a common bus.
- Minimize data access patterns that are offset by multiples of 64-KBytes in each thread.
- Adjust the private stack of each thread in an application so the spacing between these stacks is not offset by multiples of 64 KBytes or 1 MByte (prevents unnecessary cache line evictions) when targeting processors supporting HT Technology.
- Add a per-instance stack offset when two instances of the same application are executing in lock steps to avoid memory accesses that are offset by multiples of 64 KByte or 1 MByte when targeting processors supporting HT Technology.

See Section 8.6, “Memory Optimization,” for details.

8.3.4 Key Practices of Front-end Optimization

Key practices for front-end optimization on processors that support HT Technology are:

- Avoid Excessive Loop Unrolling to ensure the Trace Cache is operating efficiently.
- Optimize code size to improve locality of Trace Cache and increase delivered trace length.

See Section 8.7, “Front-end Optimization,” for details.

8.3.5 Key Practices of Execution Resource Optimization

Each physical processor has dedicated execution resources. Logical processors in physical processors supporting HT Technology share specific on-chip execution resources. Key practices for execution resource optimization include:

- Optimize each thread to achieve optimal frequency scaling first.
- Optimize multithreaded applications to achieve optimal scaling with respect to the number of physical processors.
- Use on-chip execution resources cooperatively if two threads are sharing the execution resources in the same physical processor package.
- For each processor supporting HT Technology, consider adding functionally uncorrelated threads to increase the hardware resource utilization of each physical processor package.

See Section 8.8, “Using Thread Affinities to Manage Shared Platform Resources,” for details.

8.3.6 Generality and Performance Impact

The next five sections cover the optimization techniques in detail. Recommendations discussed in each section are ranked by importance in terms of estimated local impact and generality.

Rankings are subjective and approximate. They can vary depending on coding style, application and threading domain. The purpose of including high, medium and low impact ranking with each recommendation is to provide a relative indicator as to the degree of performance gain that can be expected when a recommendation is implemented.

It is not possible to predict the likelihood of a code instance across many applications, so an impact ranking cannot be directly correlated to application-level performance gain. The ranking on generality is also subjective and approximate.

Coding recommendations that do not impact all three scaling factors are typically categorized as medium or lower.

8.4 THREAD SYNCHRONIZATION

Applications with multiple threads use synchronization techniques in order to ensure correct operation. However, thread synchronization that are improperly implemented can significantly reduce performance.

The best practice to reduce the overhead of thread synchronization is to start by reducing the application's requirements for synchronization. Intel Thread Profiler can be used to profile the execution timeline of each thread and detect situations where performance is impacted by frequent occurrences of synchronization overhead.

Several coding techniques and operating system (OS) calls are frequently used for thread synchronization. These include spin-wait loops, spin-locks, critical sections, to name a few. Choosing the optimal OS call for the circumstance and implementing synchronization code with parallelism in mind are critical in minimizing the cost of handling thread synchronization.

SSE3 provides two instructions (MONITOR/MWAIT) to help multithreaded software improve synchronization between multiple agents. In the first implementation of MONITOR and MWAIT, these instructions are available to operating system so that operating system can optimize thread synchronization in different areas. For example, an operating system can use MONITOR and MWAIT in its system idle loop (known as C0 loop) to reduce power consumption. An operating system can also use MONITOR and MWAIT to implement its C1 loop to improve the responsiveness of the C1 loop. See Chapter 7 in the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*.

8.4.1 Choice of Synchronization Primitives

Thread synchronization often involves modifying some shared data while protecting the operation using synchronization primitives. There are many primitives to choose from. Guidelines that are useful when selecting synchronization primitives are:

- Favor compiler intrinsics or an OS provided interlocked API for atomic updates of simple data operation, such as increment and compare/exchange. This will be more efficient than other more complicated synchronization primitives with higher overhead.

For more information on using different synchronization primitives, see the white paper *Developing Multi-threaded Applications: A Platform Consistent Approach*. See <http://www3.intel.com/cd/ids/developer/asm-na/eng/53797.htm>.

- When choosing between different primitives to implement a synchronization construct, using Intel Thread Checker and Thread Profiler can be very useful in dealing with multithreading functional correctness issue and performance impact under multi-threaded execution. Additional information on the capabilities of Intel Thread Checker and Thread Profiler are described in Appendix A.

Table 8-1 is useful for comparing the properties of three categories of synchronization objects available to multi-threaded applications.

Table 8-1. Properties of Synchronization Objects

Characteristics	Operating System Synchronization Objects	Light Weight User Synchronization	Synchronization Object based on MONITOR/MWAIT
Cycles to acquire and release (if there is a contention)	Thousands or Tens of thousands cycles	Hundreds of cycles	Hundreds of cycles
Power consumption	Saves power by halting the core or logical processor if idle	Some power saving if using PAUSE	Saves more power than PAUSE
Scheduling and context switching	Returns to the OS scheduler if contention exists (can be tuned with earlier spin loop count)	Does not return to OS scheduler voluntarily	Does not return to OS scheduler voluntarily
Ring level	Ring 0	Ring 3	Ring 0

Table 8-1. Properties of Synchronization Objects (Contd.)

Characteristics	Operating System Synchronization Objects	Light Weight User Synchronization	Synchronization Object based on MONITOR/MWAIT
Miscellaneous	Some objects provide intra-process synchronization and some are for inter-process communication	Must lock accesses to synchronization variable if several threads may write to it simultaneously. Otherwise can write without locks.	Same as light weight. Can be used only on systems supporting MONITOR/MWAIT
Recommended use conditions	<ul style="list-style-type: none"> ▪ Number of active threads is greater than number of cores ▪ Waiting thousands of cycles for a signal ▪ Synchronization among processes 	<ul style="list-style-type: none"> ▪ Number of active threads is less than or equal to number of cores ▪ Infrequent contention ▪ Need inter process synchronization 	<ul style="list-style-type: none"> ▪ Same as light weight objects ▪ MONITOR/MWAIT available

8.4.2 Synchronization for Short Periods

The frequency and duration that a thread needs to synchronize with other threads depends application characteristics. When a synchronization loop needs very fast response, applications may use a spin-wait loop.

A spin-wait loop is typically used when one thread needs to wait a short amount of time for another thread to reach a point of synchronization. A spin-wait loop consists of a loop that compares a synchronization variable with some pre-defined value. See Example 8-4(a).

On a modern microprocessor with a superscalar speculative execution engine, a loop like this results in the issue of multiple simultaneous read requests from the spinning thread. These requests usually execute out-of-order with each read request being allocated a buffer resource. On detection of a write by a worker thread to a load that is in progress, the processor must guarantee no violations of memory order occur. The necessity of maintaining the order of outstanding memory operations inevitably costs the processor a severe penalty that impacts all threads.

This penalty occurs on the Pentium M processor, the Intel Core Solo and Intel Core Duo processors. However, the penalty on these processors is small compared with penalties suffered on the Pentium 4 and Intel Xeon processors. There the performance penalty for exiting the loop is about 25 times more severe.

On a processor supporting HT Technology, spin-wait loops can consume a significant portion of the execution bandwidth of the processor. One logical processor executing a spin-wait loop can severely impact the performance of the other logical processor.

Example 8-4. Spin-wait Loop and PAUSE Instructions

(a) An un-optimized spin-wait loop experiences performance penalty when exiting the loop. It consumes execution resources without contributing computational work.

```
do {
    // This loop can run faster than the speed of memory access,
    // other worker threads cannot finish modifying sync_var until
    // outstanding loads from the spinning loops are resolved.
} while( sync_var != constant_value);
```

(b) Inserting the PAUSE instruction in a fast spin-wait loop prevents performance-penalty to the spinning thread and the worker thread

```
do {
    _asm pause
    // Ensure this loop is de-pipelined, i.e. preventing more than one
    // load request to sync_var to be outstanding,
    // avoiding performance penalty when the worker thread updates
    // sync_var and the spinning thread exiting the loop.
}
while( sync_var != constant_value);
```

(c) A spin-wait loop using a “test, test-and-set” technique to determine the availability of the synchronization variable. This technique is recommended when writing spin-wait loops to run on Intel 64 and IA-32 architecture processors.

```
Spin_Lock:
    CMP lockvar, 0;           // Check if lock is free.
    JE Get_lock
    PAUSE;                   // Short delay.
    JMP Spin_Lock;
Get_Lock:
    MOV EAX, 1;
    XCHG EAX, lockvar;       // Try to get lock.
    CMP EAX, 0;             // Test if successful.
    JNE Spin_Lock;
Critical_Section:
    <critical section code>
    MOV lockvar, 0;         // Release lock.
```

User/Source Coding Rule 20. (M impact, H generality) Insert the PAUSE instruction in fast spin loops and keep the number of loop repetitions to a minimum to improve overall system performance.

On processors that use the Intel NetBurst microarchitecture core, the penalty of exiting from a spin-wait loop can be avoided by inserting a PAUSE instruction in the loop. In spite of the name, the PAUSE instruction improves performance by introducing a slight delay in the loop and effectively causing the memory read requests to

be issued at a rate that allows immediate detection of any store to the synchronization variable. This prevents the occurrence of a long delay due to memory order violation.

One example of inserting the PAUSE instruction in a simplified spin-wait loop is shown in Example 8-4(b). The PAUSE instruction is compatible with all Intel 64 and IA-32 processors. On IA-32 processors prior to Intel NetBurst microarchitecture, the PAUSE instruction is essentially a NOP instruction. Additional examples of optimizing spin-wait loops using the PAUSE instruction are available in Application note AP-949, "Using Spin-Loops on Intel Pentium 4 Processor and Intel Xeon Processor." See <http://www3.intel.com/cd/ids/developer/asmo-na/eng/dc/threading/knowledge-base/19083.htm>.

Inserting the PAUSE instruction has the added benefit of significantly reducing the power consumed during the spin-wait because fewer system resources are used.

8.4.3 Optimization with Spin-Locks

Spin-locks are typically used when several threads need to modify a synchronization variable and the synchronization variable must be protected by a lock to prevent unintentional overwrites. When the lock is released, however, several threads may compete to acquire it at once. Such thread contention significantly reduces performance scaling with respect to frequency, number of discrete processors, and HT Technology.

To reduce the performance penalty, one approach is to reduce the likelihood of many threads competing to acquire the same lock. Apply a software pipelining technique to handle data that must be shared between multiple threads.

Instead of allowing multiple threads to compete for a given lock, no more than two threads should have write access to a given lock. If an application must use spin-locks, include the PAUSE instruction in the wait loop. Example 8-4(c) shows an example of the "test, test-and-set" technique for determining the availability of the lock in a spin-wait loop.

User/Source Coding Rule 21. (M impact, L generality) *Replace a spin lock that may be acquired by multiple threads with pipelined locks such that no more than two threads have write accesses to one lock. If only one thread needs to write to a variable shared by two threads, there is no need to use a lock.*

8.4.4 Synchronization for Longer Periods

When using a spin-wait loop not expected to be released quickly, an application should follow these guidelines:

- Keep the duration of the spin-wait loop to a minimum number of repetitions.
- Applications should use an OS service to block the waiting thread; this can release the processor so that other runnable threads can make use of the processor or available execution resources.

On processors supporting HT Technology, operating systems should use the HLT instruction if one logical processor is active and the other is not. HLT will allow an idle logical processor to transition to a halted state; this allows the active logical processor to use all the hardware resources in the physical package. An operating system that does not use this technique must still execute instructions on the idle logical processor that repeatedly check for work. This “idle loop” consumes execution resources that could otherwise be used to make progress on the other active logical processor.

If an application thread must remain idle for a long time, the application should use a thread blocking API or other method to release the idle processor. The techniques discussed here apply to traditional MP system, but they have an even higher impact on processors that support HT Technology.

Typically, an operating system provides timing services, for example `Sleep(dwMilliseconds)`⁶; such variables can be used to prevent frequent checking of a synchronization variable.

Another technique to synchronize between worker threads and a control loop is to use a thread-blocking API provided by the OS. Using a thread-blocking API allows the control thread to use less processor cycles for spinning and waiting. This gives the OS more time quanta to schedule the worker threads on available processors. Furthermore, using a thread-blocking API also benefits from the system idle loop optimization that OS implements using the HLT instruction.

User/Source Coding Rule 22. (H impact, M generality) Use a thread-blocking API in a long idle loop to free up the processor.

Using a spin-wait loop in a traditional MP system may be less of an issue when the number of runnable threads is less than the number of processors in the system. If the number of threads in an application is expected to be greater than the number of processors (either one processor or multiple processors), use a thread-blocking API to free up processor resources. A multithreaded application adopting one control thread to synchronize multiple worker threads may consider limiting worker threads to the number of processors in a system and use thread-blocking APIs in the control thread.

8.4.4.1 Avoid Coding Pitfalls in Thread Synchronization

Synchronization between multiple threads must be designed and implemented with care to achieve good performance scaling with respect to the number of discrete processors and the number of logical processor per physical processor. No single technique is a universal solution for every synchronization situation.

The pseudo-code example in Example 8-5(a) illustrates a polling loop implementation of a control thread. If there is only one runnable worker thread, an attempt to

6. The `Sleep()` API is not thread-blocking, because it does not guarantee the processor will be released. Example 8-5(a) shows an example of using `Sleep(0)`, which does not always realize the processor to another thread.

call a timing service API, such as `Sleep(0)`, may be ineffective in minimizing the cost of thread synchronization. Because the control thread still behaves like a fast spinning loop, the only runnable worker thread must share execution resources with the spin-wait loop if both are running on the same physical processor that supports HT Technology. If there are more than one runnable worker threads, then calling a thread blocking API, such as `Sleep(0)`, could still release the processor running the spin-wait loop, allowing the processor to be used by another worker thread instead of the spinning loop.

A control thread waiting for the completion of worker threads can usually implement thread synchronization using a thread-blocking API or a timing service, if the worker threads require significant time to complete. Example 8-5(b) shows an example that reduces the overhead of the control thread in its thread synchronization.

Example 8-5. Coding Pitfall using Spin Wait Loop

(a) A spin-wait loop attempts to release the processor incorrectly. It experiences a performance penalty if the only worker thread and the control thread runs on the same physical processor package.

```
// Only one worker thread is running,  
// the control loop waits for the worker thread to complete.
```

```
ResumeWorkThread(thread_handle);  
While (!task_not_done ) {  
    Sleep(0) // Returns immediately back to spin loop.  
    ...  
}
```

(b) A polling loop frees up the processor correctly.

```
// Let a worker thread run and wait for completion.  
ResumeWorkThread(thread_handle);  
While (!task_not_done ) {  
    Sleep(FIVE_MILLISEC)  
  
    // This processor is released for some duration, the processor  
    // can be used by other threads.  
  
    ...  
}
```

In general, OS function calls should be used with care when synchronizing threads. When using OS-supported thread synchronization objects (critical section, mutex, or semaphore), preference should be given to the OS service that has the least synchronization overhead, such as a critical section.

8.4.5 Prevent Sharing of Modified Data and False-Sharing

On an Intel Core Duo processor or a processor based on Intel Core microarchitecture, sharing of modified data incurs a performance penalty when a thread running on one core tries to read or write data that is currently present in modified state in the first level cache of the other core. This will cause eviction of the modified cache line back into memory and reading it into the first-level cache of the other core. The latency of such cache line transfer is much higher than using data in the immediate first level cache or second level cache.

False sharing applies to data used by one thread that happens to reside on the same cache line as different data used by another thread. These situations can also incur performance delay depending on the topology of the logical processors/cores in the platform.

An example of false sharing of multithreading environment using processors based on Intel NetBurst Microarchitecture is when thread-private data and a thread synchronization variable are located within the line size boundary (64 bytes) or sector boundary (128 bytes). When one thread modifies the synchronization variable, the “dirty” cache line must be written out to memory and updated for each physical processor sharing the bus. Subsequently, data is fetched into each target processor 128 bytes at a time, causing previously cached data to be evicted from its cache on each target processor.

False sharing can experience performance penalty when the threads are running on logical processors reside on different physical processors. For processors that support HT Technology, false-sharing incurs a performance penalty when two threads run on different cores, different physical processors, or on two logical processors in the physical processor package. In the first two cases, the performance penalty is due to cache evictions to maintain cache coherency. In the latter case, performance penalty is due to memory order machine clear conditions.

False sharing is not expected to have a performance impact with a single Intel Core Duo processor.

User/Source Coding Rule 23. (H impact, M generality) *Beware of false sharing within a cache line (64 bytes on Intel Pentium 4, Intel Xeon, Pentium M, Intel Core Duo processors), and within a sector (128 bytes on Pentium 4 and Intel Xeon processors).*

When a common block of parameters is passed from a parent thread to several worker threads, it is desirable for each work thread to create a private copy of frequently accessed data in the parameter block.

8.4.6 Placement of Shared Synchronization Variable

On processors based on Intel NetBurst microarchitecture, bus reads typically fetch 128 bytes into a cache, the optimal spacing to minimize eviction of cached data is 128 bytes. To prevent false-sharing, synchronization variables and system objects

(such as a critical section) should be allocated to reside alone in a 128-byte region and aligned to a 128-byte boundary.

Example 8-6 shows a way to minimize the bus traffic required to maintain cache coherency in MP systems. This technique is also applicable to MP systems using processors with or without HT Technology.

Example 8-6. Placement of Synchronization and Regular Variables

```
int regVar;  
int padding[32];  
int SynVar[32*NUM_SYNC_VARS];  
int AnotherVar;
```

On Pentium M, Intel Core Solo, Intel Core Duo processors, and processors based on Intel Core microarchitecture; a synchronization variable should be placed alone and in separate cache line to avoid false-sharing. Software must not allow a synchronization variable to span across page boundary.

User/Source Coding Rule 24. (M impact, ML generality) Place each synchronization variable alone, separated by 128 bytes or in a separate cache line.

User/Source Coding Rule 25. (H impact, L generality) Do not place any spin lock variable to span a cache line boundary.

At the code level, false sharing is a special concern in the following cases:

- Global data variables and static data variables that are placed in the same cache line and are written by different threads.
- Objects allocated dynamically by different threads may share cache lines. Make sure that the variables used locally by one thread are allocated in a manner to prevent sharing the cache line with other threads.

Another technique to enforce alignment of synchronization variables and to avoid a cacheline being shared is to use compiler directives when declaring data structures. See Example 8-7.

Example 8-7. Declaring Synchronization Variables without Sharing a Cache Line

```
__declspec(align(64)) unsigned __int64 sum;  
struct sync_struct {...};  
__declspec(align(64)) struct sync_struct sync_var;
```

Other techniques that prevent false-sharing include:

- Organize variables of different types in data structures (because the layout that compilers give to data variables might be different than their placement in the source code).
- When each thread needs to use its own copy of a set of variables, declare the variables with:
 - Directive `threadprivate`, when using OpenMP
 - Modifier `__declspec (thread)`, when using Microsoft compiler
- In managed environments that provide automatic object allocation, the object allocators and garbage collectors are responsible for layout of the objects in memory so that false sharing through two objects does not happen.
- Provide classes such that only one thread writes to each object field and close object fields, in order to avoid false sharing.

One should not equate the recommendations discussed in this section as favoring a sparsely populated data layout. The data-layout recommendations should be adopted when necessary and avoid unnecessary bloat in the size of the work set.

8.5 SYSTEM BUS OPTIMIZATION

The system bus services requests from bus agents (e.g. logical processors) to fetch data or code from the memory sub-system. The performance impact due data traffic fetched from memory depends on the characteristics of the workload, and the degree of software optimization on memory access, locality enhancements implemented in the software code. A number of techniques to characterize memory traffic of a workload is discussed in Appendix A. Optimization guidelines on locality enhancement is also discussed in Section 3.6.10, "Locality Enhancement," and Section 7.6.11, "Hardware Prefetching and Cache Blocking Techniques."

The techniques described in Chapter 3 and Chapter 7 benefit application performance in a platform where the bus system is servicing a single-threaded environment. In a multi-threaded environment, the bus system typically services many more logical processors, each of which can issue bus requests independently. Thus, techniques on locality enhancements, conserving bus bandwidth, reducing large-stride-cache-miss-delay can have strong impact on processor scaling performance.

8.5.1 Conserve Bus Bandwidth

In a multithreading environment, bus bandwidth may be shared by memory traffic originated from multiple bus agents (These agents can be several logical processors and/or several processor cores). Preserving the bus bandwidth can improve processor scaling performance. Also, effective bus bandwidth typically will decrease if there are significant large-stride cache-misses. Reducing the amount of large-stride cache misses (or reducing DTLB misses) will alleviate the problem of bandwidth reduction due to large-stride cache misses.

One way for conserving available bus command bandwidth is to improve the locality of code and data. Improving the locality of data reduces the number of cache line evictions and requests to fetch data. This technique also reduces the number of instruction fetches from system memory.

User/Source Coding Rule 26. (M impact, H generality) *Improve data and code locality to conserve bus command bandwidth.*

Using a compiler that supports profiler-guided optimization can improve code locality by keeping frequently used code paths in the cache. This reduces instruction fetches. Loop blocking can also improve the data locality. Other locality enhancement techniques can also be applied in a multithreading environment to conserve bus bandwidth (see Section 7.6, “Memory Optimization Using Prefetch”).

Because the system bus is shared between many bus agents (logical processors or processor cores), software tuning should recognize symptoms of the bus approaching saturation. One useful technique is to examine the queue depth of bus read traffic (see Appendix A.2.1.3, “Workload Characterization”). When the bus queue depth is high, locality enhancement to improve cache utilization will benefit performance more than other techniques, such as inserting more software prefetches or masking memory latency with overlapping bus reads. An approximate working guideline for software to operate below bus saturation is to check if bus read queue depth is significantly below 5.

Some MP and workstation platforms may have a chipset that provides two system buses, with each bus servicing one or more physical processors. The guidelines for conserving bus bandwidth described above also applies to each bus domain.

8.5.2 Understand the Bus and Cache Interactions

Be careful when parallelizing code sections with data sets that results in the total working set exceeding the second-level cache and /or consumed bandwidth exceeding the capacity of the bus. On an Intel Core Duo processor, if only one thread is using the second-level cache and / or bus, then it is expected to get the maximum benefit of the cache and bus systems because the other core does not interfere with the progress of the first thread. However, if two threads use the second-level cache concurrently, there may be performance degradation if one of the following conditions is true:

- Their combined working set is greater than the second-level cache size.
- Their combined bus usage is greater than the capacity of the bus.
- They both have extensive access to the same set in the second-level cache, and at least one of the threads writes to this cache line.

To avoid these pitfalls, multithreading software should try to investigate parallelism schemes in which only one of the threads access the second-level cache at a time, or where the second-level cache and the bus usage does not exceed their limits.

8.5.3 Avoid Excessive Software Prefetches

Pentium 4 and Intel Xeon Processors have an automatic hardware prefetcher. It can bring data and instructions into the unified second-level cache based on prior reference patterns. In most situations, the hardware prefetcher is likely to reduce system memory latency without explicit intervention from software prefetches. It is also preferable to adjust data access patterns in the code to take advantage of the characteristics of the automatic hardware prefetcher to improve locality or mask memory latency. Processors based on Intel Core microarchitecture also provides several advanced hardware prefetching mechanisms. Data access patterns that can take advantage of earlier generations of hardware prefetch mechanism generally can take advantage of more recent hardware prefetch implementations.

Using software prefetch instructions excessively or indiscriminately will inevitably cause performance penalties. This is because excessively or indiscriminately using software prefetch instructions wastes the command and data bandwidth of the system bus.

Using software prefetches delays the hardware prefetcher from starting to fetch data needed by the processor core. It also consumes critical execution resources and can result in stalled execution. In some cases, it may be fruitful to evaluate the reduction or removal of software prefetches to migrate towards more effective use of hardware prefetch mechanisms. The guidelines for using software prefetch instructions are described in Chapter 3. The techniques for using automatic hardware prefetcher is discussed in Chapter 7.

User/Source Coding Rule 27. (M impact, L generality) *Avoid excessive use of software prefetch instructions and allow automatic hardware prefetcher to work. Excessive use of software prefetches can significantly and unnecessarily increase bus utilization if used inappropriately.*

8.5.4 Improve Effective Latency of Cache Misses

System memory access latency due to cache misses is affected by bus traffic. This is because bus read requests must be arbitrated along with other requests for bus transactions. Reducing the number of outstanding bus transactions helps improve effective memory access latency.

One technique to improve effective latency of memory read transactions is to use multiple overlapping bus reads to reduce the latency of sparse reads. In situations where there is little locality of data or when memory reads need to be arbitrated with other bus transactions, the effective latency of scattered memory reads can be improved by issuing multiple memory reads back-to-back to overlap multiple outstanding memory read transactions. The average latency of back-to-back bus reads is likely to be lower than the average latency of scattered reads interspersed with other bus transactions. This is because only the first memory read needs to wait for the full delay of a cache miss.

User/Source Coding Rule 28. (M impact, M generality) Consider using overlapping multiple back-to-back memory reads to improve effective cache miss latencies.

Another technique to reduce effective memory latency is possible if one can adjust the data access pattern such that the access strides causing successive cache misses in the last-level cache is predominantly less than the trigger threshold distance of the automatic hardware prefetcher. See Section 7.6.3, "Example of Effective Latency Reduction with Hardware Prefetch."

User/Source Coding Rule 29. (M impact, M generality) Consider adjusting the sequencing of memory references such that the distribution of distances of successive cache misses of the last level cache peaks towards 64 bytes.

8.5.5 Use Full Write Transactions to Achieve Higher Data Rate

Write transactions across the bus can result in write to physical memory either using the full line size of 64 bytes or less than the full line size. The latter is referred to as a partial write. Typically, writes to writeback (WB) memory addresses are full-size and writes to write-combine (WC) or uncacheable (UC) type memory addresses result in partial writes. Both cached WB store operations and WC store operations utilize a set of six WC buffers (64 bytes wide) to manage the traffic of write transactions. When competing traffic closes a WC buffer before all writes to the buffer are finished, this results in a series of 8-byte partial bus transactions rather than a single 64-byte write transaction.

User/Source Coding Rule 30. (M impact, M generality) Use full write transactions to achieve higher data throughput.

Frequently, multiple partial writes to WC memory can be combined into full-sized writes using a software write-combining technique to separate WC store operations from competing with WB store traffic. To implement software write-combining, uncacheable writes to memory with the WC attribute are written to a small, temporary buffer (WB type) that fits in the first level data cache. When the temporary buffer is full, the application copies the content of the temporary buffer to the final WC destination.

When partial-writes are transacted on the bus, the effective data rate to system memory is reduced to only 1/8 of the system bus bandwidth.

8.6 MEMORY OPTIMIZATION

Efficient operation of caches is a critical aspect of memory optimization. Efficient operation of caches needs to address the following:

- Cache blocking
- Shared memory optimization
- Eliminating 64-KByte aliased data accesses

- Preventing excessive evictions in first-level cache

8.6.1 Cache Blocking Technique

Loop blocking is useful for reducing cache misses and improving memory access performance. The selection of a suitable block size is critical when applying the loop blocking technique. Loop blocking is applicable to single-threaded applications as well as to multithreaded applications running on processors with or without HT Technology. The technique transforms the memory access pattern into blocks that efficiently fit in the target cache size.

When targeting Intel processors supporting HT Technology, the loop blocking technique for a unified cache can select a block size that is no more than one half of the target cache size, if there are two logical processors sharing that cache. The upper limit of the block size for loop blocking should be determined by dividing the target cache size by the number of logical processors available in a physical processor package. Typically, some cache lines are needed to access data that are not part of the source or destination buffers used in cache blocking, so the block size can be chosen between one quarter to one half of the target cache (see Chapter 3, “General Optimization Guidelines”).

Software can use the deterministic cache parameter leaf of CPUID to discover which subset of logical processors are sharing a given cache (see Chapter 7, “Optimizing Cache Usage”). Therefore, guideline above can be extended to allow all the logical processors serviced by a given cache to use the cache simultaneously, by placing an upper limit of the block size as the total size of the cache divided by the number of logical processors serviced by that cache. This technique can also be applied to single-threaded applications that will be used as part of a multitasking workload.

User/Source Coding Rule 31. (H impact, H generality) Use cache blocking to improve locality of data access. Target one quarter to one half of the cache size when targeting Intel processors supporting HT Technology or target a block size that allow all the logical processors serviced by a cache to share that cache simultaneously.

8.6.2 Shared-Memory Optimization

Maintaining cache coherency between discrete processors frequently involves moving data across a bus that operates at a clock rate substantially slower than the processor frequency.

8.6.2.1 Minimize Sharing of Data between Physical Processors

When two threads are executing on two physical processors and sharing data, reading from or writing to shared data usually involves several bus transactions (including snooping, request for ownership changes, and sometimes fetching data

across the bus). A thread accessing a large amount of shared memory is likely to have poor processor-scaling performance.

User/Source Coding Rule 32. (H impact, M generality) *Minimize the sharing of data between threads that execute on different bus agents sharing a common bus. The situation of a platform consisting of multiple bus domains should also minimize data sharing across bus domains.*

One technique to minimize sharing of data is to copy data to local stack variables if it is to be accessed repeatedly over an extended period. If necessary, results from multiple threads can be combined later by writing them back to a shared memory location. This approach can also minimize time spent to synchronize access to shared data.

8.6.2.2 Batched Producer-Consumer Model

The key benefit of a threaded producer-consumer design, shown in Figure 8-5, is to minimize bus traffic while sharing data between the producer and the consumer using a shared second-level cache. On an Intel Core Duo processor and when the work buffers are small enough to fit within the first-level cache, re-ordering of producer and consumer tasks are necessary to achieve optimal performance. This is because fetching data from L2 to L1 is much faster than having a cache line in one core invalidated and fetched from the bus.

Figure 8-5 illustrates a batched producer-consumer model that can be used to overcome the drawback of using small work buffers in a standard producer-consumer model. In a batched producer-consumer model, each scheduling quanta batches two or more producer tasks, each producer working on a designated buffer. The number of tasks to batch is determined by the criteria that the total working set be greater than the first-level cache but smaller than the second-level cache.

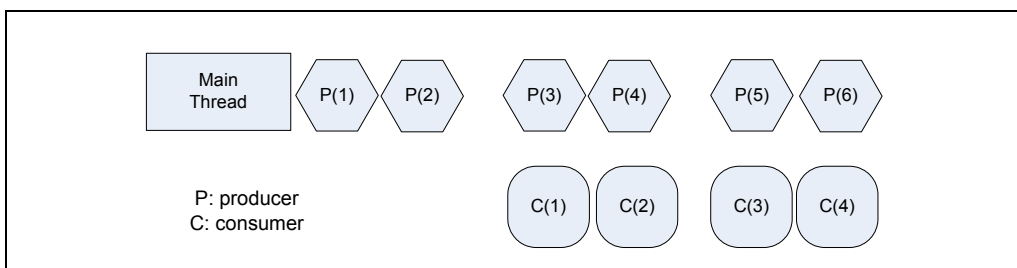


Figure 8-5. Batched Approach of Producer Consumer Model

Example 8-8 shows the batched implementation of the producer and consumer thread functions.

Example 8-8. Batched Implementation of the Producer Consumer Threads

```

void producer_thread()
{  int iter_num = workamount - batchsize;
   int mode1;
for (mode1=0; mode1 < batchsize; mode1++)
{   produce(bufs[mode1],count); }

   while (iter_num--)
   {   Signal(&signal1,1);
       produce(bufs[mode1],count); // placeholder function
       WaitForSignal(&end1);
       mode1++;
       if (mode1 > batchsize)
           mode1 = 0;
   }
}

void consumer_thread()
{  int mode2 = 0;
   int iter_num = workamount - batchsize;
   while (iter_num--)
   {   WaitForSignal(&signal1);
       consume(bufs[mode2],count); // placeholder function
       Signal(&end1,1);
       mode2++;
       if (mode2 > batchsize)
           mode2 = 0;
   }
   for (i=0;i<batchsize;i++)
   {   consume(bufs[mode2],count);
       mode2++;
       if (mode2 > batchsize)
           mode2 = 0;
   }
}

```

8.6.3 Eliminate 64-KByte Aliased Data Accesses

The 64-KByte aliasing condition is discussed in Chapter 3. Memory accesses that satisfy the 64-KByte aliasing condition can cause excessive evictions of the first-level data cache. Eliminating 64-KByte aliased data accesses originating from each thread

helps improve frequency scaling in general. Furthermore, it enables the first-level data cache to perform efficiently when HT Technology is fully utilized by software applications.

User/Source Coding Rule 33. (H impact, H generality) *Minimize data access patterns that are offset by multiples of 64 KBytes in each thread.*

The presence of 64-KByte aliased data access can be detected using Pentium 4 processor performance monitoring events. Appendix B includes an updated list of Pentium 4 processor performance metrics. These metrics are based on events accessed using the Intel VTune Performance Analyzer.

Performance penalties associated with 64-KByte aliasing are applicable mainly to current processor implementations of HT Technology or Intel NetBurst microarchitecture. The next section discusses memory optimization techniques that are applicable to multithreaded applications running on processors supporting HT Technology.

8.7 FRONT-END OPTIMIZATION

For dual-core processors where the second-level unified cache is shared by two processor cores (Intel Core Duo processor and processors based on Intel Core microarchitecture), multi-threaded software should consider the increase in code working set due to two threads fetching code from the unified cache as part of front-end and cache optimization. For quad-core processors based on Intel Core microarchitecture, the considerations that applies to Intel Core 2 Duo processors also apply to quad-core processors.

8.7.1 Avoid Excessive Loop Unrolling

Unrolling loops can reduce the number of branches and improve the branch predictability of application code. Loop unrolling is discussed in detail in Chapter 3. Loop unrolling must be used judiciously. Be sure to consider the benefit of improved branch predictability and the cost of under-utilization of the loop stream detector (LSD).

User/Source Coding Rule 34. (M impact, L generality) *Avoid excessive loop unrolling to ensure the LSD is operating efficiently..*

8.8 AFFINITIES AND MANAGING SHARED PLATFORM RESOURCES

Modern OSes provide either API and/or data constructs (e.g. affinity masks) that allow applications to manage certain shared resources , e.g. logical processors, Non-Uniform Memory Access (NUMA) memory sub-systems.

Before multithreaded software considers using affinity APIs, it should consider the recommendations in Table 8-2.

Table 8-2. Design-Time Resource Management Choices

Runtime Environment	Thread Scheduling/Processor Affinity Consideration	Memory Affinity Consideration
A single-threaded application	Support OS scheduler objectives on system response and throughput by letting OS scheduler manage scheduling. OS provides facilities for end user to optimize runtime specific environment.	Not relevant, Let OS do its job.
A multi-threaded application requiring: i) less than all processor resource in the system, ii) share system resource with other concurrent applications, iii) other concurrent applications may have higher priority.	Rely on OS default scheduler policy. Hard-coded affinity-binding will likely harm system response and throughput; and/or in some cases hurting application performance.	Rely on OS default scheduler policy. Use API that could provide transparent NUMA benefit without managing NUMA explicitly.
A multi-threaded application requiring i) foreground and higher priority, ii) uses less than all processor resource in the system, iii) share system resource with other concurrent applications, iv) but other concurrent applications have lower priority.	If application-customized thread binding policy is considered, a cooperative approach with OS scheduler should be taken instead of hard-coded thread affinity binding policy. For example, the use of <code>SetThreadIdealProcessor()</code> can provide a floating base to anchor a next-free-core binding policy for locality-optimized application binding policy, and cooperate with default OS policy.	Use API that could provide transparent NUMA benefit without managing NUMA explicitly. Use performance event to diagnose non-local memory access issue if default OS policy cause performance issue.

Table 8-2. Design-Time Resource Management Choices

Runtime Environment	Thread Scheduling/Processor Affinity Consideration	Memory Affinity Consideration
<p>A multi-threaded application runs in foreground, requiring all processor resource in the system and not sharing system resource with concurrent applications; MPI-based multi-threading.</p>	<p>Application-customized thread binding policy can be more efficient than default OS policy. Use performance event to help optimize locality and cache transfer opportunities.</p> <p>A multi-threaded application that employs its own explicit thread affinity-binding policy should deploy with some form of opt-in choice granted by the end-user or administrator. For example, permission to deploy explicit thread affinity-binding policy can be activated after permission is granted after installation.</p>	<p>Application-customized memory affinity binding policy can be more efficient than default OS policy. Use performance event to diagnose non-local memory access issues related to either OS or custom policy</p>

8.8.1 Topology Enumeration of Shared Resources

Whether multithreaded software ride on OS scheduling policy or need to use affinity APIs for customized resource management, understanding the topology of the shared platform resource is essential. The processor topology of logical processors (SMT), processor cores, and physical processors in the platform can be enumerated using information provided by CPUID. This is discussed in Chapter 8, “Multiple-Processor Management” of *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3A*. A white paper and reference code is also available from Intel.

8.8.2 Non-Uniform Memory Access

Platforms using two or more Intel Xeon processors based on Intel microarchitecture code name Nehalem support non-uniform memory access (NUMA) topology because each physical processor provides its own local memory controller. NUMA offers system memory bandwidth that can scale with the number of physical processors. System memory latency will exhibit asymmetric behavior depending on the memory transaction occurring locally in the same socket or remotely from another socket. Additionally, OS-specific construct and/or implementation behavior may present additional complexity at the API level that the multi-threaded software may need to pay attention to memory allocation/initialization in a NUMA environment.

Generally, latency sensitive workload would favor memory traffic to stay local over remote. If multiple threads shares a buffer, the programmer will need to pay attention to OS-specific behavior of memory allocation/initialization on a NUMA system.

Bandwidth sensitive workloads will find it convenient to employ a data composition threading model and aggregates application threads executing in each socket to favor local traffic on a per-socket basis to achieve overall bandwidth scalable with the number of physical processors.

The OS construct that provides the programming interface to manage local/remote NUMA traffic is referred to as memory affinity. Because OS manages the mapping between physical address (populated by system RAM) to linear address (accessed by application software); and paging allows dynamic reassignment of a physical page to map to different linear address dynamically, proper use of memory affinity will require a great deal of OS-specific knowledge.

To simplify application programming, OS may implement certain APIs and physical/linear address mapping to take advantage of NUMA characteristics transparently in certain situations. One common technique is for OS to delay commit of physical memory page assignment until the first memory reference on that physical page is accessed in the linear address space by an application thread. This means that the allocation of a memory buffer in the linear address space by an application thread does not necessarily determine which socket will service local memory traffic when the memory allocation API returns to the program. However, the memory allocation API that supports this level of NUMA transparency varies across different OSes. For example, the portable C-language API "malloc" provides some degree of transparency on Linux*, whereas the API "VirtualAlloc" behave similarly on Windows*. Different OSes may also provide memory allocation APIs that require explicit NUMA information, such that the mapping between linear address to local/remote memory traffic are fixed at allocation.

Example 8-9 shows an example that multi-threaded application could undertake the least amount of effort dealing with OS-specific APIs and to take advantage of NUMA

hardware capability. This parallel approach to memory buffer initialization is conducive to having each worker thread keep memory traffic local on NUMA systems.

Example 8-9. Parallel Memory Initialization Technique Using OpenMP and NUMA

```

#ifdef _LINUX // Linux implements malloc to commit physical page at first touch/access
    buf1 = (char *) malloc(DIM*(sizeof (double))+1024);
    buf2 = (char *) malloc(DIM*(sizeof (double))+1024);
    buf3 = (char *) malloc(DIM*(sizeof (double))+1024);
#endif
#ifdef windows
// Windows implements malloc to commit physical page at allocation, so use VirtualAlloc
    buf1 = (char *) VirtualAlloc(NULL, DIM*(sizeof (double))+1024, fAllocType, fProtect);
    buf2 = (char *) VirtualAlloc(NULL, DIM*(sizeof (double))+1024, fAllocType, fProtect);
    buf3 = (char *) VirtualAlloc(NULL, DIM*(sizeof (double))+1024, fAllocType, fProtect);
#endif
    (continue)

    a = (double *) buf1;
    b = (double *) buf2;
    c = (double *) buf3;
#pragma omp parallel
{ // use OpenMP threads to execute each iteration of the loop
// number of OpenMP threads can be specified by default or via environment variable
#pragma omp for private(num)
// each loop iteration is dispatched to execute in different OpenMP threads using private iterator
    for(num=0;num<len;num++)
        { // each thread perform first-touches to its own subset of memory address, physical pages
//            mapped to the local memory controller of the respective threads
                a[num]=10.;
                b[num]=10.;
                c[num]=10.;
        }
}
}

```

Note that the example shown in Example 8-9 implies that the memory buffers will be freed after the worker threads created by OpenMP have ended. This situation avoids a potential issue of repeated use of malloc/free across different application threads. Because if the local memory that was initialized by one thread and subsequently got

freed up by another thread, the OS may have difficulty in tracking/re-allocating memory pools in linear address space relative to NUMA topology. In Linux, another API, "numa_local_alloc" may be used.

8.9 OPTIMIZATION OF OTHER SHARED RESOURCES

Resource optimization in multi-threaded application depends on the cache topology and execution resources associated within the hierarchy of processor topology. Processor topology and an algorithm for software to identify the processor topology are discussed in Chapter 7 of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*.

In platforms with shared buses, the bus system is shared by multiple agents at the SMT level and at the processor core level of the processor topology. Thus multi-threaded application design should start with an approach to manage the bus bandwidth available to multiple processor agents sharing the same bus link in an equitable manner. This can be done by improving the data locality of an individual application thread or allowing two threads to take advantage of a shared second-level cache (where such shared cache topology is available).

In general, optimizing the building blocks of a multi-threaded application can start from an individual thread. The guidelines discussed in Chapter 3 through Chapter 9 largely apply to multi-threaded optimization.

Tuning Suggestion 3. *Optimize single threaded code to maximize execution throughput first.*

Tuning Suggestion 4. *Employ efficient threading model, leverage available tools (such as Intel Threading Building Block, Intel Thread Checker, Intel Thread Profiler) to achieve optimal processor scaling with respect to the number of physical processors or processor cores.*

8.9.1 Expanded Opportunity for HT Optimization

The Hyper-Threading Technology (HT) implementation in Intel microarchitecture code name Nehalem differs from previous generations of HT implementations. It offers broader opportunity for multi-threaded software to take advantage of HT and achieve higher system throughput over a broader range of application problems. This section provide a few heuristic recommendations and illustrates some of those situations that HT in Intel microarchitecture code name Nehalem provides more optimization opportunities.

Chapter 2, "Intel® 64 and IA-32 Architectures" covered some of the microarchitectural capability enhancement in Hyper-Threading Technology. Many of these enhancements centers around the basic needs of multi-threaded software in terms of sharing common hardware resources that may be used by more than one thread context.

Different software algorithms and workload characteristics may produce different performance characteristics due to their demands on critical microarchitectural resources that may be shared amongst several logical processors. A brief comparison of the various microarchitectural subsystem that can play a significant role in software tuning for HT is summarized in Table 8-3.

Table 8-3. Microarchitectural Resources Comparisons of HT Implementations

Microarchitectural Subsystem	Intel microarchitecture code name Nehalem	Intel NetBurst Microarchitecture
	06_1AH	0F_02H, 0F_03H, 0F_04H, 0F_06H
Issue ports, execution units	Three issue ports (0, 1, 5) distributed to handle ALU, SIMD, FP computations	Unbalanced ports, fast ALU SIMD and FP sharing the same port (port 1).
Buffering	More entries in ROB, RS, fill buffers, etc. with moderate pipeline depths	Less balance between buffer entries and pipeline depths
Branch Prediction and Misaligned memory access	More robust speculative execution with immediate reclamation after misprediction; efficient handling of cache splits	More microarchitectural hazards resulting in pipeline cleared for both threads.
Cache hierarchy	Larger and more efficient	More microarchitectural hazards to work around
Memory and bandwidth	NUMA, three channels per socket to DDR3, up to 32GB/s per socket	SMP, FSB, or dual FSB, up to 12.8 GB/s per FSB

For compute bound workloads, the HT opportunity in Intel NetBurst microarchitecture tend to favor thread contexts that executes with relatively high CPI (average cycles to retire consecutive instructions). At a hardware level, this is in part due to the issue port imbalance in the microarchitecture, as port 1 is shared by fast ALU, slow ALU (more heavy-duty integer operations), SIMD, and FP computations. At a software level, some of the cause for high CPI and may appear as benign catalyst for providing HT benefit may include: long latency instructions (port 1), some L2 hits, occasional branch mispredictions, etc. But the length of the pipeline in Intel NetBurst microarchitecture often impose additional internal hardware constraints that limits software’s ability to take advantage of HT.

The microarchitectural enhancements listed in Table 8-3 is expected to provide broader software optimization opportunities for compute-bound workloads. Whereas contention in the same execution unit by two compute-bound threads might be a concern to choose a functional-decomposition threading model over data-composition threading. Intel microarchitecture code name Nehalem will likely be more accommodating to support the programmer to choose the optimal threading decomposition models.

Memory intensive workloads can exhibit a wide range of performance characteristics, ranging from completely parallel memory traffic (saturating system memory bandwidth, as in the well-known example of Stream), memory traffic dominated by memory latency, or various mixtures of compute operations and memory traffic of either kind.

The HT implementation in Intel NetBurst microarchitecture may provide benefit to some of the latter two types of workload characteristics. The HT capability in the Intel microarchitecture code name Nehalem can broaden the operating envelop of the two latter types workload characteristics to deliver higher system throughput, due to its support for non-uniform memory access (NUMA), more efficient link protocol, and system memory bandwidth that scales with the number of physical processors.

Some cache levels of the cache hierarchy may be shared by multiple logical processors. Using the cache hierarchy is an important means for software to improve the efficiency of memory traffic and avoid saturating the system memory bandwidth. Multi-threaded applications employing cache-blocking technique may wish to partition a target cache level to take advantage of Hyper-Threading Technology. Alternatively two logical processors sharing the same L1 and L2, or logical processors sharing the L3 may wish to manage the shared resources according to their relative topological relationship. A white paper on processor topology enumeration and cache topology enumeration with companion reference code has been published (see reference at the end of chapter 1).

CHAPTER 9

64-BIT MODE CODING GUIDELINES

9.1 INTRODUCTION

This chapter describes coding guidelines for application software written to run in 64-bit mode. Some coding recommendations applicable to 64-bit mode are covered in Chapter 3. The guidelines in this chapter should be considered as an addendum to the coding guidelines described in Chapter 3 through Chapter 8.

Software that runs in either compatibility mode or legacy non-64-bit modes should follow the guidelines described in Chapter 3 through Chapter 8.

9.2 CODING RULES AFFECTING 64-BIT MODE

9.2.1 Use Legacy 32-Bit Instructions When Data Size Is 32 Bits

64-bit mode makes 16 general purpose 64-bit registers available to applications. If application data size is 32 bits, there is no need to use 64-bit registers or 64-bit arithmetic.

The default operand size for most instructions is 32 bits. The behavior of those instructions is to make the upper 32 bits all zeros. For example, when zeroing out a register, the following two instruction streams do the same thing, but the 32-bit version saves one instruction byte:

32-bit version:

```
xor eax, eax; Performs xor on lower 32bits and zeroes the upper 32 bits.
```

64-bit version:

```
xor rax, rax; Performs xor on all 64 bits.
```

This optimization holds true for the lower 8 general purpose registers: EAX, ECX, EBX, EDX, ESP, EBP, ESI, EDI. To access the data in registers R9-R15, the REX prefix is required. Using the 32-bit form there does not reduce code size.

Assembly/Compiler Coding Rule 65. (H impact, M generality) Use the 32-bit versions of instructions in 64-bit mode to reduce code size unless the 64-bit version is necessary to access 64-bit data or additional registers.

9.2.2 Use Extra Registers to Reduce Register Pressure

64-bit mode makes 8 additional 64-bit general purpose registers and 8 additional XMM registers available to applications. To access the additional registers, a single byte REX prefix is necessary. Using 8 additional registers can prevent the compiler from needing to spill values onto the stack.

Note that the potential increase in code size, due to the REX prefix, can increase cache misses. This can work against the benefit of using extra registers to access the data. When eight registers are sufficient for an algorithm, don't use the registers that require an REX prefix. This keeps the code size smaller.

Assembly/Compiler Coding Rule 66. (M impact, MH generality) When they are needed to reduce register pressure, use the 8 extra general purpose registers for integer code and 8 extra XMM registers for floating-point or SIMD code.

9.2.3 Effective Use of 64-Bit by 64-Bit Multiplies

Integer multiplies of 64-bit by 64-bit operands can produce a result that's 128-bit wide. The upper 64-bits of a 128-bit result may take a few cycles longer to be ready than the lower 64 bits. In a dependent chain of addition of integers wider than 128-bits, accessing the high 64-bits result of the multiply should be delayed relative to the low 64-bit multiply result for optimal software pipelining.

If the compiler can determine at compile time that the result of a multiply will not exceed 64 bits, then the compiler should generate the multiply instruction that produces a 64-bit result. If the compiler or assembly programmer can not determine that the result will be less than 64 bits, then a multiply that produces a 128-bit result is necessary.

Assembly/Compiler Coding Rule 67. (ML impact, M generality) Prefer 64-bit by 64-bit integer multiplies that produce 64-bit results over multiplies that produce 128-bit results.

Assembly/Compiler Coding Rule 68. (ML impact, M generality) Stagger accessing the high 64-bit result of a 128-bit multiply after accessing the low 64-bit results.

In Intel microarchitecture code name Sandy Bridge, the low 64-bit result of a 128-bit multiply is ready to use in 3 cycles and the high 64-bit result is ready one cycle after the low 64-bit result. This can speed up the calculation of integer multiply and division of large integers.

9.2.4 Replace 128-bit Integer Division with 128-bit Multiplies

Modern compiler can transform expressions of integer division in high-level language code with a constant divisor into assembly sequences that uses IMUL/MUL to replace IDIV/DIV instructions. Typically, compilers will replace divisor value that is within the range of 32-bits and if the divisor value is known at compile time. If the divisor value is not known at compile time or the divisor is greater than those represented by 32-bits, DIV or IDIV will be generated.

The latency of DIV instruction with a 128-bit dividend is quite long. For dividend value greater than 64-bits, the latency can range from 70-90 cycles.

The basic technique that compiler employs to transform integer division into 128-bit integer multiply is based on the congruence principle of modular arithmetic. It can be easily extended to handle larger divisor values to take advantage of fast 128-bit IMUL/MUL operation.

The integer equation:

Dividend = $Q * \text{divisor} + R$, or

$Q = \text{floor}(\text{Dividend}/\text{Divisor})$, $R = \text{Dividend} - Q * \text{Divisor}$

Transform to the real number domain:

$\text{floor}(\text{Dividend}/\text{Divisor}) = \text{Dividend}/\text{Divisor} - R/\text{Divisor}$, is equivalent to

$Q * C2 = \text{Dividend} * (C2 / \text{divisor}) + R * C2 / \text{Divisor}$

One can choose $C2 = 2^N$ to control the rounding of the last term, then

$Q = ((\text{Dividend} * (C2 / \text{divisor})) \gg N) + ((R * C2 / \text{Divisor}) \gg N)$.

If "divisor" is known at compile time, $(C2/\text{Divisor})$ can be pre-computed into a congruent constant $Cx = \text{Ceil}(C2/\text{divisor})$, then the quotient can computed by an integer multiple, followed by a shift:

$Q = (\text{Dividend} * Cx) \gg N$;

$R = \text{Dividend} - ((\text{Dividend} * Cx) \gg N) * \text{divisor}$;

The 128-bit IDIV/DIV instructions restrict the range of divisor, quotient, and remainder to be within 64- bits to avoid causing numerical exceptions. This presents a challenge for situations with either of the three having a value near the upper bounds of 64-bit and for dividend values nearing the upper bound of 128 bits.

This challenge can be overcome with choosing a larger shift count N , and extending the $(\text{Dividend} * Cx)$ operation from 128-bit range to the next computing-efficient range. For example, if $(\text{Dividend} * Cx)$ is greater and 128 bits and N is greater than 63 bits, one can take advantage of computing bits 191:64 of the 192-bit results using 128-bit MUL without implementing a full 192-bit multiply.

A convenient way to choose the congruent constant Cx is as follows:

- If the range of dividend is within 64 bits: $N_{\text{min}} \sim \text{BSR}(\text{Divisor}) + 63$,

64-BIT MODE CODING GUIDELINES

- In situations of disparate dynamic range of quotient/remainder relative to the range of divisor, raise N accordingly so that quotient/remainder can be computed efficiently.

Consider the computation of quotient/remainder computation for the divisor 10^{16} on unsigned dividends near the range of 64-bits. Example 9-1 illustrates using "MUL r64" instruction to handle 64-bit dividend with 64-bit divisors.

Example 9-1. Compute 64-bit Quotient and Remainder with 64-bit Divisor

```
_Cx10to16:      ; Congruent constant for 10^16 with shift count 'N' = 117
  DD      0c44de15ch      ; floor ( (2^117 / 10^16) + 1)
  DD      0e69594beh      ; Optimize length of Cx to reduce # of 128-bit multiplies
_tento16:      ; 10^16
  DD      6fc10000h
  DD      002386f2h

  mov     r9, qword ptr [rcx]      ; load 64-bit dividend value
  mov     rax, r9
  mov     rsi, _Cx10to16      ; Congruent Constant for 10^16 with shift count 117
  mul    { [rsi]                ; 128-bit multiply
  mov     r10, qword ptr 8[rsi]    ; load divisor 10^16
  shr     rdx, 53;                ;
  mov     r8, rdx

  mov     rax, r8
  mul    { r10                  ; 128-bit multiply
  sub     r9, rax;                ;
  jae    remain
  sub     r8                      ; this may be off by one due to round up
  mov     rax, r8
  mul    { r10                  ; 128-bit multiply
  sub     r9, rax;                ;
remain:
  mov     rdx, r8                ; quotient
  mov     rax, r9                ; remainder
```

Example 9-2 shows a similar technique to handle 128-bit dividend with 64-bit divisors.

Example 9-2. Quotient and Remainder of 128-bit Dividend with 64-bit Divisor

```

mov    rax, qword ptr [rcx]    ; load bits 63:0 of 128-bit dividend from memory
mov    rsi, _Cx10to16        ; Congruent Constant for 10^16 with shift count 117
mov    r9, qword ptr [rsi]    ; load Congruent Constant
mul    { r9                    ; 128-bit multiply
xor    r11, r11                ; clear accumulator
mov    rax, qword ptr 8[rcx]   ; load bits 127:64 of 128-bit dividend
shr    rdx, 53;                ;
mov    r10, rdx                ; initialize bits 127:64 of 192 bit result
mul    { r9                    ; Accumulate to bits 191:128
add    rax, r10;                ;
adc    rdx, r11;                ;
shr    rax, 53;                ;
shl    rdx, 11;                ;
or     rdx, rax;                ;
mov    r8, qword ptr 8[rsi]    ; load Divisor 10^16
mov    r9, rdx;                ; approximate quotient, may be off by 1
mov    rax, r8
mul    { r9                    ; will quotient * divisor > dividend?
sub    rax, qword ptr [rcx]    ;
sbb    rdx, qword ptr 8[rcx]   ;

jb     remain
sub    r9                    ; this may be off by one due to round up
mov    rax, r8                ; retrieve Divisor 10^16
mul    { r9                    ; final quotient * divisor
sub    rax, qword ptr [rcx]    ;
sbb    rdx, qword ptr 8[rcx]   ;
remain:
mov    rdx, r9                ; quotient
neg    rax                    ; remainder

```

The techniques illustrated in Example 9-1 and Example 9-2 can increase the speed of remainder/quotient calculation of 128-bit dividends to at or below the cost of a 32-bit integer division.

Extending the technique above to deal with divisor greater than 64-bits is relatively straightforward. One optimization worth considering is to choose shift count $N > 128$

bits. This can reduce the number of 128-bit MUL needed to compute the relevant upper bits of (Dividend * Cx).

9.2.5 Sign Extension to Full 64-Bits

When in 64-bit mode, processors based on Intel NetBurst microarchitecture can sign-extend to 64 bits in a single μ op. In 64-bit mode, when the destination is 32 bits, the upper 32 bits must be zeroed.

Zeroing the upper 32 bits requires an extra μ op and is less optimal than sign extending to 64 bits. While sign extending to 64 bits makes the instruction one byte longer, it reduces the number of μ ops that the trace cache has to store, improving performance.

For example, to sign-extend a byte into ESI, use:

```
movsx rsi, BYTE PTR[rax]
```

instead of:

```
movsx esi, BYTE PTR[rax]
```

If the next instruction uses the 32-bit form of esi register, the result will be the same. This optimization can also be used to break an unintended dependency. For example, if a program writes a 16-bit value to a register and then writes the register with an 8-bit value, if bits 15:8 of the destination are not needed, use the sign-extended version of writes when available.

For example:

```
mov r8w, r9w; Requires a merge to preserve
; bits 63:15.
mov r8b, r10b; Requires a merge to preserve bits 63:8
```

Can be replaced with:

```
movsx r8, r9w ; If bits 63:8 do not need to be
; preserved.
movsx r8, r10b ; If bits 63:8 do not need to
; be preserved.
```

In the above example, the moves to R8W and R8B both require a merge to preserve the rest of the bits in the register. There is an implicit real dependency on R8 between the 'MOV R8W, R9W' and 'MOV R8B, R10B'. Using MOVSX breaks the real dependency and leaves only the output dependency, which the processor can eliminate through renaming.

For processors based on Intel Core microarchitecture, zeroing the upper 32 bits is faster than sign-extend to 64 bits. For processors based on Intel microarchitecture code name Nehalem, zeroing or sign-extend the upper bits is single micro-op.

9.3 ALTERNATE CODING RULES FOR 64-BIT MODE

9.3.1 Use 64-Bit Registers Instead of Two 32-Bit Registers for 64-Bit Arithmetic Result

Legacy 32-bit mode offers the ability to support extended precision integer arithmetic (such as 64-bit arithmetic). However, 64-bit mode offers native support for 64-bit arithmetic. When 64-bit integers are desired, use the 64-bit forms of arithmetic instructions.

In 32-bit legacy mode, getting a 64-bit result from a 32-bit by 32-bit integer multiply requires three registers; the result is stobbed in 32-bit chunks in the EDX:EAX pair. When the instruction is available in 64-bit mode, using the 32-bit version of the instruction is not the optimal implementation if a 64-bit result is desired. Use the extended registers.

For example, the following code sequence loads the 32-bit values sign-extended into the 64-bit registers and performs a multiply:

```
movsx rax, DWORD PTR[x]
movsx rcx, DWORD PTR[y]
imul rax, rcx
```

The 64-bit version above is more efficient than using the following 32-bit version:

```
mov eax, DWORD PTR[x]
mov ecx, DWORD PTR[y]
imul ecx
```

In the 32-bit case above, EAX is required to be a source. The result ends up in the EDX:EAX pair instead of in a single 64-bit register.

Assembly/Compiler Coding Rule 69. (ML impact, M generality) Use the 64-bit versions of multiply for 32-bit integer multiplies that require a 64 bit result.

To add two 64-bit numbers in 32-bit legacy mode, the add instruction followed by the addc instruction is used. For example, to add two 64-bit variables (X and Y), the following four instructions could be used:

```
mov eax, DWORD PTR[X]
mov edx, DWORD PTR[X+4]
add eax, DWORD PTR[Y]
adc edx, DWORD PTR[Y+4]
```

The result will end up in the two-register EDX:EAX.

In 64-bit mode, the above sequence can be reduced to the following:

```
mov rax, QWORD PTR[X]
add rax, QWORD PTR[Y]
```

The result is stored in rax. One register is required instead of two.

Assembly/Compiler Coding Rule 70. (ML impact, M generality) Use the 64-bit versions of *add* for 64-bit *adds*.

9.3.2 CVTSS2SS and CVTSD2SD

In processors based on Intel Core microarchitecture and later, CVTSS2SS and CVTSD2SD are improved significantly over those in Intel NetBurst microarchitecture, in terms of latency and throughput. The improvements applies equally to 64-bit and 32-bit versions.

9.3.3 Using Software Prefetch

Intel recommends that software developers follow the recommendations in Chapter 3 and Chapter 7 when considering the choice of organizing data access patterns to take advantage of the hardware prefetcher (versus using software prefetch).

Assembly/Compiler Coding Rule 71. (L impact, L generality) If software prefetch instructions are necessary, use the prefetch instructions provided by SSE.

CHAPTER 10 SSE4.2 AND SIMD PROGRAMMING FOR TEXT-PROCESSING/LEXING/PARSING

String/text processing spans a discipline that often employs techniques different from traditional SIMD integer vector processing. Much of the traditional string/text algorithms are character based, where characters may be represented by encodings (or code points) of fixed or variable byte sizes. Textual data represents a vast amount of raw data and often carrying contextual information. The contextual information embedded in raw textual data often requires algorithmic processing dealing with a wide range of attributes, such as character values, character positions, character encoding formats, subsetting of character sets, strings of explicit or implicit lengths, tokens, delimiters; contextual objects may be represented by sequential characters within a pre-defined character subsets (e.g. decimal-valued strings); textual streams may contain embedded state transitions separating objects of different contexts (e.g. tag-delimited fields).

Traditional Integer SIMD vector instructions may, in some simpler situations, be successful to speed up simple string processing functions. SSE4.2 includes four new instructions that offer advances to computational algorithms targeting string/text processing, lexing and parsing of either unstructured or structured textual data.

10.1 SSE4.2 STRING AND TEXT INSTRUCTIONS

SSE4.2 provides four instructions, PCMPSTRI/PCMPSTRM/PCMPISTRI/PCMPISTRM that can accelerate string and text processing by combining the efficiency of SIMD programming techniques and the lexical primitives that are embedded in these 4 instructions. Simple examples of these instructions include string length determination, direct string comparison, string case handling, delimiter/token processing, locating word boundaries, locating sub-string matches in large text blocks. Sophisticated application of SSE4.2 can accelerate XML parsing and Schema validation.

Processor's support for SSE4.2 is indicated by the feature flag value returned in ECX [bit 20] after executing CPUID instruction with EAX input value of 1 (i.e. SSE4.2 is supported if CPUID.01H:ECX.SSE4_2 [bit 20] = 1). Therefore, software must verify CPUID.01H:ECX.SSE4_2 [bit 20] is set before using these 4 instructions. (Verifying CPUID.01H:ECX.SSE4_2 = 1 is also required before using PCMPGTQ or CRC32. Verifying CPUID.01H:ECX.POPCNT[Bit 23] = 1 is required before using the POPCNT instruction.)

These string/text processing instructions work by performing up to 256 comparison operations on text fragments. Each text fragment can be 16 bytes. They can handle fragments of different formats: either byte or word elements. Each of these four instructions can be configured to perform four types of parallel comparison operation on two text fragments.

The aggregated intermediate result of a parallel comparison of two text fragments become a bit patterns: 16 bits for processing byte elements or 8 bits for word elements. These instruction provide additional flexibility, using bit fields in the immediate operand of the instruction syntax, to configure an unary transformation (polarity) on the first intermediate result.

Lastly, the instruction’s immediate operand offers a output selection control to further configure the flexibility of the final result produced by the instruction. The rich configurability of these instruction is summarized in Figure 10-1.

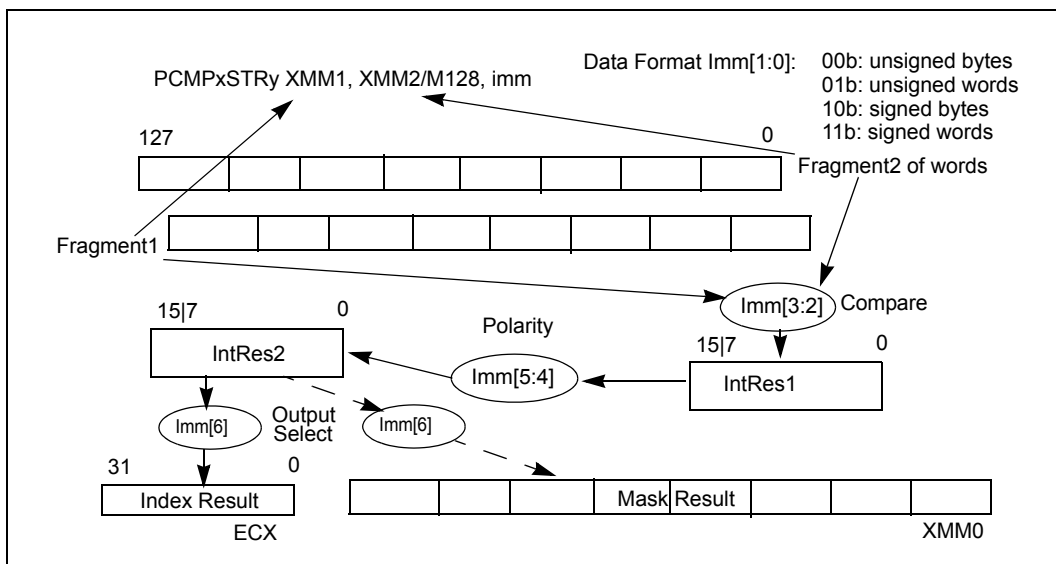


Figure 10-1. SSE4.2 String/Text Instruction Immediate Operand Control

The PCMPxSTRI instructions produce final result as an integer index in ECX, the PCMPxSTRM instructions produce final result as a bit mask in the XMM0 register. The PCMPISTRy instructions support processing string/text fragments using implicit length control via null termination for handling string/text of unknown size. the PCMPESTRy instructions support explicit length control via EDX:EAX register pair to specify the length text fragments in the source operands.

The first intermediate result, IntRes1, is an aggregated result of bit patterns from parallel comparison operations done on pairs of data elements from each text fragment, according to the imm[3:2] bit field encoding, see Table 10-1.

Table 10-1. SSE4.2 String/Text Instructions Compare Operation on N-elements

Imm[3:2]	Name	IntRes1[i] is TRUE if	Potential Usage
00B	Equal Any	Element i in fragment2 matches any element j in fragment1	Tokenization, XML parser
01B	Ranges	Element i in fragment2 is within any range pairs specified in fragment1	Subsetting, Case handling, XML parser, Schema validation
10B	Equal Each	Element i in fragment2 matches element i in fragment1	Strcmp()
11B	Equal Ordered	Element i and subsequent, consecutive valid elements in fragment2 match fully or partially with fragment1 starting from element 0	Substring Searches, KMP, Strstr()

Input data element format selection using imm[1:0] can support signed or unsigned byte/word elements.

The bit field imm[5:4] allows applying a unary transformation on IntRes1, see Table 10-2.

Table 10-2. SSE4.2 String/Text Instructions Unary Transformation on IntRes1

Imm[5:4]	Name	IntRes2[i] =	Potential Usage
00B	No Change	IntRes1[i]	
01B	Invert	-IntRes1[i]	
10B	No Change	IntRes1[i]	
11B	Mask Negative	IntRes1[i] if element i of fragment2 is invalid, otherwise -IntRes1[i]	

The output selection field, imm[6] is described in Table 10-3.

Table 10-3. SSE4.2 String/Text Instructions Output Selection Imm[6]

Imm[6]	Instruction	Final Result	Potential Usage
0B	PCMPxSTRl	ECX = offset of least significant bit set in IntRes2 if IntRes2 != 0, otherwise ECX = number of data element per 16 bytes	
0B	PCMPxSTRM	XMM0 = ZeroExtend(IntRes2);	
1B	PCMPxSTRl	ECX = offset of most significant bit set in IntRes2 if IntRes2 != 0, otherwise ECX = number of data element per 16 bytes	
1B	PCMPxSTRM	Data element i of XMM0 = SignExtend(IntRes2[i]);	

The comparison operation on each data element pair is defined in Table 10-4. Table 10-4 defines the type of comparison operation between valid data elements (last row of Table 10-4) and boundary conditions when the fragment in a source operand may contain invalid data elements (rows 1 through 3 of Table 10-4). Arithmetic comparison are performed only if both data elements are valid element in fragment1 and fragment2, as shown in row 4 of Table 10-4.

Table 10-4. SSE4.2 String/Text Instructions Element-Pair Comparison Definition

fragment1 element	fragment2 element	Imm[3:2]= 00B, Equal Any	Imm[3:2]= 01B, Ranges	Imm[3:2]= 10B, Equal Each	Imm[3:2]= 11B, Equal Ordered
invalid	invalid	Force False	Force False	Force True	Force True
invalid	valid	Force False	Force False	Force False	Force True
valid	invalid	Force False	Force False	Force False	Force False
valid	valid	Compare	Compare	Compare	Compare

The string and text processing instruction provides several aid to handle end-of-string situations, see Table 10-5. Additionally, the PCMPxSTRy instructions are designed to not require 16-byte alignment to simplify text processing requirements.

Table 10-5. SSE4.2 String/Text Instructions Eflags Behavior

EFLAGS	Description	Potential Usage
CF	Reset if IntRes2 = 0; Otherwise set	When CF=0, ECX= #of data element to scan next
ZF	Reset if entire 16-byte fragment2 is valid	likely end-of-string
SF	Reset if entire 16-byte fragment1 is valid	
OF	IntRes2[0];	

10.1.1 CRC32

CRC32 instruction computes the 32-bit cyclic redundancy checksum signature for byte/word/dword or qword stream of data. It can also be used as a hash function. For examples, a dictionary uses hash indices to dereference strings. CRC32 instruction can be easily adapted for use in this situation.

Example 10-1 shows a straight forward hash function that can be used to evaluate the hash index of a string to populate a hash table. Typically, the hash index is derived from the hash value by taking the remainder of the hash value modulo the size of a hash table.

Example 10-1. A Hash Function Examples

```

unsigned int hash_str(unsigned char* pStr)
{  unsigned int hVal = (unsigned int)(*pStr++);
   while (*pStr)
   {  hVal = (hashVal * CONST_A) + (hVal >> 24) + (unsigned int)(*pStr++);
     }
   return hVal;
}

```

CRC32 instruction can be use to derive an alternate hash function. Example 10-2 takes advantage the 32-bit granular CRC32 instruction to update signature value of the input data stream. For string of small to moderate sizes, using the hardware accelerated CRC32 can be twice as fast as Example 10-1.

Example 10-2. Hash Function Using CRC32

```

static unsigned cn_7e = 0x7efefeff, Cn_81 = 0x81010100;

unsigned int hash_str_32_crc32x(unsigned char* pStr)
{ unsigned *pDW = (unsigned *) &pStr[1];
  unsigned short *pWd = (unsigned short *) &pStr[1];
  unsigned int tmp, hVal = (unsigned int)(*pStr);
  if( !pStr[1] );
  else {
    tmp = ((pDW[0] +cn_7e ) ^ (pDW[0]^ -1)) & Cn_81;
    while ( !tmp ) // loop until there is byte in *pDW had 0x00
    {
      hVal = _mm_crc32_u32 (hVal, *pDW ++);
      tmp = ((pDW[0] +cn_7e ) ^ (pDW[0]^ -1)) & Cn_81;
    };
    if(!pDW[0]);
    else if(pDW[0] < 0x100) { // finish last byte that's non-zero
      hVal = _mm_crc32_u8 (hVal, pDW[0]);
    }
    else if(pDW[0] < 0x10000) { // finish last two byte that's non-zero
      hVal = _mm_crc32_u16 (hVal, pDW[0]);
    }
    else { // finish last three byte that's non-zero
      hVal = _mm_crc32_u32 (hVal, pDW[0]);
    }
  }
  return hVal;
}

```

10.2 USING SSE4.2 STRING AND TEXT INSTRUCTIONS

String libraries provided by high-level languages or as part of system library are used in a wide range of situations across applications and privileged system software. These situations can be accelerated using a replacement string library that implements PCMPSTR/PCMPSTRM/PCMPSTRSI/PCMPSTRM.

Although system-provided string library provides standardized string handling functionality and interfaces, most situations dealing with structured document processing requires considerable more sophistication, optimization, and services not available from system-provided string libraries. For example, structured document processing software often architect different class objects to provide building block functionality to service specific needs of the application. Often application may choose to disperse equivalent string library services into separate classes (string, lexer, parser) or integrate memory management capability into string handling/lexing/parsing objects.

PCMPSTR/PCMPSTRM/PCMPSTRSI/PCMPSTRM instructions are general-purpose primitives that software can use to build replacement string libraries or build class hierarchy to provide lexing/parsing services for structured document processing. XML parsing and schema validation are examples of the latter situations.

Unstructured, raw text/string data consist of characters, and have no natural alignment preferences. Therefore, PCMPSTR/PCMPSTRM/PCMPSTRSI/PCMPSTRM instructions are architected to not require the 16-Byte alignment restrictions of other 128-bit SIMD integer vector processing instructions.

With respect to memory alignment, PCMPSTR/PCMPSTRM/PCMPSTRSI/PCMPSTRM support unaligned memory loads like other unaligned 128-bit memory access instructions, e.g. MOVDQU.

Unaligned memory accesses may encounter special situations that require additional coding techniques, depending on the code running in ring 3 application space or in privileged space. Specifically, an unaligned 16-byte load may cross page boundary. Section 10.2.1 discusses a technique that application code can use. Section 10.2.2 discusses the situation string library functions needs to deal with. Section 10.3 gives detailed examples of using PCMPSTR/PCMPSTRM/PCMPSTRSI/PCMPSTRM instructions to implement equivalent functionality of several string library functions in situations that application code has control over memory buffer allocation.

10.2.1 Unaligned Memory Access and Buffer Size Management

In application code, the size requirements for memory buffer allocation should consider unaligned SIMD memory semantics and application usage.

For certain types of application usage, it may be desirable to make distinctions between valid buffer range limit versus valid application data size (e.g. a video frame). The former must be greater or equal to the latter.

To support algorithms requiring unaligned 128-bit SIMD memory accesses, memory buffer allocation by a caller function should consider adding some pad space so that

a callee function can safely use the address pointer safely with unaligned 128-bit SIMD memory operations.

The minimal padding size should be the width of the SIMD register that might be used in conjunction with unaligned SIMD memory access.

10.2.2 Unaligned Memory Access and String Library

String library functions may be used by application code or privileged code. String library functions must be careful not to violate memory access rights. Therefore, a replacement string library that employ SIMD unaligned access must employ special techniques to ensure no memory access violation occur.

Section 10.3.6 provides an example of a replacement string library function implemented with SSE4.2 and demonstrates a technique to use 128-bit unaligned memory access without unintentionally crossing page boundary.

10.3 SSE4.2 APPLICATION CODING GUIDELINE AND EXAMPLES

Software implementing SSE4.2 instruction must use CPUID feature flag mechanism to verify processor's support for SSE4.2. Details can be found in CHAPTER 12 of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 1* and in CPUID of CHAPTER 3 in *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A*.

In the following sections, we use several examples in string/text processing of progressive complexity to illustrates the basic techniques of adapting the SIMD approach to implement string/text processing using PCMPxSTRy instructions in SSE4.2. For simplicity, we will consider string/text in byte data format in situations that caller functions have allocated sufficient buffer size to support unaligned 128-bit SIMD loads from memory without encountering side-effects of cross page boundaries.

10.3.1 Null Character Identification (Strlen equivalent)

The most widely used string function is probably `strlen()`. One can view the lexing requirement of `strlen()` is to identify the null character in a text block of unknown size (end of string condition). Brute-force, byte-granular implementation fetches data inefficiently by loading one byte at a time.

Optimized implementation using general-purpose instructions can take advantage of dword operations in 32-bit environment (and qword operations in 64-bit environment) to reduce the number of iterations.

A 32-bit assembly implementation of `strlen()` is shown Example 10-3. The peak execution throughput of handling EOS condition is determined by eight ALU instructions in the main loop.

Example 10-3. `Strlen()` Using General-Purpose Instructions

```
int strlen_asm(const char* s1)
{int len = 0;
  _asm{
    mov ecx, s1
    test ecx, 3 ; test addr aligned to dword
    je short _main_loop1 ; dword aligned loads would be faster
  _malign_str1:
    mov al, byte ptr [ecx] ; read one byte at a time
    add ecx, 1
    test al, al ; if we find a null, go calculate the length
    je short _byte3a
    test ecx, 3; test if addr is now aligned to dword
    jne short _malign_str1; if not, repeat
  align16
  _main_loop1; read each 4-byte block and check for a NULL char in the dword
    mov eax, [ecx]; read 4 byte to reduce loop count
    mov edx, 7efefeffh
    add edx, eax
    xor eax, -1
    xor eax, edx
    add ecx, 4; increment address pointer by 4
    test eax, 81010100h ; if no null code in 4-byte stream, do the next 4 bytes
    je short _main_loop1
    ; there is a null char in the dword we just read,
    ; since we already advanced pointer ecx by 4, and the dword is lost
    mov eax, [ecx -4]; re-read the dword that contain at least a null char
    test al, al ; if byte0 is null
    je short _byte0a; the least significant byte is null
    test ah, ah ; if byte1 is null
    je short _byte1a
    test eax, 00ff0000h; if byte2 is null
    (continue)
```

Example 10-3. Strlen() Using General-Purpose Instructions

```

je    short_byte2a
test  eax, 00ff000000h; if byte3 is null
je    short_byte3a
jmp   short_main_loop1
_byte3a:
; we already found the null, but pointer already advanced by 1
lea   eax, [ecx-1]; load effective address corresponding to null code
mov   ecx, s1
sub   eax, ecx; difference between null code and start address
jmp   short_resulta
_byte2a:

lea   eax, [ecx-2]
mov   ecx, s1
sub   eax, ecx
jmp   short_resulta
_byte1a:
lea   eax, [ecx-3]
mov   ecx, s1
sub   eax, ecx
jmp   short_resulta
_byte0a:
lea   eax, [ecx-4]
mov   ecx, s1
sub   eax, ecx
_resulta:
mov   len, eax; store result
}
return len;
}

```

The equivalent functionality of EOS identification can be implemented using PCMP-ISTRI. Example 10-4 shows a simplistic SSE4.2 implementation to scan a text block by loading 16-byte text fragments and locate the null termination character. Example 10-5 shows the optimized SSE4.2 implementation that demonstrates the importance of using memory disambiguation to improve instruction-level parallelism.

Example 10-4. Sub-optimal PCMPISTRI Implementation of EOS handling

```

static char sscch2[16]= {0x1, 0xff, 0x00, }; // range values for non-null characters

int strlen_un_optimized(const char* s1)
{int len = 0;
  _asm{
    mov  eax, s1
    movdquxmm2, sscch2 ; load character pair as range (0x01 to 0xff)
    xor  ecx, ecx ; initial offset to 0
        (continue)

_loopc:
    add  eax, ecx ; update addr pointer to start of text fragment
    pcmpestri xmm2, [eax], 14h; unsigned bytes, ranges, invert, lsb index returned to ecx
        ; if there is a null char in the 16Byte fragment at [eax], zf will be set.
        ; if all 16 bytes of the fragment are non-null characters, ECX will return 16,
    jnz  short _loopc; xmm1 has no null code, ecx has 16, continue search
        ; we have a null code in xmm1, ecx has the offset of the null code i
    add  eax, ecx ; add ecx to the address of the last fragment2/xmm1
    mov  edx, s1; retrieve effective address of the input string
    sub  eax, edx; the string length
    mov  len, eax; store result
  }
  return len;
}

```

The code sequence shown in Example 10-4 has a loop consisting of three instructions. From a performance tuning perspective, the loop iteration has loop-carry dependency because address update is done using the result (ECX value) of a previous loop iteration. This loop-carry dependency deprives the out-of-order engine's capability to have multiple iterations of the instruction sequence making forward progress. The latency of memory loads, the latency of these instructions, any bypass delay could not be amortized by OOO execution in the presence of loop-carry dependency.

A simple optimization technique to eliminate loop-carry dependency is shown in Example 10-5.

Using memory disambiguation technique to eliminate loop-carry dependency, the cumulative latency exposure of the 3-instruction sequence of Example 10-5 is amortized over multiple iterations, the net cost of executing each iteration (handling 16 bytes) is less than 3 cycles. In contrast, handling 4 bytes of string data using 8 ALU instructions in Example 10-3 will also take a little less than 3 cycles per iteration.

Whereas each iteration of the code sequence in Example 10-4 will take more than 10 cycles because of loop-carry dependency.

Example 10-5. Strlen() Using PCMPISTRI without Loop-Carry Dependency

```
int strlen_sse4_2(const char* s1)
{int len = 0;
  _asm{
    mov  eax, s1
    movdquxmm2, ssch2 ; load character pair as range (0x01 to 0xff)
    xor  ecx, ecx ; initial offset to 0
    sub  eax, 16 ; address arithmetic to eliminate extra instruction and a branch
  _loopc:
    add  eax, 16 ; adjust address pointer and disambiguate load address for each iteration
    pcmpistri xmm2, [eax], 14h; unsigned bytes, ranges, invert, lsb index returned to ecx
      ; if there is a null char in [eax] fragment, zf will be set.
      ; if all 16 bytes of the fragment are non-null characters, ECX will return 16,
    jnz short _loopc ; ECX will be 16 if there is no null byte in [eax], so we disambiguate
  _endofstring:
    add  eax, ecx ; add ecx to the address of the last fragment
    mov  edx, s1; retrieve effective address of the input string
    sub  eax, edx;the string length
    mov  len, eax; store result
  }
  return len;
}
```

SSE4.2 Coding Rule 5. (H impact, H generality) Loop-carry dependency that depends on the ECX result of PCMPSTRI/PCMPSTRM/PCMPISTRI/PCMPISTRM for address adjustment must be minimized. Isolate code paths that expect ECX result will be 16 (bytes) or 8 (words), replace these values of ECX with constants in address adjustment expressions to take advantage of memory disambiguation hardware.

10.3.2 White-Space-Like Character Identification

Character-granular-based text processing algorithms have developed techniques to handle specific tasks to remedy the efficiency issue of character-granular approaches. One such technique is using look-up tables for character subset classification. For example, some application may need to separate alpha-numeric characters from white-space-like characters. More than one character may be treated as white-space characters.

Example 10-6 illustrates a simple situation of identifying white-space-like characters for the purpose of marking the beginning and end of consecutive non-white-space characters.

Example 10-6. WordCnt() Using C and Byte-Scanning Technique

```
// Counting words involves locating the boundary of contiguous non-whitespace characters.
// Different software may choose its own mapping of white space character set.
// This example employs a simple definition for tutorial purpose:
// Non-whitespace character set will consider: A-Z, a-z, 0-9, and the apostrophe mark '
// The example uses a simple technique to map characters into bit patterns of square waves
// we can simply count the number of falling edges

static char alphnrange[16]= {0x27, 0x27, 0x30, 0x39, 0x41, 0x5a, 0x61, 0x7a, 0x0};
static char alp_map8[32] = {0x0, 0x0, 0x0, 0x0, 0x80,0x0,0xff, 0x3,0xfe, 0xff, 0xff, 0x7, 0xfe,
0xff, 0xff, 0x7}; // 32 byte lookup table, 1s map to bit patterns of alpha numerics in alphnrange
int wordcnt_c(const char* s1)
{int i, j, cnt = 0;
char cc, cc2;
char flg[3]; // capture the a wavelet to locate a falling edge
cc2 = cc = s1[0];
// use the compacted bit pattern to consolidate multiple comparisons into one look up
if( alp_map8[cc>>3] & ( 1<<( cc & 7) ) )
{ flg[1] = 1; } // non-white-space char that is part of a word,
(continue)
```

Example 10-6. WordCnt() Using C and Byte-Scanning Technique

```

// we're including apostrophe in this example since counting the
// following 's' as a separate word would be kind of silly
else
{  flg[1] = 0; } // 0: whitespace, punctuations not be considered as part of a word

i = 1; // now we're ready to scan through the rest of the block
// we'll try to pick out each falling edge of the bit pattern to increment word count.
// this works with consecutive white spaces, dealing with punctuation marks, and
// treating hyphens as connecting two separate words.
while (cc2 )
{  cc2 = s1[i];
  if( alp_map8[cc2>>3] & ( 1<<( cc2 & 7) ) )
  {  flg[2] = 1; } // non-white-space
  else
  {  flg[2] = 0; } // white-space-like

  if( !flg[2] && flg[1] )
  {  cnt ++; } // found the falling edge
  flg[1] = flg[2];
  i++;
}
return cnt;
}

```

In Example 10-6, a 32-byte look-up table is constructed to represent the ascii code values 0x0-0xff, and partitioned with each bit of 1 corresponding to the specified subset of characters. While this bit-lookup technique simplifies the comparison operations, data fetching remains byte-granular.

Example 10-7 shows an equivalent implementation of counting words using PCMP-ISTRM. The loop iteration is performed at 16-byte granularity instead of byte granularity. Additionally, character set subsetting is easily expressed using range value pairs and parallel comparisons between the range values and each byte in the text fragment are performed by executing PCMPISTRI once.

Example 10-7. WordCnt() Using PCMPISTRM

```

// an SSE 4.2 example of counting words using the definition of non-whitespace character
// set of {A-Z, a-z, 0-9, '}, Each text fragment (up to 16 bytes) are mapped to a
// 16-bit pattern, which may contain one or more falling edges. Scanning bit-by-bit
// would be inefficient and goes counter to leveraging SIMD programming techniques.
// Since each falling edge must have a preceding rising edge, we take a finite
// difference approach to derive a pattern where each rising/falling edge maps to 2-bit pulse,
// count the number of bits in the 2-bit pulses using popcnt and divide by two.
int wdcnt_sse4_2(const char* s1)
{int len = 0;
  _asm{
    mov  eax, s1
    movdquxmm3, alphanrange ; load range value pairs to detect non-white-space codes
    xor  ecx, ecx
    xor  esi, esi
    xor  edx, edx
    movdquxmm1, [eax]
    pcmplstrm xmm3, xmm1, 04h ; white-space-like char becomes 0 in xmm0[15:0]
    movdqa xmm4, xmm0
    movdqa  xmm1, xmm0
    psrlq xmm4, 15 ; save MSB to use in next iteration
    movdqa xmm5, xmm1
    psllw xmm5, 1; lsb is effectively mapped to a white space
    pxor  xmm5, xmm0; the first edge is due to the artifact above
    pextrd edi, xmm5, 0
    jz   _lastfragment; if xmm1 had a null, zf would be set
    popcnt edi, edi; the first fragment will include a rising edge
    add  esi, edi
    mov  ecx, 16
        (continue)

_loopc:
    add  eax, ecx ; advance address pointer
    movdquxmm1, [eax]
    pcmplstrm xmm3, xmm1, 04h ; white-space-like char becomes 0 in xmm0[15:0]
    movdqa xmm5, xmm4 ; retrieve the MSB of the mask from last iteration
    movdqa xmm4, xmm0
    psrlq xmm4, 15 ; save mSB of this iteration for use in next iteration
    movdqa xmm1, xmm0

```

Example 10-7. WordCnt() Using PCMPISTRM

```

psllw xmm1, 1
por  xmm5, xmm1 ; combine MSB of last iter and the rest from current iter
pxor  xmm5, xmm0; differentiate binary wave form into pattern of edges
pextrdedi, xmm5, 0; the edge patterns has (1 bit from last, 15 bits from this round)
jz   _lastfragment; if xmm1 had a null, zf would be set
mov  ecx, 16; xmm1, had no null char, advance 16 bytes
popcntedi, edi; count both rising and trailing edges
add  esi, edi; keep a running count of both edges
jmp  short _loopc
_lastfragment:
popcntedi, edi; count both rising and trailing edges
add  esi, edi; keep a running count of both edges
shr  esi, 1; word count corresponds to the trailing edges
mov  len, esi
}
return len;
}

```

10.3.3 Substring Searches

Strstr() is a common function in the standard string library. Typically, A library may implement strstr(sTarg, sRef) with a brute-force, byte-granular technique of iterative comparisons between the reference string with a round of string comparison with a subset of the target string. Brute-force, byte-granular techniques provide reasonable efficiency when the first character of the target substring and the reference string are different, allowing subsequent string comparisons of target substrings to proceed forward to the next byte in the target string.

When a string comparison encounters partial matches of several characters (i.e. the sub-string search found a partial match starting from the beginning of the reference string) and determined the partial match led to a false-match. The brute-force search process need to go backward and restart string comparisons from a location that had participated in previous string comparison operations. This is referred to as re-trace inefficiency of the brute-force substring search algorithm. See Figure 10-2.

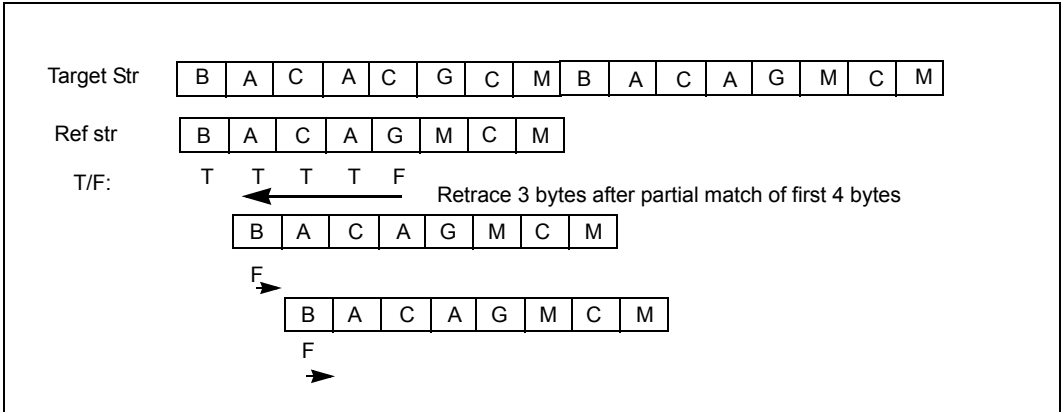


Figure 10-2. Retrace Inefficiency of Byte-Granular, Brute-Force Search

The Knuth, Morris, Pratt algorithm¹ (KMP) provides an elegant enhancement to overcome the re-trace inefficiency of brute-force substring searches. By deriving an overlap table that is used to manage retrace distance when a partial match leads to a false match, KMP algorithm is very useful for applications that search relevant articles containing keywords from a large corpus of documents.

Example 10-8 illustrates a C-code example of using KMP substring searches.

Example 10-8. KMP Substring Search in C

```
// s1 is the target string of length cnt1
// s2 is the reference string of length cnt2
// j is the offset in target string s1 to start each round of string comparison
// i is the offset in reference string s2 to perform byte granular comparison
(continue)
```

1. Donald E. Knuth, James H. Morris, and Vaughan R. Pratt; SIAM J. Comput. Volume 6, Issue 2, pp. 323-350 (1977)

Example 10-8. KMP Substring Search in C

```

int str_kmp_c(const char* s1, int cnt1, const char* s2, int cnt2)
{ int i, j;
  i = 0; j = 0;
  while ( i+j < cnt1) {
    if( s2[i] == s1[i+j]) {
      i++;
      if( i == cnt2) break; // found full match
    }
    else {
      j = j+i - overlap_tbl[i]; // update the offset in s1 to start next round of string compare
      if( i > 0) {
        i = overlap_tbl[i]; // update the offset of s2 for next string compare should start at
      }
    }
  }
};
return j;
}

void kmp_precalc(const char * s2, int cnt2)
{int i = 2;
char nch = 0;
overlap_tbl[0] = -1; overlap_tbl[1] = 0;
// pre-calculate KMP table
while(i < cnt2) {
  if( s2[i-1] == s2[nch] ) {
    overlap_tbl[i] = nch + 1;
    i++; nch++;
  }
  else if ( nch > 0) nch = overlap_tbl[nch];
  else {
    overlap_tbl[i] = 0;
    i++;
  }
};
overlap_tbl[cnt2] = 0;
}

```

Example 10-8 also includes the calculation of the KMP overlap table. Typical usage of KMP algorithm involves multiple invocation of the same reference string, so the overhead of precalculating the overlap table is easily amortized. When a false match is determined at offset *i* of the reference string, the overlap table will predict where the

next round of string comparison should start (updating the offset *j*), and the offset in the reference string that byte-granular character comparison should resume/restart.

While KMP algorithm provides efficiency improvement over brute-force byte-granular substring search, its best performance is still limited by the number of byte-granular operations. To demonstrate the versatility and built-in lexical capability of PCMP-ISTR1, we show an SSE4.2 implementation of substring search using brute-force 16-byte granular approach in Example 10-9, and combining KMP overlap table with substring search using PCMPISTR1 in Example 10-10.

Example 10-9. Brute-Force Substring Search Using PCMPISTR1 Intrinsic

```
int strsubs_sse4_2i(const char* s1, int cnt1, const char* s2, int cnt2 )
{ int kpm_i=0, idx;
  int ln1= 16, ln2=16, rcnt1 = cnt1, rcnt2= cnt2;
  __m128i *p1 = (__m128i *) s1;
  __m128i *p2 = (__m128i *) s2;
  __m128ifrag1, frag2;
  int cmp, cmp2, cmp_s;
  __m128i *pt = NULL;
  if( cnt2 > cnt1 || !cnt1) return -1;
  frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
  frag2 = _mm_loadu_si128(p2); // load up to 16 bytes of fragment

  while(rcnt1 > 0)
  {  cmp_s = _mm_cmpestrs(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1, 0x0c);
    cmp = _mm_cmpestri(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1, 0x0c);
    if( !cmp) { // we have a partial match that needs further analysis
      if( cmp_s) { // if we're done with s2
        if( pt)
          {idx = (int) ((char *) pt - (char *) s1); }
        else
          {idx = (int) ((char *) p1 - (char *) s1); }
        return idx;
      }
    }
    (continue)
  }
}
```

Example 10-9. Brute-Force Substring Search Using PCMPISTRI Intrinsic

```

// we do a round of string compare to verify full match till end of s2
if( pt == NULL) pt = p1;
cmp2 = 16;
rcnt2 = cnt2 - 16 -(int) ((char *)p2-(char *)s2);
while( cmp2 == 16 && rcnt2) { // each 16B frag matches,
    rcnt1 = cnt1 - 16 -(int) ((char *)p1-(char *)s1);
    rcnt2 = cnt2 - 16 -(int) ((char *)p2-(char *)s2);
    if( rcnt1 <=0 || rcnt2 <= 0 ) break;
    p1 = (__m128i *)(((char *)p1) + 16);
    p2 = (__m128i *)(((char *)p2) + 16);
    frag1 = _mm_loadu_si128(p1);// load up to 16 bytes of fragment
    frag2 = _mm_loadu_si128(p2);// load up to 16 bytes of fragment
    cmp2 = _mm_cmpestri(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1,
0x18); // lsb, eq each
};
if( !rcnt2 || rcnt2 == cmp2) {
    idx = (int) ((char *) pt - (char *) s1);
    return idx;
}
else if ( rcnt1 <= 0) { // also cmp2 < 16, non match
    if( cmp2 == 16 && ((rcnt1 + 16) >= (rcnt2+16) ) )
    {idx = (int) ((char *) pt - (char *) s1);
    return idx;
}
else return -1;
}
(continue)

```

Example 10-9. Brute-Force Substring Search Using PCMPISTRI Intrinsic

```

else { // in brute force, we advance fragment offset in target string s1 by 1
    p1 = (__m128i *)(((char *)pt) + 1); // we're not taking advantage of kmp
    rcnt1 = cnt1 -(int) ((char *)p1-(char *)s1);
    pt = NULL;
    p2 = (__m128i *)((char *)s2);
    rcnt2 = cnt2 -(int) ((char *)p2-(char *)s2);
    frag1 = _mm_loadu_si128(p1); // load next fragment from s1
    frag2 = _mm_loadu_si128(p2); // load up to 16 bytes of fragment
}
}
else{
    if( cmp == 16) p1 = (__m128i *)(((char *)p1) + 16);
    else p1 = (__m128i *)(((char *)p1) + cmp);
    rcnt1 = cnt1 -(int) ((char *)p1-(char *)s1);
    if( pt && cmp ) pt = NULL;
    frag1 = _mm_loadu_si128(p1); // load next fragment from s1
}
}
return idx;
}

```

In Example 10-9, address adjustment using a constant to minimize loop-carry dependency is practised in two places:

- In the inner while loop of string comparison to determine full match or false match (the result `cmp2` is not used for address adjustment to avoid dependency).
- In the last code block when the outer loop executed PCMPISTRI to compare 16 sets of ordered compare between a target fragment with the first 16-byte fragment of the reference string, and all 16 ordered compare operations produced false result (producing `cmp` with a value of 16).

Example 10-10 shows an equivalent intrinsic implementation of substring search using SSE4.2 and KMP overlap table. When the inner loop of string comparison determines a false match, the KMP overlap table is consulted to determine the address offset for the target string fragment and the reference string fragment to minimize retrace.

It should be noted that a significant portions of retrace with retrace distance less than 15 bytes are avoided even in the brute-force SSE4.2 implementation of Example 10-9. This is due to the order-compare primitive of PCMPISTRI. "Ordered compare" performs 16 sets of string fragment compare, and many false match with less than 15 bytes of partial matches can be filtered out in the same iteration that executed PCMPISTRI.

Retrace distance of greater than 15 bytes does not get filtered out by the Example 10-9. By consulting with the KMP overlap table, Example 10-10 can eliminate retraces of greater than 15 bytes.

Example 10-10. Substring Search Using PCMPISTRI and KMP Overlap Table

```
int strkmp_sse4_2(const char* s1, int cnt1, const char* s2, int cnt2 )
{ int kpm_i=0, idx;
  int ln1= 16, ln2=16, rcnt1 = cnt1, rcnt2= cnt2;
  __m128i *p1 = (__m128i *) s1;
  __m128i *p2 = (__m128i *) s2;
  __m128ifrag1, frag2;
  int cmp, cmp2, cmp_s;
  __m128i *pt = NULL;
  if( cnt2 > cnt1 || !cnt1) return -1;
  frag1 = _mm_loadu_si128(p1);// load up to 16 bytes of fragment
  frag2 = _mm_loadu_si128(p2);// load up to 16 bytes of fragment

  while(rcnt1 > 0)
  {  cmp_s = _mm_cmpestrs(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1, 0x0c);
    cmp = _mm_cmpestri(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1, 0x0c);
    if( !cmp) { // we have a partial match that needs further analysis
      if( cmp_s) { // if we've reached the end with s2
        if( pt)
          {idx = (int) ((char *) pt - (char *) s1); }
        else
          {idx = (int) ((char *) p1 - (char *) s1); }
        return idx;
      }
    }
    // we do a round of string compare to verify full match till end of s2
    if( pt == NULL) pt = p1;
    cmp2 = 16;
    rcnt2 = cnt2 - 16 -(int) ((char *)p2-(char *)s2);

    (continue)
  }
}
```

Example 10-10. Substring Search Using PCMPISTRI and KMP Overlap Table

```

while( cmp2 == 16 && rcnt2) { // each 16B frag matches
    rcnt1 = cnt1 - 16 -(int) ((char *)p1-(char *)s1);
    rcnt2 = cnt2 - 16 -(int) ((char *)p2-(char *)s2);
    if( rcnt1 <= 0 || rcnt2 <= 0 ) break;
    p1 = (__m128i *)(((char *)p1) + 16);
    p2 = (__m128i *)(((char *)p2) + 16);
    frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
    frag2 = _mm_loadu_si128(p2); // load up to 16 bytes of fragment
    cmp2 = _mm_cmpestri(frag2, (rcnt2>ln2)? ln2: rcnt2, frag1, (rcnt1>ln1)? ln1: rcnt1,
0x18); // lsb, eq each
    };
    if( !rcnt2 || rcnt2 == cmp2) {
        idx = (int) ((char *) pt - (char *) s1);
        return idx;
    }
    else if ( rcnt1 <= 0 ) { // also cmp2 < 16, non match
        return -1;
    }
    else { // a partial match led to false match, consult KMP overlap table for addr adjustment
        kpm_i = (int) ((char *)p1 - (char *)pt)+ cmp2 ;
        p1 = (__m128i *)(((char *)pt) + (kpm_i - overlap_tbl[kpm_i])); // use kmp to skip retrace
        rcnt1 = cnt1 -(int) ((char *)p1-(char *)s1);
        pt = NULL;
        p2 = (__m128i *)(((char *)s2) + (overlap_tbl[kpm_i]));
        rcnt2 = cnt2 -(int) ((char *)p2-(char *)s2);
        frag1 = _mm_loadu_si128(p1); // load next fragment from s1
        frag2 = _mm_loadu_si128(p2); // load up to 16 bytes of fragment
    }
}
}

(continue)

```

Example 10-10. Substring Search Using PCMPISTRI and KMP Overlap Table

```

else{
    if( kpm_i && overlap_tbl[kpm_i] ) {
        p2 = (__m128i *)(((char *)s2) );
        frag2 = _mm_loadu_si128(p2);// load up to 16 bytes of fragment
        //p1 = (__m128i *)(((char *)p1) );

        //rcnt1 = cnt1 -(int) ((char *)p1)-(char *)s1);
        if( pt && cmp ) pt = NULL;
        rcnt2 = cnt2 ;
        //frag1 = _mm_loadu_si128(p1);// load next fragment from s1
        frag2 = _mm_loadu_si128(p2);// load up to 16 bytes of fragment
        kpm_i = 0;
    }
    else { // equ order comp resulted in sub-frag match or non-match
        if( cmp == 16) p1 = (__m128i *)(((char *)p1) + 16);
        else p1 = (__m128i *)(((char *)p1) + cmp);
        rcnt1 = cnt1 -(int) ((char *)p1)-(char *)s1);
        if( pt && cmp ) pt = NULL;
        frag1 = _mm_loadu_si128(p1);// load next fragment from s1
    }
}
}
return idx;
}

```

The relative speed up of byte-granular KMP, brute-force SSE4.2, and SSE4.2 with KMP overlap table over byte-granular brute-force substring search is illustrated in the graph that plots relative speedup over percentage of retrace for a reference string of 55 bytes long. A retrace of 40% in the graph meant, after a partial match of the first 22 characters, a false match is determined.

So when brute-force, byte-granular code has to retrace, the other three implementation may be able to avoid the need to retrace because:

- Example 10-8 can use KMP overlap table to predict the start offset of next round of string compare operation after a partial-match/false-match, but forward movement after a first-character-false-match is still byte-granular.

- Example 10-9 can avoid retrace of shorter than 15 bytes but will be subject to retrace of 21 bytes after a partial-match/false-match at byte 22 of the reference string. Forward movement after each order-compare-false-match is 16 byte granular.
- Example 10-10 avoids retrace of 21 bytes after a partial-match/false-match, but KMP overlap table lookup incurs some overhead. Forward movement after each order-compare-false-match is 16 byte granular.

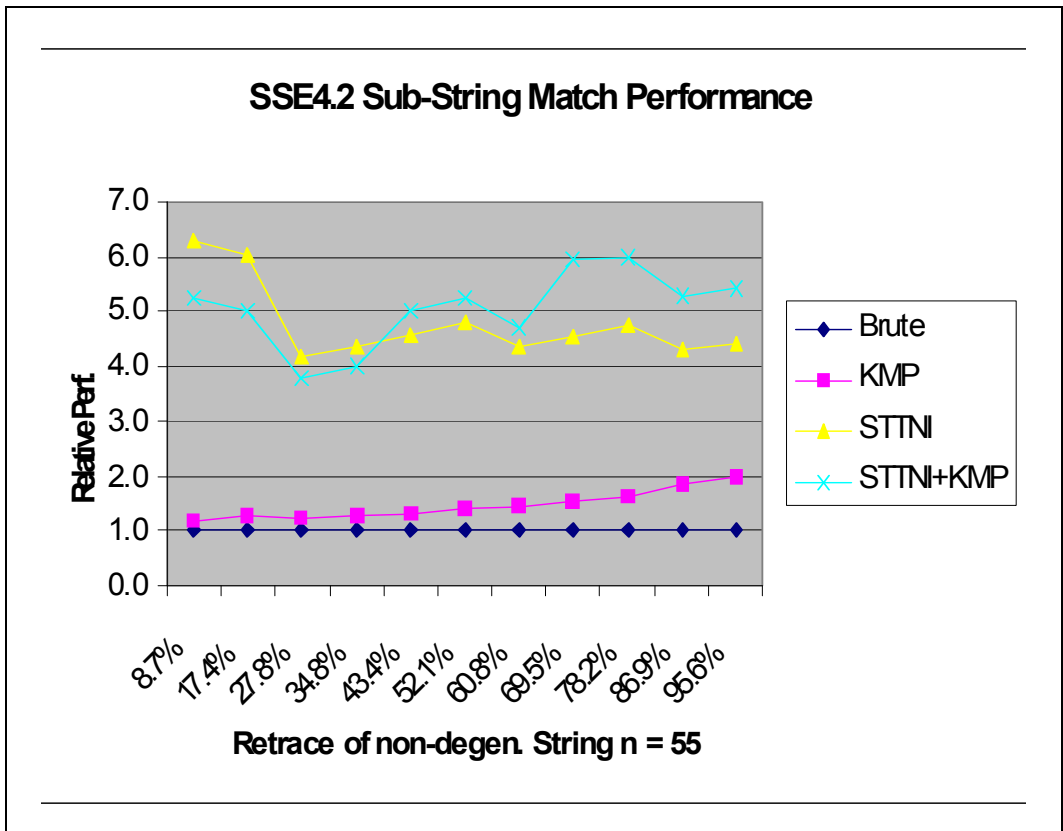


Figure 10-3. SSE4.2 Speedup of SubString Searches

10.3.4 String Token Extraction and Case Handling

Token extraction is a common task in text/string handling. It is one of the foundation of implementing lexer/parser objects of higher sophistication. Indexing services also build on tokenization primitives to sort text data from streams.

Tokenization requires the flexibility to use an array of delimiter characters.

A library implementation of `Strtok_s()` may employ a table-lookup technique to consolidate sequential comparisons of the delimiter characters into one comparison (similar to Example 10-6). An SSE4.2 implementation of the equivalent functionality of `strtok_s()` using intrinsic is shown in Example 10-11.

Example 10-11. | Equivalent `Strtok_s()` Using `PCMPISTRI` Intrinsic

```

char ws_map8[32]; // packed bit lookup table for delimiter characters

char * strtok_sse4_2i(char* s1, char *sdlm, char ** pCtxt)
{
    __m128i *p1 = (__m128i *) s1;
    __m128ifrag1, stmpz, stmp1;
    int cmp_z, jj =0;
    int start, endtok, s_idx, ldx;
    if (sdlm == NULL || pCtxt == NULL) return NULL;
    if (p1 == NULL && *pCtxt == NULL) return NULL;
    if (s1 == NULL) {
        if( *pCtxt[0] == 0 ) { return NULL; }
        p1 = (__m128i *) *pCtxt;
        s1 = *pCtxt;
    }
    else p1 = (__m128i *) s1;
    memset(&ws_map8[0], 0, 32);
    while (sdlm[jj] ) {
        ws_map8[ (sdlm[jj] >> 3) ] |= (1 << (sdlm[jj] & 7)); jj ++
    }
    frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
    stmpz = _mm_loadu_si128((__m128i *)sdelimiter);
    // if the first char is not a delimiter , proceed to check non-delimiter,
    // otherwise need to skip leading delimiter chars
    if( ws_map8[s1[0]>>3] & (1 << (s1[0]&7)) ) {
        start = s_idx = _mm_cmpistri(stmpz, frag1, 0x10); // unsigned bytes/equal any, invert, lsb
    }
    else start = s_idx = 0;
        (continue)
}

```


Example 10-11. | Equivalent Strtok_s() Using PCMPISTRI Intrinsic

```

// check if we're dealing with short input string less than 16 bytes
cmp_z = _mm_cmpistrz(stmpz, frag1, 0x10);
if( cmp_z) { // last fragment
    if( !start) {
        endtok = idx = _mm_cmpistri(stmpz, frag1, 0x00);
        if( endtok == 16) { // didn't find delimiter at the end, since it's null-terminated
            // find where is the null byte
            *pCtxt = s1 + 1 + _mm_cmpistri(frag1, frag1, 0x40);
            return s1;
        }
        else { // found a delimiter that ends this word
            s1[start+endtok] = 0;
            *pCtxt = s1+start+endtok+1;
        }
    }
    else {
        if( !s1[start]) {
            *pCtxt = s1 + start + 1;
            return NULL;
        }
        p1 = (__m128i *)(((char *)p1) + start);
        frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
        endtok = idx = _mm_cmpistri(stmpz, frag1, 0x00); // unsigned bytes/equal any, lsb
        if( endtok == 16) { // looking for delimiter, found none
            *pCtxt = (char *)p1 + 1 + _mm_cmpistri(frag1, frag1, 0x40);
            return s1+start;
        }
        else { // found delimiter before null byte
            s1[start+endtok] = 0;
            *pCtxt = s1+start+endtok+1;
        }
    }
}
}

    (continue)

```

Example 10-11. | Equivalent Strtok_s() Using PCMPISTRI Intrinsic

```

else
{ while (!cmp_z && s_idx == 16) {
    p1 = (__m128i *)(((char *)p1) + 16);
    frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
    s_idx = _mm_cmpistri(stmpz, frag1, 0x10); // unsigned bytes/equal any, invert, lsb
    cmp_z = _mm_cmpistrz(stmpz, frag1, 0x10);
}
if(s_idx != 16) start = ((char *) p1 - s1) + s_idx;
else { // corner case if we ran to the end looking for delimiter and never found a non-dilimiter
    *pCtxt = (char *)p1 + 1 + _mm_cmpistri(frag1, frag1, 0x40);
    return NULL;
}
if( !s1[start] ) { // in case a null byte follows delimiter chars
    *pCtxt = s1 + start + 1;
    return NULL;
}
// now proceed to find how many non-delimiters are there
p1 = (__m128i *)(((char *)p1) + s_idx);
frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
endtok = idx = _mm_cmpistri(stmpz, frag1, 0x00); // unsigned bytes/equal any, lsb
cmp_z = 0;
while (!cmp_z && idx == 16) {
    p1 = (__m128i *)(((char *)p1) + 16);
    frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
    idx = _mm_cmpistri(stmpz, frag1, 0x00); // unsigned bytes/equal any, lsb
    cmp_z = _mm_cmpistrz(stmpz, frag1, 0x00);
    if(cmp_z) { endtok += idx; }
}
    (continue)
}

if( cmp_z ) { // reached the end of s1
    if( idx < 16 ) // end of word found by finding a delimiter
        endtok += idx;
    else { // end of word found by finding the null
        if( s1[start+endtok] ) // ensure this frag don't start with null byte
            endtok += 1 + _mm_cmpistri(frag1, frag1, 0x40);
    }
}
*pCtxt = s1 + start + endtok + 1;
s1[start+endtok] = 0;
}
return (char *) (s1 + start);
}

```

An SSE4.2 implementation of the equivalent functionality of `strupr()` using intrinsic is shown in Example 10-12.

Example 10-12. Equivalent `strupr()` Using `PCMPISTRM` Intrinsic

```
static char uldelta[16]= {0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20, 0x20};
static char ranglc[6]= {0x61, 0x7a, 0x00, 0x00, 0x00, 0x00};
char * strup_sse4_2i( char* s1)
{int len = 0, res = 0;
 __m128i *p1 = (__m128i *) s1;
 __m128ifrag1, ranglo, rmask, stmpz, stmp1;
 int cmp_c, cmp_z, cmp_s;
 if( !s1[0]) return (char *) s1;
 frag1 = _mm_loadu_si128(p1);// load up to 16 bytes of fragment
 ranglo = _mm_loadu_si128((__m128i *)ranglc);// load up to 16 bytes of fragment
 stmpz = _mm_loadu_si128((__m128i *)uldelta);
 (continue)
```

Example 10-12. I Equivalent Strupr() Using PCMPISTRM Intrinsic

```

cmp_z = _mm_cmpistrz(ranglo, frag1, 0x44); // range compare, produce byte masks
while (!cmp_z)
{
    rmsk = _mm_cmpistrm(ranglo, frag1, 0x44); // producing byte mask
    stmp1 = _mm_blendv_epi8(stmpz, frag1, rmsk); // bytes of lc preserved, other bytes
replaced by const
    stmp1 = _mm_sub_epi8(stmp1, stmpz); // bytes of lc becomes uc, other bytes are now zero
    stmp1 = _mm_blendv_epi8(frag1, stmp1, rmsk); //bytes of lc replaced by uc, other bytes
unchanged
    _mm_storeu_si128(p1, stmp1); //
    p1 = (__m128i *)((char *)p1 + 16);
    frag1 = _mm_loadu_si128(p1); // load up to 16 bytes of fragment
    cmp_z = _mm_cmpistrz(ranglo, frag1, 0x44);
}
if (*(char *)p1 == 0) return (char *)s1;
rmsk = _mm_cmpistrm(ranglo, frag1, 0x44); // byte mask, valid lc bytes are 1, all other 0
stmp1 = _mm_blendv_epi8(stmpz, frag1, rmsk); // bytes of lc continue, other bytes replaced
by const
stmp1 = _mm_sub_epi8(stmp1, stmpz); // bytes of lc becomes uc, other bytes are now zero
stmp1 = _mm_blendv_epi8(frag1, stmp1, rmsk); //bytes of lc replaced by uc, other bytes
unchanged
rmsk = _mm_cmpistrm(frag1, frag1, 0x44); // byte mask, valid bytes are 1, invalid bytes are
zero
_mm_maskmoveu_si128(stmp1, rmsk, (char *)p1); //
return (char *)s1;
}

```

10.3.5 Unicode Processing and PCMPxSTRy

Unicode representation of string/text data is required for software localization. UTF-16 is a common encoding scheme for localized content. In UTF-16 representation, each character is represented by a code point. There are two classes of code points: 16-bit code points and 32-bit code points which consists of a pair of 16-bit code points in specified value range, the latter is also referred to as a surrogate pair.

A common technique in unicode processing uses a table-loop up method, which has the benefit of reduced branching. As a tutorial example we compare the analogous problem of determining properly-encoded UTF-16 string length using general purpose code with table-lookup vs. SSE4.2.

Example 10-13 lists the C code sequence to determine the number of properly-encoded UTF-16 code points (either 16-bit or 32-bit code points) in a unicode text block. The code also verifies if there are any improperly-encoded surrogate pairs in the text block.

Example 10-13. UTF16 VerStrlen() Using C and Table Lookup Technique

```
// This example demonstrates validation of surrogate pairs (32-bit code point) and
// tally the number of 16-bit and 32-bit code points in the text block
// Parameters: s1 is pointer to input utf-16 text block.
// pLen: store count of utf-16 code points
// return the number of 16-bit code point encoded in the surrogate range but do not form
// a properly encoded surrogate pair. if 0: s1 is a properly encoded utf-16 block,
// If return value >0 then s1 contains invalid encoding of surrogates

int u16vstrlen_c(const short* s1, unsigned * pLen)
{int i, j, cnt = 0, cnt_invl = 0, spcnt = 0;
 unsigned short cc, cc2;
 char flg[3];

 cc2 = cc = s1[0];
 // map each word in s1 into bit patterns of 0, 1 or 2 using a table lookup
 // the first half of a surrogate pair must be encoded between D800-DBFF and mapped as 2
 // the 2nd half of a surrogate pair must be encoded between DC00-DFFF and mapped as 1
 // regular 16-bit encodings are mapped to 0, except null code mapped to 3
 flg[1] = utf16map[cc];
 flg[0] = flg[1];
 if(!flg[1]) cnt ++;
 i = 1;
```

(continue)

Example 10-13. UTF16 VerStrlen() Using C and Table Lookup Technique

```

while (cc2) // examine each non-null word encoding
{ cc2 = s1[i];
  flg[2] = utf16map[cc2];
  if( (flg[1] && flg[2] && (flg[1]-flg[2] == 1)) )
  { spcnt ++; } // found a surrogate pair
  else if(flg[1] == 2 && flg[2] != 1)
  { cnt_invl += 1; } // orphaned 1st half
  else if(!flg[1] && flg[2] == 1)
  { cnt_invl += 1; } // orphaned 2nd half
  else
  { if(!flg[2]) cnt ++; // regular non-null code16-bit code point
    else ;
  }
  flg[0] = flg[1]; // save the pair sequence for next iteration
  flg[1] = flg[2];
  i++;
}
*pLen = cnt + spcnt;
return cnt_invl;
}

```

The VerStrlen() function for UTF-16 encoded text block can be implemented using SSE4.2.

Example 10-14 shows the listing of SSE4.2 assembly implementation and Example 10-15 shows the listing of SSE4.2 intrinsic listings of VerStrlen().

Example 10-14. Assembly Listings of UTF16 VerStrlen() Using PCMPISTRI

```

// complementary range values for detecting either halves of 32-bit UTF-16 code point
static short ssch0[16]= {0x1, 0xd7ff, 0xe000, 0xffff, 0, 0};
// complementary range values for detecting the 1st half of 32-bit UTF-16 code point
static short ssch1[16]= {0x1, 0xd7ff, 0xdc00, 0xffff, 0, 0};
// complementary range values for detecting the 2nd half of 32-bit UTF-16 code point
static short ssch2[16]= {0x1, 0xdbff, 0xe000, 0xffff, 0, 0};
      (continue)

```

Example 10-14. Assembly Listings of UTF16 VerStrlen() Using PCMPISTRI

```

int utf16slen_sse4_2a(const short* s1, unsigned * pLen)
{int len = 0, res = 0;
  _asm{
    mov  eax, s1
    movdquxmm2, ssch0 ; load range value to identify either halves
    movdquxmm3, ssch1 ; load range value to identify 1st half (0xd800 to 0xdbff)
    movdquxmm4, ssch2 ; load range value to identify 2nd half (0xdc00 to 0xdfff)
    xor  ecx, ecx
    xor  edx, edx; store # of 32-bit code points (surrogate pairs)
    xor  ebx, ebx; store # of non-null 16-bit code points
    xor  edi, edi ; store # of invalid word encodings
  _loopc:
    shl  ecx, 1; pcmpistri with word processing return ecx in word granularity, multiply by 2 to get
byte offset
    add  eax, ecx
    movdquxmm1, [eax] ; load a string fragment of up to 8 words
    pcmpistri xmm2, xmm1, 15h; unsigned words, ranges, invert, lsb index returned to ecx
      ; if there is a utf-16 null wchar in xmm1, zf will be set.
      ; if all 8 words in the comparison matched range,
      ; none of bits in the intermediate result will be set after polarity inversions,
      ; and ECX will return with a value of 8
    jz   short _lstfrag; if null code, handle last fragment
      ; if ecx < 8, ecx point to a word of either 1st or 2nd half of a 32-bit code point
    cmp  ecx, 8
    jne  _chksp
    add  ebx, ecx ; accumulate # of 16-bit non-null code points
    mov  ecx, 8 ; ecx must be 8 at this point, we want to avoid loop carry dependency
    jmp  _loopc
  }
}

```

Example 10-14. Assembly Listings of UTF16 VerStrlen() Using PCMPISTRI

```

_chksp; this fragment has word encodings in the surrogate value range
  add  ebx, ecx ; account for the 16-bit code points
  shl  ecx, 1; pcmpistri with word processing return ecx in word granularity, multiply by 2 to get
byte offset
  add  eax, ecx
  movdquxmm1, [eax]; ensure the fragment start with word encoding in either half
  pcmpistri xmm3, xmm1, 15h; unsigned words, ranges, invert, lsb index returned to ecx
  jz   short _lstfrag2; if null code, handle the last fragment
  cmp  ecx, 0 ; properly encoded 32-bit code point must start with 1st half
  jg   _invalidsp; some invalid s-p code point exists in the fragment
  pcmpistri xmm4, xmm1, 15h; unsigned words, ranges, invert, lsb index returned to ecx
  cmp  ecx, 1 ; the 2nd half must follow the first half
  jne  _invalidsp
  add  edx, 1; accumulate # of valid surrogate pairs
  add  ecx, 1 ; we want to advance two words
  jmp  _loopc
_invalidsp; the first word of this fragment is either the 2nd half or an un-paired 1st half
  add  edi, 1 ; we have an invalid code point (not a surrogate pair)
  mov  ecx, 1 ; advance one word and continue scan for 32-bit code points
  jmp  _loopc
_lstfrag:
  add  ebx, ecx ; account for the non-null 16-bit code points
_morept:
  shl  ecx, 1; pcmpistri with word processing return ecx in word granularity, multiply by 2 to get
byte offset
  add  eax, ecx
  mov  si, [eax]; need to check for null code
  cmp  si, 0
  je   _final
  movdquxmm1, [eax]; load remaining word elements which start with either 1st/2nd half
  pcmpistri xmm3, xmm1, 15h; unsigned words, ranges, invert, lsb index returned to ecx

```


Example 10-14. Assembly Listings of UTF16 VerStrlen() Using PCMPISTRI

```

_1stfrag2:
  cmp  ecx, 0 ; a valid 32-bit code point must start from 1st half
  jne  _invalidsp2
  pcmpestri xmm4, xmm1, 15h; unsigned words, ranges, invert, lsb index returned to ecx
  cmp  ecx, 1
  jne  _invalidsp2
  add  edx, 1
  mov  ecx, 2
  jmp  _morept
_invalidsp2:
  add  edi, 1
  mov  ecx, 1
  jmp  _morept
_final:
  add  edx, ebx; add # of 16-bit and 32-bit code points
  mov  ecx, pLen; retrieve address of pointer provided by caller
  mov  [ecx], edx; store result of string length to memory
  mov  res, edi
  }
  return res;
}

```

Example 10-15. Intrinsic Listings of UTF16 VerStrlen() Using PCMPISTRI

```

int utf16slen_i(const short* s1, unsigned * pLen)
{int len = 0, res = 0;
 __m128i *pF = (__m128i *) s1;
 __m128iu32 = _mm_loadu_si128((__m128i *)ssch0);
    (continue)
}

```

Example 10-15. Intrinsic Listings of UTF16 VerStrlen() Using PCMPISTRI

```

__m128i u32a = _mm_loadu_si128((__m128i *)ssch1);
__m128i u32b = _mm_loadu_si128((__m128i *)ssch2);
__m128ifrag1;
int offset1 = 0, cmp, cmp_1, cmp_2;
intcnt_16 = 0, cnt_sp=0, cnt_invl= 0;
short *ps;
while (1) {
    pF = (__m128i *)(((short *)pF) + offset1);
    frag1 = _mm_loadu_si128(pF); // load up to 8 words
    // does frag1 contain either halves of a 32-bit UTF-16 code point?
    cmp = _mm_cmpistri(u32, frag1, 0x15); // unsigned bytes, equal order, lsb index returned to
    ecx

    if (_mm_cmpistrz(u32, frag1, 0x15)) // there is a null code in frag1
    { cnt_16 += cmp;
      ps = (((short *)pF) + cmp);
      while (ps[0])
      { frag1 = _mm_loadu_si128( (__m128i *)ps);
        cmp_1 = _mm_cmpistri(u32a, frag1, 0x15);
        if(!cmp_1)
        { cmp_2 = _mm_cmpistri(u32b, frag1, 0x15);
          if( cmp_2 ==1) { cnt_sp++; offset1 = 2;}
          else {cnt_invl++; offset1= 1;}
        }
        else
        { cmp_2 = _mm_cmpistri(u32b, frag1, 0x15);
          if(!cmp_2) {cnt_invl ++; offset1 = 1;}
          else {cnt_16 ++; offset1 = 1;}
        }
        ps = (((short *)ps) + offset1);
      }
      break;
    }
}

```

Example 10-15. Intrinsic Listings of UTF16 VerStrlen() Using PCMPISTRI

```

if(cmp != 8) // we have at least some half of 32-bit utf-16 code points
{
    cnt_16 += cmp; // regular 16-bit UTF16 code points
    pF = (__m128i *)(((short *)pF) + cmp);
    frag1 = _mm_loadu_si128(pF);
    cmp_1 = _mm_cmpistri(u32a, frag1, 0x15);
    if(!cmp_1)
    {
        cmp_2 = _mm_cmpistri(u32b, frag1, 0x15);
        if(cmp_2 == 1) { cnt_sp++; offset1 = 2;}
        else {cnt_invl++; offset1 = 1;}
    }
    else
    {
        cnt_invl++;
        offset1 = 1;
    }
}
else {
    offset1 = 8; // increment address by 16 bytes to handle next fragment
    cnt_16 += 8;
}
};
*pLen = cnt_16 + cnt_sp;
return cnt_invl;
}

```

10.3.6 Replacement String Library Function Using SSE4.2

Unaligned 128-bit SIMD memory access can fetch data cross page boundary, since system software manages memory access rights with page granularity.

Implementing a replacement string library function using SIMD instructions must not cause memory access violation. This requirement can be met by adding a small amounts of code to check the memory address of each string fragment. If a memory address is found to be within 16 bytes of crossing over to the next page boundary, string processing algorithm can fall back to byte-granular technique.

Example 10-16 shows an SSE4.2 implementation of strcmp() that can replace byte-granular implementation supplied by standard tools.

Example 10-16. Replacement String Library Strcmp Using SSE4.2

```

// return 0 if strings are equal, 1 if greater, -1 if less
int strcmp_sse4_2(const char *src1, const char *src2)
{
    int val;
    __asm{
        mov     esi, src1 ;
        mov     edi, src2
        mov     edx, -16 ; common index relative to base of either string pointer
        xor     eax, eax
topofloop:
        add     edx, 16 ; prevent loop carry dependency
next:
        lea     ecx, [esi+edx] ; address of fragment that we want to load
        and     ecx, 0x0fff ; check least significant 12 bits of addr for page boundary
        cmp     ecx, 0x0fff
        jg     too_close_pgb ; branch to byte-granular if within 16 bytes of boundary
        lea     ecx, [edi+edx] ; do the same check for each fragment of 2nd string
        and     ecx, 0x0fff
        cmp     ecx, 0x0fff
        jg     too_close_pgb
        movdqu  xmm2, BYTE PTR[esi+edx]
        movdqu  xmm1, BYTE PTR[edi+edx]
        pcmpistri  xmm2, xmm1, 0x18 ; equal each
        ja     topofloop
        jnc    ret_tag
        add     edx, ecx ; ecx points to the byte offset that differ
not_equal:
        movzx   eax, BYTE PTR[esi+edx]
        movzx   edx, BYTE PTR[edi+edx]
        cmp     eax, edx
        cmova   eax, ONE
        cmovb   eax, NEG_ONE
        jmp     ret_tag
        (continue)
    }
}

```

Example 10-16. Replacement String Library Strcmp Using SSE4.2

```

too_close_pgb:
    add     edx, 1 ; do byte granular compare
    movzx  ecx, BYTE PTR[esi+edx-1]
    movzx  ebx, BYTE PTR[edi+edx-1]
    cmp    ecx, ebx
    jne    inequality
    add    ebx, ecx
    jnz   next
    jmp    ret_tag
inequality:
    cmovb  eax, NEG_ONE
    cmova  eax, ONE
ret_tag:
    mov    [val], eax
    }
    return(val);
}

```

In Example 10-16, 8 instructions were added following the label “next” to perform 4KByte boundary checking of address that will be used to load two string fragments into registers. If either address is found to be within 16 bytes of crossing over to the next page, the code branches to byte-granular comparison path following the label “too_close_pgb”.

The return values of Example 10-16 uses the convention of returning 0, +1, -1 using CMOV. It is straight forward to modify a few instructions to implement the convention of returning 0, positive integer, negative integer.

10.4 SSE4.2 ENABLED NUMERICAL AND LEXICAL COMPUTATION

SSE4.2 can enable SIMD programming techniques to explore byte-granular computational problems that were considered unlikely candidates for using SIMD instructions. We consider a common library function `atol()` in its full 64-bit flavor of converting a sequence of alpha numerical characters within the range representable by the data type `__int64`.

There are several attributes of this string-to-integer problem that poses as difficult challenges for using prior SIMD instruction sets (before the introduction of SSE4.2) to accelerate the numerical computation aspect of string-to-integer conversions:

- Character subset validation: Each character in the input stream must be validated with respect to the character subset definitions and conform to data representation rules of white space, signs, numerical digits. SSE4.2 provides the perfect tools for character subset validation.
- State-dependent nature of character validation: While SIMD computation instructions can expedite the arithmetic operations of "multiply by 10 and add", the arithmetic computation requires the input byte stream to consist of numerical digits only. For example, the validation of numerical digits, white-space, and the presence/absence of sign, must be validated in mid-stream. The flexibility of the SSE4.2 primitive can handle these state-dependent validation well.
- Additionally, exit condition to wrap up arithmetic computation can happen in mid-stream due to invalid characters, or due to finite representable range of the data type ($\sim 10^{19}$ for int64, no more than 10 non-zero-leading digits for int32) may lead one to believe this type data stream consisting of short bursts are not suited for exploring SIMD ISA and be content with byte-granular solutions.

Because of the character subset validation and state-dependent nature, byte-granular solutions of the standard library function tends to have a high start-up cost (for example, converting a single numerical digit to integer may take 50 or 60 cycles), and low throughput (each additional numeric digit in the input character stream may take 6-8 cycles per byte).

A high level pseudo-operation flow of implementing a library replacement of `atol()` is described in Example 10-17.

Example 10-17. High-level flow of Character Subset Validation for String Conversion

1. Check Early_Out Exit Conditions (e.g. first byte is not valid).
2. Check if 1st byte is white space and skip any additional leading white space.
3. Check for the presence of a sign byte.
4. Check the validity of the remaining byte stream if they are numeric digits.
5. If the byte stream starts with '0', skip all leading digits that are '0'.
6. Determine how many valid non-zero-leading numeric digits.
7. Convert up to 16 non-zero-leading digits to int64 value.
8. load up to the next 16 bytes safely and check for consecutive numeric digits
9. Normalize int64 value converted from the first 16 digits, according to # of remaining digits,
10. Check for out-of-bound results of normalized intermediate int64 value,
11. Convert remaining digits to int64 value and add to normalized intermediate result,
12. Check for out-of-bound final results.

Example 10-18 shows the code listing of an equivalent functionality of `atol()` capable of producing int64 output range. Auxiliary function and data constants are listed in Example 10-19.

Example 10-18. Intrinsic Listings of atol() Replacement Using PCMPISTRI

```

__int64 sse4i_atol(const char* s1)
{char *p = (char *)s1;
  int NegSgn = 0;
  __m128i mask0;
  __m128i value0, value1;
  __m128i w1, w1_l8, w1_u8, w2, w3 = _mm_setzero_si128();
  __int64 xxi;
  int index, cflag, sflag, zflag, oob=0;
  // check the first character is valid via lookup
  if ( (BtMLValDeclnt[ *p >> 3] & (1 << ((*p) & 7)) ) == 0) return 0;
  // if the first character is white space, skip remaining white spaces
  if (BtMLws[*p >> 3] & (1 << ((*p) & 7)) )
  { p++;
    value0 = _mm_loadu_si128 ((__m128i *)listws);
  skip_more_ws:
    mask0 = __m128i_strloadu_page_boundary (p);
    /* look for the 1st non-white space character */
    index = _mm_cmpistri (value0, mask0, 0x10);
    cflag = _mm_cmpistr (value0, mask0, 0x10);
    sflag = _mm_cmpistrs (value0, mask0, 0x10);
    if( !sflag && !cflag)
    { p = (char *)((size_t) p + 16);
      goto skip_more_ws;
    }
    else    p = (char *)((size_t) p + index);
  }
  if( *p == '-')
  { p++;
    NegSgn = 1;
  }
  else if( *p == '+') p++;

  /* load up to 16 byte safely and check how many valid numeric digits we can do SIMD */
  value0 = _mm_loadu_si128 ((__m128i *)rangenumint);
  mask0 = __m128i_strloadu_page_boundary (p);
  index = _mm_cmpistri (value0, mask0, 0x14);
  zflag = _mm_cmpistrz (value0, mask0, 0x14);
  (continue)

```

Example 10-18. Intrinsic Listings of atol() Replacement Using PCMPISTRI

```

/* index points to the first digit that is not a valid numeric digit */
if( !index) return 0;
else if (index == 16)
{ if( *p == '0') /* if all 16 bytes are numeric digits */
  { /* skip leading zero */
    value1 = _mm_loadu_si128 ((__m128i *) rangenumintzr);
    index = _mm_cmpistri (value1, mask0, 0x14);
    zflag = _mm_cmpistrz (value1, mask0, 0x14);
    while(index == 16 && !zflag )
    { p = ( char *) ((size_t) p + 16);
      mask0 = __m128i_strloadu_page_boundary (p);
      index = _mm_cmpistri (value1, mask0, 0x14);
      zflag = _mm_cmpistrz (value1, mask0, 0x14);
    }
    /* now the 1st digit is non-zero, load up to 16 bytes and update index */
    if( index < 16)
      p = ( char *) ((size_t) p + index);
    /* load up to 16 bytes of non-zero leading numeric digits */
    mask0 = __m128i_strloadu_page_boundary (p);
    /* update index to point to non-numeric character or indicate we may have more than 16
bytes */
    index = _mm_cmpistri (value0, mask0, 0x14);
  }
}
if( index == 0) return 0;
else if( index == 1) return (NegSgn? (long long) -(p[0]-48): (long long) (p[0]-48));
// Input digits in xmm are ordered in reverse order. the LS digit of output is next to eos
// least sig numeric digit aligned to byte 15 , and subtract 0x30 from each ascii code
mask0 = ShfLAlnLSByte( mask0, 16 -index);
w1_u8 = _mm_slli_si128 ( mask0, 1);
w1 = _mm_add_epi8( mask0, _mm_slli_epi16 (w1_u8, 3)); /* mul by 8 and add */
w1 = _mm_add_epi8( w1, _mm_slli_epi16 (w1_u8, 1)); /* 7 LS bits per byte, in bytes 0, 2, 4, 6, 8,
10, 12, 14*/
w1 = _mm_srli_epi16( w1, 8); /* clear out upper bits of each wd*/
w2 = _mm_madd_epi16(w1, _mm_loadu_si128( (__m128i *) &MultiplyPairBaseP2[0])); /*
multiply base^2, add adjacent word*/
w1_u8 = _mm_packus_epi32 ( w2, w2); /* pack 4 low word of each dword into 63:0 */
w1 = _mm_madd_epi16(w1_u8, _mm_loadu_si128( (__m128i *) &MultiplyPairBaseP4[0])); /*
multiply base^4, add adjacent word*/
w1 = _mm_cvtepu32_epi64( w1); /* converted dw was in 63:0, expand to qw */
w1_i8 = _mm_mul_epu32(w1, _mm_setr_epi32( 100000000, 0, 0, 0));
w2 = _mm_add_epi64(w1_i8, _mm_srli_si128 (w1, 8));

```

(continue)

Example 10-18. Intrinsic Listings of atol() Replacement Using PCMPISTRI

```

if( index < 16)
{ xxi = _mm_extract_epi64(w2, 0);
  return (NegSgn? (long long) -xxi: (long long) xxi);
}
/* 64-bit integer allow up to 20 non-zero-leading digits. */
/* accumulate each 16-digit fragment*/
w3 = _mm_add_epi64(w3, w2);
/* handle next batch of up to 16 digits, 64-bit integer only allow 4 more digits */
p = ( char *) ((size_t) p + 16);
if( *p == 0)
{ xxi = _mm_extract_epi64(w2, 0);
  return (NegSgn? (long long) -xxi: (long long) xxi);
}
mask0 = __m128i_strloadu_page_boundary (p);
/* index points to first non-numeric digit */
index = _mm_cmpistri (value0, mask0, 0x14);
zflag = _mm_cmpistrz (value0, mask0, 0x14);
if( index == 0) /* the first char is not valid numeric digit */
{ xxi = _mm_extract_epi64(w2, 0);
  return (NegSgn? (long long) -xxi: (long long) xxi);
}
if ( index > 3) return (NegSgn? (long long) RINT64VALNEG: (long long) RINT64VALPOS);
/* multiply low qword by base^index */
w1 = _mm_mul_epu32( _mm_shuffle_epi32( w2, 0x50), _mm_setr_epi32( MulplyByBaseExpN
[index - 1], 0, MulplyByBaseExpN[index-1], 0));
w3 = _mm_add_epi64(w1, _mm_slli_epi64 ( _mm_srli_si128(w1, 8), 32 ));
mask0 = ShfLAlnLSByte( mask0, 16 -index);
// convert upper 8 bytes of xmm: only least sig. 4 digits of output will be added to prev 16 digits
w1_u8 = _mm_cvtepi8_epi16(_mm_srli_si128 ( mask0, 8));
/* merge 2 digit at a time with multiplier into each dword*/
w1_u8 = _mm_madd_epi16(w1_u8, _mm_loadu_si128( (__m128i *) &MulplyQuadBaseExp3To0
[ 0]));
/* bits 63:0 has two dword integer, bits 63:32 is the LS dword of output; bits 127:64 is not
needed*/
w1_u8 = _mm_cvtepu32_epi64( _mm_hadd_epi32(w1_u8, w1_u8) );
w3 = _mm_add_epi64(w3, _mm_srli_si128( w1_u8, 8) );
xxi = _mm_extract_epi64(w3, 0);
if( xxi >> 63 )
    return (NegSgn? (long long) RINT64VALNEG: (long long) RINT64VALPOS);
else return (NegSgn? (long long) -xxi: (long long) xxi);
}

```

The general performance characteristics of an SSE4.2 enhanced `atol()` replacement have a start-up cost that is somewhat lower than byte-granular implementations generated from C code.

Example 10-19. Auxiliary Routines and Data Constants Used in `sse4i_atol()` listing

```
// bit lookup table of valid ascii code for decimal string conversion, white space, sign, numeric
// digits
static char BtMLValDeclnt[32] = {0x0, 0x3e, 0x0, 0x0, 0x1, 0x28, 0xff, 0x03,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0};
// bit lookup table, white space only
static char BtMLws[32] = {0x0, 0x3e, 0x0, 0x0, 0x1, 0x0, 0x0, 0x0,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0,
0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0};
// list of white space for sttni use
static char listws[16] =
    {0x20, 0x9, 0xa, 0xb, 0xc, 0xd, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0};
// list of numeric digits for sttni use
static char rangenumint[16] =
    {0x30, 0x39, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0};
static char rangenumintzr[16] =
    {0x30, 0x30, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0, 0x0};

// we use pmaddwd to merge two adjacent short integer pair, this is the second step of merging
// each pair of 2-digit integers
static short MulplyPairBaseP2[8] =
    { 100, 1, 100, 1, 100, 1, 100, 1};

// Multiplier-pair for two adjacent short integer pair, this is the third step of merging each pair of
// 4-digit integers
static short MulplyPairBaseP4[8] =
    { 10000, 1, 10000, 1, 10000, 1, 10000, 1 };

// multiplier for pmulld for normalization of > 16 digits
static int MulplyByBaseExpN[8] =
    { 10, 100, 1000, 10000, 100000, 1000000, 10000000, 100000000};

static short MulplyQuadBaseExp3To0[8] =
    { 1000, 100, 10, 1, 1000, 100, 10, 1};
    (continue)
```

Example 10-19. Auxiliary Routines and Data Constants Used in `sse4i_atol()` listing

```

__m128i __m128i_shift_right (__m128i value, int offset)
{
  switch (offset)
  {
    case 1: value = _mm_srli_si128 (value, 1); break;
    case 2: value = _mm_srli_si128 (value, 2); break;
    case 3: value = _mm_srli_si128 (value, 3); break;
    case 4: value = _mm_srli_si128 (value, 4); break;
    case 5: value = _mm_srli_si128 (value, 5); break;
    case 6: value = _mm_srli_si128 (value, 6); break;
    case 7: value = _mm_srli_si128 (value, 7); break;
    case 8: value = _mm_srli_si128 (value, 8); break;
    case 9: value = _mm_srli_si128 (value, 9); break;
    case 10: value = _mm_srli_si128 (value, 10); break;
    case 11: value = _mm_srli_si128 (value, 11); break;
    case 12: value = _mm_srli_si128 (value, 12); break;
    case 13: value = _mm_srli_si128 (value, 13); break;
    case 14: value = _mm_srli_si128 (value, 14); break;
    case 15: value = _mm_srli_si128 (value, 15); break;
  }
  return value;
}
/* Load string at S near page boundary safely. */

__m128i __m128i_strloadu_page_boundary (const char *s)
{
  int offset = ((size_t) s & (16 - 1));
  if (offset)
  {
    __m128i v = _mm_load_si128 ((__m128i *) (s - offset));
    __m128i zero = _mm_setzero_si128 ();
    int bmsk = _mm_movemask_epi8 (_mm_cmpeq_epi8 (v, zero));
    if ((bmsk >> offset) != 0) return __m128i_shift_right (v, offset);
  }
  return _mm_loadu_si128 ((__m128i *) s);
}

```

(continue)

Example 10-19. Auxiliary Routines and Data Constants Used in sse4i_atol() listing

```

__m128i ShfLAlnLSByte( __m128i value, int offset)
{
    /*now remove constant bias, so each byte element are unsigned byte int */
    value = _mm_sub_epi8(value, _mm_setr_epi32(0x30303030, 0x30303030, 0x30303030,
0x30303030));
    switch (offset)
    {
        case 1:
            value = _mm_slli_si128 (value, 1); break;
        case 2:
            value = _mm_slli_si128 (value, 2); break;
        case 3:
            value = _mm_slli_si128 (value, 3); break;
        case 4:
            value = _mm_slli_si128 (value, 4); break;
        case 5:
            value = _mm_slli_si128 (value, 5); break;
        case 6:
            value = _mm_slli_si128 (value, 6); break;
        case 7:
            value = _mm_slli_si128 (value, 7); break;
        case 8:
            value = _mm_slli_si128 (value, 8); break;
        case 9:
            value = _mm_slli_si128 (value, 9); break;
        case 10:
            value = _mm_slli_si128 (value, 10); break;
        case 11:
            value = _mm_slli_si128 (value, 11); break;
        case 12:
            value = _mm_slli_si128 (value, 12); break;
        case 13:
            value = _mm_slli_si128 (value, 13); break;
        case 14:
            value = _mm_slli_si128 (value, 14); break;
        case 15:
            value = _mm_slli_si128 (value, 15); break;
    }
    return value;
}

```

With an input byte stream no more than 16 non-zero-leading digits, it has a constant performance. An input string consisting of more than 16 bytes of non-zero-leading

digits can be processed in about 100 cycles or less, compared byte-granular solution needing around 200 cycles. Even for shorter input strings of 9 non-zero-leading digits, SSE4.2 enhanced replacement can also achieve $\sim 2X$ performance of byte-granular solutions.

CHAPTER 11

OPTIMIZATIONS FOR INTEL® ADVANCED VECTOR EXTENSION

Intel® Advanced Vector Extension (Intel® AVX), is a major enhancement to Intel Architecture. It extends the functionality of previous generations of 128-bit SSE vector instructions and increased the vector register width to support 256-bit operations. The Intel AVX ISA enhancement is focused on float-point instructions. Some 256-bit integer vectors are supported via floating-point to integer and integer to floating-point conversions.

Intel microarchitecture code name Sandy Bridge implements the Intel AVX instructions, in most cases, on 256-bit hardware. Thus, each core has 256-bit floating-point Add and Multiply units. The Divide and Square-root units are not enhanced to 256-bits. Thus, Intel AVX instructions use the 128-bit hardware in two steps to complete these 256-bit operations.

Prior generations of Intel® Streaming SIMD Extensions (Intel® SSE) instructions generally are two-operand syntax, where one of the operands serves both as source and as destination. Intel AVX instructions are encoded with a VEX prefix, which includes a bit field to encode vector lengths and support three-operand syntax. A typical instruction has two sources and one destination. Four operand instructions such as VBLENDVPS and VBLENDVPD exist as well. The added operand enables non destructive source (NDS) and it eliminates the need for register duplication using MOVAPS operations.

With the exception of MMX instructions, almost all legacy 128-bit SSE instructions have AVX equivalents that support three operand syntax. 256-bit AVX instructions employ three-operand syntax and some with 4-operand syntax.

The 256-bit vector register, YMM, extends the 128-bit XMM register to 256 bits. Thus the lower 128-bits of YMM is aliased to the legacy XMM registers.

While 256-bit AVX instructions writes 256 bits of results to YMM, 128-bit AVX instructions writes 128-bits of results into the XMM register and zeros the upper bits above bit 128 of the corresponding YMM. 16 vector registers are available in 64-bit mode. Only the lower 8 vector registers are available in non-64-bit modes.

Software can continue to use any mixture of legacy SSE code, 128-bit AVX code and 256-bit AVX code. Section covers guidelines to deliver optimal performance across mixed-vector-length code modules without experiencing transition delays between legacy SSE and AVX code. There are no transition delays of mixing 128-bit AVX code and 256-bit AVX code.

The optimal memory alignment of an Intel AVX 256-bit vector, stored in memory, is 32 bytes. Some data-movement 256-bit Intel AVX instructions enforce 32-byte alignment and will signal #GP fault if memory operand is not properly aligned. The majority of 256-bit Intel AVX instructions do not require address alignment. These

instructions generally combine load and compute operations, so any non-aligned memory address can be used in these instructions.

For best performance, software should pay attention to align the load and store addresses to 32 bytes whenever possible.

The major differences between using AVX instructions and legacy SSE instructions are summarized in Table 11-1 below:

Table 11-1. Features between 256-bit AVX, 128-bit AVX and Legacy SSE Extensions

Features	256-bit AVX	128-bit AVX	Legacy SSE-AESNI
Functionality Scope	Floating-point operation, Data Movement	Matches legacy SIMD ISA (except MMX)	128-bit FP and integer SIMD ISA
Register Operand	YMM	XMM	XMM
Operand Syntax	Up to 4; non-destructive source	Up to 4; non-destructive source	2 operand syntax; destructive source
Memory alignment	Load-Op semantics do not require alignment.	Load-Op semantics do not require alignment.	Always enforce 16B alignment
Aligned Move Instructions	32 byte alignment	16 byte alignment	16 byte alignment
Non-destructive source operand	Yes	Yes	No
Register State Handling	Updates bits 255:0	Updates 127:0; Zeroes bits above 128.	Updates 127:0; Bits above 128 unmodified.
Intrinsic Support	<ul style="list-style-type: none"> ▪ New 256-bit data types; ▪ <code>_mm256</code> prefix for promoted functionality ▪ new intrinsics for new functionalities 	<ul style="list-style-type: none"> ▪ Existing data types ▪ Inherit same prototype for exiting functionalities ▪ Use “<code>_mm</code>” prefix for new VEX-128 functionalities. 	baseline datatypes and prototype definitions
128-bit Lanes	Applies to most 256-bit operations	One 128-bit lane	One 128-bit lane
Mixed Code Handling	Use <code>VZEROUPPER</code> to avoid transition penalty	No transition penalty	Transition penalty after executing 256-bit AVX code.

11.1 INTEL® AVX INTRINSICS CODING

256-bit AVX instructions have new intrinsics. Specifically, 256-bit AVX instructions that are promoted to 256-bit vector length from existing SSE functionality are gener-

ally prototyped with a “_mm256” prefix instead of the “_mm” prefix and using new data types defined for 256-bit operation. New functionality in 256-bit AVX instructions have brand new prototype.

The 128-bit AVX instruction that were promoted from legacy SIMD ISA uses the same prototype as before. Newer functionality common in 256-bit and 128-bit AVX instructions are prototyped with “_mm256” and “_mm” prefixes respectively.

Thus porting from legacy SIMD code written in intrinsic can be ported to 256-bit AVX code with a modest effort.

The following guidelines show how to convert a simple intrinsic from Intel SSE code sequence to Intel AVX:

- Align statically and dynamically allocated buffers to 32-bytes
- May need to double supplemental buffer size
- Change __mm_ intrinsic name prefix with __mm256_
- Change variable data types names from __m128 to __m256
- Divide by 2 iteration count (or double stride length)

This example below on Cartesian coordinate transformation demonstrates the Intel AVX Instruction format, 32 byte YMM registers, dynamic and static memory allocation with data alignment of 32bytes, and the C data type representing 8 floating point elements in a YMM register.

Example 11-1. Cartesian Coordinate Transformation with Intrinsics

<pre>//Use SSE intrinsic #include "wmmmintrin.h" int main() { int len = 3200; //Dynamic memory allocation with 16byte //alignment float* plnVector = (float*) _mm_malloc(len*sizeof(float), 16); float* pOutVector = (float*) _mm_malloc(len*sizeof(float), 16); //init data for(int i=0; i<len; i++) plnVector[i] = 1; float cos_teta = 0.8660254037; float sin_teta = 0.5; //Static memory allocation of 4 floats with 16byte alignment __declspec(align(16)) float cos_sin_teta_vec[4] = {cos_teta, sin_teta, cos_teta, sin_teta};</pre>	<pre>// Use Intel AVX intrinsic #include "immintrin.h" int main() { int len = 3200; //Dynamic memory allocation with 16byte //alignment float* plnVector = (float*) _mm_malloc(len*sizeof(float), 32); float* pOutVector = (float*) _mm_malloc(len*sizeof(float), 32); //init data for(int i=0; i<len; i++) plnVector[i] = 1; float cos_teta = 0.8660254037; float sin_teta = 0.5; //Static memory allocation of 8 floats with 32byte alignment __declspec(align(32)) float cos_sin_teta_vec[8] = {cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta};</pre>
--	--

Example 11-1. Cartesian Coordinate Transformation with Intrinsics (Contd.)

<pre> __declspec(align(16)) float sin_cos_teta_vec[4] = {sin_teta, cos_teta, sin_teta, cos_teta}; //__m128 data type represents an xmm //register with 4 float elements __m128 Xmm_cos_sin = _mm_load_ps(cos_sin_teta_vec); //SSE 128bit packed single load __m128 Xmm_sin_cos = _mm_load_ps(sin_cos_teta_vec); __m128 Xmm0, Xmm1, Xmm2, Xmm3 //processing 8 elements in an unrolled //twice loop for(int i=0; i<len; i+=8) { Xmm0 = _mm_load_ps(plnVector+i); Xmm1 = _mm_moveldup_ps(Xmm0); Xmm2 = _mm_movehdup_ps(Xmm0); Xmm1 = _mm_mul_ps(Xmm1,Xmm_cos_sin); Xmm2 = _mm_mul_ps(Xmm2,Xmm_sin_cos); Xmm3 = _mm_addsub_ps(Xmm1, Xmm2); _mm_store_ps(pOutVector + i, Xmm3); Xmm0 = _mm_load_ps(plnVector+i+4); Xmm1 = _mm_moveldup_ps(Xmm0); Xmm2 = _mm_movehdup_ps(Xmm0); Xmm1 = _mm_mul_ps(Xmm1,Xmm_cos_sin); Xmm2 = _mm_mul_ps(Xmm2,Xmm_sin_cos); Xmm3 = _mm_addsub_ps(Xmm1, Xmm2); _mm_store_ps(pOutVector+i+4, Xmm3); } _mm_free(plnVector); _mm_free(pOutVector); return 0; } </pre>	<pre> __declspec(align(32)) float sin_cos_teta_vec[8] = {sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta }; //__m256 data type holds 8 float elements __m256 Ymm_cos_sin = _mm256_load_ps(cos_sin_teta_vec); //AVX 256bit packed single load __m256 Ymm_sin_cos = _mm256_load_ps(sin_cos_teta_vec); __m256 Ymm0, Ymm1, Ymm2, Ymm3; //processing 8 elements in an unrolled //twice loop for(int i=0; i<len; i+=16) { Ymm0 = _mm256_load_ps(plnVector+i); Ymm1 = _mm256_moveldup_ps(Ymm0); Ymm2 = _mm256_movehdup_ps(Ymm0); Ymm1 = _mm256_mul_ps(Ymm1,Ymm_cos_sin); Ymm2 = _mm256_mul_ps(Ymm2,Ymm_sin_cos); Ymm3 = _mm256_addsub_ps(Ymm1, Ymm2); _mm256_store_ps(pOutVector + i, Ymm3); Ymm0 = _mm256_load_ps(plnVector+i+8); Ymm1 = _mm256_moveldup_ps(Ymm0); Ymm2 = _mm256_movehdup_ps(Ymm0); Ymm1 = _mm256_mul_ps(Ymm1,Ymm_cos_sin); Ymm2 = _mm256_mul_ps(Ymm2,Ymm_sin_cos); Ymm3 = _mm256_addsub_ps(Ymm1, Ymm2); _mm256_store_ps(pOutVector+i+8, Ymm3); } _mm_free(plnVector); _mm_free(pOutVector); return 0; } </pre>
--	--

11.1.1 Intel® AVX Assembly Coding

Similar to the intrinsic porting guidelines, assembly porting guidelines is listed below:

- Align statically and dynamically allocated buffers to 32-bytes
- Double the supplemental buffer sizes if needed
- Add a "v" prefix to instruction names
- Change register names from xmm to ymm
- Add destination registers to computational Intel AVX instructions
- Divide the iteration count by two (or double stride length)

Example 11-2. Cartesian Coordinate Transformation with Assembly

<pre>//Use SSE Assembly int main() { int len = 3200; //Dynamic memory allocation with 16byte //alignment float* plnVector = (float*) _mm_malloc(len*sizeof(float), 16); float* pOutVector = (float*) _mm_malloc(len*sizeof(float), 16); //init data for(int i=0; i<len; i++) plnVector[i] = 1; //Static memory allocation of 4 floats //with 16byte alignment float cos_teta = 0.8660254037; float sin_teta = 0.5; __declspec(align(16)) float cos_sin_teta_vec[4] = {cos_teta, sin_teta, cos_teta, sin_teta}; __declspec(align(16)) float sin_cos_teta_vec[4] = {sin_teta, cos_teta, sin_teta, cos_teta}; //processing 8 elements in an unrolled //twice loop</pre>	<pre>// Use Intel AVX assembly int main() { int len = 3200; //Dynamic memory allocation with 32byte //alignment float* plnVector = (float*) _mm_malloc(len*sizeof(float), 32); float* pOutVector = (float*) _mm_malloc(len*sizeof(float), 32); //init data for(int i=0; i<len; i++) plnVector[i] = 1; //Static memory allocation of 8 floats //with 32byte alignment float cos_teta = 0.8660254037; float sin_teta = 0.5; __declspec(align(32)) float cos_sin_teta_vec[8] = {cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta}; __declspec(align(32)) float sin_cos_teta_vec[8] = {sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta, sin_teta, cos_teta}; //processing 16 elements in an unrolled //twice loop</pre>
--	--

Example 11-2. Cartesian Coordinate Transformation with Assembly (Contd.)

<pre> __asm { mov eax, plnVector mov ebx, pOutVector // Load into an xmm register of 16 bytes movups xmm3, xmmword ptr[cos_sin_teta_vec] movups xmm4, xmmword ptr[sin_cos_teta_vec] mov edx, len xor ecx, ecx loop1: movsldup xmm0, [eax+ecx] movshdup xmm1, [eax+ecx] //example: mulps has 2 operands mulps xmm0, xmm3 mulps xmm1, xmm4 addsubps xmm0, xmm1 // 16 byte store from an xmm register movaps [ebx+ecx], xmm0 movsldup xmm0, [eax+ecx+16] movshdup xmm1, [eax+ecx+16] mulps xmm0, xmm3 mulps xmm1, xmm4 addsubps xmm0, xmm1 // offset of 16 bytes from previous store movaps [ebx+ecx+16], xmm0 // Processed 32bytes in this loop //(The code is unrolled twice) add ecx, 32 cmp ecx, edx jl loop1 } _mm_free(plnVector); _mm_free(pOutVector); return 0; } </pre>	<pre> __asm { mov eax, plnVector mov ebx, pOutVector // Load into an ymm register of 32 bytes vmovups ymm3, ymmword ptr[cos_sin_teta_vec] vmovups ymm4, ymmword ptr[sin_cos_teta_vec] mov edx, len xor ecx, ecx loop1: vmovsldup ymm0, [eax+ecx] vmovshdup ymm1, [eax+ecx] //example: vmulps has 3 operands vmulps ymm0, ymm0, ymm3 vmulps ymm1, ymm1, ymm4 vaddsubps ymm0, ymm0, ymm1 // 32 byte store from an ymm register vmovaps [ebx+ecx], ymm0 vmovsldup ymm0, [eax+ecx+32] vmovshdup ymm1, [eax+ecx+32] vmulps ymm0, ymm0, ymm3 vmulps ymm1, ymm1, ymm4 vaddsubps ymm0, ymm0, ymm1 // offset of 32 bytes from previous store vmovaps [ebx+ecx+32], ymm0 // Processed 64bytes in this loop //(The code is unrolled twice) add ecx, 64 cmp ecx, edx jl loop1 } _mm_free(plnVector); _mm_free(pOutVector); return 0; } </pre>
---	--

11.2 NON-DESTRUCTIVE SOURCE (NDS)

Most Intel AVX instructions have three operands. A typical instruction has two sources and one destination, with both source operands unmodified by the instruction. This section describes how using the NDS feature to save register copies, reduce amount of instructions, reduce amount of micro-ops and improve performance. In this example, the Intel AVX code is more than 2x faster than the Intel SSE code.

The following example uses a vectorized calculation of the polynomial A^3+A^2+A . The polynomial calculation pseudo code is:

```
While (i<len)
{
    B[i] := A[i]3 + A[i]2 + A[i]
    i++
}
```

The left side of Example 11-3 shows the vectorized implementation using SSE assembly. In this code, A is copied by an additional load from memory to a register and A2 is copied using a register to register assignment. The code uses 10 micro-ops to process four elements.

The middle example uses 128-bit AVX and takes advantage of NDS. The additional load and register copies are eliminated. This code uses 8 micro-ops to process 4 elements and is about 30% faster than the baseline above.

The right-most example uses 256-bit Intel AVX instructions. It uses 8 micro-ops to process 8 elements. Combining the NDS feature with doubling of vector width, this speeds up over the baseline by more than 2x.

Example 11-3. Direct Polynomial Calculation

SSE code	128-bit AVX code	256-bit AVX code
<pre>float* pA = InputBuffer; float* pB = OutputBuffer; int len = miBufferWidth-4; __asm { mov rax, pA mov rbx, pB movsxd r8, len</pre>	<pre>float* pA = InputBuffer1; float* pB = OutputBuffer1; int len = miBufferWidth-4; __asm { mov rax, pA mov rbx, pB movsxd r8, len</pre>	<pre>float* pA = InputBuffer1; float* pB = OutputBuffer1; int len = miBufferWidth-8; __asm { mov rax, pA mov rbx, pB movsxd r8, len</pre>

Example 11-3. Direct Polynomial Calculation (Contd.)

SSE code	128-bit AVX code	256-bit AVX code
<pre> loop1: //Load A movupsxmm0, [rax+r8*4] //Copy A movupsxmm1, [rax+r8*4] //A^2 mulpsxmm1, xmm1 //Copy A^2 movupsxmm2, xmm1 //A^3 mulps xmm2, xmm0 //A + A^2 addps xmm0, xmm1 //A + A^2 + A^3 addps xmm0, xmm2 //Store result movups[rbx+r8*4], xmm0 sub r8, 4 jge loop1 } </pre>	<pre> loop1: //Load A vmovups xmm0, [rax+r8*4] //A^2 vmulps xmm1, xmm0, xmm0 //A^3 vmulps xmm2, xmm1, xmm0 //A+A^2 vaddps xmm0, xmm0, xmm1 //A+A^2+A^3 vaddps xmm0, xmm0, xmm2 //Store result vmovups[rbx+r8*4], xmm0 sub r8, 4 jge loop1 } </pre>	<pre> loop1: //Load A vmovups ymm0, [rax+r8*4] //A^2 vmulpsymm1, ymm0, ymm0 //A^3 vmulps ymm2, ymm1, ymm0 //A+A^2 vaddps ymm0, ymm0, ymm1 //A+A^2+A^3 vaddps ymm0, ymm0, ymm2 //Store result vmovups [rbx+r8*4], ymm0 sub r8, 8 jge loop1 } </pre>

11.3 MIXING AVX CODE WITH SSE CODE

The Intel AVX architecture allows programmers to port a large code base gradually, resulting in mixed AVX code and SSE code. If your code includes both Intel AVX and Intel SSE, consider the following:

- Recompilation of Intel SSE code with the Intel compiler and the option `"/QxAVX"` in Windows or `"-xAVX"` in Linux. This transforms all SSE instructions to 128-bit Intel AVX instructions automatically. This refers to inline assembly and intrinsic code. `"GCC -c -mAVX"` will generate AVX code, including assembly files. GCC assembler also supports `"-msse2AVX"` switch to generate AVX code from SSE.
- Intel AVX and Intel SSE code can co-exist and execute in the same run. This can happen if your application includes third party libraries with Intel SSE code, a new DLL using AVX code is deployed that calls other modules running SSE code, or you cannot recompile all your application at once. In these cases, the AVX code must use the `VZERoupper` instruction to avoid AVX/SSE transition penalty.

AVX instruction always modifies the upper bits of YMM registers and SSE instructions do not modify the upper bits. From a hardware perspective, the upper bits of the YMM register collection can be considered to be in one of three states:

- Clean: All upper bits of YMM are zero. This is the state when processor starts from RESET.
- Modified and saved to XSAVE region The content of the upper bits of YMM registers matches saved data in XSAVE region. This happens when after XSAVE/XRSTOR executes.
- Modified and Unsaved: The execution of one AVX instruction (either 256-bit or 128-bit) modifies the upper bits of the destination YMM.

The AVX/SSE transition penalty applies whenever the processor states is "Modified and Unsaved". Using VZEROUPPER move the processor states to "Clean" and avoid the transition penalty.

Table 11-2. State Transitions of Mixing AVX and SSE Code

	Execute SSE			Execute AVX-128			Execute AVX-256		
	Clean	M/U	Saved	Clean	M/U	Saved	Clean	M/U	Saved
Begin	Clean	M/U	Saved	Clean	M/U	Saved	Clean	M/U	Saved
End	Clean	Saved	Saved	Clean	Clean	M/U	M/U	M/U	M/U
Penalty	No	XSAVE	No	No	No	Yes	No	No	Yes

Table 11-2 lists the effect of mixing Intel AVX and Intel SSE code. Initially the processor is in clean state (1), where Intel SSE and Intel AVX instructions are executed with no penalty. When a 256-bit Intel AVX instruction is executed, the processor marks that it is in the Dirty Upper state (2). While in this state, executing an Intel SSE instruction saves the upper 128 bits of all YMM registers and the state changes to Saved Dirty Upper state (3). Next time an Intel AVX instruction is executed the upper 128 bits of all YMM registers are restored and the processor is back at state (2). These save and restore operations have a high penalty. Frequent execution of these transitions causes significant performance loss.

To enable fast transition between 256-bit Intel AVX and Intel SSE code, use the VZEROUPPER instruction. The VZEROUPPER instruction resets the upper 128 bits of all Intel AVX registers. This instruction has zero latency. In addition, the processor changes back to Clean state (1), after which execution of Intel SSE instructions or Intel AVX instructions has no transition penalty.

128-bit Intel AVX instructions zero the upper 128-bits of the destination registers. Therefore, 128-bit and 256-bit Intel AVX instructions can be mixed with no penalty.

Assembly/Compiler Coding Rule 72. (H impact, H generality) *Whenever 256-bit AVX code and 128-bit SSE code might execute together, use the VZEROUPPER instruction whenever a transition from "Modified/Unsaved" state is expected.*

11.3.1 Mixing Intel® AVX and Intel SSE in Function Calls

Intel AVX to Intel SSE transitions can occur unexpectedly when calling functions or returning from functions. For example, if a function that uses 256-bit Intel AVX, calls another function, the callee might be using SSE code. Similarly, after a 256-bit Intel AVX function returns, the caller might be executing Intel SSE code

Assembly/Compiler Coding Rule 73. (H impact, H generality) Add *VZERoupper* instruction after 256-bit AVX instructions are executed and before any function call that might execute SSE code. Add *VZERoupper* at the end of any function that uses 256-bit AVX instructions.

Example 11-4. Function Calls and AVX/SSE transitions

<pre> __attribute__((noinline)) void SSE_function() { __asm addps xmm1, xmm2 __asm xorps xmm3, xmm4 } __attribute__((noinline)) void AVX_function_no_zeroupper() { __asm vaddps ymm1, ymm2, ymm3 __asm vxorps ymm4, ymm5, ymm6 } __attribute__((noinline)) void AVX_function_with_zeroupper() { __asm vaddps ymm1, ymm2, ymm3 __asm vxorps ymm4, ymm5, ymm6 //add vzeroupper when returning from an AVX function __asm vzeroupper } // Code encounter transition penalty __asm vaddps ymm1, ymm2, ymm3 .. //penalty SSE_function(); AVX_function_no_zeroupper(); //penalty __asm addps xmm1, xmm2 </pre>	<pre> // Code mitigated transition penalty __asm vaddps ymm1, ymm2, ymm3 //add vzeroupper before //calling SSE function from AVX code __asm vzeroupper //no penalty SSE_function(); AVX_function_with_zeroupper(); //no penalty __asm addps xmm1, xmm2 </pre>
---	---

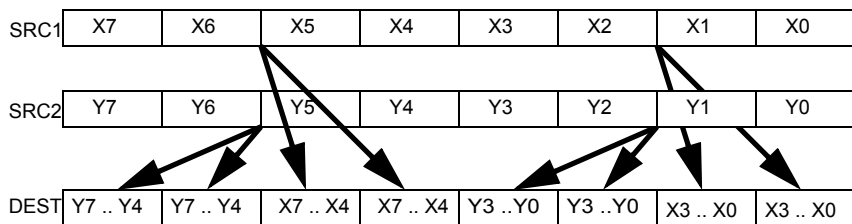
11.4 128-BIT LANE OPERATION AND AVX

256-bit operations in Intel AVX are generally performed in two halves of 128-bit lanes. Most of the 256-bit Intel AVX instructions are defined as in-lane: the destination elements in each lane are calculated using source elements only from the same lane. There are only a few cross-lane instructions, which are described below.

The majority of SSE computational instructions perform computation along vertical slots with each data elements. The 128-bit lanes does not affect porting 128-bit code into 256-bit AVX code. VADDPS is one example of this.

Many 128-bit SSE instruction moves data elements horizontally, e.g. SHUFPS uses an imm8 byte to control the horizontal movement of data elements.

Intel AVX promotes these horizontal 128-bit SIMD instruction in-lane into 256-bit operation by using the same control field within the low 128-bit lane and the high 128-bit lane. For example, the 256-bit VSHUFPS instruction uses a control byte containing 4 control values to select the source location of each destination element in a 128-bit lane. This is shown below.



Imm8: Imm[7:6] Imm[5:4] Imm[3:2] Imm[1:0] Imm[7:6] Imm[5:4] Imm[3:2] Imm[1:0]

Control Values 00b: X0/Y0 (Low lane), X4/Y4 (high lane)

Control Values 01b: X1/Y1 (Low lane), X5/Y5 (high lane)

Control Values 10b: X2/Y2 (Low lane), X6/Y6 (high lane)

Control Values 11b: X3/Y3 (Low lane), X7/Y7 (high lane)

11.4.1 Programming With the Lane Concept

Using the lane concept, algorithms implemented with SSE instruction set can be easily converted to use 256-bit Intel AVX. An SSE algorithm that executes iterations 0 to n can be converted such that the calculation of iteration i is done in the low lane and the calculation of iteration i+k is done in the high lane. For consecutive iterations k equals one.

Some vectorized algorithms implemented with SSE instructions cannot use a simple conversion described above. For example, shuffles that move elements within 16 bytes cannot be naturally converted to shuffles with 32 byte since 32 byte shuffles can't cross lanes.

You can use the following instructions as building blocks for working with lanes.

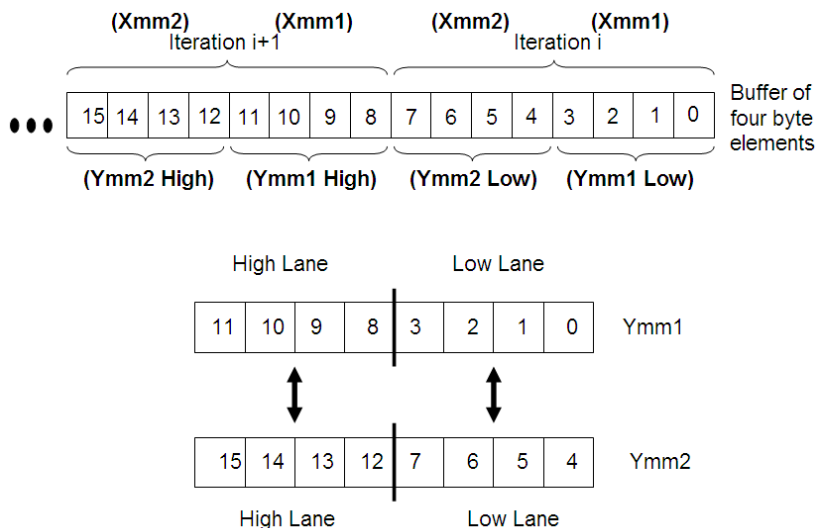
- VINSERTF128 - insert packed floating-point values
- VEXTRACTF128 - extract packed floating-point values
- VPERM2F128 - permute floating-point values
- VBROADCAST - load with broadcast

The sections below describe two techniques: the strided loads and the cross register overlap. These methods implement the in lane data arrangement described above and are useful in many algorithms that initially seem to require cross lane calculations.

11.4.2 Strided Load Technique

The strided load technique is a programming method that uses Intel AVX instructions and is useful for algorithms that involve unsupported cross-lane shuffles.

The method describes how to arrange data in a way that avoids cross-lane shuffles. The main idea is to use 128-bit loads in a way that mimics the corresponding Intel SSE algorithm, and enables the 256 Intel AVX instructions to execute iterations i of the loop in the low lanes and the iteration $i+k$ in the high lanes. In the following example, k equals one.



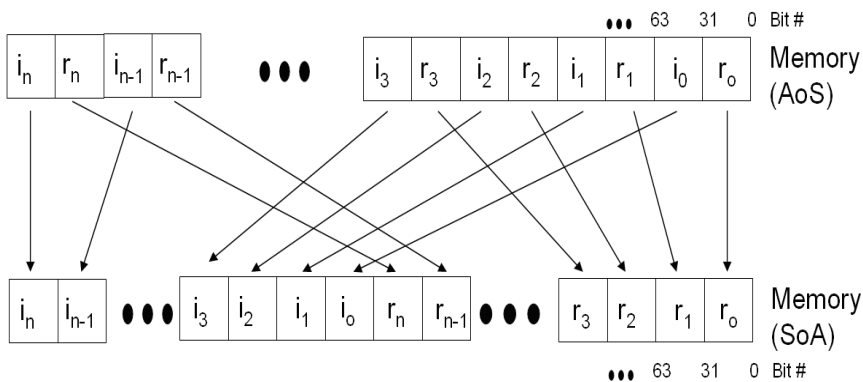
The values in the low lanes of Ymm1 and Ymm2 in the figure above correspond to iteration *i* in the SSE implementation. Similarly, the values in the high lanes of Ymm1 and Ymm2 correspond to iteration *i*+1.

The following example demonstrates the strided load method in a conversion of an Array of Structures (AoS) to a Structure of Arrays (SoA). In this example, the input buffer contains complex numbers in an AoS format. Each complex number is made of a real and an imaginary float values. The output buffer is arranged as SoA. All the real components of the complex numbers are located at the first half of the output buffer and all the imaginary components are located at the second half of the buffer. The following pseudo code and figure illustrate the conversion:

Example 11-5. AoS to SoA Conversion of Complex Numbers in C Code

```

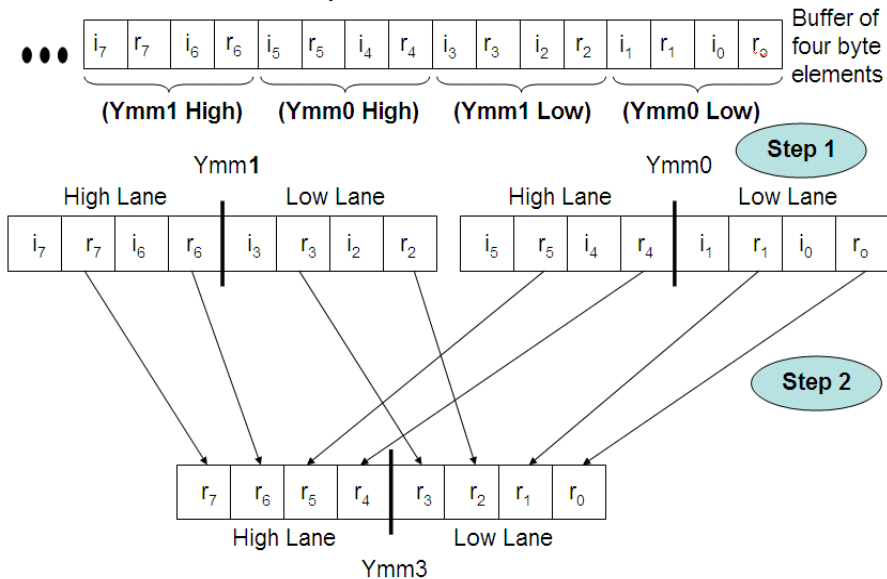
lfor (i=0 | < N)
{
  Real[i] =Complex[i].Real
  Imaginary[i] =Complex[i].Imaginary
}
    
```



A simple extension of the Intel SSE algorithm from 16-byte to 32-byte operations would require cross-lane data transition, as shown in the following figure. However, this is not possible with Intel AVX architecture and a different technique is required

The challenge of cross-lane shuffle can be overcome with Intel AVX for AoS to SoA conversion. Using VINSERTF128 to load 16 bytes to the appropriate lane in the YMM registers obviates the need for cross-lane shuffle. Once the data is organized prop-

erly in the YMM registers for step 1, 32-byte VSHUFPS can be used to move the data in lanes, as shown in step 2.



The following code compares the Intel SSE implementation of AoS to SoA with the 256-bit Intel AVX implementation and demonstrates the performance gained.

Example 11-6. AOS to SoA Conversion of Complex Numbers Using AVX

SSE Code	AVX Code
<pre>xor rbx, rbx xor rdx, rdx mov rcx, len mov rdi, inPtr mov rsi, outPtr1 mov rax, outPtr2 loop1: movups xmm0, [rdi+rbx] //i1 r1 i0 r0 movups xmm1, [rdi+rbx+16] // i3 r3 i2 r2 movups xmm2, xmm0</pre>	<pre>xor rbx, rbx xor rdx, rdx mov rcx, len mov rdi, inPtr mov rsi, outPtr1 mov rax, outPtr2 loop1: vmovups xmm0, [rdi+rbx] //i1 r1 i0 r0 vmovups xmm1, [rdi+rbx+16] // i3 r3 i2 r2 vinsertf128 ymm0, ymm0, [rdi+rbx+32], 1 //i5 r5 i4 r4; i1 r1 i0 r0</pre>

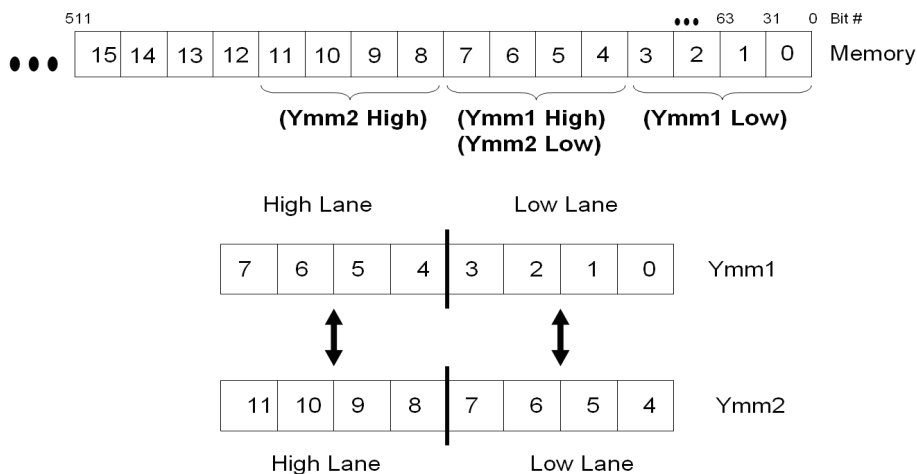
Example 11-6. AoS to SoA Conversion of Complex Numbers Using AVX (Contd.)

SSE Code	AVX Code
<pre>shufps xmm0, xmm1, 0xdd //i3 i2 i1 i0 shufps xmm2, xmm1, 0x88 //r3 r2 r1 r0 movups [rax+rdx], xmm0 movups [rsi+rdx], xmm2 add rdx, 16 add rbx, 32 cmp rcx, rbx jnz loop1</pre>	<pre>vinsertf128 ymm1, ymm1, [rdi+rbx+48], 1 //i7 r7 i6 r6; i3 r3 i2 r2 vshufps ymm2, ymm0, ymm1, 0xdd //i7 i6 i5 i4; i3 i2 i1 i0 vshufps ymm3, ymm0, ymm1, 0x88 //r7 r6 r5 r4; r3 r2 r1 r0 vmovups [rax+rdx], ymm2 vmovups [rsi+rdx], ymm3 add rdx, 32 add rbx, 64 cmp rcx, rbx jnz loop1</pre>

11.4.3 The Register Overlap Technique

The register overlap technique is useful for algorithms that use shuffling. Similar to the strided load technique, the register overlap technique arranges data to avoid cross-lane shuffles.

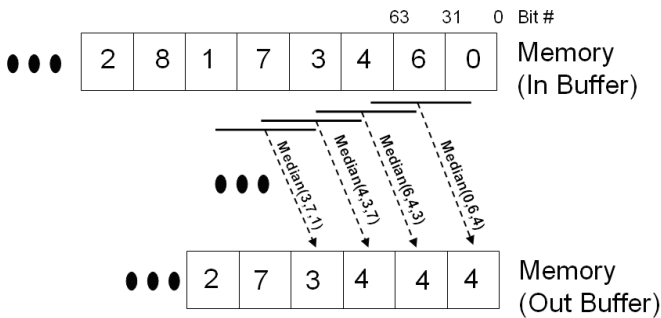
This technique is useful for algorithm that process continues data, which is partially shared by sequential iterations. The following figure illustrates the desired data layout. This is enabled by using overlapping 256-bit loads, or by using the VPERM2F128 instruction



The Median3 code sample below demonstrates the register overlap technique. The median3 technique calculates the median of every three consecutive elements in a vector.

$$Y[i] = \text{Median}(X[i], X[i+1], X[i+2])$$

Where Y is the output vector and X is the input vector. The following figure illustrates the calculation done by the median algorithm.

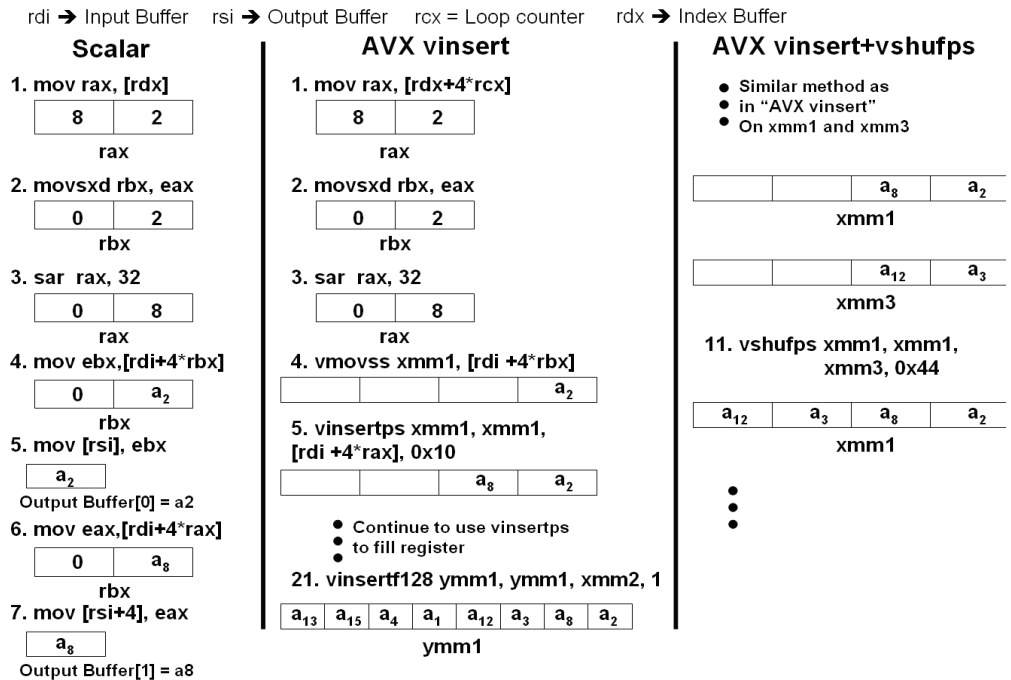


Following are three implementations of the Median3 algorithm. Alternative 1 is the Intel SSE implementation. Alternatives 2 and 3 implement the register overlap technique in two ways. Alternative 2 loads the data from the input buffer into the YMM registers using overlapping 256-bit load operations. Alternative 3 loads the data from the input buffer into the YMM registers using a 256-bit load operation and VPERM2F128. Alternatives 2 and 3 gain performance by using wider vectors.

Example 11-7. Register Overlap Method for Median of 3 Numbers

1: SSE Code	2: 256-bit AVX with Overlapping Loads	3: 256-bit AVX with VPERM2F128
<pre>xor ebx, ebx mov rcx, len mov rdi, inPtr mov rsi, outPtr movaps xmm0, [rdi] loop_start: movaps xmm4, [rdi+16] movaps xmm2, [rdi] movaps xmm1, [rdi] movaps xmm3, [rdi]</pre>	<pre>xor ebx, ebx mov rcx, len mov rdi, inPtr mov rsi, outPtr vmovaps ymm0, [rdi] loop_start: vshufps ymm2, ymm0, [rdi+16], 0x4E vshufps ymm1, ymm0, ymm2, 0x99</pre>	<pre>xor ebx, ebx mov rcx, len mov rdi, inPtr mov rsi, outPtr vmovaps ymm0, [rdi] loop_start: add rdi, 32 vmovaps ymm6, [rdi] vperm2f128 ymm1, ymm0, ymm6, 0x21 vshufps ymm3, ymm0, ymm1, 0x4E</pre>

Following are 3 implementations for the gather operation from an array of 4 byte elements. Alternative 1 is a scalar implementation using general purpose registers. Alternative 2 and 3 use Intel AVX instructions. The figure below shows code snippets from Example 11-8 assuming that it runs the first iteration on data from the previous figure.



Performance of the Intel AVX examples is similar to the performance of a corresponding Intel SSE implementation. The table below shows the three gather implementations

Example 11-8. Data Gather - AVX versus Scalar Code

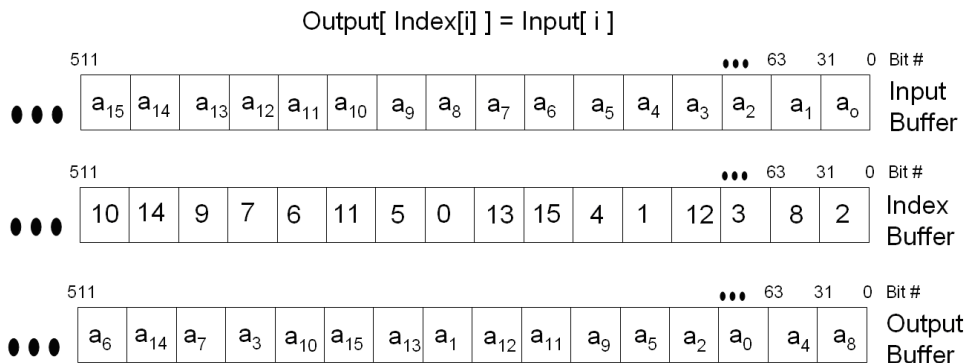
1: Scalar Code	2: AVX w/ VINSRT	3: VINSRT+VSHUFPS
movr di, InBuf	movr di, InBuf	movr di, InBuf
mov rsi, OutBuf	mov rsi, OutBuf	mov rsi, OutBuf
mov rdx, Index	mov rdx, Index	mov rdx, Index
xor rcx, rcx	xor rcx, rcx	xor rcx, rcx

Example 11-8. Data Gather - AVX versus Scalar Code (Contd.)

1: Scalar Code	2: AVX w/ VINSRT	3: VINSRT+VSHUFPS
<pre> loop1: mov rax, [rdx] movsxd rbx, eax sar rax, 32 mov ebx,[rdi+4*rbx] mov [rsi], ebx mov eax,[rdi+4*rax] mov [rsi+4], eax mov rax, [rdx+8] movsxd rbx, eax sar rax, 32 mov ebx, [rdi+4*rbx] mov [rsi+8], ebx mov eax,[rdi+4*rax] mov [rsi+12], eax mov rax, [rdx+16] movsxd rbx, eax sar rax, 32 mov ebx, [rdi+4*rbx] mov [rsi+16], ebx mov eax, [rdi+4*rax] mov [rsi+20], eax mov rax, [rdx+24] movsxd rbx, eax sar rax, 32 mov ebx, [rdi+4*rbx] mov [rsi+24], ebx mov eax, [rdi+4*rax] mov [rsi+28], eax addrsi, 32 addrdx, 32 addrcx, 8 cmprcx, len jl loop1 </pre>	<pre> loop1: mov rax, [rdx+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm1, [rdi +4*rbx] vinsertps xmm1, xmm1, [rdi +4*rax], 0x10 mov rax, [rdx + 8+4*rcx] movsxd rbx, eax sar rax, 32 vinsertps xmm1, xmm1, [rdi +4*rbx], 0x20 vinsertps xmm1, xmm1, [rdi +4*rax], 0x30 mov rax, [rdx + 16+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm2, [rdi +4*rbx] vinsertps xmm2, xmm2, [rdi +4*rax], 0x10 mov rax,[rdx + 24+4*rcx] movsxd rbx, eax sar rax, 32 vinsertps xmm2, xmm2, [rdi +4*rbx], 0x20 vinsertps xmm2, xmm2, [rdi +4*rax], 0x30 vinsertf128 ymm1, ymm1, xmm2, 1 vmovaps [rsi+4*rcx], ymm1 add rcx, 8 cmp rcx, len jl loop1 </pre>	<pre> loop1: mov rax, [rdx+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm1, [rdi +4*rbx] vinsertps xmm1, xmm1, [rdi +4*rax], 0x10 mov rax, [rdx + 8+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm3, [rdi +4*rbx] vinsertps xmm3, xmm3, [rdi +4*rax], 0x10 vshufps xmm1, xmm1, xmm3, 0x44 mov rax, [rdx + 16+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm2, [rdi +4*rbx] vinsertps xmm2, xmm2, [rdi +4*rax], 0x10 Mov rax,[rdx + 24+4*rcx] movsxd rbx, eax sar rax, 32 vmovss xmm4, [rdi +4*rbx] vinsertps xmm4, xmm4, [rdi +4*rax], 0x10 vshufpsxmm2, xmm2, xmm4, 0x44 vinsertf128 ymm1, ymm1, xmm2, 1 vmovaps [rsi+4*rcx], ymm1 add rcx, 8 cmp rcx, len jl loop1 </pre>

11.5.2 Data Scatter

The scatter operation reads elements from an input buffer sequentially. It then writes them to an output buffer based on indexes specified in an index buffer. The following figure illustrates an example for a scatter operation.



The following table includes a scalar implementation and an Intel AVX implementation of a scatter sequence. The Intel AVX examples consist mainly of 128-bit Intel AVX instructions. Performance of the Intel AVX examples is similar to the performance of corresponding Intel SSE implementation.

Example 11-9. Scatter Operation Using AVX

Scalar Code	AVX Code
<pre> mov rdi, InBuf mov rsi, OutBuf mov rdx, Index xorrr cx, rcx loop1: movsxd rax, [rdx] mov ebx, [rdi] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 4] mov ebx, [rdi + 4] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 8] </pre>	<pre> mov rdi, InBuf mov rsi, OutBuf mov rdx, Index xorrr cx, rcx loop1: vmovaps ymm0, [rdi + 4*rcx] movsxd rax, [rdx + 4*rcx] movsxd rbx, [rdx + 4*rcx + 4] vmovss [rsi + 4*rax], xmm0 movsxd rax, [rdx + 4*rcx + 8] vpsllq xmm1, xmm0, 4 </pre>

Example 11-9. Scatter Operation Using AVX (Contd.)

Scalar Code	AVX Code
<pre> mov ebx, [rdi + 8] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 12] mov ebx, [rdi + 12] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 16] mov ebx, [rdi + 16] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 20] mov ebx, [rdi + 20] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 24] mov ebx, [rdi + 24] mov [rsi + 4*rax], ebx movsxd rax, [rdx + 28] mov ebx, [rdi + 28] mov [rsi + 4*rax], ebx add rdi, 32 add rdx, 32 add rcx, 8 cmp rcx, len jl loop1 </pre>	<pre> vmovss [rsi + 4*rbx], xmm1 movsxd rbx, [rdx + 4*rcx + 12] vpsalignr xmm2, xmm0, xmm0, 8 vmovss [rsi + 4*rax], xmm2 movsxd rax, [rdx + 4*rcx + 16] vpsalignr xmm3, xmm0, xmm0, 12 vmovss [rsi + 4*rbx], xmm3 movsxd rbx, [rdx + 4*rcx + 20] vextractf128 xmm0, ymm0, 1 vmovss [rsi + 4*rax], xmm0 movsxd rax, [rdx + 4*rcx + 24] vpsalignr xmm1, xmm0, xmm0, 4 vmovss [rsi + 4*rbx], xmm1 movsxd rbx, [rdx + 4*rcx + 28] vpsalignr xmm2, xmm0, xmm0, 8 vmovss [rsi + 4*rax], xmm2 vpsalignr xmm3, xmm0, xmm0, 12 vmovss [rsi + 4*rbx], xmm3 add rcx, 8 cmp rcx, len jl loop1 </pre>

11.6 DATA ALIGNMENT FOR INTEL® AVX

This section explains the benefit of aligning data that is used by Intel AVX instructions and proposes some methods to improve performance when such alignment is not possible. Most examples in this section are variations of the SAXPY kernel. SAXPY is the Scalar Alpha * X + Y algorithm.

The C code below is a C implementation of SAXPY.

```

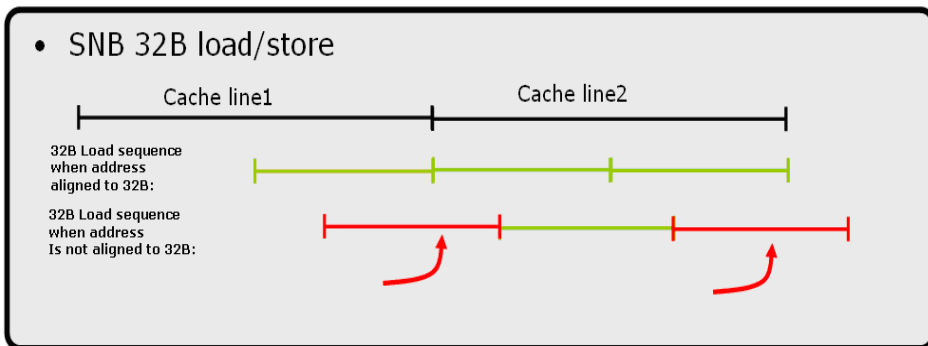
for (int i = 0; i < n; i++)
{ c[i] = alpha * a[i] + b[i]; }

```

11.6.1 Align Data to 32 Bytes

Aligning data to vector length is recommended. When using 16-byte SIMD instructions, loaded data should be aligned to 16 bytes. Similarly, for best results when using Intel AVX instructions with 32-byte registers align the data to 32-bytes.

When using Intel AVX with unaligned 32-byte vectors, every second load is a cache-line split, since the cache-line is 64 bytes. This doubles the cache line split rate compared to Intel SSE code that uses 16-byte vectors. Even though split line access penalties were reduced significantly since Intel microarchitecture code name Nehalem, a high cache-line split rate in memory-intensive code may cause performance degradation.



Example 11-10. SAXPY using Intel AVX

```

mov    rax, src1
mov    rbx, src2
mov    rcx, dst
mov    rdx, len
xor    rdi, rdi
vbroadcastss    ymm0, alpha
start_loop:
vmovups    ymm1, [ rax + rdi ]
vmulps    ymm1, ymm1, ymm0
vmovups    ymm2, [ rbx + rdi ]
vaddps    ymm1, ymm1, ymm2
vmovups    [ rcx + rdi ], ymm1

```

Example 11-10. SAXPY using Intel AVX (Contd.)

```

vmovups    ymm1, [ rax + rdi + 32 ]
vmulps     ymm1, ymm1, ymm0
vmovups    ymm2, [ rbx + rdi + 32]
vaddps     ymm1, ymm1, ymm2
vmovups    [ rcx + rdi + 32 ], ymm1

add        rdi, 64
cmp        rdi, rdx
jl         start_loop

```

SAXPY is a memory intensive kernel that emphasizes the importance of data alignment. Optimal performance requires both data source address to be 32-byte aligned and destination address also 32-byte aligned. If only one of the three address is not aligned to 32-byte boundary, the performance may be halved. If all three addresses are mis-aligned relative to 32 byte, the performance degrades further. In some cases, unaligned accesses may result in lower performance for Intel AVX code compared to Intel SSE code. Other Intel AVX kernels typically have more computation which can reduce the effect of the data alignment penalty.

Assembly/Compiler Coding Rule 74. (H impact, M generality) *Align data to 32-byte boundary when possible. Prefer store alignment over load alignment.*

You can use dynamic data alignment using the `_mm_malloc` intrinsic instruction with the Intel® Compiler, or `_aligned_malloc` of the Microsoft* Compiler. For example:

```
//dynamically allocating 32byte aligned buffer with 2048 float elements.
```

```
InputBuffer = (float*) _mm_malloc (2048*sizeof(float), 32);
```

You can use static data alignment using `__declspec(align(32))`. For example:

```
//Statically allocating 32byte aligned buffer with 2048 float elements.
```

```
__declspec(align(32)) float InputBuffer[2048];
```

11.6.2 Consider 16-Byte Memory Access when Memory is Unaligned

For best results use Intel AVX 32-byte loads and align data to 32-bytes. However, there are cases where you cannot align the data, or data alignment is unknown. This can happen when you are writing a library function and the input data alignment is unknown. In these cases, using 16-Byte memory accesses may be the best alternative. The following method uses 16-byte loads while still benefiting from the 32-byte YMM registers.

Consider replacing unaligned 32-byte memory accesses using a combination of `VMOVUPS`, `VINSERTF128`, and `VEXTRACTF128` instructions.

Example 11-11. Using 16-Byte Memory Operations for Unaligned 32-Byte Memory Operation

Convert 32-byte loads as follows:

```
vmovups   ymm0, mem  -> vmovups xmm0, mem
                               vinsertf128 ymm0, ymm0, mem+16, 1
```

Convert 32-byte stores as follows:

```
vmovups mem, ymm0   -> vmovups mem, xmm0
                               vextractf128 mem+16, ymm0, 1
```

The following intrinsics are available to handle unaligned 32-byte memory operating using 16-byte memory accesses:

```
_mm256_loadu2_m128 ( float const * addr_hi, float const * addr_lo);
_mm256_loadu2_m128d ( double const * addr_hi, double const * addr_lo);
_mm256_loadu2_m128i ( __m128i const * addr_hi, __m128i const * addr_lo);
_mm256_storeu2_m128 ( float * addr_hi, float * addr_lo, __m256 a);
_mm256_storeu2_m128d ( double * addr_hi, double * addr_lo, __m256d a);
_mm256_storeu2_m128i ( __m128i * addr_hi, __m128i * addr_lo, __m256i a);
```

Example 11-12 shows two implementations for SAXPY with unaligned addresses. Alternative 1 uses 32 byte loads and alternative 2 uses 16 byte loads. These code samples are executed with two source buffers, src1, src2, at 4 byte offset from 32-Byte alignment, and a destination buffer, DST, that is 32-Byte aligned. Using two 16-byte memory operations in lieu of 32-byte memory access performs faster.

Example 11-12. SAXPY Implementations for Unaligned Data Addresses

AVX with 32-byte memory operation	AVX using two 16-byte memory operations
<pre>mov rax, src1 mov rbx, src2 mov rcx, dst mov rdx, len xor rdi, rdi vbroadcastss ymm0, alpha start_loop: vmovups ymm1, [rax + rdi] vmulps ymm1, ymm1, ymm0 vmovups ymm2, [rbx + rdi] vaddps ymm1, ymm1, ymm2 vmovups [rcx + rdi], ymm1</pre>	<pre>mov rax, src1 mov rbx, src2 mov rcx, dst mov rdx, len xor rdi, rdi vbroadcastss ymm0, alpha start_loop: vmovups xmm2, [rax+rdi] vinsertf128 ymm2, ymm2, [rax+rdi+16], 1 vmulps ymm1, ymm0, ymm2 vmovups xmm2, [rbx + rdi] vinsertf128 ymm2, ymm2, [rbx+rdi+16], 1 vaddps ymm1, ymm1, ymm2</pre>

Example 11-12. SAXPY Implementations for Unaligned Data Addresses (Contd.)

AVX with 32-byte memory operation	AVX using two 16-byte memory operations
<pre> vmovups ymm1, [rax+rdi+32] vmulps ymm1, ymm1, ymm0 vmovups ymm2, [rbx+rdi+32] vaddps ymm1, ymm1, ymm2 vmovups [rcx+rdi+32], ymm1 add rdi, 64 cmp rdi, rdx jl start_loop </pre>	<pre> vaddps ymm1, ymm1, ymm2 vmovaps [rcx+rdi], ymm1 vmovups xmm2, [rax+rdi+32] vinsertf128 ymm2, ymm2, [rax+rdi+48], 1 vmulps ymm1, ymm0, ymm2 vmovups xmm2, [rbx+rdi+32] vinsertf128 ymm2, ymm2, [rbx+rdi+48], 1 vaddps ymm1, ymm1, ymm2 vmovups [rcx+rdi+32], ymm1 add rdi, 64 cmp rdi, rdx jl start_loop </pre>

Assembly/Compiler Coding Rule 75. (M impact, H generality) *Align data to 32-byte boundary when possible. Prefer store alignment over load alignment.*

11.6.3 Prefer Aligned Stores Over Aligned Loads

There are cases where it is possible to align only a subset of the processed data buffers. In these cases, aligning data buffers used for store operations usually yields better performance than aligning data buffers used for load operations.

Unaligned stores are likely to cause greater performance degradation than unaligned loads, since there is a very high penalty on stores to a split cache-line that crosses pages. This penalty is estimated at 150 cycles. Loads that cross a page boundary are executed at retirement. In Example 11-12, unaligned store address can affect SAXPY performance for 3 unaligned addresses to about one quarter of the aligned case.

11.7 L1D CACHE LINE REPLACEMENTS

When a load misses the L1D Cache, a cache line with the requested data is brought from a higher memory hierarchy level. In memory intensive code where the L1 DCache is always active, replacing a cache line in the L1 DCache may delay other loads. The penalty for 32-Byte loads may be higher than the penalty for 16-Byte loads. Therefore, memory intensive Intel AVX code with 32-Byte loads and with data set larger than the L1 DCache may be slower than similar code with 16-Byte loads.

When Example 11-12 is run with a data set that resides in the L2 Cache, the 16-byte memory access implementation is slightly faster than the 32-byte memory operation.

Be aware that the relative merit of 16-Byte memory accesses versus 32-byte memory access is implementation specific across generations of microarchitectures.

11.8 4K ALIASING

4-KByte memory aliasing occurs when the code stores to one memory location and shortly after that it loads from a different memory location with a 4-KByte offset between them. For example, a load to linear address 0x400020 follows a store to linear address 0x401020.

The load and store have the same value for bits 5 - 11 of their addresses and the accessed byte offsets should have partial or complete overlap.

4K aliasing may have a five-cycle penalty on the load latency. This penalty may be significant when 4K aliasing happens repeatedly and the loads are on the critical path. If the load spans two cache lines it might be delayed until the conflicting store is committed to the cache. Therefore 4K aliasing that happens on repeated unaligned Intel AVX loads incurs a higher performance penalty.

To detect 4K aliasing, use the `LD_BLOCKS_PARTIAL.ADDRESS_ALIAS` event that counts the number of times Intel AVX loads were blocked due to 4K aliasing.

To resolve 4K aliasing, try the following methods in the following order:

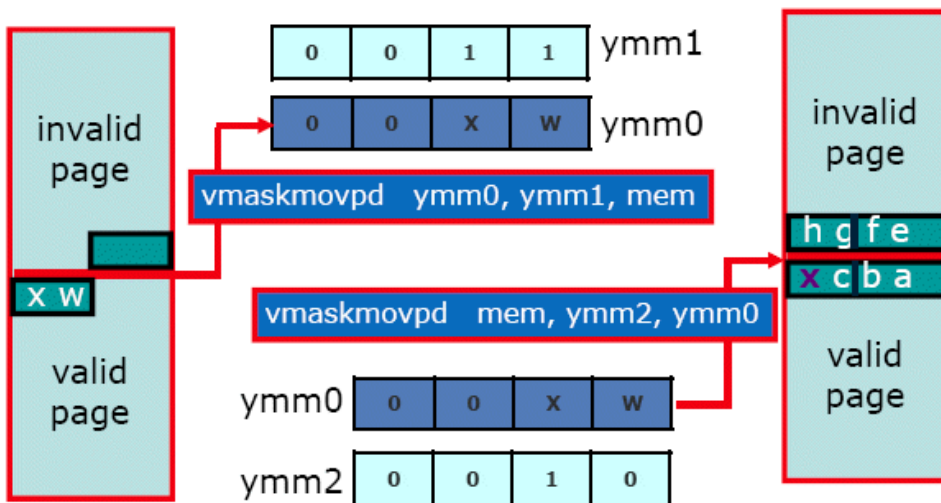
- Align data to 32 Bytes.
- Change offsets between input and output buffers if possible.
- Use 16-Byte memory accesses on memory which is not 32-Byte aligned.

11.9 CONDITIONAL SIMD PACKED LOADS AND STORES

The `VMASKMOV` instruction conditionally moves packed data elements to/from memory, depending on the mask bits associated with each data element. The mask bit for each data element is the most significant bit of the corresponding element in the mask register.

When performing a mask load, the returned value is 0 for elements which have a corresponding mask value of 0. The mask store instruction writes to memory only the elements with a corresponding mask value of 1, while preserving memory values for elements with a corresponding mask value of 0. Faults can occur only for memory accesses that are required by the mask. Faults do not occur due to referencing any memory location if the corresponding mask bit value for that memory location is zero. For example, no faults are detected if the mask bits are all zero.

The following figure shows an example for a mask load and a mask store which does not cause a fault. In this example, the mask register for the load operation is `ymm1` and the mask register for the store operation is `ymm2`.



When using masked load or store consider the following:

- The address of a VMASKMOV store is considered as resolved only after the mask is known. Loads that follow a masked store can be blocked until the mask value is known (unless relieved by the memory disambiguator).
- If the mask is not all 1 or all 0, loads that depend on the masked store have to wait until the store data is written to the cache. If the mask is all 1 the data can be forwarded from the masked store to the dependent loads. If the mask is all 0 the loads do not depend on the masked store.
- Masked loads including an illegal address range do not result in an exception if the range is under a zero mask value. However, the processor may take a multi-cycle "assist" to determine that no part of the illegal range have a one mask value. This assist may occur even when the mask is "zero" and it seems obvious to the programmer that the load should not be executed.

When using VMASKMOV, consider the following:

- Use VMASKMOV only in cases where VMOVUPS cannot be used.
- Use VMASKMOV on 32Byte aligned addresses if possible.
- If possible use valid address range for masked loads, even if the illegal part is masked with zeros.
- Determine the mask as early as possible.
- Avoid store-forwarding issues by performing loads prior to a VMASKMOV store if possible.

11.9.1 Conditional Loops

VMASKMOV enables vectorization of loops that contain conditional code. There are two main benefits in using VMASKMOV over the scalar implementation in these cases.

- VMASKMOV code is vectorized
- Branch mispredictions are eliminated

Below is a conditional loop C code:

Example 11-13. Loop with Conditional Expression

```
for(int i = 0; i < miBufferWidth; i++)
{
    if(A[i]>0)
    {
        B[i] = (E[i]*C[i]);
    }
    else
    {
        B[i] = (E[i]*D[i]);
    }
}
```

Example 11-14. Handling Loop Conditional with VMASKMOV

Scalar	AVX using VMASKMOV
<pre>float* pA = A; float* pB = B; float* pC = C; float* pD = D; float* pE = E; uint64 len = (uint64) (miBuffer- Width)*sizeof(float); __asm { mov rax, pA mov rbx, pB mov rcx, pC mov rdx, pD mov rsi, pE mov r8, len</pre>	<pre>float* pA = A; float* pB = B; float* pC = C; float* pD = D; float* pE = E; uint64 len = (uint64) (miBufferWidth)*sizeof(float); __asm { mov rax, pA mov rbx, pB mov rcx, pC mov rdx, pD mov rsi, pE mov r8, len</pre>

Example 11-14. Handling Loop Conditional with VMASKMOV (Contd.)

Scalar	AVX using VMASKMOV
<pre>//xmm8 all zeros vxorps xmm8, xmm8, xmm8 xor r9,r9 loop1: vmovss xmm1, [rax+r9] vcomiss xmm1, xmm8 jbe a_le a_gt: vmovss xmm4, [rcx+r9] jmp mul a_le: vmovss xmm4, [rdx+r9] mul: vmulss xmm4, xmm4, [rsi+r9] vmovss [rbx+r9], xmm4 add r9, 4 cmp r9, r8 jl loop1 }</pre>	<pre>//ymm8 all zeros vxorps ymm8, ymm8, ymm8 //ymm9 all ones vcmppps ymm9, ymm8, ymm8, 0 xor r9,r9 loop1: vmovups ymm1, [rax+r9] vcmppps ymm2, ymm8, ymm1, 1 vmaskmovps ymm4, ymm2, [rcx+r9] vxorps ymm2, ymm2, ymm9 vmaskmovps ymm5, ymm2, [rdx+r9] vorps ymm4, ymm4, ymm5 vmulps ymm4, ymm4, [rsi+r9] vmovups [rbx+r9], ymm4 add r9, 32 cmp r9, r8 jl loop1 }</pre>

The performance of the left side of Example 11-14 is sensitive to branch mis-predictions and can be an order of magnitude slower than the VMASKMOV example which has no data-dependent branches.

11.10 MIXING INTEGER AND FLOATING-POINT CODE

Integer SIMD functionalities in Intel AVX instructions are limited to 128-bit. There are some algorithm that uses mixed integer SIMD and floating-point SIMD instructions. Therefore, porting such legacy 128-bit code into 256-bit AVX code requires special attention.

For example, PALINGR (Packed Align Right) is an integer SIMD instruction that is useful arranging data elements for integer and floating-point code. But VPALINGR instruction does not have a corresponding 256-bit instruction in AVX.

There are two approaches to consider when porting legacy code consisting of mostly floating point with some integer operations into 256-bit AVX code:.

- Locate a 256-bit AVX alternative to replace the critical 128-bit Integer SIMD instructions if such an AVX instructions exist. This is more likely to be true with integer SIMD instruction that re-arranges data elements.
- Mix 128-bit AVX and 256-bit AVX instructions

The performance gain from these two approaches may vary. Where possible, use method (1), since this method utilizes the full 256-bit vector width.

In case the code is mostly integer, convert the code from 128-bit SSE to 128 bit AVX instructions and gain from the Non destructive Source (NDS) feature.

Example 11-15. Three-Tap Filter in C Code

```
for(int i = 0; i < len -2; i++)
{
    pOut[i] = A[i]*coeff[0]+A[i+1]*coeff[1]+A[i+2]*coeff[2];B[i] = (E[i]*D[i]);
}
```

Example 11-16. Three-Tap Filter with 128-bit Mixed Integer and FP SIMD

```
xor ebx, ebx
mov rcx, len
mov rdi, inPtr
mov rsi, outPtr
mov r15, coeffs
movss xmm2, [r15] //load coeff 0
shufps xmm2, xmm2, 0 //broadcast coeff 0
movss xmm1, [r15+4] //load coeff 1
shufps xmm1, xmm1, 0 //broadcast coeff 1
movss xmm0, [r15+8] //coeff 2
shufps xmm0, xmm0, 0 //broadcast coeff 2
movaps xmm5, [rdi] //xmm5={A[n+3],A[n+2],A[n+1],A[n]}
```

loop_start:

```
movaps xmm6, [rdi+16] //xmm6={A[n+7],A[n+6],A[n+5],A[n+4]}
movaps xmm7, xmm6
movaps xmm8, xmm6
add rdi, 16 //inPtr+=32
add rbx, 4 //loop counter
palignr xmm7, xmm5, 4 //xmm7={A[n+4],A[n+3],A[n+2],A[n+1]}
palignr xmm8, xmm5, 8 //xmm8={A[n+5],A[n+4],A[n+3],A[n+2]}
mulps xmm5, xmm2 //xmm5={C0*A[n+3],C0*A[n+2],C0*A[n+1], C0*A[n]}
```

Example 11-16. Three-Tap Filter with 128-bit Mixed Integer and FP SIMD (Contd.)

```

mulps xmm7,xmm1 //xmm7={C1*A[n+4],C1*A[n+3],C1*A[n+2],C1*A[n+1]}
mulps xmm8,xmm0 //xmm8={C2*A[n+5],C2*A[n+4],C2*A[n+3],C2*A[n+2]}
addps xmm7,xmm5
addps xmm7,xmm8
movaps [rsi],xmm7
movaps xmm5,xmm6
add rsi,16 //outPtr+=16
cmp rbx,rcx
jl loop_start

```

Example 11-17. 256-bit AVX Three-Tap Filter Code with VSHUFPS

```

xor ebx,ebx
mov rcx,len
mov rdi,inPtr
mov rsi,outPtr
mov r15,coeffs
vbroadcastss ymm2,[r15] //load and broadcast coeff 0
vbroadcastss ymm1,[r15+4] //load and broadcast coeff 1
vbroadcastss ymm0,[r15+8] //load and broadcast coeff 2
loop_start:
vmovaps ymm5,[rdi] // Ymm5={A[n+7],A[n+6],A[n+5],A[n+4];
// A[n+3],A[n+2],A[n+1],A[n]}
vshufps ymm6,ymm5,[rdi+16],0x4e // ymm6={A[n+9],A[n+8],A[n+7],A[n+6];
// A[n+5],A[n+4],A[n+3],A[n+2]}
vshufps ymm7,ymm5,ymm6,0x99 // ymm7={A[n+8],A[n+7],A[n+6],A[n+5];
// A[n+4],A[n+3],A[n+2],A[n+1]}
vmulps ymm3,ymm5,ymm2 // ymm3={C0*A[n+7],C0*A[n+6],C0*A[n+5],C0*A[n+4];
// C0*A[n+3],C0*A[n+2],C0*A[n+1],C0*A[n]}
vmulps ymm9,ymm7,ymm1 // ymm9={C1*A[n+8],C1*A[n+7],C1*A[n+6],C1*A[n+5];
// C1*A[n+4],C1*A[n+3],C1*A[n+2],C1*A[n+1]}
vmulps ymm4,ymm6,ymm0 // ymm4={C2*A[n+9],C2*A[n+8],C2*A[n+7],C2*A[n+6];
// C2*A[n+5],C2*A[n+4],C2*A[n+3],C2*A[n+2]}
vaddps ymm8,ymm3,ymm4
vaddps ymm10,ymm8,ymm9
vmovaps [rsi],ymm10

```

Example 11-17. 256-bit AVX Three-Tap Filter Code with VSHUFPS (Contd.)

```

add   rdi, 32 //inPtr+=32
add   rbx, 8 //loop counter
add   rsi, 32 //outPtr+=32
cmp   rbx, rcx
jl    loop_start

```

Example 11-18. Three-Tap Filter Code with Mixed 256-bit AVX and 128-bit AVX Code

```

xor   ebx, ebx
mov   rcx, len
mov   rdi, inPtr
mov   rsi, outPtr
mov   r15, coeffs
vbroadcastss ymm2, [r15] //load and broadcast coeff 0
vbroadcastss ymm1, [r15+4] //load and broadcast coeff 1
vbroadcastss ymm0, [r15+8] //load and broadcast coeff 2
vmovaps xmm3, [rdi] //xmm3={A[n+3],A[n+2],A[n+1],A[n]}
vmovaps xmm4, [rdi+16] //xmm4={A[n+7],A[n+6],A[n+5],A[n+4]}
vmovaps xmm5, [rdi+32] //xmm5={A[n+11], A[n+10],A[n+9],A[n+8]}
loop_start:
vinsertf128 ymm3, ymm3, xmm4, 1 // ymm3={A[n+7],A[n+6],A[n+5],A[n+4];
// A[n+3], A[n+2],A[n+1],A[n]}
vpsalignr xmm6, xmm4, xmm3, 4 // xmm6={A[n+4],A[n+3],A[n+2],A[n+1]}
vpsalignr xmm7, xmm5, xmm4, 4 // xmm7={A[n+8],A[n+7],A[n+6],A[n+5]}
vinsertf128 ymm6, ymm6, xmm7, 1 // ymm6={A[n+8],A[n+7],A[n+6],A[n+5];
// A[n+4],A[n+3],A[n+2],A[n+1]}
vpsalignr xmm8, xmm4, xmm3, 8 // xmm8={A[n+5],A[n+4],A[n+3],A[n+2]}
vpsalignr xmm9, xmm5, xmm4, 8 // xmm9={A[n+9],A[n+8],A[n+7],A[n+6]}
vinsertf128 ymm8, ymm8, xmm9, 1 // ymm8={A[n+9],A[n+8],A[n+7],A[n+6];
// A[n+5],A[n+4],A[n+3],A[n+2]}
vmulps ymm3, ymm3, ymm2 // Ymm3={C0*A[n+7],C0*A[n+6],C0*A[n+5], C0*A[n+4];
// C0*A[n+3],C0*A[n+2],C0*A[n+1],C0*A[n]}
vmulps ymm6, ymm6, ymm1 // Ymm9={C1*A[n+8],C1*A[n+7],C1*A[n+6],C1*A[n+5];
// C1*A[n+4],C1*A[n+3],C1*A[n+2],C1*A[n+1]}
vmulps ymm8, ymm8, ymm0 // Ymm4={C2*A[n+9],C2*A[n+8],C2*A[n+7],C2*A[n+6];
// C2*A[n+5],C2*A[n+4],C2*A[n+3],C2*A[n+2]}
vaddps ymm3, ymm3, ymm6
vaddps ymm3, ymm3, ymm8
vmovaps [rsi], ymm3

```

Example 11-18. Three-Tap Filter Code with Mixed 256-bit AVX and 128-bit AVX Code (Contd.)

```

vmovaps xmm3, xmm5
add    rdi, 32 //inPtr+=32
add    rbx, 8 //loop counter
add    rsi, 32 //outPtr+=32
cmp    rbx, rcx
jl     loop_start

```

Example 11-17 uses 256-bit VSHUFPS to replace the PALIGNR in 128-bit mixed SSE code. This speeds up almost 70% over the 128-bit mixed SSE code of Example 11-16 and slightly ahead of Example 11-18.

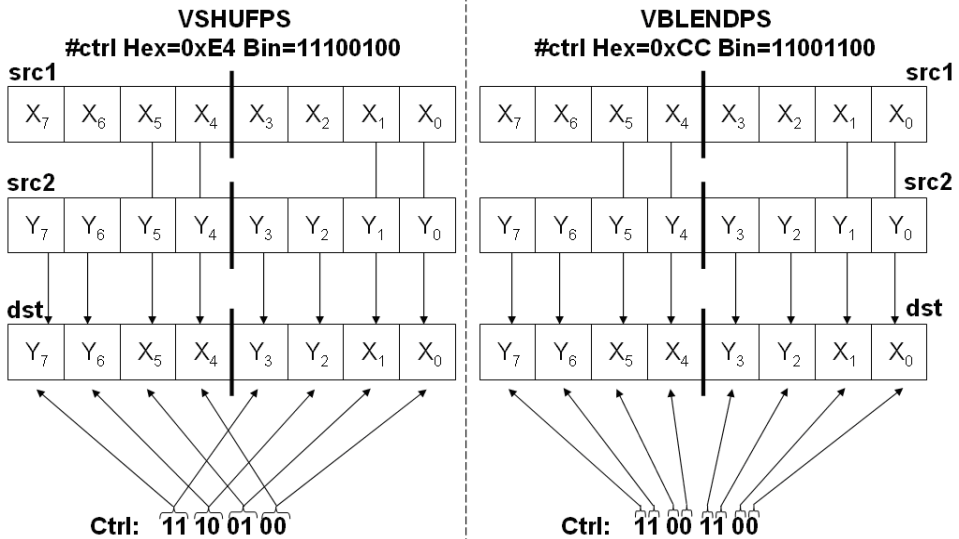
For code that includes integer instructions and is written with 256-bit Intel AVX instructions, replace the integer instruction with floating-point instructions that have similar functionality and performance. If there is no similar floating-point instruction, consider using a 128-bit Intel AVX instruction to perform the required integer operation.

11.11 HANDLING PORT 5 PRESSURE

Port 5 in Intel microarchitecture code name Sandy Bridge includes shuffle units and it frequently becomes a performance bottleneck. Sometimes it is possible to replace shuffle instructions that dispatch only on port 5, with different instructions and improve performance by reducing port 5 pressure. For more information, see Table 2-4.

11.11.1 Replace Shuffles with Blends

There are a few cases where shuffles such as VSHUFPS or VPERM2F128 can be replaced by blend instructions. Intel AVX shuffles are executed only on port 5, while blends are also executed on port 0. Therefore, replacing shuffles with blends could reduce port 5 pressure. The following figure shows how a VSHUFPS is implemented using VBLENDPS.



The following example shows two implementations of an 8x8 Matrix transpose. In both cases, the bottleneck is Port 5 pressure. Alternative 1 uses 12 vshufps instructions that are executed only on port 5. Alternative 2 replaces eight of the vshufps instructions with the vblendps instruction which can be executed on Port 0.

Example 11-19. 8x8 Matrix Transpose - Replace Shuffles with Blends

256-bit AVX using VSHUFPS	AVX replacing VSHUFPS with VBLENDPS
<pre>movrcx, inpBuf movrdx, outBuf movr10, NumOfLoops movrbx, rdx loop1: vmovaps ymm9, [rcx] vmovaps ymm10, [rcx+32] vmovaps ymm11, [rcx+64] vmovaps ymm12, [rcx+96] vmovaps ymm13, [rcx+128] vmovaps ymm14, [rcx+160] vmovaps ymm15, [rcx+192] vmovaps ymm2, [rcx+224]</pre>	<pre>movrcx, inpBuf movrdx, outBuf movr10, NumOfLoops movrbx, rdx loop1: vmovaps ymm9, [rcx] vmovaps ymm10, [rcx+32] vmovaps ymm11, [rcx+64] vmovaps ymm12, [rcx+96] vmovaps ymm13, [rcx+128] vmovaps ymm14, [rcx+160] vmovaps ymm15, [rcx+192] vmovaps ymm2, [rcx+224]</pre>

Example 11-19. 8x8 Matrix Transpose - Replace Shuffles with Blends (Contd.)

256-bit AVX using VSHUFPS	AVX replacing VSHUFPS with VBLENDPS
vunpcklps ymm6, ymm9, ymm10	vunpcklps ymm6, ymm9, ymm10
vunpcklps ymm1, ymm11, ymm12	vunpcklps ymm1, ymm11, ymm12
vunpckhps ymm8, ymm9, ymm10	vunpckhps ymm8, ymm9, ymm10
vunpcklps ymm0, ymm13, ymm14	vunpcklps ymm0, ymm13, ymm14
vunpcklps ymm9, ymm15, ymm2	vunpcklps ymm9, ymm15, ymm2
vshufps ymm3, ymm6, ymm1, 0x4E	vshufps ymm3, ymm6, ymm1, 0x4E
vshufps ymm10, ymm6, ymm3, 0xE4	vblendps ymm10, ymm6, ymm3, 0xCC
vshufps ymm6, ymm0, ymm9, 0x4E	vshufps ymm6, ymm0, ymm9, 0x4E
vunpckhps ymm7, ymm11, ymm12	vunpckhps ymm7, ymm11, ymm12
vshufps ymm11, ymm0, ymm6, 0xE4	vblendps ymm11, ymm0, ymm6, 0xCC
vshufps ymm12, ymm3, ymm1, 0xE4	vblendps ymm12, ymm3, ymm1, 0xCC
vperm2f128 ymm3, ymm10, ymm11, 0x20	vperm2f128 ymm3, ymm10, ymm11, 0x20
vmovaps [rdx], ymm3	vmovaps [rdx], ymm3
vunpckhps ymm5, ymm13, ymm14	vunpckhps ymm5, ymm13, ymm14
vshufps ymm13, ymm6, ymm9, 0xE4	vblendps ymm13, ymm6, ymm9, 0xCC
vunpckhps ymm4, ymm15, ymm2	vunpckhps ymm4, ymm15, ymm2
vperm2f128 ymm2, ymm12, ymm13, 0x20	vperm2f128 ymm2, ymm12, ymm13, 0x20
vmovaps 32[rdx], ymm2	vmovaps 32[rdx], ymm2
vshufps ymm14, ymm8, ymm7, 0x4	vshufps ymm14, ymm8, ymm7, 0x4E
vshufps ymm15, ymm14, ymm7, 0xE4	vblendps ymm15, ymm14, ymm7, 0xCC
vshufps ymm7, ymm5, ymm4, 0x4E	vshufps ymm7, ymm5, ymm4, 0x4E
vshufps ymm8, ymm8, ymm14, 0xE4	vblendps ymm8, ymm8, ymm14, 0xCC
vshufps ymm5, ymm5, ymm7, 0xE4	vblendps ymm5, ymm5, ymm7, 0xCC
vperm2f128 ymm6, ymm8, ymm5, 0x20	vperm2f128 ymm6, ymm8, ymm5, 0x20
vmovaps 64[rdx], ymm6	vmovaps 64[rdx], ymm6
vshufps ymm4, ymm7, ymm4, 0xE4	vblendps ymm4, ymm7, ymm4, 0xCC
vperm2f128 ymm7, ymm15, ymm4, 0x20	vperm2f128 ymm7, ymm15, ymm4, 0x20
vmovaps 96[rdx], ymm7	vmovaps 96[rdx], ymm7
vperm2f128 ymm1, ymm10, ymm11, 0x31	vperm2f128 ymm1, ymm10, ymm11, 0x31
vperm2f128 ymm0, ymm12, ymm13, 0x31	vperm2f128 ymm0, ymm12, ymm13, 0x31
vmovaps 128[rdx], ymm1	vmovaps 128[rdx], ymm1
vperm2f128 ymm5, ymm8, ymm5, 0x31	vperm2f128 ymm5, ymm8, ymm5, 0x31
vperm2f128 ymm4, ymm15, ymm4, 0x31	vperm2f128 ymm4, ymm15, ymm4, 0x31
vmovaps 160[rdx], ymm0	vmovaps 160[rdx], ymm0
vmovaps 192[rdx], ymm5	vmovaps 192[rdx], ymm5
vmovaps 224[rdx], ymm4	vmovaps 224[rdx], ymm4
decr 10	dec r10
jnz loop1	jnz loop1

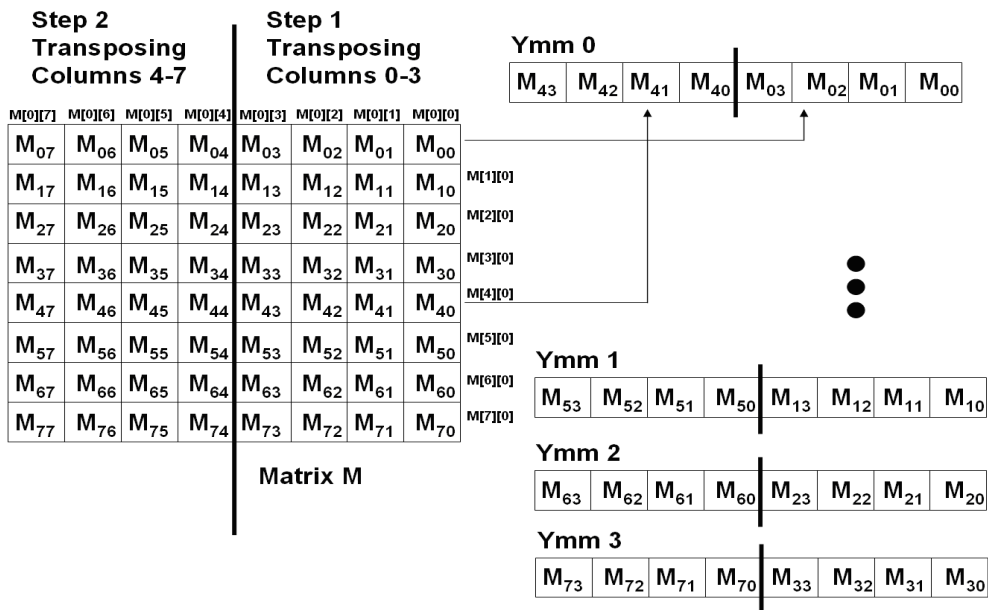
In Example 11-19, replacing VSHUFPS with VBLENDPS relieved port 5 pressure and can gain almost 40% speedup.

Assembly/Compiler Coding Rule 76. (M impact, M generality) Use Blend instructions in lieu of shuffle instruction in AVX whenever possible.

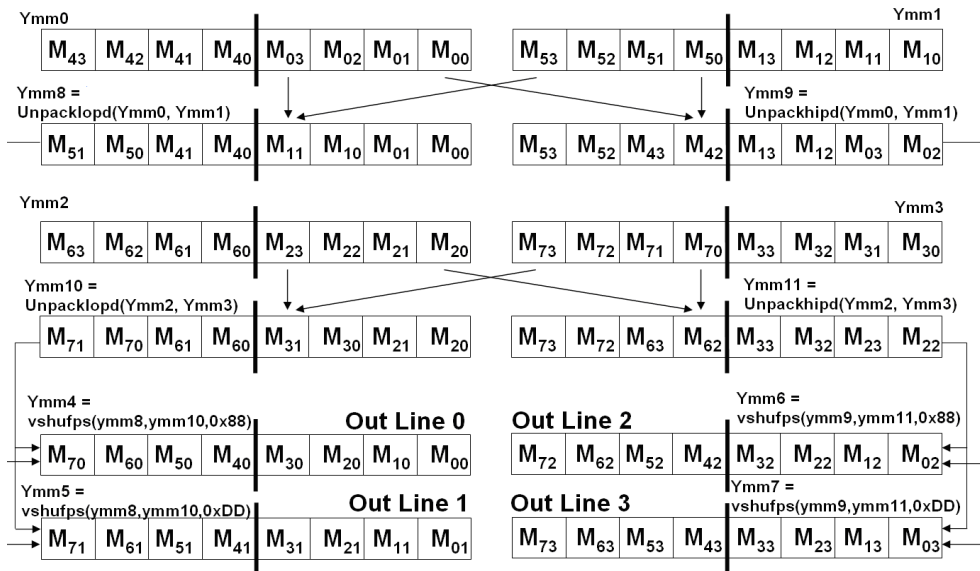
11.11.2 Design Algorithm With Fewer Shuffles

In some cases you can reduce port 5 pressure by changing the algorithm to use less shuffles. The figure below shows that the transpose moved all the elements in rows 0-4 to the low lanes, and all the elements in rows 4-7 to the high lanes. Therefore, using 256-bit loads in the beginning of the algorithm requires using VPERM2F128 in order to swap elements between the lanes. The processor executes the VPERM2F128 instruction only on port 5.

Example 11-19 used eight 256-bit loads and eight VPERM2F128 instructions. You can implement the same 8x8 Matrix Transpose using VINSERTF128 instead of the 256-bit loads and the eight VPERM2F128. Using VINSERTF128 from memory is executed in the load ports and on port 0 or 5. The original method required loads that are performed on the load ports and VPERM2F128 that is only performed on port 5. Therefore redesigning the algorithm to use VINSERTF128 reduces port 5 pressure and improves performance.



The following figure describes step 1 of the 8x8 matrix transpose with vinsertf128. Step 2 performs the same operations on different columns.



Example 11-20. 8x8 Matrix Transpose Using VINSRTPS

```

mov    rcx, inpBuf
mov    rdx, outBuf
mov    r8, iLineSize
mov    r10, NumOfLoops
loop1:
vmovaps  xmm0, [rcx]
vinsertf128 ymm0, ymm0, [rcx + 128], 1
vmovaps  xmm1, [rcx + 32]
vinsertf128 ymm1, ymm1, [rcx + 160], 1

vunpcklpd  ymm8, ymm0, ymm1
vunpckhpd  ymm9, ymm0, ymm1
vmovaps  xmm2, [rcx+64]
vinsertf128 ymm2, ymm2, [rcx + 192], 1
    
```

Example 11-20. 8x8 Matrix Transpose Using VINSRTPS (Contd.)

```

vmovaps    xmm3, [rcx+96]
vinsertf128 ymm3, ymm3, [rcx + 224], 1

vunpcklpd  ymm10, ymm2, ymm3
vunpckhpd  ymm11, ymm2, ymm3

vshufps    ymm4, ymm8, ymm10, 0x88
vmovaps    [rdx], ymm4
vshufps    ymm5, ymm8, ymm10, 0xDD
vmovaps    [rdx+32], ymm5
vshufps    ymm6, ymm9, ymm11, 0x88
vmovaps    [rdx+64], ymm6
vshufps    ymm7, ymm9, ymm11, 0xDD
vmovaps    [rdx+96], ymm7
vmovaps    xmm0, [rcx+16]
vinsertf128 ymm0, ymm0, [rcx + 144], 1
vmovaps    xmm1, [rcx + 48]
vinsertf128 ymm1, ymm1, [rcx + 176], 1

vunpcklpd  ymm8, ymm0, ymm1
vunpckhpd  ymm9, ymm0, ymm1

vmovaps    xmm2, [rcx+80]
vinsertf128 ymm2, ymm2, [rcx + 208], 1
vmovaps    xmm3, [rcx+112]
vinsertf128 ymm3, ymm3, [rcx + 240], 1

vunpcklpd  ymm10, ymm2, ymm3
vunpckhpd  ymm11, ymm2, ymm3

vshufps    ymm4, ymm8, ymm10, 0x88
vmovaps    [rdx+128], ymm4
vshufps    ymm5, ymm8, ymm10, 0xDD
vmovaps    [rdx+160], ymm5
vshufps    ymm6, ymm9, ymm11, 0x88
vmovaps    [rdx+192], ymm6
vshufps    ymm7, ymm9, ymm11, 0xDD
vmovaps    [rdx+224], ymm7
dec        r10
jnz        loop1

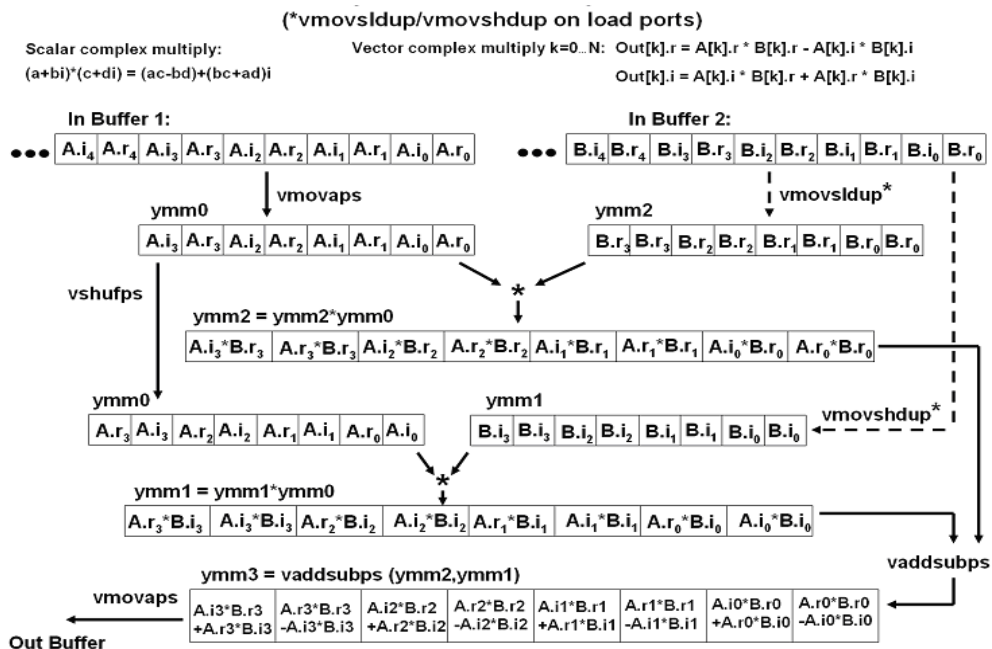
```

In Example 11-20, this reduced port 5 pressure further than the combination of VSHUFPS with VBLENDPS in Example 11-19. It can gain 70% speedup relative to relying on VSHUFPS alone in Example 11-19.

11.11.3 Perform Basic Shuffles on Load Ports

Some shuffles can be executed in the load ports (ports 2, 3) if the source is from memory. The following example shows how moving some shuffles (vmovsldup/vmovshdup) from Port 5 to the load ports improves performance significantly.

The following figure describes an Intel AVX implementation of the complex multiply algorithm with vmovsldup/vmovshdup on the load ports.



Example 11-21 includes two versions of the complex multiply. Both versions are unrolled twice. Alternative 1 shuffles all the data in registers. Alternative 2 shuffles data while it is loaded from memory.

Example 11-21. Port 5 versus Load Port Shuffles

Shuffles data in registers	Shuffling loaded data
<pre> mov rax, inPtr1 mov rbx, inPtr2 mov rdx, outPtr mov r8, len xor rcx, rcx loop1: vmovaps ymm0, [rax +8*rcx] vmovaps ymm4, [rax +8*rcx +32] vmovaps ymm3, [rbx +8*rcx] vmovsldup ymm2, ymm3 vmulps ymm2, ymm2, ymm0 vshufps ymm0, ymm0, ymm0, 177 vmovshdup ymm1, ymm3 vmulps ymm1, ymm1, ymm0 vmovaps ymm7, [rbx +8*rcx +32] vmovsldup ymm6, ymm7 vmulps ymm6, ymm6, ymm4 vaddsubps ymm2, ymm2, ymm1 vmovshdup ymm5, ymm7 vmovaps [rdx+8*rcx], ymm2 vshufps ymm4, ymm4, ymm4, 177 vmulps ymm5, ymm5, ymm4 vaddsubps ymm6, ymm6, ymm5 vmovaps [rdx+8*rcx+32], ymm6 addrcx, 8 cmprcx, r8 jl loop1 </pre>	<pre> mov rax, inPtr1 mov rbx, inPtr2 mov rdx, outPtr mov r8, len xor rcx, rcx loop1: vmovaps ymm0, [rax +8*rcx] vmovaps ymm4, [rax +8*rcx +32] vmovsldup ymm2, [rbx +8*rcx] vmulps ymm2, ymm2, ymm0 vshufps ymm0, ymm0, ymm0, 177 vmovshdup ymm1, [rbx +8*rcx] vmulps ymm1, ymm1, ymm0 vmovsldup ymm6, [rbx +8*rcx +32] vmulps ymm6, ymm6, ymm4 vaddsubps ymm3, ymm2, ymm1 vmovshdup ymm5, [rbx +8*rcx +32] vmovaps [rdx +8*rcx], ymm3 vshufps ymm4, ymm4, ymm4, 177 vmulps ymm5, ymm5, ymm4 vaddsubps ymm7, ymm6, ymm5 vmovaps [rdx +8*rcx +32], ymm7 addrcx, 8 cmprcx, r8 jl loop1 </pre>

11.12 DIVIDE AND SQUARE ROOT OPERATIONS

The Intel SSE divide and square root instructions `DIVPS` and `SQRTPS` have a latency of 14 cycles and they are not pipelined. This means that the throughput of these instructions is one in every 14 cycles. The 256-bit Intel AVX instructions `VDIVPS` and `VSQRTPS` have a latency of 28 cycles and they are not pipelined as well. Therefore, the performance of the Intel SSE divide and square root instructions is similar to the Intel AVX 256-bit instructions on Intel microarchitecture code name Sandy Bridge.

It is possible to speed up single-precision divide and square root calculations using the (V)RSQRTPS and (V)RCPPS instructions. The Intel SSE implementation of these instructions has a latency of 5 cycles and it is pipelined with a throughput of one instruction every cycle. The 256-bit Intel AVX implementation of these instructions has a latency of 7 cycles and throughput of one instruction every two cycles. These instructions can be used with a single Newton-Raphson iteration or Taylor approximation to achieve almost the same precision as the (V)DIVPS and (V)SQRTPS instructions. See Intel® 64 and IA-32 Architectures Software Developer's Manual for more information on these instructions.

In some cases, when the divide or square root operations are part of a larger algorithm that hides some of the latency of these operations, the approximation with Newton-Raphson can slow down execution, because more micro-ops, coming from the additional instructions, fill the pipe.

The following table describes three single precision operations with recommended calculation methods depending on the desired accuracy level.

Table 11-3. Single-Precision Divide and Square Root Alternatives

Operation	Accuracy	Recommendation
Divide	24 bits (IEEE)	DIVPS
	~ 22 bits	RCPPS + 1 Newton-Raphson Iteration + MULPS
	~ 11 bits	RCPPS + MULPS
Reciprocal square root	24 bits (IEEE)	SQRTPS + DIVPS
	~ 22 bits	RSQRTPS + 1 Newton-Raphson Iteration
	~ 11 bits	RSQRTPS
Square root	24 bits (IEEE)	SQRTPS
	~ 22 bits	RSQRTPS + 1 Newton-Raphson Iteration + MULPS
	~ 11 bits	RSQRTPS + RCPPS

11.12.1 Single-Precision Divide

To compute:

$$Z[i]=A[i]/B[i]$$

On a large vector of single-precision numbers, $Z[i]$ can be calculated by a divide operation, or by multiplying $1/B[i]$ by $A[i]$.

Denoting $B[i]$ by N , it is possible to calculate $1/N$ using the (V)RCPPS instruction, achieving approximately 11-bit precision.

For better accuracy you can use the one Newton-Raphson iteration:

$X_{(0)} \sim 1/N$; Initial estimation, $\text{rcp}(N)$

$X_{(0)} = 1/N * (1 - E)$

$E = 1 - N * X_{(0)}$; $E \sim 2^{(-11)}$

$X_{(1)} = X_{(0)} * (1 + E) = 1/N * (1 - E^2)$; $E^2 \sim 2^{(-22)}$

$X_{(1)} = X_{(0)} * (1 + 1 - N * X_{(0)}) = 2 * X_{(0)} - N * X_{(0)}^2$

$X_{(1)}$ is an approximation of $1/N$ with approximately 22-bit precision.

Example 11-22. Divide Using DIVPS for 24-bit Accuracy

SSE code using DIVPS	Using VDIVPS
<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx loop1: movups xmm0, [rax+rdx*1] movups xmm1, [rbx+rdx*1] divps xmm0, xmm1 movups [rcx+rdx*1], xmm0 add rdx, 0x10 cmp rdx, rsi jl loop1 </pre>	<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx loop1: vmovups ymm0, [rax+rdx*1] vmovups ymm1, [rbx+rdx*1] vdivps ymm0, ymm0, ymm1 vmovups [rcx+rdx*1], ymm0 add rdx, 0x20 cmp rdx, rsi jl loop1 </pre>

Example 11-23. Divide Using RCPPS 11-bit Approximation

SSE code using RCPPS	Using VRCPPS
<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx </pre>	<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx </pre>

Example 11-23. Divide Using RCPPS 11-bit Approximation (Contd.)

SSE code using RCPPS	Using VRCPPS
<pre> loop1: movups xmm0,[rax+rdx*1] movups xmm1,[rbx+rdx*1] rcpps xmm1,xmm1 mulps xmm0,xmm1 movups [rcx+rdx*1],xmm0 add rdx, 16 cmp rdx, rsi jl loop1 </pre>	<pre> loop1: vmovups ymm0,[rax+rdx] vmovups ymm1,[rbx+rdx] vrcpps ymm1, ymm1 vmulps ymm0, ymm0, ymm1 vmovups [rcx+rdx], ymm0 add rdx, 32 cmp rdx, rsi jl loop1 </pre>

Example 11-24. Divide Using RCPPS and Newton-Raphson Iteration

RCPPS + MULPS ~ 22 bit accuracy	VRCPPS + VMULPS ~ 22 bit accuracy
<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx loop1: movups xmm0,[rax+rdx*1] movups xmm1,[rbx+rdx*1] rcpps xmm3, xmm1 movaps xmm2, xmm3 addps xmm3, xmm2 mulps xmm2, xmm2 mulps xmm2, xmm1 subps xmm3, xmm2 mulps xmm0, xmm3 movups xmmword ptr [rcx+rdx*1], xmm0 add rdx, 0x10 cmp rdx, rsi jl loop1 </pre>	<pre> mov rax, pln1 mov rbx, pln2 mov rcx, pOut mov rsi, iLen xor rdx, rdx loop1: vmovups ymm0,[rax+rdx] vmovups ymm1,[rbx+rdx] vrcpps ymm3, ymm1 vaddps ymm2, ymm3, ymm3 vmulps ymm3, ymm3, ymm3 vmulps ymm3, ymm3, ymm1 vsubps ymm2, ymm2, ymm3 vmulps ymm0, ymm0, ymm2 vmovups [rcx+rdx], ymm0 add rdx, 32 cmp rdx, rsi jl loop1 </pre>

Table 11-4. Comparison of Single-Precision Divide Alternatives

Accuracy	Method	SSE Performance	SSE Performance
24 bits	(V)DIVPS	Baseline	1X
~ 22 bits	(V)RCPPS + Newton-Raphson	2.7X	4.5X
~ 11 bits	(V)RCPPS	6X	8X

11.12.2 Single-Precision Reciprocal Square Root

To compute $Z[i]=1/(A[i])^{0.5}$ on a large vector of single-precision numbers, denoting $A[i]$ by N , it is possible to calculate $1/N$ using the (V)RSQRTPS instruction.

For better accuracy you can use one Newton-Raphson iteration:

$$X_0 \sim 1/N$$

$$E=1-N*X_0^2$$

$$X_0 = (1/N)^{0.5} * ((1-E)^{0.5}) = (1/N)^{0.5} * (1-E/2) ; E/2 \sim 2^{(-11)}$$

$$X_1 = X_0 * (1+E/2) \sim (1/N)^{0.5} * (1-E^2/4) ; E^2/4 \sim 2^{(-22)}$$

$$X_1 = X_0 * (1 + 1/2 * 1/2 * N * X_0^2) = 1/2 * X_0 * (3 - N * X_0^2)$$

X_1 is an approximation of $(1/N)^{0.5}$ with approximately 22-bit precision

Example 11-25. Reciprocal Square Root Using DIVPS+SQRTPS for 24-bit Accuracy

Using SQRTPS, DIVPS	Using VSQRTPS, VDIVPS
<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: movups xmm1, [rax+rdx] sqrtps xmm0, xmm1 divps xmm0, xmm1 movups [rbx+rdx], xmm0 add rdx, 16 cmp rdx, rcx jl loop1 </pre>	<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: vmovups ymm1, [rax+rdx] vsqrtps ymm0, ymm1 vdivps ymm0, ymm0, ymm1 vmovups [rbx+rdx], ymm0 add rdx, 32 cmp rdx, rcx jl loop1 </pre>

Example 11-26. Reciprocal Square Root Using RCPPS 11-bit Approximation

SSE code using RCPPS	Using VRCPPS
<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: rsqrtps xmm0, [rax+rdx] movups [rbx+rdx], xmm0 add rdx, 16 cmp rdx, rcx jl loop1 </pre>	<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: vrsqrtps ymm0, [rax+rdx] vmovups [rbx+rdx], ymm0 add rdx, 32 cmp rdx, rcx jl loop1 </pre>

Example 11-27. Reciprocal Square Root Using RCPPS and Newton-Raphson Iteration

RCPPS + MULPS ~ 22 bit accuracy	VRCPPS + VMULPS ~ 22 bit accuracy
<pre> __declspec(align(16)) float minus_half[4] = {-0.5, -0.5, -0.5, -0.5}; __declspec(align(16)) float three[4] = {3.0, 3.0, 3.0, 3.0}; __asm { mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx movups xmm3, [three] movups xmm4, [minus_half] loop1: movups xmm5, [rax+rdx] rsqrtps xmm0, xmm5 movaps xmm2, xmm0 mulps xmm0, xmm0 mulps xmm0, xmm5 subps xmm0, xmm3 mulps xmm0, xmm2 mulps xmm0, xmm4 movups [rbx+rdx], xmm0 </pre>	<pre> __declspec(align(32)) float half[8] = {0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5}; __declspec(align(32)) float three[8] = {3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0}; __asm { mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx vmovups ymm3, [three] vmovups ymm4, [half] loop1: vmovups ymm5, [rax+rdx] vrsqrtps ymm0, ymm5 vmulps ymm2, ymm0, ymm0 vmulps ymm2, ymm2, ymm5 vsubps ymm2, ymm3, ymm2 vmulps ymm0, ymm0, ymm2 vmulps ymm0, ymm0, ymm4 </pre>

Example 11-27. Reciprocal Square Root Using RCPPS and Newton-Raphson Iteration (Contd.)

RCPPS + MULPS ~ 22 bit accuracy	VRCPPS + VMULPS ~ 22 bit accuracy
<pre> add rdx, 16 cmp rdx, rcx jl loop1 } </pre>	<pre> vmovups [rbx+rdx], ymm0 add rdx, 32 cmp rdx, rcx jl loop1 } </pre>

Table 11-5. Comparison of Single-Precision Reciprocal Square Root Operation

Accuracy	Method	SSE Performance	SSE Performance
24 bits	(V)SQRTPS + (V)DIVPS	Baseline	1X
~ 22 bits	(V)RCPPS + Newton-Raphson	5.2X	9.1X
~ 11 bits	(V)RCPPS	13.5X	17.5X

11.12.3 Single-Precision Square Root

To compute $Z[i] = (A[i])^{0.5}$ on a large vector of single-precision numbers, denoting $A[i]$ by N , the approximation for $N^{0.5}$ is N multiplied by $(1/N)^{0.5}$, where the approximation for $(1/N)^{0.5}$ is described in the previous section.

To get approximately 22-bit precision of $N^{0.5}$, use the following calculation:

$$N^{0.5} = X_1 * N = 1/2 * N * X_0 * (3 - N * X_0^2)$$

Example 11-28. Square Root Using SQRTPS for 24-bit Accuracy

Using SQRTPS	Using VSQRTPS
<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: movups xmm1, [rax+rdx] sqrtps xmm1, xmm1 movups [rbx+rdx], xmm1 add rdx, 16 cmp rdx, rcx jl loop1 </pre>	<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: vmovups ymm1, [rax+rdx] vsqrtps ymm1, ymm1 vmovups [rbx+rdx], ymm1 add rdx, 32 cmp rdx, rcx jl loop1 </pre>

Example 11-29. Square Root Using RCPPS 11-bit Approximation

SSE code using RCPPS	Using VRCPPS
<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx loop1: movups xmm1, [rax+rdx] xorps xmm8, xmm8 cmpneqps xmm8, xmm1 rsqrtps xmm1, xmm1 rcpps xmm1, xmm1 andps xmm1, xmm8 movups [rbx+rdx], xmm1 add rdx, 16 cmp rdx, rcx jl loop1 </pre>	<pre> mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx vxorps ymm8, ymm8, ymm8 loop1: vmovups ymm1, [rax+rdx] vcmpneqps ymm9, ymm8, ymm1 vrsqrtps ymm1, ymm1 vrcpps ymm1, ymm1 vandps ymm1, ymm1, ymm9 vmovups [rbx+rdx], ymm1 add rdx, 32 cmp rdx, rcx jl loop1 </pre>

Example 11-30. Square Root Using RCPPS and One Taylor Series Expansion

RCPPS + Taylor ~ 22 bit accuracy	VRCPPS + Taylor ~ 22 bit accuracy
<pre> __declspec(align(16)) float minus_half[4] = {-0.5, -0.5, -0.5, -0.5}; __declspec(align(16)) float three[4] = {3.0, 3.0, 3.0, 3.0}; __asm { mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx movups xmm6, [three] movups xmm7, [minus_half] </pre>	<pre> __declspec(align(32)) float three[8] = {3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0}; __declspec(align(32)) float minus_half[8] = {-0.5, -0.5, -0.5, -0.5, -0.5, -0.5, -0.5, -0.5}; __asm { mov rax, pln mov rbx, pOut mov rcx, iLen xor rdx, rdx vmovups ymm6, [three] vmovups ymm7, [minus_half] vxorps ymm8, ymm8, ymm8 </pre>

Example 11-30. Square Root Using RCPPS and One Taylor Series Expansion (Contd.)

RCPPS + Taylor ~ 22 bit accuracy	VRCPPS + Taylor ~ 22 bit accuracy
<pre> loop1: movups xmm3, [rax+rdx] rsqrtps xmm1, xmm3 xorps xmm8, xmm8 cmpneqps xmm8, xmm3 andps xmm1, xmm8 movaps xmm4, xmm1 mulps xmm1, xmm3 movaps xmm5, xmm1 mulps xmm1, xmm4 subps xmm1, xmm6 mulps xmm1, xmm5 mulps xmm1, xmm7 movups [rbx+rdx], xmm1 add rdx, 16 cmp rdx, rcx jl loop1 } </pre>	<pre> loop1: vmovups ymm3, [rax+rdx] vrsqrtps ymm4, ymm3 vcmpneqps ymm9, ymm8, ymm3 vandps ymm4, ymm4, ymm9 vmulps ymm1, ymm4, ymm3 vmulps ymm2, ymm1, ymm4 vsubps ymm2, ymm2, ymm6 vmulps ymm1, ymm1, ymm2 vmulps ymm1, ymm1, ymm7 vmovups [rbx+rdx], ymm1 add rdx, 32 cmp rdx, rcx jl loop1 } </pre>

Table 11-6. Comparison of Single-Precision Square Root Operation

Accuracy	Method	SSE Performance	SSE Performance
24 bits	(V)SQRTPS	Baseline	1X
~ 22 bits	(V)RCPPS + Taylor-Expansion	2.3X	4.3X
~ 11 bits	(V)RCPPS	4.7X	5.9X

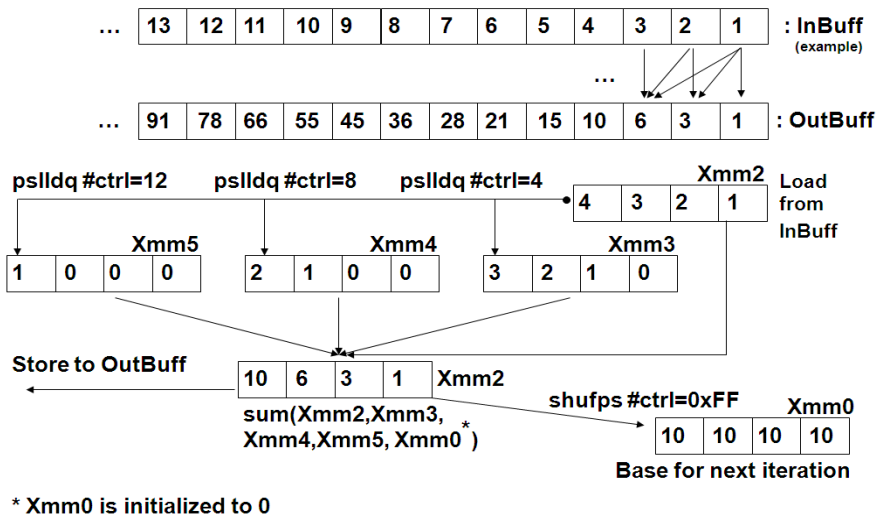
11.13 OPTIMIZATION OF ARRAY SUB SUM EXAMPLE

This section shows the transformation of SSE implementation of Array Sub Sum algorithm to Intel AVX implementation.

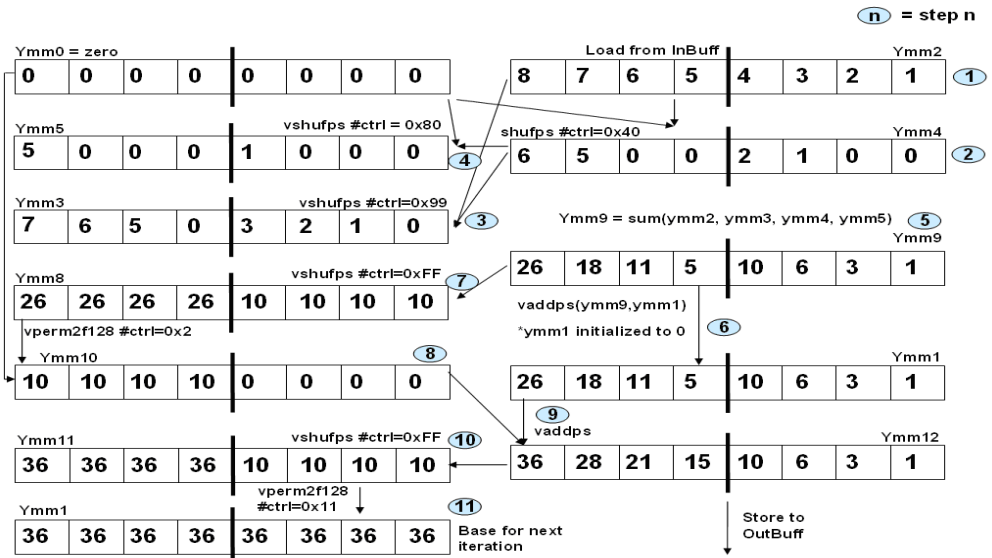
The Array Sub Sum algorithm is:

$$Y_{[i]} = \text{Sum of } k \text{ from } 0 \text{ to } i (X_{[k]}) = X_{[0]} + X_{[1]} + \dots + X_{[i]}$$

The following figure describes the SSE implementation



The figure below describes the Intel AVX implementation of the Array Sub Sums algorithm. The PSLLDQ is an integer SIMD instruction which does not have a 256-bit equivalent. It is replaced by VSHUFPS.



Example 11-31. Array Sub Sums Algorithm

SSE code	AVX code
<pre> mov rax, InBuff mov rbx, OutBuff mov rdx, len xor rcx, rcx xorps xmm0, xmm0 loop1: movaps xmm2, [rax+4*rcx] movaps xmm3, [rax+4*rcx] movaps xmm4, [rax+4*rcx] movaps xmm5, [rax+4*rcx] pslldq xmm3, 4 pslldq xmm4, 8 pslldq xmm5, 12 addps xmm2, xmm3 addps xmm4, xmm5 addps xmm2, xmm4 addps xmm2, xmm0 movaps xmm0, xmm2 shufps xmm0, xmm2, 0xFF movaps [rbx+4*rcx], xmm2 add rcx, 4 cmp rcx, rdx jl loop1 </pre>	<pre> mov rax, InBuff mov rbx, OutBuff mov rdx, len xor rcx, rcx vxorps ymm0, ymm0, ymm0 vxorps ymm1, ymm1, ymm1 loop1: vmovaps ymm2, [rax+4*rcx] vshufps ymm4, ymm0, ymm2, 0x40 vshufps ymm3, ymm4, ymm2, 0x99 vshufps ymm5, ymm0, ymm4, 0x80 vaddps ymm6, ymm2, ymm3 vaddps ymm7, ymm4, ymm5 vaddps ymm9, ymm6, ymm7 vaddps ymm1, ymm9, ymm1 vshufps ymm8, ymm9, ymm9, 0xFF vperm2f128 ymm10, ymm8, ymm0, 0x2 vaddps ymm12, ymm1, ymm10 vshufps ymm11, ymm12, ymm12, 0xFF vperm2f128 ymm1, ymm11, ymm11, 0x11 vmovaps [rbx+4*rcx], ymm12 add rcx, 8 cmp rcx, rdx jl loop1 </pre>

Example 11-31 shows SSE implementation of array sub summ and AVX implementation. The AVX code is about 40% faster.

CHAPTER 12

POWER OPTIMIZATION FOR MOBILE USAGES

12.1 OVERVIEW

Mobile computing allows computers to operate anywhere, anytime. Battery life is a key factor in delivering this benefit. Mobile applications require software optimization that considers both performance and power consumption. This chapter provides background on power saving techniques in mobile processors¹ and makes recommendations that developers can leverage to provide longer battery life.

A microprocessor consumes power while actively executing instructions and doing useful work. It also consumes power in inactive states (when halted). When a processor is active, its power consumption is referred to as active power. When a processor is halted, its power consumption is referred to as static power.

ACPI 3.0 (ACPI stands for Advanced Configuration and Power Interface) provides a standard that enables intelligent power management and consumption. It does this by allowing devices to be turned on when they are needed and by allowing control of processor speed (depending on application requirements). The standard defines a number of P-states to facilitate management of active power consumption; and several C-state types² to facilitate management of static power consumption.

Pentium M, Intel Core Solo, Intel Core Duo processors, and processors based on Intel Core microarchitecture implement features designed to enable the reduction of active power and static power consumption. These include:

- Enhanced Intel SpeedStep[®] Technology enables operating system (OS) to program a processor to transition to lower frequency and/or voltage levels while executing a workload.
- Support for various activity states (for example: Sleep states, ACPI C-states) to reduce static power consumption by turning off power to sub-systems in the processor.

Enhanced Intel SpeedStep Technology provides low-latency transitions between operating points that support P-state usages. In general, a high-numbered P-state operates at a lower frequency to reduce active power consumption. High-numbered C-state types correspond to more aggressive static power reduction. The trade-off is that transitions out of higher-numbered C-states have longer latency.

-
1. For Intel[®] Centrino[®] mobile technology and Intel[®] Centrino[®] Duo mobile technology, only processor-related techniques are covered in this manual.
 2. ACPI 3.0 specification defines four C-state types, known as C0, C1, C2, C3. Microprocessors supporting the ACPI standard implement processor-specific states that map to each ACPI C-state type.

12.2 MOBILE USAGE SCENARIOS

In mobile usage models, heavy loads occur in bursts while working on battery power. Most productivity, web, and streaming workloads require modest performance investments. Enhanced Intel SpeedStep Technology provides an opportunity for an OS to implement policies that track the level of performance history and adapt the processor's frequency and voltage. If demand changes in the last 300 ms³, the technology allows the OS to optimize the target P-state by selecting the lowest possible frequency to meet demand.

Consider, for example, an application that changes processor utilization from 100% to a lower utilization and then jumps back to 100%. The diagram in Figure 12-1 shows how the OS changes processor frequency to accommodate demand and adapt power consumption. The interaction between the OS power management policy and performance history is described below:

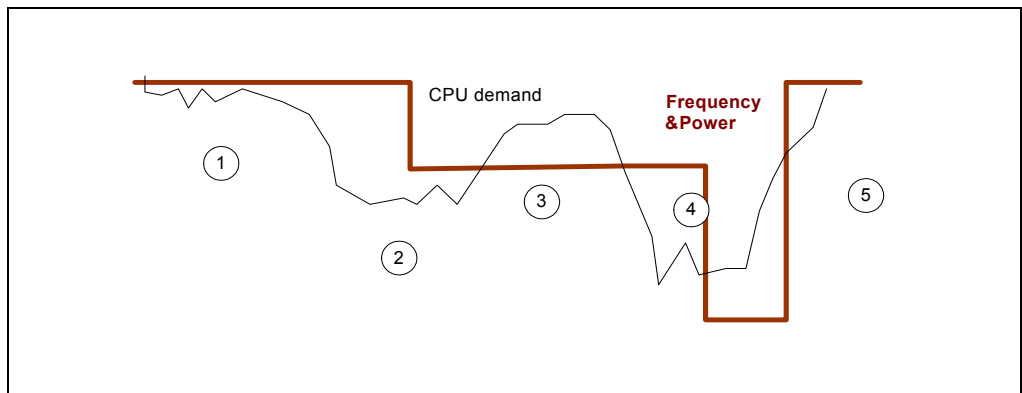


Figure 12-1. Performance History and State Transitions

1. Demand is high and the processor works at its highest possible frequency (P0).
2. Demand decreases, which the OS recognizes after some delay; the OS sets the processor to a lower frequency (P1).
3. The processor decreases frequency and processor utilization increases to the most effective level, 80-90% of the highest possible frequency. The same amount of work is performed at a lower frequency.

3. This chapter uses numerical values representing time constants (300 ms, 100 ms, etc.) on power management decisions as examples to illustrate the order of magnitude or relative magnitude. Actual values vary by implementation and may vary between product releases from the same vendor.

4. Demand decreases and the OS sets the processor to the lowest frequency, sometimes called Low Frequency Mode (LFM).
5. Demand increases and the OS restores the processor to the highest frequency.

12.2.1 Intelligent Energy Efficient Software

With recent advances in power technology and wide range of computing scenarios demanded by end users, intelligent balance between power consumption and performance becomes more and more important. Energy efficient software plays a key role in exploring the latest hardware power savings offered by current generation architecture. Poorly-written code can prevent a system from taking advantage of new hardware features and serving the dynamic needs of end users.

A mobile platform consists of various components such as a CPU, LCD, HDD, DVD, and chipsets, which individually contribute to the power drain of the notebook. Understanding the power contribution of each major component in the platform provides a better view on the total power usage, provides guidance on optimizing power consumption, and may help software to adjust dynamic balance of power budgets between some components.

The following are a few general guidelines for energy efficient software:

- Application should leverage modern OS facility to select appropriate operating frequency instead of setting processor frequency by itself. The latter is likely to have negative impact on both power consumption and performance.
- When your application is waiting for user input or another event to happen, let your application use services that are optimized to go to idle mode quickly. The idle behavior can have a big impact on power consumption. When an application knows it will be operating in a mostly idle context, reduce the frequency of application events that wake up the processor, avoid periodic polling, and reduce the number of services that are active in memory.
- Build context awareness into applications to extend battery life further and optimal user experience.
- Architect an awareness of power consumption and/or dynamic power policy into your application for contextual usages and optimal end user experiences. Specific detail may vary across OSES. For Microsoft Windows OS consult http://www.microsoft.com/whdc/system/pnppwr/powermgmt/PMpolicy_Windows.msp#. For Linux, consult <http://www.lesswatts.org>.
- Characterize your application's power consumption. There are various techniques available to measure the power consumption of a platform:
 - Use hardware instrumentation such as Fluke NetDAQ*. This provides power measurements for each component such as CPU, HDD, and memory.
 - Use C-state residency counters. See Appendix B, "Model-Specific Registers (MSRs)" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

- Study parameters such as CPU usage, kernel time, and time interrupt rate, to gain insight into the behavior of the software, which can then be related to platform power consumption if the hardware instrumentation is not available.

Section 12.5 provide some examples on how to relate performance with power consumption and techniques for optimizing software.

12.3 ACPI C-STATES

When computational demands are less than 100%, part of the time the processor is doing useful work and the rest of the time it is idle. For example, the processor could be waiting on an application time-out set by a Sleep() function, waiting for a web server response, or waiting for a user mouse click. Figure 12-2 illustrates the relationship between active and idle time.

When an application moves to a wait state, the OS issues a HLT instruction and the processor enters a halted state in which it waits for the next interrupt. The interrupt may be a periodic timer interrupt or an interrupt that signals an event.

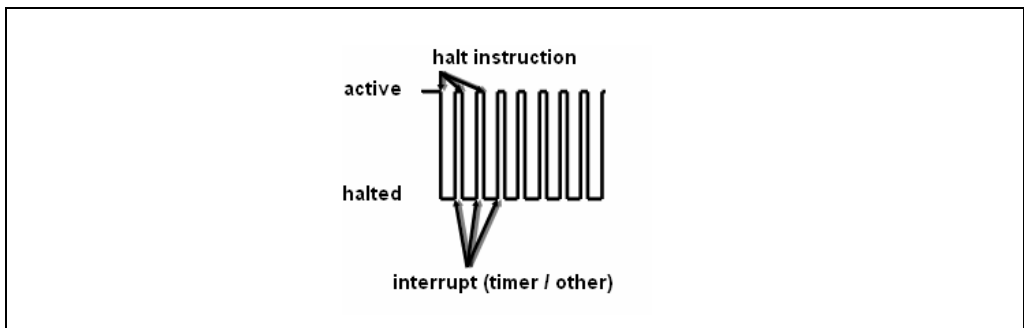


Figure 12-2. Active Time Versus Halted Time of a Processor

As shown in the illustration of Figure 12-2, the processor is in either active or idle (halted) state. ACPI defines four C-state types (C0, C1, C2 and C3). Processor-specific C states can be mapped to an ACPI C-state type via ACPI standard mechanisms. The C-state types are divided into two categories: active (C0), in which the processor consumes full power; and idle (C1-3), in which the processor is idle and may consume significantly less power.

The index of a C-state type designates the depth of sleep. Higher numbers indicate a deeper sleep state and lower power consumption. They also require more time to wake up (higher exit latency).

C-state types are described below:

- **C0** — The processor is active and performing computations and executing instructions.

- **C1** — This is the lowest-latency idle state, which has very low exit latency. In the C1 power state, the processor is able to maintain the context of the system caches.
- **C2** — This level has improved power savings over the C1 state. The main improvements are provided at the platform level.
- **C3** — This level provides greater power savings than C1 or C2. In C3, the processor stops clock generating and snooping activity. It also allows system memory to enter self-refresh mode.

The basic technique to implement OS power management policy to reduce static power consumption is by evaluating processor idle durations and initiating transitions to higher-numbered C-state types. This is similar to the technique of reducing active power consumption by evaluating processor utilization and initiating P-state transitions. The OS looks at history within a time window and then sets a target C-state type for the next time window, as illustrated in Figure 12-3:

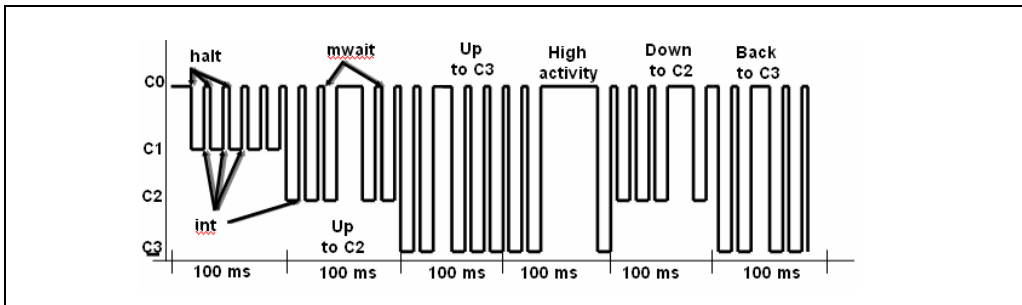


Figure 12-3. Application of C-states to Idle Time

Consider that a processor is in lowest frequency (LFM- low frequency mode) and utilization is low. During the first time slice window (Figure 12-3 shows an example that uses 100 ms time slice for C-state decisions), processor utilization is low and the OS decides to go to C2 for the next time slice. After the second time slice, processor utilization is still low and the OS decides to go into C3.

12.3.1 Processor-Specific C4 and Deep C4 States

The Pentium M, Intel Core Solo, Intel Core Duo processors, and processors based on Intel Core microarchitecture⁴ provide additional processor-specific C-states (and associated sub C-states) that can be mapped to ACPI C3 state type. The processor-

-
4. Pentium M processor can be detected by CPUID signature with family 6, model 9 or 13; Intel Core Solo and Intel Core Duo processor has CPUID signature with family 6, model 14; processors based on Intel Core microarchitecture has CPUID signature with family 6, model 15.

specific C states and sub C-states are accessible using MWAIT extensions and can be discovered using CPUID. One of the processor-specific state to reduce static power consumption is referred to as C4 state. C4 provides power savings in the following manner:

- The voltage of the processor is reduced to the lowest possible level that still allows the L2 cache to maintain its state.
- In an Intel Core Solo, Intel Core Duo processor or a processor based on Intel Core microarchitecture, after staying in C4 for an extended time, the processor may enter into a Deep C4 state to save additional static power.

The processor reduces voltage to the minimum level required to safely maintain processor context. Although exiting from a deep C4 state may require warming the cache, the performance penalty may be low enough such that the benefit of longer battery life outweighs the latency of the deep C4 state.

12.3.2 Processor-Specific Deep C-States and Intel® Turbo Boost Technology

Processors based on Intel microarchitecture code name Nehalem implements several processor-specific C-states.

Table 12-1. ACPI C-State Type Mappings to Processor Specific C-State for Mobile Processors Based on Intel microarchitecture code name Nehalem

ACPI C-State Type	Processor-Specific C-State
C0	C0
C1	C1
C2	C3
C3	C7

The processor-specific deep C-states are implementation dependent. Generally, the low power C-states (higher numbered C-states) have higher exit latencies. For example, when the cores are already in C7, the last level cache (L3) is flushed. The processor support auto-demotion of OS request to deep C-states (C3/C7) and demote to C1/C3 state to support flexible power-performance settings.

In addition to low-power, deep C-states, Intel Turbo Boost Technology can opportunistically boost performance in normal state (C0) by mapping p1 state to the processor’s qualified high-frequency mode operation. Headroom in the system’s TDP can be converted to an even higher frequency than P1 state target. When the operating system requests P0 state, the processor sets core frequencies between P1 to P0 range. A P0 state with only one core busy, achieves the maximum possible Intel Turbo Boost Technology frequency, whereas when the processor is running two to

four cores the frequency is constrained by processor limitations. Under normal conditions the frequency does not go below P1, even when all cores are running.

12.3.3 Processor-Specific Deep C-States for Intel® microarchitecture code name Sandy Bridge

Processors based on Intel microarchitecture code name Sandy Bridge implements several processor-specific C-states.

Table 12-2. ACPI C-State Type Mappings to Processor Specific C-State of Intel microarchitecture code name Sandy Bridge

ACPI C-State Type	Processor-Specific C-State
C0	C0
C1	C1
C2	C3
C3	C6/C7

The microarchitectural behavior of processor-specific deep C-states are implementation dependent. The following summarizes some of their key power-saving and intelligent responsive characteristics:

- For mobile platforms, while the cores are already in C7, the last level cache (L3) is flushed.
- Auto-demotion: The processor can demote OS requests to a target C-state (core C6/C7 or C3 state) to a numerically lower C-state (core C3 or C1 state) in the following cases:
 - When history indicates that C6/C7 or C3 states are less energy efficient than C3 or C1 states.
 - When history indicates that a deeper sleep state may impact performance.
 - Energy inefficiency or performance degradation can occur due to the deeper C-state transition overhead occurring too frequently. Intel microarchitecture code name Sandy Bridge has an enhanced algorithm that improves power gain from this feature.
- Un-demotion: An OS request to a deeper C-state can be demoted by auto-demotion, resulting in C1 or C3 states. After long residency in the demoted state, the hardware returns control back to the OS. The expectation is that in this case, the OS will repeat the deeper C-state request and hardware un-demotion will enter into the OS-requested deeper C state.

12.3.4 Intel® Turbo Boost Technology 2.0

Intel® Turbo Boost Technology 2.0 is a second generation enhancement of Intel® Turbo Boost Technology. The latter can opportunistically boost the processor core's frequency to a higher frequency above the qualified frequency depending on the TDP headroom.

The TDP of Intel Core processors based on Intel microarchitecture code name Sandy Bridge include budgets for the processor core and processor graphic sub-system. Intel® Turbo Boost Technology 2.0 allows more opportunity to convert the thermal and power budget headroom to boost the processor core frequency and/or operating frequency of the processor graphic sub-system.

Energy consumption by the processor cores and/or by the processor graphic unit can be measured using a set of MSR interface⁵. Operating system requirements to support Intel Turbo Boost Technology, to use hints to optimize performance and energy bias in turbo mode operation, are described in Chapter 14, "Power and Thermal Management" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*.

12.4 GUIDELINES FOR EXTENDING BATTERY LIFE

Follow the guidelines below to optimize to conserve battery life and adapt for mobile computing usage:

- Adopt a power management scheme to provide just-enough (not the highest) performance to achieve desired features or experiences.
- Avoid using spin loops.
- Reduce the amount of work the application performs while operating on a battery.
- Take advantage of hardware power conservation features using ACPI C3 state type and coordinate processor cores in the same physical processor.
- Implement transitions to and from system sleep states (S1-S4) correctly.
- Allow the processor to operate at a higher-numbered P-state (lower frequency but higher efficiency in performance-per-watt) when demand for processor performance is low.
- Allow the processor to enter higher-numbered ACPI C-state type (deeper, low-power states) when user demand for processor activity is infrequent.

5. Generally, energy measurements and power management decisions based on these MSR interfaces should operate within the same processor family/model and refrain from extrapolating across different family/models or unsupported environmental conditions.

12.4.1 Adjust Performance to Meet Quality of Features

When a system is battery powered, applications can extend battery life by reducing the performance or quality of features, turning off background activities, or both. Implementing such options in an application increases the processor idle time. Processor power consumption when idle is significantly lower than when active, resulting in longer battery life.

Example of techniques to use are:

- Reducing the quality/color depth/resolution of video and audio playback.
- Turning off automatic spell check and grammar correction.
- Turning off or reducing the frequency of logging activities.
- Consolidating disk operations over time to prevent unnecessary spin-up of the hard drive.
- Reducing the amount or quality of visual animations.
- Turning off, or significantly reducing file scanning or indexing activities.
- Postponing possible activities until AC power is present.

Performance/quality/battery life trade-offs may vary during a single session, which makes implementation more complex. An application may need to implement an option page to enable the user to optimize settings for user's needs (see Figure 12-4).

To be battery-power-aware, an application may use appropriate OS APIs. For Windows XP, these include:

- **GetSystemPowerStatus** — Retrieves system power information. This status indicates whether the system is running on AC or DC (battery) power, whether the battery is currently charging, and how much battery life remains.
- **GetActivePwrScheme** — Retrieves the active power scheme (current system power scheme) index. An application can use this API to ensure that system is running best power scheme. Avoid Using Spin Loops.

Spin loops are used to wait for short intervals of time or for synchronization. The main advantage of a spin loop is immediate response time. Using the `PeekMessage()` in Windows API has the same advantage for immediate response (but is rarely needed in current multitasking operating systems).

However, spin loops and `PeekMessage()` in message loops require the constant attention of the processor, preventing it from entering lower power states. Use them sparingly and replace them with the appropriate API when possible. For example:

- When an application needs to wait for more than a few milliseconds, it should avoid using spin loops and use the Windows synchronization APIs, such as `WaitForSingleObject()`.
- When an immediate response is not necessary, an application should avoid using `PeekMessage()`. Use `WaitMessage()` to suspend the thread until a message is in the queue.

Intel[®] Mobile Platform Software Development Kit⁶ provides a rich set of APIs for mobile software to manage and optimize power consumption of mobile processor and other components in the platform.

12.4.2 Reducing Amount of Work

When a processor is in the C0 state, the amount of energy a processor consumes from the battery is proportional to the amount of time the processor executes an active workload. The most obvious technique to conserve power is to reduce the number of cycles it takes to complete a workload (usually that equates to reducing the number of instructions that the processor needs to execute, or optimizing application performance).

Optimizing an application starts with having efficient algorithms and then improving them using Intel software development tools, such as Intel VTune Performance Analyzers, Intel compilers, and Intel Performance Libraries.

See Chapter 3 through Chapter 7 for more information about performance optimization to reduce the time to complete application workloads.

12.4.3 Platform-Level Optimizations

Applications can save power at the platform level by using devices appropriately and redistributing the workload. The following techniques do not impact performance and may provide additional power conservation:

- Read ahead from CD/DVD data and cache it in memory or hard disk to allow the DVD drive to stop spinning.
- Switch off unused devices.
- When developing a network-intensive application, take advantage of opportunities to conserve power. For example, switch to LAN from WLAN whenever both are connected.
- Send data over WLAN in large chunks to allow the WiFi card to enter low power mode in between consecutive packets. The saving is based on the fact that after every send/receive operation, the WiFi card remains in high power mode for up to several seconds, depending on the power saving mode. (Although the purpose keeping the WiFi in high power mode is to enable a quick wake up).
- Avoid frequent disk access. Each disk access forces the device to spin up and stay in high power mode for some period after the last access. Buffer small disk reads and writes to RAM to consolidate disk operations over time. Use the `GetDevicePowerState()` Windows API to test disk state and delay the disk access if it is not spinning.

6. Evaluation copy may be downloaded at <http://www.intel.com/cd/software/products/asmo-na/eng/219691.htm>

12.4.4 Handling Sleep State Transitions

In some cases, transitioning to a sleep state may harm an application. For example, suppose an application is in the middle of using a file on the network when the system enters suspend mode. Upon resuming, the network connection may not be available and information could be lost.

An application may improve its behavior in such situations by becoming aware of sleep state transitions. It can do this by using the WM_POWERBROADCAST message. This message contains all the necessary information for an application to react appropriately.

Here are some examples of an application reaction to sleep mode transitions:

- Saving state/data prior to the sleep transition and restoring state/data after the wake up transition.
- Closing all open system resource handles such as files and I/O devices (this should include duplicated handles).
- Disconnecting all communication links prior to the sleep transition and re-establishing all communication links upon waking up.
- Synchronizing all remote activity, such as like writing back to remote files or to remote databases, upon waking up.
- Stopping any ongoing user activity, such as streaming video, or a file download, prior to the sleep transition and resuming the user activity after the wake up transition.

Recommendation: *Appropriately handling the suspend event enables more robust, better performing applications.*

12.4.5 Using Enhanced Intel SpeedStep® Technology

Use Enhanced Intel SpeedStep Technology to adjust the processor to operate at a lower frequency and save energy. The basic idea is to divide computations into smaller pieces and use OS power management policy to effect a transition to higher P-states.

Typically, an OS uses a time constant on the order of 10s to 100s of milliseconds⁷ to detect demand on processor workload. For example, consider an application that requires only 50% of processor resources to reach a required quality of service (QoS). The scheduling of tasks occurs in such a way that the processor needs to stay in P0 state (highest frequency to deliver highest performance) for 0.5 seconds and may then go to sleep for 0.5 seconds. The demand pattern then alternates.

Thus the processor demand switches between 0 and 100% every 0.5 seconds, resulting in an average of 50% of processor resources. As a result, the frequency switches accordingly between the highest and lowest frequency. The power

7. The actual number may vary by OS and by OS release.

consumption also switches in the same manner, resulting in an average power usage represented by the equation $P_{average} = (P_{max} + P_{min}) / 2$.

Figure 12-4 illustrates the chronological profiles of coarse-grain (> 300 ms) task scheduling and its effect on operating frequency and power consumption.

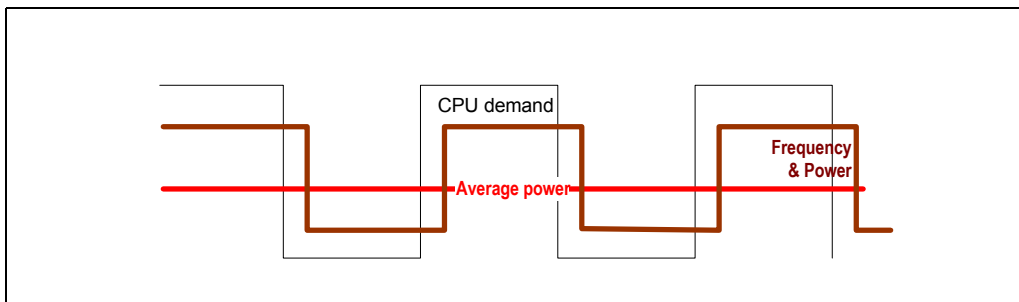


Figure 12-4. Profiles of Coarse Task Scheduling and Power Consumption

The same application can be written in such a way that work units are divided into smaller granularity, but scheduling of each work unit and Sleep() occurring at more frequent intervals (e.g. 100 ms) to deliver the same QOS (operating at full performance 50% of the time). In this scenario, the OS observes that the workload does not require full performance for each 300 ms sampling. Its power management policy may then commence to lower the processor’s frequency and voltage while maintaining the level of QOS.

The relationship between active power consumption, frequency and voltage is expressed by the equation:

$$Power = \alpha * C * V^2 * F$$

In the equation: ‘V’ is core voltage, ‘F’ is operating frequency, and ‘α’ is the activity factor. Typically, the quality of service for 100% performance at 50% duty cycle can be met by 50% performance at 100% duty cycle. Because the slope of frequency scaling efficiency of most workloads will be less than one, reducing the core frequency to 50% can achieve more than 50% of the original performance level. At the same time, reducing the core frequency to 50% allows for a significant reduction of the core voltage.

Because executing instructions at higher P-state (lower power state) takes less energy per instruction than at P0 state, Energy savings relative to the half of the duty cycle in P0 state ($P_{max} / 2$) more than compensate for the increase of the half of the duty cycle relative to inactive power consumption ($P_{min} / 2$). The non-linear relationship between power consumption to frequency and voltage means that changing the task unit to finer granularity will deliver substantial energy savings. This optimization is possible when processor demand is low (such as with media streaming, playing a DVD, or running less resource intensive applications like a word processor, email or web browsing).

An additional positive effect of continuously operating at a lower frequency is that frequent changes in power draw (from low to high in our case) and battery current eventually harm the battery. They accelerate its deterioration.

When the lowest possible operating point (highest P-state) is reached, there is no need for dividing computations. Instead, use longer idle periods to allow the processor to enter a deeper low power mode.

12.4.6 Enabling Intel® Enhanced Deeper Sleep

In typical mobile computing usages, the processor is idle most of the time. Conserving battery life must address reducing static power consumption.

Typical OS power management policy periodically evaluates opportunities to reduce static power consumption by moving to lower-power C-states. Generally, the longer a processor stays idle, OS power management policy directs the processor into deeper low-power C-states.

After an application reaches the lowest possible P-state, it should consolidate computations in larger chunks to enable the processor to enter deeper C-States between computations. This technique utilizes the fact that the decision to change frequency is made based on a larger window of time than the period to decide to enter deep sleep. If the processor is to enter a processor-specific C4 state to take advantage of aggressive static power reduction features, the decision should be based on:

- Whether the QOS can be maintained in spite of the fact that the processor will be in a low-power, long-exit-latency state for a long period.
- Whether the interval in which the processor stays in C4 is long enough to amortize the longer exit latency of this low-power C state.

Eventually, if the interval is large enough, the processor will be able to enter deeper sleep and save a considerable amount of power. The following guidelines can help applications take advantage of Intel® Enhanced Deeper Sleep:

- Avoid setting higher interrupt rates. Shorter periods between interrupts may keep OSEs from entering lower power states. This is because transition to/from a deep C-state consumes power, in addition to a latency penalty. In some cases, the overhead may outweigh power savings.
- Avoid polling hardware. In a ACPI C3 type state, the processor may stop snooping and each bus activity (including DMA and bus mastering) requires moving the processor to a lower-numbered C-state type. The lower-numbered state type is usually C2, but may even be C0. The situation is significantly improved in the Intel Core Solo processor (compared to previous generations of the Pentium M processors), but polling will likely prevent the processor from entering into highest-numbered, processor-specific C-state.

12.4.7 Multicore Considerations

Multicore processors deserves some special considerations when planning power savings. The dual-core architecture in Intel Core Duo processor and mobile processors based on Intel Core microarchitecture provide additional potential for power savings for multi-threaded applications.

12.4.7.1 Enhanced Intel SpeedStep® Technology

Using domain-composition, a single-threaded application can be transformed to take advantage of multicore processors. A transformation into two domain threads means that each thread will execute roughly half of the original number of instructions. Dual core architecture enables running two threads simultaneously, each thread using dedicated resources in the processor core. In an application that is targeted for the mobile usages, this instruction count reduction for each thread enables the physical processor to operate at lower frequency relative to a single-threaded version. This in turn enables the processor to operate at a lower voltage, saving battery life.

Note that the OS views each logical processor or core in a physical processor as a separate entity and computes CPU utilization independently for each logical processor or core. On demand, the OS will choose to run at the highest frequency available in a physical package. As a result, a physical processor with two cores will often work at a higher frequency than it needs to satisfy the target QOS.

For example if one thread requires 60% of single-threaded execution cycles and the other thread requires 40% of the cycles, the OS power management may direct the physical processor to run at 60% of its maximum frequency.

However, it may be possible to divide work equally between threads so that each of them require 50% of execution cycles. As a result, both cores should be able to operate at 50% of the maximum frequency (as opposed to 60%). This will allow the physical processor to work at a lower voltage, saving power.

So, while planning and tuning your application, make threads as symmetric as possible in order to operate at the lowest possible frequency-voltage point.

12.4.7.2 Thread Migration Considerations

Interaction of OS scheduling and multicore unaware power management policy may create some situations of performance anomaly for multi-threaded applications. The problem can arise for multithreading application that allow threads to migrate freely.

When one full-speed thread is migrated from one core to another core that has idled for a period of time, an OS without a multicore-aware P-state coordination policy may mistakenly decide that each core demands only 50% of processor resources (based on idle history). The processor frequency may be reduced by such multicore unaware P-state coordination, resulting in a performance anomaly. See Figure 12-5.

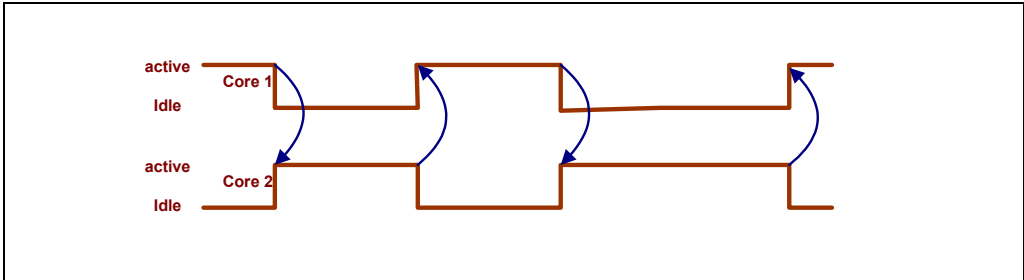


Figure 12-5. Thread Migration in a Multicore Processor

Software applications have a couple of choices to prevent this from happening:

- **Thread affinity management** — A multi-threaded application can enumerate processor topology and assign processor affinity to application threads to prevent thread migration. This can work around the issue of OS lacking multicore aware P-state coordination policy.
- **Upgrade to an OS with multicore aware P-state coordination policy** — Some newer OS releases may include multicore aware P-state coordination policy. The reader should consult with specific OS vendors.

12.4.7.3 Multicore Considerations for C-States

There are two issues that impact C-states on multicore processors.

Multicore-unaware C-state Coordination May Not Fully Realize Power Savings

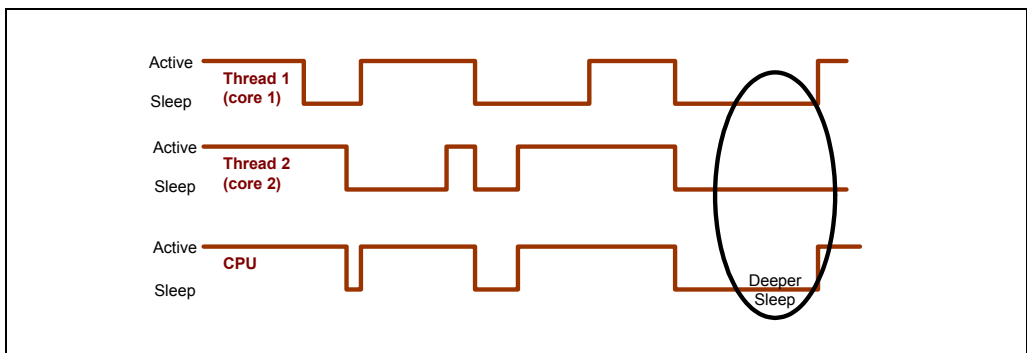


Figure 12-6. Progression to Deeper Sleep

When each core in a multicore processor meets the requirements necessary to enter a different C-state type, multicore-unaware hardware coordination causes the physical processor to enter the lowest possible C-state type (lower-numbered C state has less power saving). For example, if Core 1 meets the requirement to be in ACPI C1 and Core 2 meets requirement for ACPI C3, multicore-unaware OS coordination takes the physical processor to ACPI C1. See Figure 12-6.

Enabling Both Cores to Take Advantage of Intel Enhanced Deeper Sleep.

To best utilize processor-specific C-state (e.g., Intel Enhanced Deeper Sleep) to conserve battery life in multithreaded applications, a multi-threaded application should synchronize threads to work simultaneously and sleep simultaneously using OS synchronization primitives. By keeping the package in a fully idle state longer (satisfying ACPI C3 requirement), the physical processor can transparently take advantage of processor-specific Deep C4 state if it is available.

Multi-threaded applications need to identify and correct load-imbalances of its threaded execution before implementing coordinated thread synchronization. Identifying thread imbalance can be accomplished using performance monitoring events. Intel Core Duo processor provides an event for this purpose. The event (Serial_Execution_Cycle) increments under the following conditions:

- Core actively executing code in C0 state
- Second core in physical processor in idle state (C1-C4)

This event enables software developers to find code that is executing serially, by comparing Serial_Execution_Cycle and Unhalted_Ref_Cycles. Changing sections of serialized code to execute into two parallel threads enables coordinated thread synchronization to achieve better power savings.

Although Serial_Execution_Cycle is available only on Intel Core Duo processors, application thread with load-imbalance situations usually remains the same for symmetric application threads and on symmetrically configured multicore processors, irrespective of differences in their underlying microarchitecture. For this reason, the technique to identify load-imbalance situations can be applied to multi-threaded applications in general, and not specific to Intel Core Duo processors.

12.5 TUNING SOFTWARE FOR INTELLIGENT POWER CONSUMPTION

This section describes some techniques for tuning software for balance of both power and performance. Most of the power optimization techniques are generic. The last sub section (Section 12.5.8) describes features specific to Intel microarchitecture code name Sandy Bridge. Explore these features to optimize software for performance and corresponding power benefits.

12.5.1 Reduction of Active Cycles

Finishing the task quicker by reducing the amount of active cycles, then transfer control to the system idle loop will take advantage of modern operating system's power saving optimizations.

Reduction of active cycles can be achieved in several ways, from applying performance-oriented coding techniques discussed in Chapter 3, vectorization using SSE and/or AVX, to multi-threading.

12.5.1.1 Multi-threading to reduce Active Cycles

If given a task of some fixed amount of computational work that has thread-level parallelism, one can apply data-decomposition for multi-threading. The amount of reduction in active cycles will depend on the degree of parallelism. Similar principle can also apply to function-decomposition situations.

A balanced multi-threading implementation is more likely to achieve more optimal results in intelligent efficient performance and power saving benefits. Choosing the right synchronization primitives also has significant impact on both power and performance.

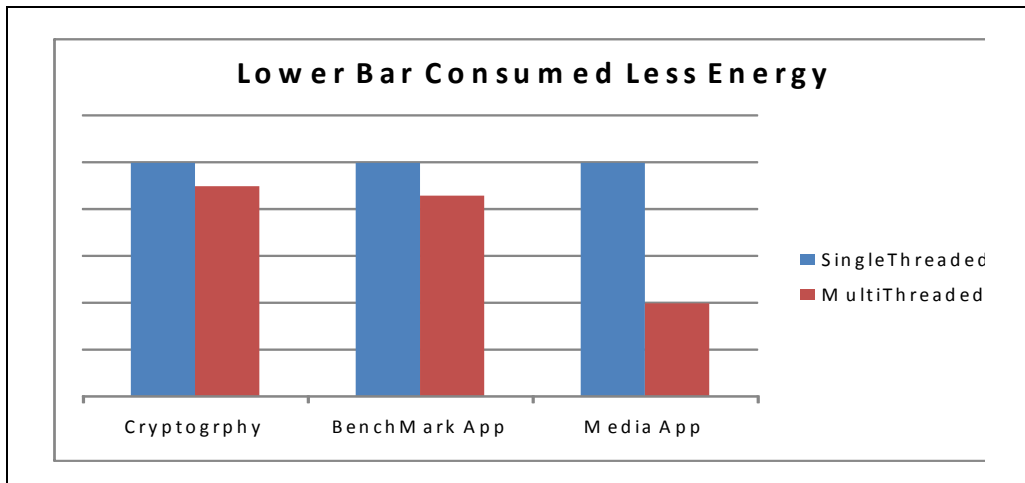


Figure 12-7. Energy Saving due to Performance Optimization

Figure 12-7 above shows the result of a study that compares processor energy consumption of single threaded workloads with their corresponding performance-optimized implementations, using three sets of applications across different application domains. In this particular study, optimization effort in application 1 (Cryptography) achieved 2X gain in performance alone. At the same time, its energy consumption reduced about 12%. In application 3 (a media application), perfor-

mance optimization efforts including multi-threading and other techniques achieved 10X performance gain. Its energy consumption reduced about 60%.

12.5.1.2 Vectorization

Use SIMD instructions can reduce the path length of completing a given computational task, often reducing active cycles. Code that performs the same operation on multiple independent data elements is a good candidate for vectorization. Vectorization techniques are typically applied to applications with loops with elements that can be processed in single instruction. Typically, the slight power increase per unit time of using SIMD instructions are compensated by much greater reduction of active cycles. The net effect is improved energy consumption.

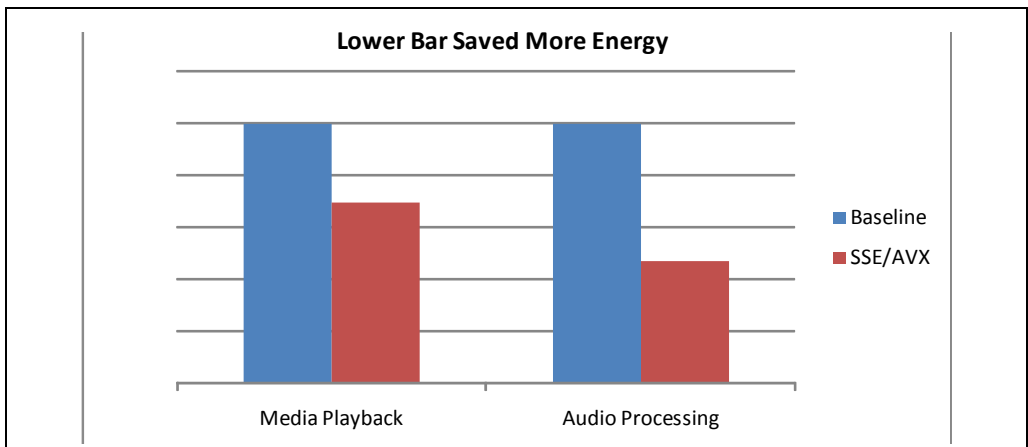


Figure 12-8. Energy Saving due to Vectorization

Figure 12-7 shows the result of a study on the energy saving effect due to vectorization. A media playback workload achieved 2.15X speedup due to using SSE2 and SSE4 instruction sets. Another audio processing workload increased performance to ~5X by using Intel AVX instruction sets. At the same time, the latter also had better energy saving.

12.5.2 PAUSE and Sleep(0) Loop Optimization

In multi-threading implementation, a popular construct in thread synchronization and for yielding scheduling quanta to another thread waiting to carry out its task is to sit in a loop and issuing SLEEP(0).

These are typically called "sleep loops", see Example 12-1. It should be noted that a `SwitchToThread` call can also be used. The "sleep loop" is common in locking algorithms and thread pools as the threads are waiting on work.

Example 12-1. Unoptimized Sleep Loop

```
while(!acquire_lock())
{ Sleep( 0); }
do_work();
release_lock();
```

This construct of sitting in a tight loop and calling `Sleep()` service with a parameter of 0 is actually a polling loop with side effects:

- Each call to `Sleep()` experiences the expensive cost of a context switch, which can be 10000+ cycles.
- It also suffers the cost of ring 3 to ring 0 transitions, which can be 1000+ cycles.
- When there is no other thread waiting to take possession of control, this sleep loop behaves to the OS as a highly active task demanding CPU resource, preventing the OS to put the CPU into a low-power state

Example 12-2. Power Consumption Friendly Sleep Loop Using PAUSE

```
if (acquire_lock())
{ /* Spin on pause max_spin_count times before backing off to sleep */
  for(int j = 0; j < max_spin_count; ++j)
  { /* intrinsic for PAUSE instruction*/
    _mm_pause();
    if (read_volatile_lock())
    {
      if (acquire_lock()) goto PROTECTED_CODE;
    }
  }
  /* Pause loop didn't work, sleep now */
  Sleep(0);
  goto ATTEMPT_AGAIN;
}
PROTECTED_CODE:
do_work();
release_lock();
```

Example 12-2 shows the technique of using PAUSE instruction to make the sleep loop power friendly.

By slowing down the "spin-wait" with the PAUSE instruction, the multi-threading software gains

- performance by facilitating the waiting tasks to acquire resources more easily from a busy wait,
- power-savings by both using fewer parts of the pipeline while spinning and
- elimination of great majority of unnecessarily executed instructions caused by the overhead of a Sleep(0) call.

In one case study, this technique achieved 4.3x of performance gain, which translated to 21% power savings at the processor and 13% power savings at platform level.

12.5.3 Spin-Wait Loops

Use the PAUSE instruction in all spin wait loops. The PAUSE instruction de-pipelines the spin-wait loop to prevent it from consuming execution resources excessively and consuming power needlessly.

When executing a spin-wait loop, the processor can suffer a severe performance penalty when exiting the loop because it detects a possible memory order violation and flushes the core processor's pipeline.

The PAUSE instruction provides a hint to the processor that the code sequence is a spin-wait loop. The processor uses this hint to avoid the memory order violation and prevent the pipeline flush. However, you should try to keep spin-wait loops with PAUSE short.

12.5.4 Using Event Driven Service Instead of Polling in Code

Consistently polling for devices or state changes can cause the platform to wake up and consume more power. Minimize polling whenever possible and use an event driven framework if available. If an OS provides notification services for various device state changes, such as transition from AC to battery, use them instead of polling for device state changes. Using this new event notification framework reduces the overhead for the code to poll the status of the power source, because the code can get notifications asynchronously when status changes happen.

12.5.5 Reducing Interrupt Rate

High interrupt rate may have two consequences that impact processor power and performance:

- It prevents the processor package and its cores from going into deeper sleep states (C-states), which means that the system does not enable the hardware to utilize the power saving features.
- It limits the frequency to which Intel Turbo Boost Technology 2.0 can reach, and therefore the performance of other applications running on the processor degrades.

If a user session and/or an application experiences a rate of thousands of interrupts per second, it would have inhibited the processor to achieve intelligent balance between performance and saving power.

To avoid this situation minimize sporadic wakeups. Schedule all periodic activities of an application or driver into one wakeup period and reduce the interrupt rate to the minimum required.

Many media applications set a very high timer tick rate (1ms). Where possible, use the operating system default timer tick rate. If high granularity is absolutely necessary make sure the software resets the timer tick rate when the task finishes.

For more information, see the article at <http://software.intel.com/en-us/articles/cpu-power-utilization-on-intel-architectures/>.

12.5.6 Reducing Privileged Time

Applications spending significant time in privileged mode lead to excessive energy use due to various reasons. Some examples are: high system call rate and IO bottlenecks. You can use Windows Perfmon to get an estimate of privileged mode time.

A high system call rate, as measured by system calls per second, indicates that the software is causing frequent kernel mode transitions. That is, the application jumps from Ring3 - user mode to Ring0 - kernel mode, frequently. A very common example of this is using an expensive synchronization call such as the Win32 API `WaitForSingleObject()`. This is a very important synchronization API, especially for inter-process communication. However, it enters kernel mode irrespective of whether the lock is achieved or not. For multi-threaded code with no or a short period contention on the lock, you can use `EnterCriticalSection` with a spin count. The advantage of this API over `WaitForSingleObject()` is that it does not enter kernel mode unless there is a contention on the lock. Hence, when there is no contention, `EnterCriticalSection` with spin count is much cheaper to use and reduces the time spent in privilege mode.

Studies were done by taking a small test application which has four active threads on an Intel microarchitecture code name Sandy Bridge-based system. The locks in the test case were implemented by using `WaitForSingleObject` and `EnterCriticalSection`. There was no contention on the lock, so each thread achieved the lock at the first attempt. As shown in the graph below, when there is no contention, using `WaitForSingleObject()` has negative impact on both power and performance as compared to using `EnterCriticalSection()`.

As indicated in the following graph, using WaitForSingleObject() on an un-contended lock uses more power. Using EnterCriticalSection() provides a 10x performance gain and 60% energy reduction.

For more information see: <http://software.intel.com/en-us/articles/implementing-scalable-atomic-locks-for-multi-core-intel-em64t-and-ia32-architectures/>

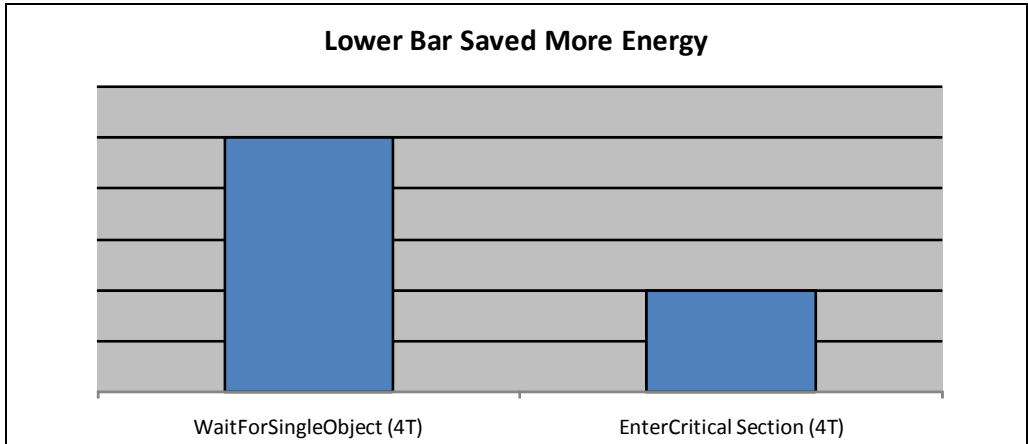


Figure 12-9. Energy Saving Comparison of Synchronization Primitives

12.5.7 Setting Context Awareness in the Code

Context awareness provides a way for software to save power when energy resources are limited. The following are some examples of context awareness that you can implement to conserve power:

- When running a game on laptop on battery power, change the frame rate to 60FPS instead of running uncapped.
- Dim the display when running on battery and not in use.
- Provide easy options for the end user to turn off devices such as wireless when not connected to network

Applications can do these changes transparently when running the game on battery power, or provide hints to users how to extend battery life. For either case, the application needs to be context aware to identify when battery power is in use as opposed to AC power.

A study done with two games running at different frame rates. The blue bar represents baseline default frame rate (uncapped) for these games. The brown line represents games running at 60FPS and the yellow line represents games running at 30FPS. This study shows that capping frame rate can help reduce power consumption

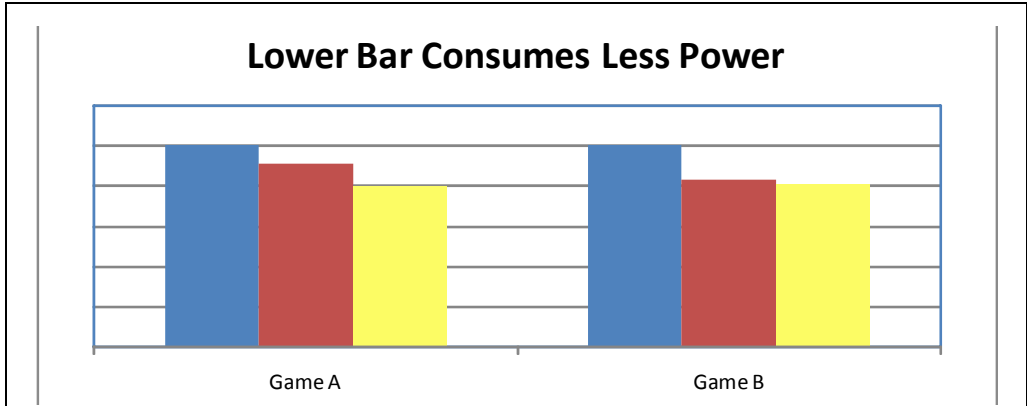


Figure 12-10. Power Saving Comparison of Power-Source-Aware Frame Rate Configurations

12.5.8 Saving Energy by Optimizing for Performance

As general rule, performance optimizations that reduce the number of cycles the CPU is busy running code also save energy consumption. Here are some examples of optimizing performance for specific microarchitectural features that produced energy savings.

The additional load port feature of Intel microarchitecture code name Sandy Bridge can save cycles, as demonstrated in Section 3.6.1. As a result the load port feature also saves power. For example, in an experiment with the kernel in Section 3.6.1.3, a sample application running on an engineering system with and without a bank conflict, the version without the bank conflict utilizes the second load port and provided a performance improvement along with 25mWhr energy savings.

Another features, the Decoded ICache and the LSD, cache micro-ops, hence eliminating power consumption by the decoders. For example, using the code alignment technique of arranging no more than three unconditional branches within an aligned 32 byte chunk (see Section 3.4.2.5) for switch statement, where we see 1.23x speedup, helps provide 1.12x power saving as well compared to the aligned version of dense unconditional branches in 32-byte chunk that can not fit the decoded ICache.

Using vectorization with Intel AVX allows processing of several elements in one cycle, hence reducing the overall processing cycles and saving energy. An example for the energy saving is shown in Figure 12-8.

12.6 PROCESSOR SPECIFIC POWER MANAGEMENT OPTIMIZATION FOR SYSTEM SOFTWARE

This section covers power management information that are processor specific. Specifically, they apply to second generation Intel Core processors based on Intel microarchitecture code name Sandy Bridge. These processors have CPUID DisplayFamily_DisplayModel signature of 06_2AH. These processor-specific capability may help system software to optimize/balance responsiveness and power consumption of the processor.

12.6.1 Power Management Recommendation of Processor-Specific Inactive State Configurations

Programming ACPI's `_CST` object with exit latency values appropriate to various inactive states will help OS power management to deliver optimal power reduction. Intel recommended values are model-specific.

Table 12-3 and Table 12-4 list Package C-State entry/exit latency for processors with CPUID DisplayFamily_DisplayModel signature of 06_2AH, and for two configurations of voltage regulator slew rate capabilities. Table 12-3 applies to slow VR configuration, and Table 12-4 applies to fast VR configuration. For each configuration, the VR device can operate in either a fast interrupt break mode enabled or slow interrupt break mode, depending on the setting of `MSR_POWER_CTL[bit 4]`. These C-State entry/exit latency are not processor specifications but estimates derived from empirical measurements. There may be some situations exit latency from a core is higher than those listed in Table 12-3 and Table 12-4.

Table 12-3. C-State Total Processor Exit Latency for Client Systems (Core+ Package Exit Latency) with Slow VR

C-State ¹	Typical Exit Latency ²	Worst Case Exit Latency
	MSR_POWER_CTL MSR.[4] =0	MSR_POWER_CTL MSR.[4] =1
C1	1 μ s	1 μ s
C3	156 μ s	80 μ s
C6	181 μ s	104 μ s
C7	199 μ s	109 μ s

NOTES:

1. These C-State Entry/Exit Latencies are Intel estimates only and not processor specifications.
2. It is assumed that package is in C0 when one of the core is active.
3. Fast interrupt break mode is enabled if `MSR_POWER_CTL[4] = 1`.
4. A device that connect to PCH may result in latencies equivalent to that of a slow interrupt break mode.

Table 12-4. C-State Total Processor Exit Latency for Client Systems (Core+ Package Exit Latency) with Fast VR

C-State ¹	Typical Worst Case Exit Latency Time (All Skus) ²	
	MSR_POWER_CTL MSR.[4] =0	MSR_POWER_CTL MSR.[4] =1
C1	1 μ s	1 μ s
C3	156 μ s	80 μ s
C6	181 μ s	104 μ s
C7	199 μ s	109 μ s

NOTES:

1. These C-State Entry/Exit Latencies are Intel estimates only and not processor specifications.
2. It is assumed that package is in C0 when one of the core is active.
3. If the package is in a deeper C-states, the exit latency of Local APIC timer wake up depends on the typical core level exit latency; If the package is in C0, it may vary between typical or worst case of the respective core-level exit latency.

Table 12-5 lists Core-only C-State entry/exit latency for processors with CPUID DisplayFamily_DisplayModel signature of 06_2AH, and for two configurations of voltage regulator slew rate capabilities. Core-only exit latency is not affected by MSR_POWER_CTL.[4].

Table 12-5. C-State Core-Only Exit Latency for Client Systems with Slow VR

C-State ¹	Typical Worst Case Exit Latency Time (All Skus) ²	
	MSR_POWER_CTL MSR.[4] =0	MSR_POWER_CTL MSR.[4] =1
C1	1 μ s	1 μ s
C3	21 μ s	240 μ s
C6	46 μ s	250 μ s
C7	46 μ s	250 μ s

NOTES:

1. These C-State Entry/Exit Latencies are Intel estimates only and not processor specifications.
2. A slow VR device refers to a device with ramp time of 10 mv/ μ s in fast mode and 2.5 mv/ μ s in slow mode .

12.6.1.1 Balancing Power Management and Responsiveness of Inactive To Active State Transitions

MSR_PKGC3_IRTL, MSR_PKGC6_IRTL, MSR_PKGC7_IRTL provide processor-specific interfaces for system software to balance power consumption and system responsiveness. System software may change budgeted values from a package inactive

POWER OPTIMIZATION FOR MOBILE USAGES

states to C0 during runtime to accommodate system specific requirements. For example, more aggressive timings when on battery vs. on AC power.

The exit latency is greatly impacted by the VR swing rate. Table 12-5 specifies the total interrupt response times per state (including the core component) for a “fast” exit rate (the default recommended configuration in the PCH and CPU for all events except the internal HPET and CPU timers).

Selecting a “slow” rate by the BIOS (disable the fast break event method in the POWER_CTL MSR bit 4) for various events from the PCIE will require extending the above budget respectively. Otherwise CPU may select a shallower PKG_Cstate to still meet the budget at a much slower VR swing rate.

Selecting a “slow” exit rate for various PCH-connected devices (PCH BIOS setting) will not be visible to the above latency calculation mechanism and thus result in the CPU not meeting the required latency goals.

Table 12-6. POWER_CTL MSR in Next Generation Intel Processor (Intel® microarchitecture code name Sandy Bridge)

Register Address		Register Name	Scope	Bit Description
Hex	Dec			
1FCH	508	MSR_POWER_CTL	Core	Power Control Register
		3:0		Reserved.
		4		FAST_Brk_Int_En. When set to 1, enables the voltage regulator for fast slope for PCI-E interrupt wakeup events.
		63:5		Reserved.

13.1 OVERVIEW

This chapter covers a brief overview the Intel Atom microarchitecture, and specific coding techniques for software whose primary targets are processors based on the Intel Atom microarchitecture. The key features of Intel Atom processors to support low power consumption and efficient performance include:

- Enhanced Intel SpeedStep® Technology enables operating system (OS) to program a processor to transition to lower frequency and/or voltage levels while executing a workload.
- Support deep power down technology to reduces static power consumption by turning off power to cache and other sub-systems in the processor.
- Intel Hyper-Threading Technology providing two logical processor for multi-tasking and multi-threading workloads
- Support Single-instruction multiple-data extensions up to SSE3 and SSSE3.
- Support for Intel 64 and IA-32 architecture.

The Intel Atom microarchitecture is designed to support the general performance requirements of modern workloads within the power-consumption envelop of small form-factor and/or thermally-constrained environments.

13.2 INTEL® ATOM™ MICROARCHITECTURE

Intel Atom microarchitecture achieves efficient performance and low power operation with a two-issue wide, in-order pipeline that support Hyper-Threading Technology. The in-order pipeline differs from out-of-order pipelines by treating an IA-32 instruction with a memory operand as a single pipeline operation instead of multiple micro-operations.

The basic block diagram of the Intel Atom microarchitecture pipeline is shown in Figure 13-1.

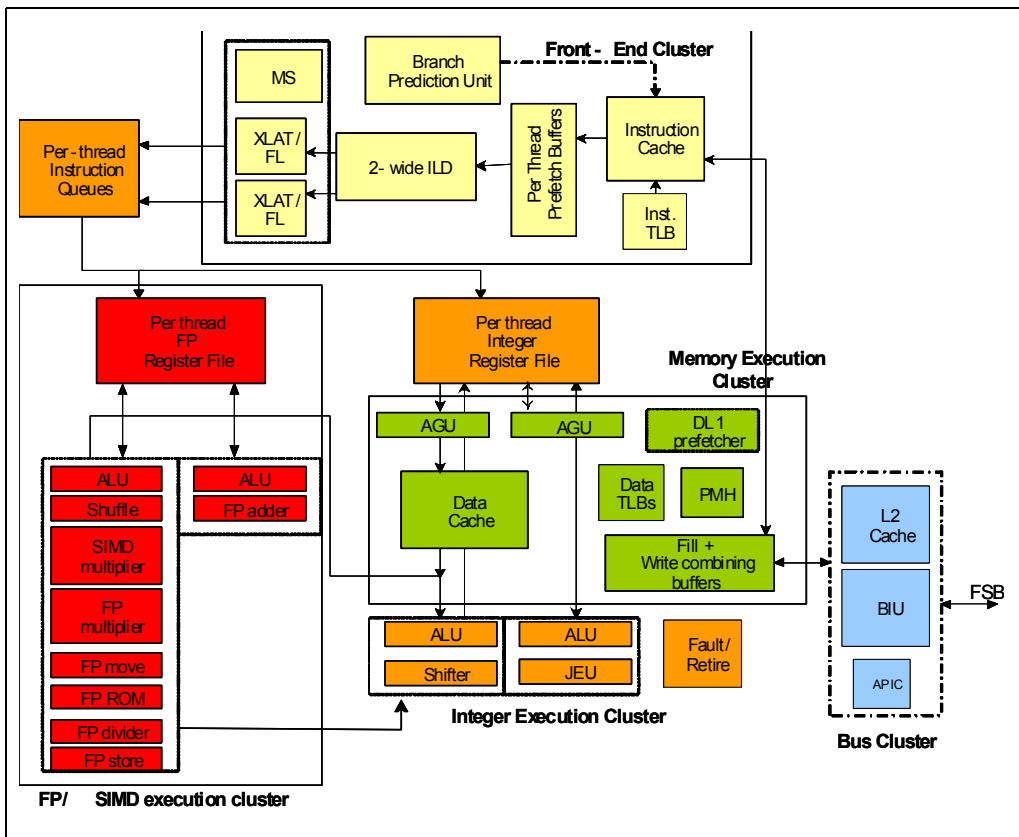


Figure 13-1. Intel Atom Microarchitecture Pipeline

The front end features a power-optimized pipeline, including

- 32KB, 8-way set associative, first-level instruction cache,
- Branch prediction units and ITLB,
- Two instruction decoders, each can decode up to one instruction per cycle.

The front end can deliver up to two instructions per cycle to the instruction queue for scheduling. The scheduler can issue up to two instructions per cycle to the integer or SIMD/FP execution clusters via two issue ports.

Each of the two issue ports can dispatch an instruction per cycle to the integer cluster or the SIMD/FP cluster to execute. The port-bindings of the integer and SIMD/FP clusters have the following features:

- Integer execution cluster:
 - Port 0: ALU0, Shift/Rotate unit, Load/Store,

- Port 1: ALU1, Bit processing unit, jump unit and LEA,
- Effective “load-to-use” latency of 0 cycle
- SIMD/FP execution cluster:
 - Port 0: SIMD ALU, Shuffle unit, SIMD/FP multiply unit, Divide unit, (support IMUL, IDIV)
 - Port 1: SIMD ALU, FP Adder,
 - The two SIMD ALUs and the shuffle unit in the SIMD/FP cluster are 128-bit wide, but 64-bit integer SIMD computation is restricted to port 0 only.
 - FP adder can execute ADDPS/SUBPS in 128-bit datapath, data path for other FP add operations are 64-bit wide,
 - Safe Instruction Recognition algorithm for FP/SIMD execution allow younger, short-latency integer instruction to execute without being blocked by older FP/SIMD instruction that might cause exception,
 - FP multiply pipe also supports memory loads
 - FP ADD instructions with memory load reference can use both ports to dispatch

The memory execution sub-system (MEU) can support 48-bit linear address for Intel 64 Architecture, either 32-bit or 36-bit physical addressing modes. The MEU provides:

- 24KB first level data cache,
- Hardware prefetching for L1 data cache,
- Two levels of DTLB for 4KByte and larger paging structure.
- Hardware pagewalker to service DTLB and ITLB misses.
- Two address generation units (port 0 supports loads and stores, port 1 supports LEA and stack operations)
- Store-forwarding support for integer operations
- 8 write combining buffers.

The bus logic sub-system provides

- 512KB, 8-way set associative, unified L2 cache,
- Hardware prefetching for L2 and interface logic to the front side bus.

13.2.1 Hyper-Threading Technology Support in Intel® Atom™ Microarchitecture

The instruction queue is statically partitioned for scheduling instruction execution from two threads. The scheduler is able to pick one instruction from either thread and dispatch to either of port 0 or port 1 for execution. The hardware makes selection

choice on fetching/decoding/dispatching instructions between two threads based on criteria of fairness as well as each thread's readiness to make forward progress.

13.3 CODING RECOMMENDATIONS FOR INTEL® ATOM™ MICROARCHITECTURE

Instruction scheduling heuristics and coding techniques that apply to out-of-order microarchitectures may not deliver optimal performance on an in-order microarchitecture. Likewise instruction scheduling heuristics and coding techniques for an in-order pipeline like Intel Atom microarchitecture may not achieve optimal performance on out-of-order microarchitectures. This section covers specific coding recommendations for software whose primary deployment targets are processors based on Intel Atom microarchitecture.

13.3.1 Optimization for Front End of Intel® Atom™ Microarchitecture

The two decoders in the front end of Intel Atom microarchitecture can handle most instructions in the Intel 64 and IA-32 architecture. Some instructions dealing with complicated operations require the use of an MSROM in the front end. Instructions that go through the two decoders generally can be decoded by either decoder unit of the front end in most cases. Instructions that must use the MSROM or conditions that cause the front end to re-arrange decoder assignments will experience a delay in the front end.

Software can use specific performance monitoring events to detect instruction sequences and/or conditions that cause front end to re-arrange decoder assignment.

Assembly/Compiler Coding Rule 1. (MH impact, ML generality) For Intel Atom processors, minimize the presence of complex instructions requiring MSROM to take advantage the optimal decode bandwidth provided by the two decode units.

Using the performance monitoring events "MACRO_INSTS.NON_CISC_DECODED" and "MACRO_INSTS.CISC_DECODED" can be used to evaluate the percentage instructions in a workload that required MSROM.

Assembly/Compiler Coding Rule 2. (M impact, H generality) For Intel Atom processors, keeping the instruction working set footprint small will help the front end to take advantage the optimal decode bandwidth provided by the two decode units.

Assembly/Compiler Coding Rule 3. (MH impact, ML generality) For Intel Atom processors, avoiding back-to-back X87 instructions will help the front end to take advantage the optimal decode bandwidth provided by the two decode units.

Using the performance monitoring events "DECODE_RESTRICTION" can count the number of occurrences in a workload that encountered delays causing reduction of decode throughput.

In general the front end restrictions are not typical a performance limiter until the retired “cycle per instruction” becomes less than unity (maximum theoretical retirement throughput corresponds to CPI of 0.5). To reach CPI below unity, it is important to generate instruction sequences that go through the front end as instruction pairs decodes in parallel by the two decoders. After the front end, the scheduler and execution hardware do not need to dispatch the decode pairings through port 0 and port 1 in the same order.

The decoders cannot decode past a jump instruction, so jumps should be paired as the second instruction in a decoder-optimized pairing. The front end can only handle one X87 instruction per cycle, and only decoder unit 0 can request a transfer to use MSROM. Instructions that are longer than 8 bytes or having more than three prefixes will result in a MSROM transfer, experiencing two cycles of delay in the front end.

Instruction lengths and alignment can impact decode throughput. The prefetching buffers inside the front end imposes a throughput limit that if the number of bytes being decoded in any 7-cycle window exceeds 48 bytes, the front end will experience a delay to wait for a buffer. Additionally, every time an instruction pair crosses 16 byte boundary, it requires the front end buffer to be held on for at least one more cycle. So instruction alignment crossing 16 byte boundary is highly problematic.

Instruction alignment can be improved using a combination of an ignore prefix and an instruction.

Example 13-1. Instruction Pairing and Alignment to Optimize Decode Throughput on Intel® Atom™ Microarchitecture

Address	Instruction Bytes	Disassembly
7FFFFDF0	0F594301	mulps xmm0, [ebx+ 01h]
7FFFFDF4	8341FFFF	add dword ptr [ecx-01h], -1
7FFFFDF8	83C2FF	add edx, -1
7FFFFDFB	64	; FS prefix override is ignored, improves code alignment
7FFFFDFC	F20f58E4	add xmm4, xmm4
7FFFFE00	0F594B11	mulps xmm1, [ebx+ 11h]
7FFFFE04	8369EFFF	sub dword ptr [ecx- 11h], -1
7FFFFE08	83EAFB	sub edx, -1
7FFFFE0B	64	; FS prefix override is ignored, improves code alignment
7FFFFE0C	F20F58ED	addsd xmm5, xmm5
7FFFFE10	0F595301	mulps xmm2, [ebx +1]

Example 13-1. Instruction Pairing and Alignment to Optimize Decode Throughput on Intel® Atom™ Microarchitecture

7FFFFFFE14	8341DFFF	add dword ptr [ecx-21H], -1
7FFFFFFE18	83C2FF	add edx, -1
7FFFFFFE1B	64	; FS prefix override is ignored, improves code alignment
7FFFFFFE1C	F20F58F6	addssd xmm6, xmm6
7FFFFFFE20	0F595B11	mulps xmm3, [ebx+ 11h]
7FFFFFFE24	8369CFFF	sub dword ptr [ecx- 31h], -1
7FFFFFFE28	83EAFB	sub edx, -1

When a small loop contains some long-latency operation inside, loop unrolling may be considered as a technique to find adjacent instruction that could be paired with the long-latency instruction to enable that adjacent instruction to make forward progress. However, loop unrolling must also be evaluated on its impact to increased code size and pressure to the branch target buffer.

The performance monitoring event "BACLEARs" can provide a means to evaluate whether loop unrolling is helping or hurting front end performance. Another event "ICACHE_MISSES" can help evaluate if loop unrolling is increasing the instruction footprint.

Branch predictors in Intel Atom processor do not distinguish different branch types. Sometimes mixing different branch types can cause confusion in the branch prediction hardware.

The performance monitoring event "BR_MISSP_TYPE_RETIRED" can provide a means to evaluate branch prediction issues due to branch types.

13.3.2 Optimizing the Execution Core

This section covers several items that can help software use the two-issue-wide execution core to make forward progress with two instructions more frequently.

13.3.2.1 Integer Instruction Selection

In an in-order machine, instruction selection and pairing can have an impact on the machine's ability to discover instruction-level-parallelism for instructions that have data ready to execute. Some examples are:

- **EFLAG:** The consumer instruction of any EFLAG flag bit can not be issued in the same cycle as the producer instruction of the EFLAG register. For example, ADD could modify the carry bit, so it is a producer; JC (or ADC) reads the carry bit and is a consumer.

- Conditional jumps are able to issue in the following cycle after the consumer.
- A consumer instruction of other EFLAG bits must wait one cycle to issue after the producer (two cycle delay).

Assembly/Compiler Coding Rule 4. (M impact, H generality) *For Intel Atom processors, place a MOV instruction between a flag producer instruction and a flag consumer instruction that would have incurred a two-cycle delay. This will prevent partial flag dependency.*

- **Long-latency Integer Instructions:** They will block shorter latency instruction on the same thread from issuing (required by program order). Additionally, they will also block shorter-latency instruction on both threads for one cycle to resolve writeback resource.
- **Common Destination:** Two instructions that produce results to the same destination can not issue in the same cycle.
- **Expensive Instructions:** Some instructions have special requirements and become expensive in consuming hardware resources for an extended period during execution. It may be delayed in execution until it is the oldest in the instruction queue; it may delay the issuing of other younger instructions. Examples of these include FDIV, instructions requiring execution units from both ports, etc.

13.3.2.2 Address Generation

The hardware optimizes the general case of instruction ready to execute must have data ready, and address generation precedes data being ready. If address generation encounters a dependency that needs data from another instruction, this dependency in address generation will incur a delay of 3 cycles.

The address generation unit (AGU) may be used directly in three situations that affect execution throughput of the two-wide machine. The situations are:

- **Implicit ESP updates:** When the ESP register is not used as the destination of an instruction (explicit ESP updates), an implicit ESP update will occur with instructions like PUSH, POP, CALL, RETURN. Mixing explicit ESP updates and implicit ESP updates will also lead to dependency between address generation and data execution.
- **LEA:** The LEA instruction uses the AGU instead of the ALU. If one of the source register of LEA must come from an execution unit. This dependency will also cause a 3 cycle delay. Thus, LEA should not be used in the technique of adding two values and produce the result in a third register. LEA should be used for address computation.
- **Integer-FP/SIMD transfer:** Instructions that transfer integer data to the FP/SIMD side of the machine also uses AGU. Examples of these instructions include MOVD, PINSRW. If one of the source register of these instructions depends on the result of an execution unit, this dependency will also cause a delay of 3 cycles.

Example 13-2. Alternative to Prevent AGU and Execution Unit Dependency

- | |
|---|
| <p>a) Three cycle delay when using LEA in ternary operations</p> <pre>mov eax, 0x01 lea eax, 0x8000[eax+ebp]; values in eax comes from execution of previous instruction ; 3 cycle delay due to lea and execution dependency</pre> <p>b) Dependency handled in execution, avoiding AGU and execution dependency</p> <pre>mov eax, 0x01 add eax, 0x8000 add eax, ebp</pre> |
|---|

Assembly/Compiler Coding Rule 5. (MH impact, H generality) For Intel Atom processors, LEA should be used for address manipulation; but software should avoid the following situations which creates dependencies from ALU to AGU: an ALU instruction (instead of LEA) for address manipulation or ESP updates; a LEA for ternary addition or non-destructive writes which do not feed address generation. Alternatively, hoist producer instruction more than 3 cycles above the consumer instruction that uses the AGU.

13.3.2.3 Integer Multiply

Integer multiply instruction takes several cycles to execute. They are pipelined such that an integer multiply instruction and another long-latency instruction can make forward progress in the execution phase. However, integer multiply instructions will block other single-cycle integer instructions from issuing due to requirement of program order.

Assembly/Compiler Coding Rule 6. (M impact, M generality) For Intel Atom processors, sequence an independent FP or integer multiply after an integer multiply instruction to take advantage of pipelined IMUL execution.

Example 13-3. Pipelining Instruction Execution in Integer Computation

- a) Multi-cycle Imul instruction can block 1-cycle integer instruction
- ```
imul eax, eax
add ecx, ecx ; 1 cycle int instruction blocked by imul for 4 cycles
imul ebx, ebx ; instruction blocked by in-order issue
```
- b) Back-to-back issue of independent imul are pipelined
- ```
imul eax, eax
imul ebx, ebx ; 2nd imul can issue 1 cycle later
add ecx, ecx ; 1 cycle int instruction blocked by imul
```

13.3.2.4 Integer Shift Instructions

Integer shift instructions that encodes shift count in the immediate byte have one-cycle latency. In contrast, shift instructions using shift count in the ECX register may need to wait for the register count are updated. Thus shift instruction using register count has 3-cycle latency.

Assembly/Compiler Coding Rule 7. (M impact, M generality) For Intel Atom processors, hoist the producer instruction for the implicit register count of an integer shift instruction before the shift instruction by at least two cycles.

13.3.2.5 Partial Register Access

Although partial register access does not cause additional delay, the in-order hardware tracks dependency on the full register. Thus 8-bit registers like AL and AH are not treated as independent registers. Additionally some instructions like LEA, vanilla loads, and pop are slower when the input is smaller than 4 bytes.

Assembly/Compiler Coding Rule 8. (M impact, MH generality) For Intel Atom processors, LEA, simple loads and POP are slower if the input is smaller than 4 bytes.

13.3.2.6 FP/SIMD Instruction Selection

Table 13-1 summarizes the characteristics of various execution units in Intel Atom microarchitecture that are likely used most frequently by software.

Table 13-1. Instruction Latency/Throughput Summary of Intel® Atom™ Microarchitecture

Instruction Category	Latency (cycles)	Throughput	# of Execution Unit
SIMD Integer ALU			
128-bit ALU/logical/move	1	1	2
64-bit ALU/logical/move	1	1	2
SIMD Integer Shift			
128-bit	1	1	1
64-bit	1	1	1
SIMD Shuffle			
128-bit	1	1	1
64-bit	1	1	1
SIMD Integer Multiply			
128-bit	5	2	1
64-bit	4	1	1
FP Adder			
X87 Ops (FADD)	5	1	1
Scalar SIMD (addsd, addss)	5	1	1
Packed single (addps)	5	1	1
Packed double (addpd)	6	5	1
FP Multiplier			
X87 Ops (FMUL)	5	2	1
Scalar single (mulss)	4	1	1

Table 13-1. Instruction Latency/Throughput Summary of Intel® Atom™ Microarchitecture

IMUL	Scalar double (mulsd)	5	2	1
	Packed single (mulps)	5	2	1
	Packed double (mulpd)	9	9	1
	IMUL r32, r/m32	5	1	1
	IMUL r12, r/m16	6	1	1

SIMD/FP instruction selection generally should favor shorter latency first, then favor faster throughput alternatives whenever possible. Note that packed double-precision instructions are not pipelined, using two scalar double-precision instead can achieve higher performance in the execution cluster.

Assembly/Compiler Coding Rule 9. (MH impact, H generality) For Intel Atom processors, prefer SIMD instructions operating on XMM register over X87 instructions using FP stack. Use Packed single-precision instructions where possible. Replace packed double-precision instruction with scalar double-precision instructions.

Assembly/Compiler Coding Rule 10. (M impact, ML generality) For Intel Atom processors, library software performing sophisticated math operations like transcendental functions should use SIMD instructions operating on XMM register instead of native X87 instructions.

Assembly/Compiler Coding Rule 11. (M impact, M generality) For Intel Atom processors, enable DAZ and FTZ whenever possible.

Several performance monitoring events may be useful for SIMD/FP instruction selection tuning: "SIMD_INST_RETIRED.{PACKED_SINGLE, SCALAR_SINGLE, PACKED_DOUBLE, SCALAR_DOUBLE}" can be used to determine the instruction selection in the program. "FP_ASSIST" and "SIR" can be used to see if floating exceptions (or false alarms) are impacting program performance.

The latency and throughput of divide instructions vary with input values and data size. Intel Atom microarchitecture implements a radix-2 based divider unit. So, divide/sqrt latency will be significantly longer than other FP operations. The issue throughput rate of divide/sqrt will be correspondingly lower. The divide unit is shared between two logical processors, so software should consider all alternatives to using the divide instructions.

Assembly/Compiler Coding Rule 12. (H impact, L generality) For Intel Atom processors, use divide instruction only when it is absolutely necessary, and pay attention to use the smallest data size operand.

The performance monitoring events "DIV" and "CYCLES_DIV_BUSY" can be used to see if the divides are a bottleneck in the program.

FP operations generally have longer latency than integer instructions. Writeback of results from FP operation generally occur later in the pipe stages than integer pipeline. Consequently, if an instruction has dependency on the result of some FP operation, there will be a two-cycle delay. Examples of these type of instructions are FP-to-integer conversions CVTxx2xx, MOVD from XMM to general purpose registers.

In situations where software needs to do computation with consecutive groups 4 single-precision data elements, PALIGNR+MOVAPS is preferred over MOVUPS. Loading 4 data elements with unconstrained array index k , such as MOVUPS xmm1, _pArray[k], where the memory address _pArray is aligned on 16-byte boundary, will periodically causing cache line split, incurring a 14-cycle delay.

The optimal approach is for each k that is not a multiple of 4, round down k to multiples of 4 with $j = 4*(k/4)$, do a MOVAPS MOVAPS xmm1, _pArray[j] and MOVAPS xmm1, _pArray[j+4], and use PALIGNR to splice together the four data elements needed for computation.

Assembly/Compiler Coding Rule 13. (MH impact, M generality) For Intel Atom processors, prefer a sequence MOVAPS+PALIGN over MOVUPS. Similarly, MOVDQA+PALIGNR is preferred over MOVDQU.

13.3.3 Optimizing Memory Access

This section covers several items that can help software optimize the performance of the memory sub-system.

Memory access to system memory or cache access that encounter certain hazards can cause the memory access to become an expensive operation, blocking short-latency instructions to issue even when they have data ready to execute.

The performance monitoring events "REISSUE" can be used to assess the impact of re-issued memory instructions in the program.

13.3.3.1 Store Forwarding

In a few limited situations, Intel Atom microarchitecture can forward data from a preceding store operation to a subsequent load instruction. The situations are:

- Store-forwarding is supported only in the integer pipe line, and does not apply to FP nor SIMD data. Furthermore, the following conditions must be met:
- The store and load operations must be of the same size and to the same address. Data size larger than 8 bytes do not forward from a store operation.

- When data forwarding proceeds, data is forwarded based on the least significant 12 bits of the address. So software must avoid the address aliasing situation of storing to an address and then loading from another address that aliases in the lowest 12-bits with the store address.

13.3.3.2 First-level Data Cache

Intel Atom microarchitecture handles each 64-byte cache line of the first-level data cache in 16 4-byte chunks. This implementation characteristic has a performance impact to data alignment and some data access patterns.

Assembly/Compiler Coding Rule 14. (MH impact, H generality) *For Intel Atom processors, ensure data are aligned in memory to its natural size. For example, 4-byte data should be aligned to 4-byte boundary, etc. Additionally, smaller access (less than 4 bytes) within a chunk may experience delay if they touch different bytes.*

13.3.3.3 Segment Base

In Intel Atom microarchitecture, the address generation unit assumes that the segment base will be 0 by default. Non-zero segment base will cause load and store operations to experience a delay.

- If the segment base isn't aligned to a cache line boundary, the max throughput of memory operations is reduced to one every 9 cycles.

If the segment base is non-zero but cache line aligned the penalty varies by segment base.

- DS will have a max throughput of one every two cycles.
- FS, and GS will have a max throughput of one every two cycles. However, FS and GS are anticipated to be used only with non-zero bases and therefore have a max throughput of one every two cycles even if the segment base is zero.
- ES,
 - If used as the implicit segment base for the destination of string operation, will have a max throughput of one every two cycles for non-zero but cacheline aligned bases,
 - Otherwise, only do one operation every nine cycles.
- CS, and SS will always have a max throughput of one every nine cycles if its segment base is non-zero but cache line aligned.

Assembly/Compiler Coding Rule 15. (H impact, ML generality) For Intel Atom processors, use segments with base set to 0 whenever possible; avoid non-zero segment base address that is not aligned to cache line boundary at all cost.

Assembly/Compiler Coding Rule 16. (H impact, L generality) For Intel Atom processors, when using non-zero segment bases, Use DS, FS, GS; string operation should use implicit ES.

Assembly/Compiler Coding Rule 17. (M impact, ML generality) For Intel Atom processors, favor using ES, DS, SS over FS, GS with zero segment base.

13.3.3.4 String Moves

Using MOVS/STOS instruction and REP prefix on Intel Atom processor should recognize the following items:

- For small count values, using REP prefix is less efficient than not using REP prefix. This is because the hardware does not have small REP count optimization.
- For small count values, using REP prefix is less efficient than not using REP prefix. This is because the hardware does not have small REP count optimization.
- For large count values, using REP prefix will be less efficient than using 16-byte SIMD instructions.
- Incrementing address in loop iterations should favor LEA instruction over explicit ADD instruction.
- If data footprint is such that memory operation is accessing L2, use of software prefetch to bring data to L1 can avoid memory operation from being re-issued.
- If string/memory operation is accessing system memory, using non-temporal hints of streaming store instructions can avoid cache pollution.

Example 13-4. Memory Copy of 64-byte

```
T1:   prefetch0 [eax+edx+0x80] ; prefetch ahead by two iterations
      movdqa  xmm0, [eax+ edx] ; load data from source (in L1 by prefetch)
      movdqa  xmm1, [eax+ edx+0x10]
      movdqa  xmm2, [eax+ edx+0x20]
      movdqa  xmm3, [eax+ edx+0x30]
      movdqa  [ebx+ edx], xmm0; store data to destination
      movdqa  [ebx+ edx+0x10], xmm1
      movdqa  [ebx+ edx+0x30], xmm2
      movdqa  [ebx+ edx+0x30], xmm3
      lea    edx, 0x40 ; use LEA to adjust offset address for next iteration
      dec    ecx
      jnz    T1
```

13.3.3.5 Parameter Passing

Due to the limited situations of load-to-store forwarding support in Intel Atom microarchitecture, parameter passing via the stack places restrictions on optimal usage by the callee function. For example, "bool" and "char" data usually are pushed onto the stack as 32-bit data, a callee function that reads "bool" or "char" data off the stack will face store-forwarding delay and causing the memory operation to be re-issued.

Compiler should recognize this limitation and generate prolog for callee function to read 32-bit data instead of smaller sizes.

Assembly/Compiler Coding Rule 18. (MH impact, M generality) *For Intel Atom processors, "bool" and "char" value should be passed onto and read off the stack as 32-bit data.*

13.3.3.6 Function Calls

In Intel Atom microarchitecture, using PUSH/POP instructions to manage stack space and address adjustment between function calls/returns will be more optimal than using ENTER/LEAVE alternatives. This is because PUSH/POP will not need MSROM flows and stack pointer address update is done at AGU.

When a callee function need to return to the caller, the callee could issue POP instruction to restore data and restore the stack pointer from the EBP.

Assembly/Compiler Coding Rule 19. (MH impact, M generality) *For Intel Atom processors, favor register form of PUSH/POP and avoid using LEAVE; Use LEA to adjust ESP instead of ADD/SUB.*

13.3.3.7 Optimization of Multiply/Add Dependent Chains

Computations of dependent multiply and add operations can illustrate the usage of several coding techniques to optimize for the front end and in-order execution pipeline of the Intel Atom microarchitecture.

Example 13-5a shows a code sequence that may be used on out-of-order microarchitectures. This sequence is far from optimal on Intel Atom microarchitecture. The full latency of multiply and add operations are exposed and it is not very successful at taking advantage of the two-issue pipeline.

Example 13-5b shows an improved code sequence that takes advantage of the two-issue in-order pipeline of Intel Atom microarchitecture. Because the dependency between multiply and add operations are present, the exposure of latency are only partially covered.

Example 13-5. Examples of Dependent Multiply and Add Computation

a) Instruction sequence that encounters stalls

; accumulator xmm2 initialized

```
Top:  movaps xmm0, [esi] ; vector stored in 16-byte aligned memory
      movaps xmm1, [edi] ; vector stored in 16-byte aligned memory
      mulps xmm0, xmm1
      addps xmm2, xmm0 ; dependency and branch exposes latency of mul and add
      add esi, 16 ;
      add edi, 16
      sub ecx, 1
      jnz top
```

b) Improved instruction sequence to increase execution throughput

; accumulator xmm4 initialized

```
Top:  movaps xmm0, [esi] ; vector stored in 16-byte aligned memory
      lea esi, [esi+16] ; can schedule in parallel with load
      mulps xmm0, [edi] ;
      lea edi, [edi+16] ; can schedule in parallel with multiply
      addps xmm4, xmm0 ; latency exposures partially covered by independent instructions
      dec ecx ;
      jnz top
```

c) Improving instruction sequence further by unrolling and interleaving

; accumulator xmm0, xmm1, xmm2, xmm3 initialized

```
Top:  movaps xmm0, [esi] ; vector stored in 16-byte aligned memory
      lea esi, [esi+16] ; can schedule in parallel with load
      mulps xmm0, [edi] ;
      lea edi, [edi+16] ; can schedule in parallel with multiply
      addps xmm5, xmm1 ; dependent multiply hoisted by unrolling and interleaving
      movaps xmm1, [esi] ; vector stored in 16-byte aligned memory
      lea esi, [esi+16] ; can schedule in parallel with load
      mulps xmm1, [edi] ;
      lea edi, [edi+16] ; can schedule in parallel with multiply
      addps xmm6, xmm2 ; dependent multiply hoisted by unrolling and interleaving
      (continue)
```

Example 13-5. Examples of Dependent Multiply and Add Computation

```

movaps xmm2, [esi] ; vector stored in 16-byte aligned memory
lea esi, [esi+16] ; can schedule in parallel with load
mulps xmm2, [edi] ;
lea edi, [edi+16] ; can schedule in parallel with multiply
addps xmm7, xmm3 ; dependent multiply hoisted by unrolling and interleaving
movaps xmm3, [esi] ; vector stored in 16-byte aligned memory
lea esi, [esi+16] ; can schedule in parallel with load
mulps xmm3, [edi] ;
lea edi, [edi+16] ; can schedule in parallel with multiply
addps xmm4, xmm0 ; dependent multiply hoisted by unrolling and interleaving
sub ecx, 4;
jnz top
; sum up accumulators xmm0, xmm1, xmm2, xmm3 to reduce dependency inside the loop

```

Example 13-5c illustrates a technique that increases instruction-level parallelism and further reduces latency exposures of the multiply and add operations. By unrolling four times, each ADDPS instruction can be hoisted far from its dependent producer instruction MULPS. Using an interleaving technique, non-dependent ADDPS and MULPS can be placed in close proximity. Because the hardware that executes MULPS and ADDPS is pipelined, the associated latency can be covered much more effectively by this technique relative to Example 13-5b.

13.3.3.8 Position Independent Code

Position independent code often needs to obtain the value of the instruction pointer. Example 13-5a show one technique to put the value of IP into the ECX register by issuing a CALL without a matching RET. Example 13-5b show an alternative technique to put the value of IP into the ECX register using a matched pair of CALL/RET.

Example 13-6. Instruction Pointer Query Techniques

```

a) Using call without return to obtain IP
   call _label; return address pushed is the IP of next instruction
_label:
   pop ECX; IP of this instruction is now put into ECX

```

Example 13-6. Instruction Pointer Query Techniques

```

b) Using matched call/ret pair

    call _blcx;
    ... ; ECX now contains IP of this instruction
    ...
_blcx
    mov ecx, [esp];
    ret

```

13.4 INSTRUCTION LATENCY

This section lists the port-binding and latency information of Intel Atom microarchitecture. The port-binding information for each instruction may show one of 3 situations:

- ‘single digit’ - the specific port that must be issued,
- (0, 1) - either port 0 or port 1,
- ‘B’ - both ports are required.

In the “Instruction” column:

- if different operand syntax of the same instruction have the same port-binding and latency, operand syntax is omitted.
- when different operand syntax may produce different latency or port binding, the operand syntax is listed; but instruction syntax of different operand sizes may be compacted and abbreviated with a footnote.

Instruction that required decoder assistance from MSROM are marked in the “Comment” column (should be used minimally if more decode-efficient alternatives are available).

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
ADD/AND/CMP/OR/SUB/XOR/TEST ¹ (E)AX/AL, imm;	(0, 1)	1	0.5
ADD/AND/CMP/OR/SUB/XOR ² mem, imm8; ADD/AND/CMP/OR/SUB/XOR/TEST ⁴ mem, imm; TEST m8, imm8	0	1	1
ADD/AND/CMP/OR/SUB/XOR/TEST ² mem, reg; ADD/AND/CMP/OR/SUB/XOR ² reg, mem;	0	1	1
ADD/AND/CMP/OR/SUB/XOR ² reg, imm8; ADD/AND/CMP/OR/SUB/XOR ⁴ reg, imm	(0, 1)	1	0.5
ADDPD/ADDSUBPD/MAXPD/MAXPS/MINPD/MINPS/SUBPD xmm, mem	B	7	6
ADDPD/ADDSUBPD/MAXPD/MAXPS/MINPD/MINPS/SUBPD xmm, xmm	B	6	5
ADDPS/ADDS/ADDSS/ADDSUBPS/SUBPS/SUBSD/SUBSS xmm, mem	B	5	1
ADDPS/ADDS/ADDSS/ADDSUBPS/SUBPS/SUBSD/SUBSS xmm, xmm	1	5	1
ANDNP/ANDNPS/ANDPD/ANDPS/ORPD/ORPS/XORPD/XORPS xmm, mem	0	1	1
ANDNP/ANDNPS/ANDPD/ANDPS/ORPD/ORPS/XORPD/XORPS xmm, xmm	(0, 1)	1	1
BSF/BSR r16, m16	B	17	16
BSF/BSR ³ reg, mem	B	16	15
BSF/BSR ⁴ reg, reg	B	16	15
BT m16, imm8; BT ³ mem, imm8	(0, 1)	2; 1	1
BT m16, r16; BT ³ mem, reg	B	10, 9	8
BT ⁴ reg, imm8; BT ⁴ reg, reg	1	1	1
BTC m16, imm8; BTC ³ mem, imm8	B	3; 2	2
BTC/BTR/BTS m16; r16	B	12	11
BTC/BTR/BTS ³ mem, reg	B	11	10
BTC/BTR/BTS ⁴ reg, imm8; BTC/BTR/BTS ⁴ reg, reg	1	1	1
CALL mem	(0, 1)	2	2
CALL reg; CALL rel16; CALL rel32	B	1	1

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
CMOV ⁴ reg, mem; MOV ¹ (E)AX/AL, MOFFS; MOV ² mem, imm	0	1	1
CMOV ⁴ reg, reg; MOV ² reg, imm; MOV ² reg, reg; ; SETcc r8	(0, 1)	1	0.5
CMPPD/CMPPS xmm, mem, imm; CVTTPS2DQ xmm, mem	B	7	6
CMPPD/CMPPS xmm, xmm, imm; CVTTPS2DQ xmm, xmm	B	6	5
CMPSSD/CMPSS xmm, mem, imm	B	5	1
CMPSSD/CMPSS xmm, xmm, imm	1	5	1
(U)COMISD/(U)COMISS xmm, mem;	B	10	9
(U)COMISD/(U)COMISS xmm, xmm;	B	9	8
CVTDQ2PD/CVTPD2DQ/CVTPD2PS xmm, mem	B	8	7
CVTDQ2PD/CVTPD2DQ/CVTPD2PS xmm, xmm	B	7	6
CVTDQ2PS/CVTS2SS/CVTSI2SS/CVTSS2SD xmm, mem	B	7	6
CVTDQ2PS/CVTS2SS/CVTSS2SD xmm, xmm	B	6	5
CVT(T)PD2PI mm, mem; CVTPI2PD xmm, mem	B	8	7
CVT(T)PD2PI mm, xmm; CVTPI2PD xmm, mm	B	7	6
CVTPI2PS/CVTSI2SD xmm, mem;	B	5	4
CVTPI2PS xmm, mm;	1	5	1
CVT(T)PS2PI mm, mem;	B	5	5
CVT(T)PS2PI mm, xmm;	1	5	1
CVT(T)SD2SI ³ reg, mem; CVT(T)SS2SI r32, mem	B	9	8
CVT(T)SD2SI ³ reg, xmm; CVT(T)SS2SI r32, xmm	B	8	7
CVTSI2SD xmm, r32; CVTSI2SS xmm, r32	B	7; 6	5
CVTSI2SD xmm, r64; CVTSI2SS xmm, r64	B	6; 7	5
CVT(T)SS2SI r64, mem; RCPPS xmm, mem	B	10	9
CVT(T)SS2SI r64, xmm; RCPPS xmm, xmm	B	9	8
CVTTPD2DQ xmm, mem	B	8	7
CVTTPD2DQ xmm, xmm	B	7	6
DEC/INC ² mem; MASKMOVQ; MOVAPD/MOVAPS mem, xmm	0	1	1
DEC/INC ² reg; FLD ST; FST/FSTP ST; MOVDQ2Q mm, xmm	(0, 1)	1	0.5
DIVPD; DIVPS	B	125; 70	124; 69

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
DIVSD; DIVSS	B	62; 34	61; 33
EMMS; LDMXCSR	B	5	4
FABS/FCHS/FXCH; MOVQ2DQ xmm, mm; MOVSX/MOVZX r16, r16	(0, 1)	1	0.5
FADD/FSUB/FSUBR ³ mem	B	5	4
FADD/FADDP/FSUB/FSUBP/FSUBR/FSUBRP ST;	1	5	1
FCMOV	B	6	5
FCOM/FCOMP ³ mem	B	1	1
FCOM/FCOMP/FCOMPP/FUCOM/FUCOMP ST; FTST	1	1	1
FCOMI/FCOMIP/FUCOMI/FUCOMIP ST	B	9	8
FDIV/FSQRT ³ mem; FDIV/FSQRT ST	0	25-65	24-64
FIADD/FIMUL ⁵ mem	B	11	10
FICOM/FICOMP mem	B	7	6
FILD ⁴ mem	B	5	4
FLD ³ mem; FXAM; MOVAPD/MOVAPS/MOVD xmm, mem	0	1	1
FLDCW	B	5	4
FMUL/FMULP ST; FMUL ³ mem	0	5	1
FNSTSW AX; FNSTSW m16	B	10; 14	9; 13
FST/FSTP ³ mem	B	2	1
HADDPD/HADDP/HSUBPD/HSUBPS xmm, mem	B	9	8
HADDPD/HADDP/HSUBPD/HSUBPS xmm, xmm	B	8	7
IDIV r/m8; IDIV r/m16; IDIV r/m32; IDIV r/m64;	B	33;42;57; 197	32;41;5 6;196
IMUL/MUL ⁶ EAX/AL, mem; IMUL/MUL AX, m16	B	7; 8	6; 7
IMUL/MUL ⁷ AX/AL, reg; IMUL/MUL EAX, r32	B	7; 6	6; 5
IMUL m16, imm8/imm16; IMUL r16, m16	B	7;	6
IMUL r/m32, imm8/imm32; IMUL r32, r/m32	0	5	1
IMUL r/m64, imm8/imm32;	B	14	13
IMUL r16, r16; IMUL r16, imm8/imm16	B	6	5
IMUL r64, r/m64; IMUL/MUL RAX, r/m64	B	11; 12	10; 11

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
JCC ¹ ; JMP ⁴ reg; JMP ¹	1	1	1
JCXZ; JECXZ; JRCXZ	B	4	1
JMP mem ⁴ ;	B	2	1
LDDQU; MOVDQU/MOVUPD/MOVUPS xmm, mem;	B	3	2
LEA r16, mem; MASKMOVDQU; SETcc m8	(0, 1)	2	1
LEA, reg, mem	1	1	1
LEAVE;	B	2;	2
MAXSD/MAXSS/MINSD/MINSS xmm, mem	B	5	1
MAXSD/MAXSS/MINSD/MINSS xmm, xmm	1	5	1
MOV ² MOFFS, (E)AX/AL; MOV ² reg, mem; MOV ² mem, reg	0	1	1
MOVD mem ³ , mm; MOVD xmm, reg ³ ; MOVD mm, mem ³	0	1	1
MOVD reg ³ , mm; MOVD reg ³ , xmm; PMOVMASK reg ³ , mm	0	3	1
MOVDQA/MOVQ xmm, mem; MOVDQA/MOVD mem, xmm;	0	1	1
MOVDQA/MOVDQU/MOVUPD xmm, xmm; MOVQ mm, mm	(0, 1)	1	0.5
MOVDQU/MOVUPD/MOVUPS mem, xmm;	B	2	2
MOVHPS; MOVLHPS; MOVHPD; MOVHPS; MOVLPD; MOVLPS	0	1	1
MOVMSKPD/MOVSKPS/PMOVMASKB reg ³ , xmm	0	3	1
MOVNTI ³ mem, reg; MOVNTPD/MOVNTPS; MOVNTQ	0	1	1
MOVQ mem, mm; MOVQ mm, mem; MOVDUP	0	1	1
MOVSD/MOVSS xmm, xmm; MOVXSD ⁵ reg, reg	(0, 1)	1	0.5
MOVSD/MOVSS xmm, mem; PALIGNR	0	1	1
MOVSD/MOVSS mem, xmm; PINSRW	0	1	1
MOVSHDUP/MOVSLDUP xmm, mem	0	1	1
MOVSHDUP/MOVSLDUP/MOVUPS xmm, xmm	(0, 1)	1	0.5
MOVSB/MOVB r16, m8; MOVSB/MOVB r16, r8	0	3; 2	1
MOVSB/MOVB reg ³ , r/m8; MOVSB/MOVB reg ³ , r/m16	0	1	1
MOVXSD ⁵ reg, mem; MOVXSD r64, r/m32	0	1	1
MULPS/MULSD xmm, mem; MULSS xmm, mem;	0	5; 4	2
MULPS/MULSD xmm, xmm; MULSS xmm, xmm	0	5; 4	2

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
MULPD	B	5; 4	2
NEG/NOT ² mem; PREFETCHNTA; PREFETCHTx	0	10	9
NEG/NOT ² reg; NOP	(0, 1)	1	0.5
PABSB/D/W mm, mem; PABSB/D/W xmm, mem	0	1	1
PABSB/D/W mm, mm; PABSB/D/W xmm, xmm	(0, 1)	1	0.5
PACKSSDw/WB mm, mem; PACKSSDw/WB xmm, mem	0	1	1
PACKSSDw/WB mm, mm; PACKSSDw/WB xmm, xmm	0	1	1
PACKUSWB mm, mem; PACKUSWB xmm, mem	0	1	1
PACKUSWB mm, mm; PACKUSWB xmm, xmm	0	1	1
PADDB/D/W/Q mm, mem; PADDB/D/W/Q xmm, mem	0	1	1
PADDB/D/W/Q mm, mm; PADDB/D/W/Q xmm, xmm	(0, 1)	1	0.5
PADDSB/W mm, mem; PADDSB/W xmm, mem	0	1	1
PADDSB/W mm, mm; PADDSB/W xmm, xmm	(0, 1)	1	0.5
PADDUSB/W mm, mem; PADDUSB/W xmm, mem	0	1	1
PADDUSB/W mm, mm; PADDUSB/W xmm, xmm	(0, 1)	1	0.5
PAND/PANDN/POR/PXOR mm, mem; PAND/PANDN/POR/PXOR xmm, mem	0	1	1
PAND/PANDN/POR/PXOR mm, mm; PAND/PANDN/POR/PXOR xmm, xmm	(0, 1)	1	0.5
PAVGB/W mm, mem; PAVGB/W xmm, mem	0	1	1
PAVGB/W mm, mm; PAVGB/W xmm, xmm	(0, 1)	1	0.5
PCMPEQB/D/W mm, mem; PCMPEQB/D/W xmm, mem	0	1	1
PCMPEQB/D/W mm, mm; PCMPEQB/D/W xmm, xmm	(0, 1)	1	0.5
PCMPGTB/D/W mm, mem; PCMPGTB/D/W xmm, mem	0	1	1
PCMPGTB/D/W mm, mm; PCMPGTB/D/W xmm, xmm	(0, 1)	1	0.5
PEXTRw;	B	4	1
PHADDD/PHSUBD mm, mem; PHADDD/PHSUBD xmm, mem	B	4	3
PHADDD/PHSUBD mm, mm; PHADDD/PHSUBD xmm, xmm	B	3	2
PHADDW/PHADDSw mm, mem; PHADDW/PHADDSw xmm, mem	B	6; 8	5; 7
PHADDW/PHADDSw mm, mm; PHADDW/PHADDSw xmm, xmm	B	5; 7	M

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
PHSUBW/PHSUBSW mm, mem; PHSUBW/PHSUBSW xmm, mem	B	6; 8	M
PHSUBW/PHSUBSW mm, mm; PHSUBW/PHSUBSW xmm, xmm	B	5; 7	M
PMADDUBSW/PMADDWD/PMULHRSW/PSADBW mm, mm; PMADDUBSW/PMADDWD/PMULHRSW/PSADBW mm, mem	0	4	1
PMADDUBSW/PMADDWD/PMULHRSW/PSADBW xmm, xmm; PMADDUBSW/PMADDWD/PMULHRSW/PSADBW xmm, mem	0	5	1
PMAXSW/UB mm, mem; PMAXSW/UB xmm, mem	0	1	1
PMAXSW/UB mm, mm; PMAXSW/UB xmm, xmm	(0, 1)	1	0.5
PMINSW/UB mm, mem; PMINSW/UB xmm, mem	0	1	1
PMINSW/UB mm, mm; PMINSW/UB xmm, xmm	(0, 1)	1	0.5
PMULHUW/PMULHW/PMULLW/PMULUDQ mm, mm; PMULHUW/PMULHW/PMULLW/PMULUDQ mm, mem	0	4	1
PMULHUW/PMULHW/PMULLW/PMULUDQ xmm, xmm; PMULHUW/PMULHW/PMULLW/PMULUDQ xmm, mem	0	5	1
POP mem ⁵ ; PSLLD/Q/W mm, mem; PSLLD/Q/W xmm, mem	B	3	2
POP r16; PUSH mem ⁴ ; PSLLD/Q/W mm, mm; PSLLD/Q/W xmm, xmm	B	2	1
POP reg ³ ; PUSH reg ⁴ ; PUSH imm	B	1	1
POPA ; POPAD	B	9	8
PSHUFB mm, mem; PSHUFD; PSHUFHW; PSHUFLW; PSHUFW	0	1	1
PSHUFB mm, mm; PSLLD/Q/W mm, imm; PSLLD/Q/W xmm, imm	0	1	1
PSHUFB xmm, mem	B	5	4
PSHUFB xmm, xmm	B	4	3
PSIGNB/D/W mm, mem; PSIGNB/D/W xmm, mem	0	1	1
PSIGNB/D/W mm, mm; PSIGNB/D/W xmm, xmm	(0, 1)	1	0.5
PSRAD/W mm, imm; PSRAD/W xmm, imm;	0	1	1
PSRLD/Q/W mm, mem; PSRLD/Q/W xmm, mem	B	3	2
PSRLD/Q/W mm, mm; PSRLD/Q/W xmm, xmm	B	2	1
PSRLD/Q/W mm, imm; PSRLD/Q/W xmm, imm;	0	1	1
PSLLDQ/PSRLDQ xmm, imm; SHUFPD/SHUFPS	0	1	1
PSUBB/D/W/Q mm, mem; PSUBB/D/W/Q xmm, mem	0	1	1

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
PSUBB/D/W/Q mm, mm; PSUBB/D/W/Q xmm, xmm	(0, 1)	1	0.5
PSUBSB/W mm, mem; PSUBSB/W xmm, mem	0	1	1
PSUBSB/W mm, mm; PSUBSB/W xmm, xmm	(0, 1)	1	0.5
PSUBUSB/W mm, mem; PSUBUSB/W xmm, mem	0	1	1
PSUBUSB/W mm, mm; PSUBUSB/W xmm, xmm	(0, 1)	1	0.5
PUNPCKHBW/DQ/WD; PUNPCKLBW/DQ/WD	0	1	1
PUNPCKHQDQ; PUNPCKLQDQ	0	1	1
PUSHA ; PUSHAD	B	8	7
RCL mem ² , 1; RCL reg ² , 1	0	1	1
RCL m8, CL; RCL m16, CL; RCL mem ³ , CL;	B	18;16; 14	17;15;13
RCL m8, imm; RCL m16, imm; RCL mem ³ , imm;	B	18; 17; 14	17;16;13
RCL r8, CL; RCL r16, CL; RCL reg ³ , CL;	B	17; 16; 14	16;15;14
RCL r8, imm; RCL r16, imm; RCL reg ³ , imm;	B	18;16; 14	17;15;13
RCPSS	0	4	1
RCR mem ² , 1; RCR reg ² , 1	B	7; 5	6;4
RCR m8, CL; RCR m16, CL; RCR mem ³ , CL;	B	15; 13; 12	14;12;11
RCR m8, imm; RCR m16, imm; RCR mem ³ , imm;	B	16;14; 12	15;13;11
RCR r8, CL; RCR r16, CL; RCR reg ³ , CL;	B	14; 13; 12	13;12;11
RCR r8, imm; RCR r16, imm; RCR reg ³ , imm;	B	15, 14, 12	14;13;11
RET imm16	B	1	1
RET (far)	B	79	
ROL; ROR; SAL; SAR; SHL; SHR	0	1	1
SETcc		1	1
SHLD ⁸ mem, reg, imm; SHLD r64, r64, imm; SHLD m64, r64, CL	B	11	10

Table 13-2. Intel® Atom™ Microarchitecture Instructions Latency Data (Contd.)

Instruction	Ports	Latency	Throughput
DisplayFamily_DisplayModel	06_1CH	06_1CH	06_1CH
SHLD m32, r32; SHLD r32, r32	B	4; 2	3; 1
SHLD m16, r16, CL; SHLD r16, r16, imm; SHLD r64, r64, CL	B	10	9
SHLD r16, r16, CL; SHRD m64, r64; SHRD r64, r64, imm	B	9	8
SHRD m32, r32; SHRD r32, r32	B	4; 2	3; 1
SHRD m16, r16; SHRD r16, r16	B	6	5
SHRD r64, r64, CL	B	8	7
STMXCSR	B	15	14
TEST ² reg, reg; TEST ⁴ reg, imm	(0, 1)	1	0.5
UNPCKHPD; UNPCKHPS; UNPCKLPD, UNPCKLPS	0	1	1

Notes on operand size (osize) and address size (asize):

1. osize = 8, 16, 32 or asize = 8, 16, 32
2. osize = 8, 16, 32, 64
3. osize = 32, 64
4. osize = 16, 32, 64 or asize = 16, 32, 64
5. osize = 16, 32
6. osize = 8, 32
7. osize = 8, 16
8. osize = 16, 64

APPENDIX A

APPLICATION PERFORMANCE TOOLS

Intel offers an array of application performance tools that are optimized to take advantage of the Intel architecture (IA)-based processors. This appendix introduces these tools and explains their capabilities for developing the most efficient programs without having to write assembly code.

The following performance tools are available:

- **Intel® C++ Compiler and Intel® Fortran Compiler** — Intel compilers generate highly optimized executable code for Intel 64 and IA-32 processors. The compilers support advanced optimizations that include auto-vectorization for MMX technology, and the Streaming SIMD Extensions (SSE) instruction set architectures (SSE, SSE2, SSE3, SSSE3, and SSE4) of our latest processors.
- **VTune Performance Analyzer** — The VTune analyzer collects, analyzes, and displays Intel architecture-specific software performance data from the system-wide view down to a specific line of code.
- **Intel® Performance Libraries** — The Intel Performance Library family consists of a set of software libraries optimized for Intel architecture processors. The library family includes the following:
 - Intel® Math Kernel Library (Intel® MKL)
 - Intel® Integrated Performance Primitives (Intel® IPP)
- **Intel® Threading Tools** — Intel Threading Tools consist of the following:
 - Intel® Thread Checker
 - Intel® Thread Profiler
- **Intel® Cluster Tools** - The Intel® Cluster Toolkit 3.1 helps you develop, analyze and optimize performance of parallel applications for clusters using IA-32, IA-64, and Intel® 64 architectures. Intel Cluster Tools consist of the following:
 - Intel® Cluster Tool Kit
 - Intel® MPI Library
 - Intel® Trace Analyzer and Collector
- **Intel® XML Products** - Intel XML Products consist of the following:
 - Intel® XML Software Suite 1.0 Beta
 - Intel® SOA Security Toolkit 1.0 Beta for Axis2
 - Intel® XSLT Accelerator 1.1 for Java* Environments on Linux* and Windows* Operating Systems

A.1 COMPILERS

Intel compilers support several general optimization settings, including `/O1`, `/O2`, `/O3`, and `/fast`. Each of them enables a number of specific optimization options. In most cases, `/O2` is recommended over `/O1` because the `/O2` option enables function expansion, which helps programs that have many calls to small functions. The `/O1` may sometimes be preferred when code size is a concern. The `/O2` option is on by default.

The `/Od` (`-O0` on Linux) option disables all optimizations. The `/O3` option enables more aggressive optimizations, most of which are effective only in conjunction with processor-specific optimizations described below.

The `/fast` option maximizes speed across the entire program. For most Intel 64 and IA-32 processors, the `/fast` option is equivalent to `/O3 /Qipo /QxP` (`-Q3 -ipo -static -xP` on Linux). For Mac OS, the `-fast` option is equivalent to `-O3 -ipo`.

All the command-line options are described in Intel® C++ Compiler documentation.

A.1.1 Recommended Optimization Settings for Intel® 64 and IA-32 Processors

64-bit addressable code can only run in 64-bit mode of processors that support Intel 64 architecture. The optimal compiler settings for 64-bit code generation is different from 32-bit code generation. Table A-1 lists recommended compiler options for generating 32-bit code for Intel 64 and IA-32 processors. Table A-1 also applies to code targeted to run in compatibility mode on an Intel 64 processor, but does not apply to running in 64-bit mode. Table A-2 lists recommended compiler options for generating 64-bit code for Intel 64 processors, it only applies to code target to run in 64-bit mode. Intel compilers provide separate compiler binary to generate 64-bit code versus 32-bit code. The 64-bit compiler binary generates only 64-bit addressable code.

Table A-1. Recommended IA-32 Processor Optimization Options

Need	Recommendation	Comments
Best performance on Intel Core 2 processor family and Intel Xeon processor 5400 series, utilizing SSE4 instructions	<ul style="list-style-type: none"> <code>/QxS</code> (<code>-xS</code> on Linux and Mac OS) 	<ul style="list-style-type: none"> Single code path

Table A-1. Recommended IA-32 Processor Optimization Options

Need	Recommendation	Comments
Best performance on Intel Core 2 processor family and Intel Xeon processor 5400 series, utilizing SSE4 instructions	<ul style="list-style-type: none"> • /QaxS (-axS on Linux and Mac OS) 	<ul style="list-style-type: none"> • Multiple code path are generated • Be sure to validate your application on all systems where it may be deployed.
Best performance on Intel Core 2 processor family and Intel Xeon processor 3000 and 5100 series, utilizing SSSE3 and other processor-specific instructions	<ul style="list-style-type: none"> • /QxT (-xT on Linux) 	<ul style="list-style-type: none"> • Single code path • Will not run on earlier processors that do not support SSSE3
Best performance on Intel Core 2 processor family and Intel Xeon processor 3000 and 5100 series, utilizing SSSE3; runs on non-Intel processor supporting SSE2	<ul style="list-style-type: none"> • /QaxT /QxW (-axT -xW on Linux) 	<ul style="list-style-type: none"> • Multiple code path are generated • Be sure to validate your application on all systems where it may be deployed.
Best performance on IA-32 processors with SSE3 instruction support	/QxP (-xP on Linux)	<ul style="list-style-type: none"> • Single code path • Will not run on earlier processors.that do not support SSE3
Best performance on IA-32 processors with SSE2 instruction support	/QaxN (-axN on Linux) Optimized for Pentium 4 and Pentium M processors, and an optimized, generic code-path to run on other processors	<ul style="list-style-type: none"> • Multiple code paths are generated. • Use /QxN (-xN for Linux) if you know your application will not be run on processors older than the Pentium 4 or Pentium M processors.

Table A-1. Recommended IA-32 Processor Optimization Options

Need	Recommendation	Comments
Best performance on IA-32 processors with SSE3 instruction support for multiple code paths	<ul style="list-style-type: none"> • /QaxP /QxW (-axP -xW on Linux) • Optimized for Pentium 4 processor and Pentium 4 processor with SSE3 instruction support 	Generates two code paths: <ul style="list-style-type: none"> • one for the Pentium 4 processor • one for the Pentium 4 processor or non-Intel processors with SSE3 instruction support.

Table A-2. Recommended Processor Optimization Options for 64-bit Code

Need	Recommendation	Comments
Best performance on Intel Core 2 processor family and Intel Xeon processor 3000 and 5100 series, utilizing SSSE3 and other processor-specific instructions	<ul style="list-style-type: none"> • /QxT (-xT on Linux) 	<ul style="list-style-type: none"> • Single code path • Will not run on earlier processors that do not support SSSE3
Best performance on Intel Core 2 processor family and Intel Xeon processor 3000 and 5100 series, utilizing SSSE3; runs on non-Intel processor supporting SSE2	<ul style="list-style-type: none"> • /QaxT /QxW (-axT -xW on Linux) 	<ul style="list-style-type: none"> • Multiple code path are generated • Be sure to validate your application on all systems where it may be deployed.
Best performance on other processors supporting Intel 64 architecture, utilizing SSE3 where possible	<ul style="list-style-type: none"> • /QxP (-xP on Linux) 	<ul style="list-style-type: none"> • Single code path are generated • Will not run on processors that do not support Intel 64 architecture and SSE3.
Best performance on other processors supporting Intel 64 architecture, utilizing SSE3 where possible, while still running on older Intel as well as non-Intel x86-64 processors supporting SSE2	<ul style="list-style-type: none"> • /QaxP /QxW (-axP -xW on Linux) 	<ul style="list-style-type: none"> • Multiple code path are generated • Be sure to validate your application on all systems where it may be deployed.

A.1.2 Vectorization and Loop Optimization

The Intel C++ and Fortran Compiler's vectorization feature can detect sequential data access by the same instruction and transforms the code to use SSE, SSE2, SSE3, SSSE3 and SSE4, depending on the target processor platform. The vectorizer supports the following features:

- Multiple data types: Float/double, char/short/int/long (both signed and unsigned), _Complex float/double are supported.
- Step by step diagnostics: Through the /Qvec-report[n] (-vec-report[n] on Linux and Mac OS) switch (see Table A-3), the vectorizer can identify, line-by-line and variable-by-variable, what code was vectorized, what code was not vectorized, and more importantly, why it was not vectorized. This feedback gives the developer the information necessary to slightly adjust or restructure code, with dependency directives and restrict keywords, to allow vectorization to occur.
- Advanced dynamic data-alignment strategies: Alignment strategies include loop peeling and loop unrolling. Loop peeling can generate aligned loads, enabling faster application performance. Loop unrolling matches the prefetch of a full cache line and allows better scheduling.
- Portable code: By using appropriate Intel compiler switches to take advantage new processor features, developers can avoid the need to rewrite source code.

The processor-specific vectorizer switch options are: -Qx[K,W,N,P,S,T] and -Qax[K,W,N,P,S,T]. The compiler provides a number of other vectorizer switch options that allow you to control vectorization. The latter switches require the -Qx[,K,W,N,P,T] or -Qax[,K,W,N,P,T] switch to be on. The default is off.

Table A-3. Vectorization Control Switch Options

-Qvec_report[n]	Controls the vectorizer's diagnostic levels, where n is either 0, 1, 2, or 3.
-Qrestrict	Enables pointer disambiguation with the restrict qualifier.

A.1.2.1 Multithreading with OpenMP*

Both the Intel C++ and Fortran Compilers support shared memory parallelism using OpenMP compiler directives, library functions and environment variables. OpenMP directives are activated by the compiler switch /Qopenmp (-openmp on Linux and Mac OS). The available directives are described in the Compiler User's Guides available with the Intel C++ and Fortran Compilers. For information about the OpenMP standard, see <http://www.openmp.org>.

A.1.2.2 Automatic Multithreading

Both the Intel C++ and Fortran Compilers can generate multithreaded code automatically for simple loops with no dependencies. This is activated by the compiler switch /Qparallel (-parallel in Linux and Mac OS).

A.1.3 Inline Expansion of Library Functions (/Oi, /Oi-)

The compiler inlines a number of standard C, C++, and math library functions by default. This usually results in faster execution. Sometimes, however, inline expansion of library functions can cause unexpected results. For explanation, see the Intel C++ Compiler documentation.

A.1.4 Floating-point Arithmetic Precision (/Op, /Op-, /Qprec, /Qprec_div, /Qpc, /Qlong_double)

These options provide a means of controlling optimization that might result in a small change in the precision of floating-point arithmetic.

A.1.5 Rounding Control Option (/Qrcr, /Qrcd)

The compiler uses the -Qrcd option to improve the performance of code that requires floating-point calculations. The optimization is obtained by controlling the change of the rounding mode.

The -Qrcd option disables the change to truncation of the rounding mode in floating-point-to-integer conversions.

For more on code optimization options, see the Intel C++ Compiler documentation.

A.1.6 Interprocedural and Profile-Guided Optimizations

The following are two methods to improve the performance of your code based on its unique profile and procedural dependencies.

A.1.6.1 Interprocedural Optimization (IPO)

You can use the /Qip (-ip in Linux and Mac OS) option to analyze your code and apply optimizations between procedures within each source file. Use multifile IPO with /Qipo (-ipo in Linux and Mac OS) to enable the optimizations between procedures in separate source files.

A.1.6.2 Profile-Guided Optimization (PGO)

Creates an instrumented program from your source code and special code from the compiler. Each time this instrumented code is executed, the compiler generates a dynamic information file. When you compile a second time, the dynamic information files are merged into a summary file. Using the profile information in this file, the compiler attempts to optimize the execution of the most heavily travelled paths in the program.

Profile-guided optimization is particularly beneficial for the Pentium 4 and Intel Xeon processor family. It greatly enhances the optimization decisions the compiler makes regarding instruction cache utilization and memory paging. Also, because PGO uses execution-time information to guide the optimizations, branch-prediction can be significantly enhanced by reordering branches and basic blocks to keep the most commonly used paths in the microarchitecture pipeline, as well as generating the appropriate branch-hints for the processor.

When you use PGO, consider the following guidelines:

- Minimize the changes to your program after instrumented execution and before feedback compilation. During feedback compilation, the compiler ignores dynamic information for functions modified after that information was generated.

NOTE

The compiler issues a warning that the dynamic information corresponds to a modified function.

- Repeat the instrumentation compilation if you make many changes to your source files after execution and before feedback compilation.

For more on code optimization options, see the Intel C++ Compiler documentation.

A.1.7 Auto-Generation of Vectorized Code

This section covers several high-level language examples that programmers can use Intel Compiler to generate vectorized code automatically.

Example 13-7. Storing Absolute Values

```
int dst[1024], src[1024]
for (i = 0; i < 1024; i++) {
    dst[i] = (src[i] >= 0) ? src[i] : -src[i];
}
```

The following examples are illustrative of the likely differences of two compiler switches.

Example 13-8. Auto-Generated Code of Storing Absolutes

Compiler Switch QxW	Compiler Switch QxT
<pre>movdqa xmm1, _src[eax*4] pxor xmm0, xmm0 pcmpgtd xmm0, xmm1 pxor xmm1, xmm0 psubd xmm1, xmm0 movdqa _dst[eax*4], xmm1 add eax, 4 cmp eax, 1024 jb \$B1\$3</pre>	<pre>pabsd xmm0, _src[eax*4] movdqa _dst[eax*4], xmm0 add eax, 4 cmp eax, 1024 jb \$B1\$3</pre>

Example 13-9. Changes Signs

```
int dst[NUM], src[1024];
for (i = 0; i < 1024; i++) {
    if (src[i] == 0)
        { dst[i] = 0; }
    else if (src[i] < 0)
        { dst[i] = -dst[i]; }
}
```

Example 13-10. Auto-Generated Code of Sign Conversion

Compiler Switch Qxw	Compiler Switch QxT
<pre> \$B1\$3: mov edx, _src[eax*4] add eax, 1 test edx, edx jne \$B1\$5 \$B1\$4: mov _dst[eax*4], 0 jmp \$B1\$7 ALIGN 4 \$B1\$5: jge \$B1\$7 \$B1\$6: mov edx, _dst[eax*4] neg edx mov _dst[eax*4], edx \$B1\$7: cmp eax, 1024 jl \$B1\$3 </pre>	<pre> \$B1\$3: movdqa xmm0, _dst[eax*4] psignd xmm0, _src[eax*4] movdqa _dst[eax*4], xmm0 add eax, 4 cmp eax, 1024 jb \$B1\$3 </pre>

Example 13-11. Data Conversion

```

int dst[1024];
unsigned char src[1024];
for (i = 0; i < 1024; i++) {
    dst[i] = src[i]
}

```

Example 13-12. Auto-Generated Code of Data Conversion

Compiler Switch QxW	Compiler Switch QxT
<pre> \$B1\$2: xor eax, eax pxor xmm0, xmm0 \$B1\$3: movd xmm1, _src[eax] punpcklbw xmm1, xmm0 punpcklwd xmm1, xmm0 movdqa _dst[eax*4], xmm1 add eax, 4 cmp eax, 1024 jb \$B1\$3 </pre>	<pre> \$B1\$2: movdqa xmm0, _2i10f12t\$1DD xor eax, eax \$B1\$3: movd xmm1, _src[eax] pshufb xmm1, xmm0 movdqa _dst[eax*4], xmm1 add eax, 4 cmp eax, 1024 jb \$B1\$3 ... _2i10f12t\$1DD 0ffffff00H,0ffffff01H,0ffffff02H,0ffffff03H </pre>

Example 13-13. Un-aligned Data Operation

<pre> __declspec(align(16)) float src[1024], dst[1024]; for(i = 2; i < 1024-2; i++) dst[i] = src[i-2] - src[i-1] - src[i+2]; </pre>
--

Intel Compiler can use PALIGNR to generate code to avoid penalties associated with unaligned loads.

Example 13-14. Auto-Generated Code to Avoid Unaligned Loads

Compiler Switch Qxw	Compiler Switch QxT
<pre> \$B2\$2 movups xmm0, _src[eax+4] movaps xmm1, _src[eax] movaps xmm4, _src[eax+16] movsd xmm3, _src[eax+20] subps xmm1, xmm0 subps xmm1, _src[eax+16] movss xmm2, _src[eax+28] movhps xmm2, _src[eax+32] movups _dst[eax+8], xmm1 shufps xmm3, xmm2, 132 subps xmm4, xmm3 subps xmm4, _src[eax+32] movlps _dst[eax+24], xmm4 movhps _dst[eax+32], xmm4 add eax, 32 cmp eax, 4064 jb \$B2\$2 </pre>	<pre> \$B2\$2: movaps xmm2, _src[eax+16] movaps xmm0, _src[eax] movdqa xmm3, _src[eax+32] movdqa xmm1, xmm2 palignr xmm3, xmm2, 4 palignr xmm1, xmm0, 4 subps xmm0, xmm1 subps xmm0, _src[eax+16] movups _dst[eax+8], xmm0 subps xmm2, xmm3 subps xmm2, _src[eax+32] movlps _dst[eax+24], xmm2 movhps _dst[eax+32], xmm2 add eax, 32 cmp eax, 4064 jb \$B2\$2 </pre>

A.2 INTEL® VTUNE™ PERFORMANCE ANALYZER

The Intel VTune Performance Analyzer is a powerful software-profiling tool for Microsoft Windows and Linux. The VTune analyzer helps you understand the performance characteristics of your software at all levels: system, application, microarchitecture.

The sections that follow describe the major features of the VTune analyzer and briefly explain how to use them. For more details on these features, run the VTune analyzer and see the online documentation.

All features are available for Microsoft Windows. On Linux, sampling and call graph are available.

A.2.1 Sampling

Sampling allows you to profile all active software on your system, including operating system, device driver, and application software. It works by occasionally interrupting the processor and collecting the instruction address, process ID, and thread ID. After the sampling activity completes, the VTune analyzer displays the data by process, thread, software module, function, or line of source. There are two methods for generating samples: Time-based sampling and Event-based sampling.

A.2.1.1 Time-based Sampling

Time-based sampling (TBS) uses an operating system's (OS) timer to periodically interrupt the processor to collect samples. The sampling interval is user definable. TBS is useful for identifying the software on your computer that is taking the most CPU time. This feature is only available in the Windows version of the VTune Analyzer

A.2.1.2 Event-based Sampling

Event-based sampling (EBS) can be used to provide detailed information on the behavior of the microprocessor as it executes software. Some of the events that can be used to trigger sampling include clockticks, cache misses, and branch mispredictions. The VTune analyzer indicates where micro architectural events, specific to the Intel Core microarchitecture, Pentium 4, Pentium M and Intel Xeon processors, occur the most often. On processors based on Intel Core microarchitecture, it is possible to collect up to 5 events (three events using fixed-function counters, two events using general-purpose counters) at a time from a list of over 400 events (see Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*). On Pentium M processors, the VTune analyzer can collect two different events at a time. The number of the events that the VTune analyzer can collect at once on the Pentium 4 and Intel Xeon processor depends on the events selected.

Event-based samples are collected periodically after a specific number of processor events have occurred while the program is running. The program is interrupted, allowing the interrupt handling driver to collect the Instruction Pointer (IP), load module, thread and process ID's. The instruction pointer is then used to derive the function and source line number from the debug information created at compile time. The Data can be displayed as horizontal bar charts or in more detail as spread sheets that can be exported for further manipulation and easy dissemination.

A.2.1.3 Workload Characterization

Using event-based sampling and processor-specific events can provide useful insights into the nature of the interaction between a workload and the microarchitecture. A few metrics useful for workload characterization are discussed in Appendix B. The event lists available on various Intel processors can be found in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

A.2.2 Call Graph

Call graph helps you understand the relationships between the functions in your application by providing timing and caller/callee (functions called) information. Call graph works by instrumenting the functions in your application. Instrumentation is the process of modifying a function so that performance data can be captured when the function is executed. Instrumentation does not change the functionality of the

program. However, it can reduce performance. The VTune analyzer can detect modules as they are loaded by the operating system, and instrument them at run-time. Call graph can be used to profile Win32*, Java*, and Microsoft.NET* applications. Call graph only works for application (ring 3) software.

Call graph profiling provides the following information on the functions called by your application: total time, self-time, total wait time, wait time, callers, callees, and the number of calls. This data is displayed using three different views: function summary, call graph, and call list. These views are all synchronized.

The Function Summary View can be used to focus the data displayed in the call graph and call list views. This view displays all the information about the functions called by your application in a sortable table format. However, it does not provide callee and caller information. It just provides timing information and number of times a function is called.

The Call Graph View depicts the caller/callee relationships. Each thread in the application is the root of a call tree. Each node (box) in the call tree represents a function. Each edge (line with an arrow) connecting two nodes represents the call from the parent to the child function. If the mouse pointer is hovered over a node, a tool tip will pop up displaying the function's timing information.

The Call List View is useful for analyzing programs with large, complex call trees. This view displays only the caller and callee information for the single function that you select in the Function Summary View. The data is displayed in a table format.

A.2.3 Counter Monitor

Counter monitor helps you identify system level performance bottlenecks. It periodically polls software and hardware performance counters. The performance counter data can help you understand how your application is impacting the performance of the computer's various subsystems. Counter monitor data can be displayed in real-time and logged to a file. The VTune analyzer can also correlate performance counter data with sampling data. This feature is only available in the Windows version of the VTune Analyzer

A.3 INTEL® PERFORMANCE LIBRARIES

The Intel Performance Library family contains a variety of specialized libraries which has been optimized for performance on Intel processors. These optimizations take advantage of appropriate architectural features, including MMX technology, Streaming SIMD Extensions (SSE), Streaming SIMD Extensions 2 (SSE2) and Streaming SIMD Extensions 3 (SSE3). The library set includes the Intel Math Kernel Library (MKL) and the Intel Integrated Performance Primitives (IPP).

- The Intel Math Kernel Library for Linux and Windows: MKL is composed of highly optimized mathematical functions for engineering, scientific and financial applications requiring high performance on Intel platforms. The functional areas of the

library include linear algebra consisting of LAPACK and BLAS, Discrete Fourier Transforms (DFT), vector transcendental functions (vector math library/VML) and vector statistical functions (VSL). Intel MKL is optimized for the latest features and capabilities of the Intel[®] Itanium[®], Intel[®] Xeon[®], Intel[®] Pentium[®] 4, and Intel[®] Core2 Duo processor-based systems. Special attention has been paid to optimizing multi-threaded performance for the new Quad-Core Intel[®] Xeon[®] processor 5300 series.

- Intel[®] Integrated Performance Primitives for Linux* and Windows*: IPP is a cross-platform software library which provides a range of library functions for video decode/encode, audio decode/encode, image color conversion, computer vision, data compression, string processing, signal processing, image processing, JPEG decode/encode, speech recognition, speech decode/encode, cryptography plus math support routines for such processing capabilities.

Intel IPP is optimized for the broad range of Intel microprocessors: Intel Core 2 processor family, Dual-core Intel Xeon processors, Intel Pentium 4 processor, Pentium M processor, Intel Xeon processors, the Intel Itanium architecture, Intel[®] SA-1110 and Intel[®] PCA application processors based on the Intel XScale[®] microarchitecture. With a single API across the range of platforms, the users can have platform compatibility and reduced cost of development.

A.3.1 Benefits Summary

The overall benefits the libraries provide to the application developers are as follows:

- **Time-to-Market** — Low-level building block functions that support rapid application development, improving time to market.
- **Performance** — Highly-optimized routines with a C interface that give Assembly-level performance in a C/C++ development environment (MKL also supports a Fortran interface).
- **Platform tuned** — Processor-specific optimizations that yield the best performance for each Intel processor.
- **Compatibility** — Processor-specific optimizations with a single application programming interface (API) to reduce development costs while providing optimum performance.
- **Threaded application support** — Applications can be threaded with the assurance that the MKL and IPP functions are safe for use in a threaded environment.

A.3.2 Optimizations with the Intel[®] Performance Libraries

The Intel Performance Libraries implement a number of optimizations that are discussed throughout this manual. Examples include architecture-specific tuning such as loop unrolling, instruction pairing and scheduling; and memory management with explicit and implicit data prefetching and cache tuning.

The Libraries take advantage of the parallelism in the SIMD instructions using MMX technology, Streaming SIMD Extensions (SSE), Streaming SIMD Extensions 2 (SSE2), and Streaming SIMD Extensions 3 (SSE3). These techniques improve the performance of computationally intensive algorithms and deliver hand coded performance in a high level language development environment.

For performance sensitive applications, the Intel Performance Libraries free the application developer from the time consuming task of assembly-level programming for a multitude of frequently used functions. The time required for prototyping and implementing new application features is substantially reduced and most important, the time to market is substantially improved. Finally, applications developed with the Intel Performance Libraries benefit from new architectural features of future generations of Intel processors simply by relinking the application with upgraded versions of the libraries.

A.4 INTEL® THREADING ANALYSIS TOOLS

The Intel® Threading Analysis Tools consist of the Intel Thread Checker 3.0, the Thread Profiler 3.0, and the Intel Threading Building Blocks 1.0 (1). The Intel Thread Checker and Thread Profiler supports Windows and Linux. The Intel Threading Building Blocks 1.0 supports Windows, Linux, and Mac OS.

A.4.1 Intel® Thread Checker 3.0

The Intel Thread Checker locates programming errors (for example: data races, stalls and deadlocks) in threaded applications. Use the Intel Thread Checker to find threading errors and reduce the amount of time you spend debugging your threaded application.

The Intel Thread Checker product is an Intel VTune Performance Analyzer plug-in data collector that executes your program and automatically locates threading errors. As your program runs, the Intel Thread Checker monitors memory accesses and other events and automatically detects situations which could cause unpredictable threading-related results. The Intel Thread Checker detects thread deadlocks, stalls, data race conditions and more.

A.4.2 Intel® Thread Profiler 3.0

The thread profiler is a plug-in data collector for the Intel VTune Performance Analyzer. Use it to analyze threading performance and identify parallel performance problems. The thread profiler graphically illustrates what each thread is doing at various levels of detail using a hierarchical summary. It can identify inactive threads,

1 For additional threading resources, visit <http://www3.intel.com/cd/software/products/asmo-na/eng/index.htm>

critical paths and imbalances in thread execution, etc. Mountains of data are collapsed into relevant summaries, sorted to identify parallel regions or loops that require attention. Its intuitive, color-coded displays make it easy to assess your application's performance.

Figure A-1 shows the execution timeline of a multi-threaded application when run in (a) a single-threaded environment, (b) a multi-threaded environment capable of executing two threads simultaneously, (c) a multi-threaded environment capable of executing four threads simultaneously. In Figure A-1, the color-coded timeline of three hardware configurations are super-imposed together to compare processor scaling performance and illustrate the imbalance of thread execution.

Load imbalance problem is visually identified in the two-way platform by noting that there is a significant portion of the timeline, during which one logical processor had no task to execute. In the four-way platform, one can easily identify those portions of the timeline of three logical processors, each having no task to execute.

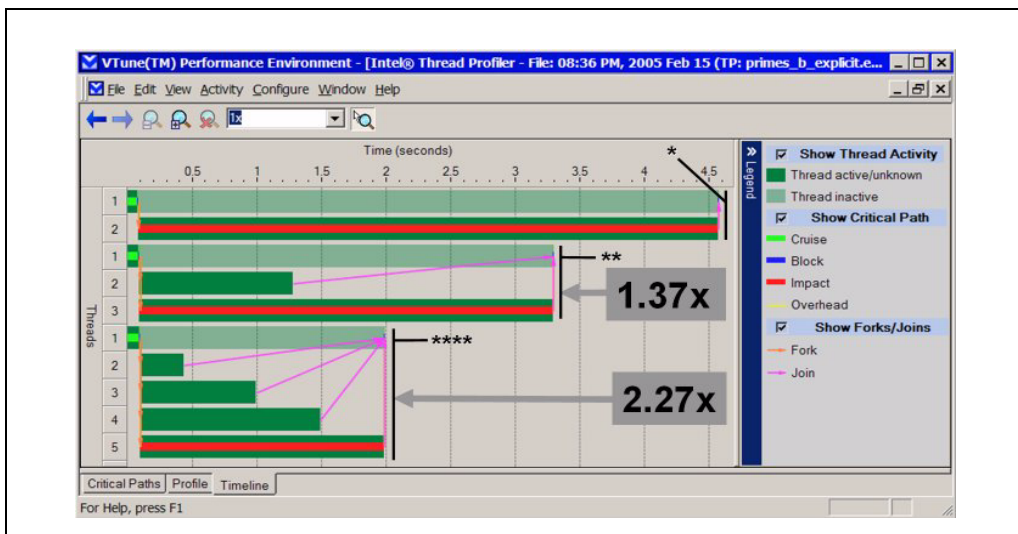


Figure A-1. Intel Thread Profiler Showing Critical Paths of Threaded Execution Timelines

A.4.3 Intel® Threading Building Blocks 1.0

The Intel Threading Building Blocks is a C++ template-based runtime library that simplifies threading for scalable, multi-core performance. It can help avoid re-writing, re-testing, re-tuning common parallel data structures and algorithms.

A.5 INTEL® CLUSTER TOOLS

The Intel® Cluster Toolkit 3.1 consists of the Intel® Trace Analyzer and Collector, Intel® Math Kernel Library (Intel® MKL), Intel® MPI Library, and Intel® MPI Benchmarks in a single package. The Intel® Cluster Toolkit 3.1 helps you develop, analyze and optimize performance of parallel applications for clusters using IA-32, IA-64, and Intel® 64 architectures. The Intel® Cluster Toolkit 3.1 supports Windows, Linux and SGI ProPack.

A.5.1 Intel® MPI Library 3.1

The Intel® MPI Library 3.1 is a multi-fabric message passing library that implements the Message Passing Interface, v2 (MPI-2) specification. It provides a standard library across Intel® platforms. The Intel® MPI Library supports multiple hardware fabrics including InfiniBand, Myrinet*, and Quadrics. Intel® MPI Library covers all your configurations by providing an accelerated universal, multi-fabric layer for fast interconnects via the Direct Access Programming Library (DAPL) methodology. Develop MPI code independent of the fabric, knowing it will run efficiently on whatever fabric is chosen by the user at runtime.

Intel MPI Library dynamically establishes the connection, but only when needed, which reduces the memory footprint. It also automatically chooses the fastest transport available. The fallback to sockets at job startup avoids the chance of execution failure even if the interconnect selection fails. This is especially helpful for batch computing. Any products developed with Intel MPI Library are assured run time compatibility since your users can download Intel's free runtime environment kit. Application performance can also be increased via the large message bandwidth advantage from the optional use of DAPL inside a multi-core or SMP node.

A.5.2 Intel® Trace Analyzer and Collector 7.1

The Intel® Trace Analyzer and Collector 7.1 helps to provide information critical to understanding and optimizing application performance on clusters by quickly finding performance bottlenecks in MPI communication. Version 7.1 includes trace file comparison, counter data displays, and an MPI correctness checking library which can detect deadlocks, data corruption, or errors with MPI parameters, data types, buffers, communicators, point-to-point messages and collective operations.

A.5.3 Intel® MPI Benchmarks 3.1

The Intel MPI Benchmarks will help enable an easy performance comparison of MPI functions and patterns, the benchmark features improvements in usability, application performance, and interoperability.

A.5.4 Benefits Summary

The overall benefits the improved MPI Benchmarks provide are as follows:

A.5.4.1 Multiple usability improvements

- New benchmarks (Gather(v), Scatter(v))

A.5.4.2 Improved application performance

- New Command line flags to control cache reuse and to limit memory usage
- Options for cold cache operation mode, maximum buffer size setting and dynamic iteration count determination
- Run time improvements for collectives like Alltoall(v) on large clusters

A.5.4.3 Extended interoperability

- Support for Windows Compute Cluster Server

A.6 INTEL® XML PRODUCTS

Intel® XML Software products deliver outstanding performance for XML processing including: XSLT, Parsing, XPath and Schema Validation. The XML Software suites offer an enterprise solution for both C/C++ and Java environments running in Linux and Windows operating systems.

A.6.1 Intel® XML Software Suite 1.0

The Intel® XML Software Suite is a comprehensive suite of high-performance C++ and Java* software-based runtime libraries for Linux* and Windows* operating systems. Intel® XML Software Suite is standards compliant, to allow for easy integration into existing XML environments and is optimized to support complex and large-size XML document processing. The key functional components of the software suite are: XML parsing, XML schema validation, XML transformation, and XML XPath navigation.

A.6.1.1 Intel® XSLT Accelerator

XSLT (eXtensible Stylesheet Language Transformation) is an XML-based language used to transform XML documents into other XML or human readable documents. Intel® XSLT Accelerator facilitates efficient XML transformations in a variety of formats and can be applied to a full range of XML documents such as a tree (the DOM tree model) or a series of events (the SAX model). Intel® XSLT Accelerator supports

the following groups of XSLT extension functions: Common operations, Math computations, String manipulations, Sets handling, and Date-and-Time functions. User Defined Java extension functions are supported allowing developers to access Java class functions (static or non-static methods) from an XSLT stylesheet to augment native XSLT transformations.

A.6.1.2 Intel® XPath Accelerator

XPath is a language that enables the navigation and data manipulation of XML documents. Intel® XPath Accelerator evaluates an XML Path (XPath) expression over an XML document DOM tree or a derived instance of Source (StreamSource, DOMSource, SAXSource or XMLDocSource) and returns a node, node set, string, number or Boolean value. Intel® XPath Accelerator supports and resolves user-defined namespace context, variables and functions. Optionally, XPath expressions can be compiled to further enhance XML processing performance.

A.6.1.3 Intel® XML Schema Accelerator

XML schema validation compares an XML document against a document that contains a set of rules and constraints specific to the XML application environment adherent to W3C specifications. Validation ensures that an XML document meets application and environment requirements for processing as described by the schema document. Intel® XML Schema Accelerator quickly and efficiently validates XML documents in Stream, SAX, or DOM mode against an XML Schema document.

A.6.1.4 Intel® XML Parsing Accelerator

The XML parser reads an XML file and makes the data in the file available for manipulation and processing to applications and programming languages. The parser is also responsible for testing if a document is well-formed. Intel® XML Parsing Accelerator parses data by following specific models: Simple API for XML (SAX) model as a sequence of events; Document Object Model (DOM) as a tree node structure; and an internal storage data-stream model for effective XML processing between Intel XML Software Suite components. Intel® XML Parsing Accelerator can enable document validation Intel® XML Schema Accelerator before passing data to the application.

A.6.2 Intel® SOA Security Toolkit 1.0 Beta for Axis2

This toolkit provides XML Digital Signature and XML Encryption following the WS-Security standard. Low cost of ownership and easy integration with consistent behavior are facilitated via a simple integrated Axis2* interface and Apache Rampart* configuration files. Key Features include the following:

A.6.2.1 High Performance

The Intel® SOA Security Toolkit achieves high performance for XML security processing. The toolkit's efficient design provides more than three times the performance compared to competitive Open Source solutions, enabling fast throughput for business processes.

A.6.2.2 Standards Compliant

A standards compliance design allows for functional interoperability with existing code and XML based applications. Intel® SOA Security Toolkit implements the following standards:

- WS-Security 1.1
- SOAP v1.1, v1.2

A.6.2.3 Easy Integration

A simple interface allows drop-in compatibility and functional interoperability for the following environments:

- Apache Rampart*
- Axis2*

A.6.3 Intel® XSLT Accelerator 1.1 for Java* Environments on Linux* and Windows* Operating Systems

Intel® XSLT Accelerator is a standards compliant software-based runtime library delivering high performance eXtensible Stylesheet Language Transformations (XSLT) processing. Intel® XSLT Accelerator is optimized to support complex and large-size XML document transformations. Main features include:

A.6.3.1 High Performance Transformations

Fast transformations enable fast throughput for business processes.

- 2X over Apache* Xalan* XSLTC* processor
- 4X over Apache Xalan-J* processor

A.6.3.2 Large XML File Transformations

Large file support facilitates application scalability, data growth, and application reliability.

- Process large XML documents
- Sustained workload support

A.6.3.3 Standards Compliant

The standards compliance design allows for functional interoperability with existing code and applications. Intel® XSLT Accelerator complies with the following standards:

- W3C XML 1.0
- W3C XSLT 1.0
- JAXP 1.3 (TrAX API)
- SAX
- DOM

A.6.3.4 Thread-Safe

Intel® XSLT Accelerator is thread-safe supporting multi-threaded applications and designed for optimal performance on Intel® Core microarchitecture.

A.7 INTEL® SOFTWARE COLLEGE

You can find information on classroom training offered by the Intel Software College at <http://developer.intel.com/software/college>. Find general information for developers at <http://softwarecommunity.intel.com/isn/home/>.

APPENDIX B

USING PERFORMANCE MONITORING EVENTS

Performance monitoring events provide facilities to characterize the interaction between programmed sequences of instructions and microarchitectural sub-systems. Performance monitoring events are described in Chapter 18 and Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

The first part of this chapter provides information on how to use performance events specific to Intel Xeon processor 5500 series. Section B.3 describe performance tuning techniques that apply to processors based on Intel microarchitecture code name Sandy Bridge Section B.5 discusses similar topics for performance events available on Intel Core Solo and Intel Core Duo processors.

B.1 INTEL® XEON® PROCESSOR 5500 SERIES

Intel Xeon processor 5500 series are based on the same microarchitecture as Intel Core i7 processors, see Section 2.3, "Intel® Microarchitecture code name Nehalem". In addition, Intel Xeon processor 5500 series support non-uniform memory access (NUMA) in platforms that have two physical processors, see Figure B-1. Figure B-1 illustrates 4 processor cores and an uncore sub-system in each physical processor. The uncore sub-system consists of L3, an integrated memory controller (IMC), and Intel QuickPath Interconnect (QPI) interfaces. The memory sub-system consists of three channels of DDR3 memory locally connected to each IMC. Access to physical memory connected to a non-local IMC is often described as a remote memory access.

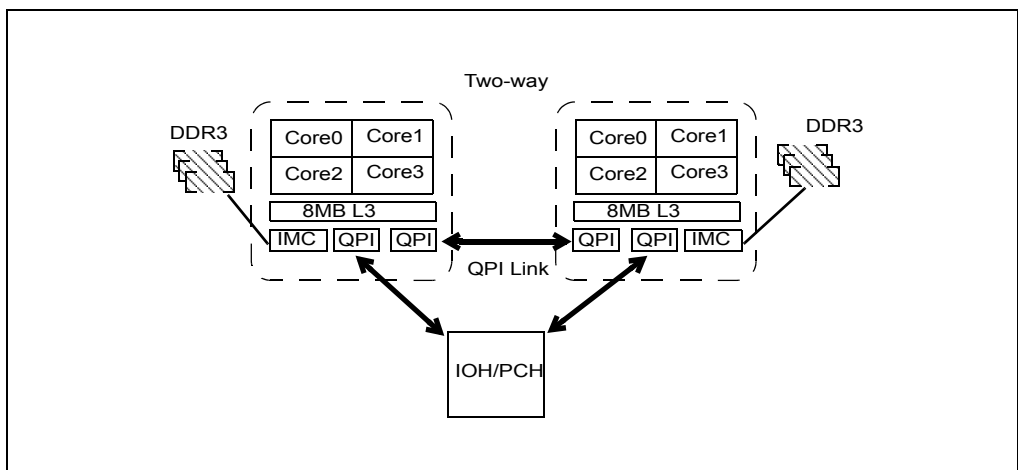


Figure B-1. System Topology Supported by Intel® Xeon® Processor 5500 Series

The performance monitoring events on Intel Xeon processor 5500 series can be used to analyze the interaction between software (code and data) and microarchitectural units hierarchically:

- **Per-core PMU:** Each processor core provides 4 programmable counters and 3 fixed counters. The programmable per-core counters can be configured to investigate front-end/micro-op flow issues, stalls inside a processor core. Additionally, a subset of per-core PMU events support precise event-based sampling (PEBS). Load latency measurement facility is new in Intel Core i7 processor and Intel Xeon processor 5500.
- **Uncore PMU:** The uncore PMU provides 8 programmable counters and 1 fixed counter. The programmable per-core counters can be configured to characterize L3 and Intel QPI operations, local and remote data memory accesses.

The number and variety of performance counters and the breadth of programmable performance events available in Intel Xeon processor 5500 offer software tuning engineers the ability to analyze performance issues and achieve higher performance. Using performance events to analyze performance issues can be grouped into the following subjects:

- Cycle Accounting and Uop Flow
- Stall Decomposition and Core Memory Access Events (non-PEBS)
- Precise Memory Access Events (PEBS)
- Precise Branch Events (PEBS, LBR)
- Core Memory Access Events (non-PEBS)
- Other Core Events (non-PEBS)
- Front End Issues
- Uncore Events

B.2 PERFORMANCE ANALYSIS TECHNIQUES FOR INTEL® XEON® PROCESSOR 5500 SERIES

The techniques covered in this chapter focuses on identifying opportunity to remove/reduce performance bottlenecks that are measurable at runtime. Compile-time and source-code level techniques are covered in other chapters in this document. Individual sub-sections describe specific techniques to identify tuning opportunity by examining various metrics that can be measured or derived directly from performance monitoring events.

B.2.1 Cycle Accounting and Uop Flow Analysis

The objectives, performance metrics and component events of the basic cycle accounting technique is summarized in Table B-1.

Table B-1. Cycle Accounting and Micro-ops Flow Recipe

Summary	
Objective	Identify code/basic block that had significant stalls
Method	Binary decomposition of cycles into “productive“ and “unproductive“ parts
PMU-Pipeline Focus	Micro-ops issued to execute
Event code/Umask	Event code B1H, Umask= 3FH for micro-op execution; Event code 3CH, Umask= 1, CMask=2 for counting total cycles
EvtSelc	Use CMask, Invert, Edge fields to count cycles and separate stalled vs. active cycles
Basic Equation	“Total Cycles“ = UOPS_EXECUTED.CORE_STALLS_CYCLES + UOPS_EXECUTED.CORE_ACTIVE_CYCLES
Metric	UOPS_EXECUTED.CORE_STALLS_CYCLES / UOPS_EXECUTED.CORE_STALLS_COUNT
Drill-down scope	Counting: Workload; Sampling: basic block
Variations	Port 0,1, 5 cycle counting for computational micro-ops execution.

Cycle accounting of executed micro-ops is an effective technique to identify stalled cycles for performance tuning. Within the microarchitecture pipeline, the meaning of micro-ops being “issued”, “dispatched”, “executed”, “retired” has precise meaning. This is illustrated in Figure B-2.

Cycles are divided into those where micro-ops are dispatched to the execution units and those where no micro-ops are dispatched, which are thought of as execution stalls.

“Total cycles” of execution for the code under test can be directly measured with CPU_CLK_UNHALTED.THREAD (event code 3CH, Umask= 1) and setting CMask = 2 and INV=1 in IA32_PERFEVTSELn.

The signals used to count the memory access uops executed (ports 2, 3 and 4) are the only core events which cannot be counted per-logical processor. Thus, Event code B1H with Umask=3FH only counts on a per-core basis, and the total execution stall

USING PERFORMANCE MONITORING EVENTS

cycles can only be evaluated on a per core basis. If HT is disabled, this presents no difficulty to conduct per-thread analysis of micro-op flow cycle accounting.

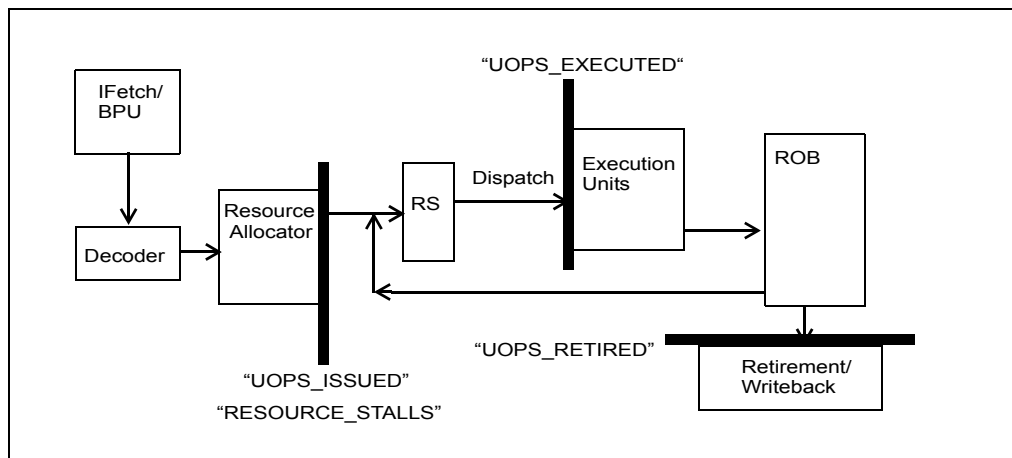


Figure B-2. PMU Specific Event Logic Within the Pipeline

The PMU signals to count `uops_executed` in port 0, 1, 5 can count on a per-thread basis even when HT is active. This provides an alternate cycle accounting technique when the workload under test interacts with HT.

The alternate metric is built from `UOPS_EXECUTED.PORT015_STALL_CYCLES`, using appropriate `Cmask`, `Inv`, and `Edge` settings. Details of performance events are shown in Table B-2.

Table B-2. CMask/Inv/Edge/Thread Granularity of Events for Micro-op Flow

Event Name	Umask	Event Code	Cmask	Inv	Edge	All Thread
<code>CPU_CLK_UNHALTED.TOTAL_CYCLES</code>	0H	3CH	2	1	0	0
<code>UOPS_EXECUTED.CORE_STALLS_CYCLES</code>	3FH	B1H	1	1	0	1
<code>UOPS_EXECUTED.CORE_STALLS_COUNT</code>	3FH	B1H	1	1	!	1
<code>UOPS_EXECUTED.CORE_ACTIVE_CYCLES</code>	3FH	B1H	1	0	0	1

Table B-2. CMask/Inv/Edge/Thread Granularity of Events for Micro-op Flow

Event Name	Umask	Event Code	Cmask	Inv	Edge	All Thread
UOPS_EXECUTED.PORT015_S TALLS_CYCLES	40H	B1H	1	1	0	0
UOPS_RETIRED.STALL_CYCLE S	1H	C2H	1	1	0	0
UOPS_RETIRED.ACTIVE_CYCL ES	1H	C2H	1	0	0	0

B.2.1.1 Cycle Drill Down and Branch Mispredictions

While executed micro-ops are considered productive from the perspective of execution units being subscribed, not all such micro-ops contribute to forward progress of the program. Branch mispredictions can introduce execution inefficiencies in OOO processor that are typically decomposed into three components.

- Wasted work associated with executing the uops of the incorrectly predicted path
- Cycles lost when the pipeline is flushed of the incorrect uops
- Cycles lost while waiting for the correct uops to arrive at the execution units

In processors based on Intel microarchitecture code name Nehalem, there are no execution stalls associated with clearing the pipeline of mispredicted uops (component 2). These uops are simply removed from the pipeline without stalling executions or dispatch. This typically lowers the penalty for mispredicted branches. Further, the penalty associated with instruction starvation (component 3) can be measured.

The wasted work within executed uops are those uops that will never be retired. This is part of the cost associated with mispredicted branches. It can be found through monitoring the flow of uops through the pipeline. The uop flow can be measured at 3 points in Figure B-2, going into the RS with the event UOPS_ISSUED, going into the execution units with UOPS_EXECUTED and at retirement with UOPS_RETIRED. The differences of between the upstream measurements and at retirement measure the wasted work associated with these mispredicted uops.

As UOPS_EXECUTED must be measured per core, rather than per thread, the wasted work per core is evaluated as

$$\text{Wasted Work} = \text{UOPS_EXECUTED.PORT234_CORE} + \text{UOPS_EXECUTED.PORT015_All_Thread} - \text{UOPS_RETIRED.ANY_ALL_THREAD}$$

The ratio above can be converted to cycles by dividing the average issue rate of uops. The events above were designed to be used in this manner without corrections for micro fusion or macro fusion.

A "per thread" measurement can be made from the difference between the uops issued and uops retired as the latter two of the above events can be counted per thread. It over counts slightly, by the mispredicted uops that are eliminated in the RS before they can waste cycles being executed, but this is usually a small correction.

$$\text{Wasted Work/thread} = (\text{UOPS_ISSUED.ANY} + \text{UOPS_ISSUED.FUSED}) - \text{UOPS_RETIRED.ANY}$$

Table B-3. Cycle Accounting of Wasted Work Due to Misprediction

Summary	
Objective	Evaluate uops that executed but not retired due to misprediction
Method	Examine uop flow differences between execution and retirement
PMU-Pipeline Focus	Micro-ops execute and retirement
Event code/Umask	Event code B1H, Umask= 3FH for micro-op execution; Event code C2H, Umask= 1, AllThread=1 for per-core counting
EvtSelc	Zero CMask, Invert, Edge fields to count uops
Basic Equation	"Wasted work" = UOPS_EXECUTED.PORT234_CORE + UOPS_EXECUTED.PORT015_ALL_THREAD - UOPS_RETIRED.ANY_ALL_THREAD
Drill-down scope	Counting: Branch misprediction cost
Variations	Divide by average uop issue rate for cycle accounting. Set AllThread=0 to estimate per-thread cost.

The third component of the misprediction penalty, instruction starvation, occurs when the instructions associated with the correct path are far away from the core and execution is stalled due to lack of uops in the RAT. Because the two primary cause of uops not being issued are either front-end starvation or resource not available in the back end. So we can explicitly measured at the output of the resource allocation as follows:

- Count the total number of cycles where no uops were issued to the OOO engine.
- Count the cycles where resources (RS, ROB entries, load buffer, store buffer, etc.) are not available for allocation.

If HT is not active, instruction starvation is simply the difference:

$$\text{Instruction Starvation} = \text{UOPS_ISSUED.STALL_CYCLES} - \text{RESOURCE_STALLS.ANY}$$

When HT is enabled, the uop delivery to the RS alternates between the two threads. In an ideal case the above condition would then over count, as 50% of the issuing stall cycles may be delivering uops for the other thread. We can modify the expression by subtracting the cycles that the other thread is having uops issued.

$$\text{Instruction Starvation (per thread)} = \text{UOPS_ISSUED.STALL_CYCLES} - \text{RESOURCE_STALLS.ANY} - \text{UOPS_ISSUED.ACTIVE_CYCLES_OTHER_THREAD.}$$

The per-thread expression above will over count somewhat because the resource_stall condition could exist on "this" thread while the other thread in the same core was issuing uops. An alternative might be

$$\text{CPU_CLK_UNHALTED.THREAD} - \text{UOPS_ISSUED.CORE_CYCLES_ACTIVE-RESOURCE_STALLS.ANY}$$

The above technique is summarized in Table B-4.

Table B-4. Cycle Accounting of Instruction Starvation

Summary	
Objective	Evaluate cycles that uops issuing is starved after misprediction
Method	Examine cycle differences between uops issuing and resource allocation
PMU-Pipeline Focus	Micro-ops issue and resource allocation
Event code/Umask	Event code 0EH, Umask= 1, for uops issued. Event code A2H, Umask=1, for Resource allocation stall cycles
EvtSelc	Set CMask=1, Inv=1, fields to count uops issue stall cycles. Set CMask=1, Inv=0, fields to count uops issue active cycles. Use AllThread = 0 and AllThread=1 on two counter to evaluate contribution from the other thread for UOPS_ISSUED.ACTIVE_CYCLES_OTHER_THREAD
Basic Equation	“Instruction Starvation“ (HT off) = UOPS_ISSUED.STALL_CYCLES - RESOURCE_STALLS.ANY;
Drill-down scope	Counting: Branch misprediction cost
Variations	Evaluate per-thread contribution with Instruction Starvation = UOPS_ISSUED.STALL_CYCLES - RESOURCE_STALLS.ANY - UOPS_ISSUED.ACTIVE_CYCLES_OTHER_THREAD

Details of performance events are shown in Table B-5.

Table B-5. CMask/Inv/Edge/Thread Granularity of Events for Micro-op Flow

Event Name	Umask	Event Code	Cmask	Inv	Edge	All Thread
UOPS_EXECUTED.PORT234_CORE	80H	B1H	0	0	0	1
UOPS_EXECUTED.PORT015_ALL_THREAD	40H	B1H	0	0	0	1
UOPS_RETIRED.ANY_ALL_THREAD_READ	1H	C2H	0	0	0	1
RESOURCE_STALLS.ANY	1H	A2H	0	0	0	0
UOPS_ISSUED.ANY	1H	0EH	0	0	0	0
UOPS_ISSUED.STALL_CYCLES	1H	0EH	1	1	0	0
UOPS_ISSUED.ACTIVE_CYCLES	1H	0EH	1	0	0	0
UOPS_ISSUED.CORE_CYCLES_ACTIVE	1H	0EH	1	0	0	1

B.2.1.2 Basic Block Drill Down

The event INST_RETIRED.ANY (instructions retired) is commonly used to evaluate a cycles/instruction ratio (CPI). Another important usage is to determine the performance-critical basic blocks by evaluating basic block execution counts.

In a sampling tool (such as VTune Analyzer), the samples tend to cluster around certain IP values. This is true when using INST_RETIRED.ANY or cycle counting events. Disassembly listing based on the hot samples may associate some instructions with high sample counts and adjacent instructions with no samples.

Because all instructions within a basic block are retired exactly the same number of times by the very definition of a basic block. Drilling down the hot basic blocks will be more accurate by averaging the sample counts over the instructions of the basic block.

$$\text{Basic Block Execution Count} = \text{Sum (Sample counts of instructions within basic block)} * \text{Sample_after_value} / (\text{number of instructions in basic block})$$

Inspection of disassembly listing to identify basic blocks associated with loop structure being a hot loop or not can be done systematically by adapting the technique

above to evaluate the trip count of each loop construct. For a simple loop with no conditional branches, the trip count ends up being the ratio of the basic block execution count of the loop block to the basic block execution count of the block immediately before and/or after the loop block. Judicious use of averaging over multiple blocks can be used to improve the accuracy.

This will allow the user to identify loops with high trip counts to focus on tuning efforts. This technique can be implemented using fixed counters.

Chains of dependent long-latency instructions (fmul, fadd, imul, etc) can result in the dispatch being stalled while the outputs of the long latency instructions become available. In general there are no events that assist in counting such stalls with the exception of instructions using the divide/sqrt execution unit. In such cases, the event ARITH can be used to count both the occurrences of these instructions and the duration in cycles that they kept their execution units occupied. The event ARITH.CYCLES_DIV_BUSY counts the cycles that either the divide/sqrt execution unit was occupied.

B.2.2 Stall Cycle Decomposition and Core Memory Accesses

The decomposition of the stall cycles is accomplished through a standard approximation. It is assumed that the penalties occur sequentially for each performance impacting event. Consequently, the total loss of cycles available for useful work is then the number of events, N_i , times the average penalty for each type of event, P_i

$$\text{Counted_Stall_Cycles} = \text{Sum} (N_i * P_i)$$

This only accounts for the performance impacting events that are or can be counted with a PMU event. Ultimately there will be several sources of stalls that cannot be counted, however their total contribution can be estimated:

$$\text{Unaccounted stall cycles} = \text{Stall_Cycles} - \text{Counted_Stall_Cycles} = \\ \text{UOPS_EXECUTED.CORE_STALLS_CYCLES} - \text{Sum} (N_i * P_i)_{\text{both_threads}}$$

The unaccounted component can become negative as the sequential penalty model is overly simple and usually over counts the contributions of the individual microarchitectural issues.

As noted in Section B.2.1.1, UOPS_EXECUTED.CORE_STALL_CYCLES counts on a per core basis rather than on a per thread basis, the over counting can become severe. In such cases it may be preferable to use the port 0,1,5 uop stalls, as that can be done on a per thread basis:

$$\text{Unaccounted stall cycles (per thread) } = \\ \text{UOPS_EXECUTED.PORT015_THREADED_STALLS_CYCLES} - \text{Sum} (N_i * P_i)$$

This unaccounted component is meant to represent the components that were either not counted due to lack of performance events or simply neglected during the data collection.

One can also choose to use the "retirement" point as the basis for stalls. The PEBS event, UOPS_RETIRED.STALL_CYCLES, has the advantage of being evaluated on a

per thread basis and being having the HW capture the IP associated with the retiring uop. This means that the IP distribution will not be effected by STI/CLI deferral of interrupts in critical sections of OS kernels, thus producing a more accurate profile of OS activity.

B.2.2.1 Measuring Costs of Microarchitectural Conditions

Decomposition of stalled cycles in this manner should start by first focusing on conditions that carry large performance penalty, for example, events with penalties of greater than 10 cycles. Short penalty events ($P < 5$ cycles) can frequently be hidden by the combined actions of the OOO execution and the compiler. The OOO engine manages both types of situations in the instruction stream and strive to keep the execution units busy during stalls of either type due to instruction dependencies. Usually, the large penalty operations are dominated by memory access and the very long latency instructions for divide and sqrt.

The largest penalty events are associated with load operations that require a cache-line which is not in L1 or L2 of the cache hierarchy. Not only must we count how many occur, but we need to know what penalty to assign.

The standard approach to measuring latency is to measure the average number of cycles a request is in a queue:

$$\text{Latency} = \text{Sum}(\text{CYCLES_Queue_entries_outstanding}) / \text{Queue_inserts}$$

where "queue_inserts" refers to the total number of entries that caused the outstanding cycles in that queue. However, the penalty associated with each queue insert (i.e. cachemiss), is the latency divided by the average queue occupancy. This correction is needed to avoid over counting associated with overlapping penalties.

$$\text{Avg_Queue_Depth} = \text{Sum}(\text{CYCLES_Queue_entries_outstanding}) / \text{Cycles_Queue_not_empty}$$

The the penalty (cost) of each occurrence is

$$\text{Penalty} = \text{Latency} / \text{Avg_Queue_Depth} = \text{Cycles_Queue_not_empty} / \text{Queue_inserts}$$

An alternative way of thinking about this is to realize that the sum of all the penalties, for an event that occupies a queue for its duration, cannot exceed the time that the queue is not empty

$$\text{Cycles_Queue_not_empty} = \text{Events} * \langle \text{Penalty} \rangle$$

The standard techniques described above are simple conceptually. In practice, the large amount of memory references in the workload and wide range of varying state/location-specific latencies made standard sampling techniques less practical. Using precise-event-based sampling (PEBS) is the preferred technique on processors based on Intel microarchitecture code name Nehalem.

The profiling the penalty by sampling (to localize the measurement in IP) is likely to have accuracy difficulties. Since the latencies for L2 misses can vary from 40 to 400 cycles, collecting the number of required samples will tend to be invasive.

The use of the precise latency event, that will be discussed later, provides a more accurate and flexible measurement technique when sampling is used. As each sample records both a load to use latency and a data source, the average latency per data source can be evaluated. Further as the PEBS hardware supports buffering the events without generating a PMI until the buffer is full, it is possible to make such an evaluation efficient without perturbing the workload intrusively.

A number of performance events in core PMU can be used to measure the costs of memory accesses that originated in the core and experienced delays due to various conditions, locality, or traffic due to cache coherence requirements. The latency of memory accesses vary, depending on locality of L3, DRAM attached to the local memory controller or remote controller, and cache coherency factors. Some examples of the approximate latency values are shown in Table B-6.

Table B-6. Approximate Latency of L2 Misses of Intel Xeon Processor 5500

Data Source	Latency
L3 hit, Line exclusive	~ 42 cycles
L3 Hit, Line shared	~ 63 cycles
L3 Hit, modified in another core	~ 73 cycles
Remote L3	100 - 150 cycles
Local DRAM	~ 50 ns
Remote DRAM	~ 90 ns

B.2.3 Core PMU Precise Events

The Precise Event Based Sampling (PEBS) mechanism enables the PMU to capture the architectural state and IP at the completion of the instruction that caused the event. This provides two significant benefit for profiling and tuning:

- The location of the eventing condition in the instruction space can be accurate profiled,
- Instruction arguments can be reconstructed in a post processing phase, using captured PEBS records of the register states.

The PEBS capability has been greatly expanded in processors based on Intel microarchitecture code name Nehalem, covering a large number of and more types of precise events.

The mechanism works by using the counter overflow to arm the PEBS data acquisition. Then on the next event, the data is captured and the interrupt is raised.

The captured IP value is sometimes referred to as IP +1, because at the completion of the instruction, the IP value is that of the next instruction.

By their very nature precise events must be “at-retirement” events. For the purposes of this discussion the precise events are divided into Memory Access events, associated with the retirement of loads and stores, and Execution Events, associated with the retirement of all instructions or specific non memory instructions (branches, FP assists, SSE uops).

B.2.3.1 Precise Memory Access Events

There are two important common properties to all precise memory access events:

- The exact instruction can be identified because the hardware captures the IP of the offending instruction. Of course the captured IP is that of the following instruction but one simply moves the samples up one instruction. This works even when the recorded IP points to the first instruction of a basic block because in such a case the offending instruction has to be the last instruction of the previous basic block, as branch instructions never load or store data, instruction arguments can be reconstructed in a post processing phase, using captured PEBS records of the register states.
- The PEBS buffer contains the values of all 16 general registers, R1-R16, where R1 is also called RAX. When coupled with the disassembly the address of the load or store can be reconstructed and used for data access profiling. The Intel® Performance Tuning Utility does exactly this, providing a wide variety of powerful analysis techniques

Precise memory access events mainly focus on loads as those are the events typically responsible for the very long duration execution stalls. They are broken down by the data source, thereby indicating the typical latency and the data locality in the intrinsically NUMA configurations. These precise load events are the only L2, L3 and DRAM access events that only count loads. All others will also include the L1D and/or L2 hardware prefetch requests. Many will also include RFO requests, both due to stores and to the hardware prefetchers.

All four general counters can be programmed to collect data for precise events. The ability to reconstruct the virtual addresses of the load and store instructions allows an analysis of the cacheline and page usage efficiency. Even though cachelines and pages are defined by physical address the lower order bits are identical, so the virtual address can be used.

As the PEBS mechanism captures the values of the register at completion of the instruction, one should be aware that pointer-chasing type of load operation will not be captured because it is not possible to infer the load instruction from the dereferenced address.

The basic PEBS memory access events falls into the following categories:

- `MEM_INST_RETIRED`: This category counts instruction retired which contain a load operation, it is selected by event code 0BH.

- **MEM_LOAD_RETIRED:** This category counts retired load instructions that experienced specific condition selected by the Umask value, the event code is 0CBH.
- **MEM_UNCORE_RETIRED:** This category counts memory instructions retired and received data from the uncore sub-system, it is selected by event code 0FH.
- **MEM_STORE_RETIRED:** This category counts instruction retired which contain a store operation, it is selected by event code 0CH.
- **ITL_MISS_RETIRED:** This counts instruction retired which missed the ITLB, it is selected by event code 0C8H

Umask values and associated name suffixes for the above PEBS memory events are listed under the in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

The precise events listed above allow load driven cache misses to be identified by data source. This does not identify the "home" location of the cachelines with respect to the NUMA configuration. The exceptions to this statement are the events

MEM_UNCORE_RETIRED.LOCAL_DRAM and **MEM_UNCORE_RETIRED.NON_LOCAL_DRAM**. These can be used in conjunction with instrumented malloc invocations to identify the NUMA "home" for the critical contiguous buffers used in an application.

The sum of all the **MEM_LOAD_RETIRED** events will equal the **MEM_INST_RETIRED.LOADS** count.

A count of L1D misses can be achieved with the use of all the **MEM_LOAD_RETIRED** events, except **MEM_LOAD_RETIRED.L1D_HIT**. It is better to use all of the individual **MEM_LOAD_RETIRED** events to do this, rather than the difference of **MEM_INST_RETIRED.LOADS-MEM_LOAD_RETIRED.L1D_HIT** because while the total counts of precise events will be correct, and they will correctly identify instructions that caused the event in question, the distribution of the events may not be correct due to PEBS SHADOWING, discussed later in this section.

$$\text{L1D_MISSES} = \text{MEM_LOAD_RETIRED.HIT_LFB} + \text{MEM_LOAD_RETIRED.L2_HIT} + \text{MEM_LOAD_RETIRED.L3_UNSHARED_HIT} + \text{MEM_LOAD_RETIRED.OTHER_CORE_HIT_HITM} + \text{MEM_LOAD_RETIRED.L3_MISS}$$

The **MEM_LOAD_RETIRED.L3_UNSHARED_HIT** event merits some explanation. The inclusive L3 has a bit pattern to identify which core has a copy of the line. If the only bit set is for the requesting core (unshared hit) then the line can be returned from the L3 with no snooping of the other cores. If multiple bits are set, then the line is in a shared state and the copy in the L3 is current and can also be returned without snooping the other cores.

If the line is read for ownership (RFO) by another core, this will put the copy in the L3 into an exclusive state. If the line is then modified by that core and later evicted, the written back copy in the L3 will be in a modified state and snooping will not be required. **MEM_LOAD_RETIRED.L3_UNSHARED_HIT** counts all of these. The event should really have been called **MEM_LOAD_RETIRED.L3_HIT_NO_SNOOP**.

The event MEM_LOAD_RETIRED.L3_HIT_OTHER_CORE_HIT_HITM could have been named as MEM_LOAD_RETIRED.L3_HIT_SNOOP intuitively for similar reason.

When a modified line is retrieved from another socket it is also written back to memory. This causes remote HITM access to appear as coming from the home dram. The MEM_UNCORE_RETIRED.LOCAL_DRAM and MEM_UNCORE_RETIRED.REMOTE_DRAM events thus also count the L3 misses satisfied by modified lines in the caches of the remote socket.

There is a difference in the behavior of MEM_LOAD_RETIRED.DTLB_MISSES with respect to that on Intel® Core™2 processors. Previously the event only counted the first miss to the page, as do the imprecise events. The event now counts all loads that result in a miss, thus it includes the secondary misses as well.

B.2.3.2 Load Latency Event

Intel Processors based on the Intel microarchitecture code name Nehalem provide support for “load-latency event”, MEM_INST_RETIRED with event code 0BH and Umask value of 10H (LATENCY_ABOVE_THRESHOLD). This event samples loads, recording the number of cycles between the execution of the instruction and actual deliver of the data. If the measured latency is larger than the minimum latency programmed into MSR 0x3f6, bits 15:0, then the counter is incremented.

Counter overflow arms the PEBS mechanism and on the next event satisfying the latency threshold, the PMU writes the measured latency, the virtual or linear address, and the data source into a PEBS record format in the PEBS buffer. Because the virtual address is captured into a known location, the sampling driver could also execute a virtual to physical translation and capture the physical address. The physical address identifies the NUMA home location and in principle allows an analysis of the details of the cache occupancies.

Further, as the address is captured before retirement even the pointer chasing encoding “MOV RAX, [RAX+const]” have their addresses captured. Because the MSR_PEBS_LD_LAT_THRESHOLD MSR is required to specify the latency threshold value, only one minimum latency value can be sampled on a core during a given period. To enable this, the Intel performance tools restrict the programming of this event to counter 4 to simplify the scheduling. Table B-7 lists a few examples of event programming configurations used by the Intel® PTU and Vtune™ Performance Analyzer for the load latency events. Different threshold values for the minimum latencies are specified in MSR_PEBS_LD_LAT_THRESHOLD (address 0x3f6).

Table B-7. Load Latency Event Programming

Load Latency Precise Events	MSR 0x3F6	Umask	Event Code
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_4	4	10H	0BH

Table B-7. Load Latency Event Programming

Load Latency Precise Events	MSR 0x3F6	Umask	Event Code
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_8	8	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_10	16	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_20	32	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_40	64	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_80	128	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_100	256	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_200	512	10H	0BH
MEM_INST_RETIRED.LATENCY_ABOVE_THRESHOLD_800	32768	10H	0BH

One of the three fields written to each PEBS record by the PEBS assist mechanism of the load latency event, encodes the data source locality information.

Table B-8. Data Source Encoding for Load Latency PEBS Record

Encoding	Description
0x0	Unknown L3 cache miss
0x1	Minimal latency core cache hit. This request was satisfied by the L1 data cache.
0x2	Pending core cache HIT. Outstanding core cache miss to same cache-line address was already underway. The data is not yet in the data cache, but is located in a fill buffer that will soon be committed to cache.
0x3	This data request was satisfied by the L2.
0x4	L3 HIT. Local or Remote home requests that hit L3 cache in the uncore with no coherency actions required (snooping).
0x5	L3 HIT (other core hit snoop). Local or Remote home requests that hit the L3 cache and was serviced by another processor core with a cross core snoop where no modified copies were found. (clean).
0x6	L3 HIT (other core HITM). Local or Remote home requests that hit the L3 cache and was serviced by another processor core with a cross core snoop where modified copies were found. (HITM).
0x7	Reserved

Table B-8. Data Source Encoding for Load Latency PEBS Record (Contd.)

Encoding	Description
0x8	L3 MISS (remote cache forwarding). Local homed requests that missed the L3 cache and was serviced by forwarded data following a cross package snoop where no modified copies found. (Remote home requests are not counted).
0x9	Reserved
0xA	L3 MISS (local DRMA go to S). Local home requests that missed the L3 cache and was serviced by local DRAM (go to shared state).
0xB	L3 MISS (remote DRMA go to S). Remote home requests that missed the L3 cache and was serviced by remote DRAM (go to shared state).
0xC	L3 MISS (local DRMA go to E). Local home requests that missed the L3 cache and was serviced by local DRAM (go to exclusive state).
0xD	L3 MISS (remote DRMA go to E). Remote home requests that missed the L3 cache and was serviced by remote DRAM (go to exclusive state).
0xE	I/O, Request of input/output operation
0xF	The request was to un-cacheable memory.

The latency event is the recommended method to measure the penalties for a cycle accounting decomposition. Each time a PMI is raised by this PEBS event a load to use latency and a data source for the cacheline is recorded in the PEBS buffer. The data source for the cacheline can be deduced from the low order 4 bits of the data source field and the table shown above. Thus an average latency for each of the 16 sources can be evaluated from the collected data. As only one minimum latency at a time can be collected it may be awkward to evaluate the latency for an MLC hit and a remote socket dram. A minimum latency of 32 cycles should give a reasonable distribution for all the offcore sources however. The Intel® PTU version 3.2 performance tool can display the latency distribution in the data profiling mode and allows sophisticated event filtering capabilities for this event.

B.2.3.3 Precise Execution Events

PEBS capability in core PMU goes beyond load and store instructions. Branches, near calls and conditional branches can all be counted with precise events, for both retired and mispredicted (and retired) branches of the type selected. For these events, the PEBS buffer will contain the target of the branch. If the Last Branch Record (LBR) is also captured then the location of the branch instruction can also be determined.

When the branch is taken the IP value in the PEBS buffer will also appear as the last target in the LBR. If the branch was not taken (conditional branches only) then it won't and the branch that was not taken and retired is the instruction before the IP in the PEBS buffer.

In the case of near calls retired, this means that Event Based Sampling (EBS) can be used to collect accurate function call counts. As this is the primary measurement for

driving the decision to inline a function, this is an important improvement. In order to measure call counts, you must sample on calls. Any other trigger introduces a bias that cannot be guaranteed to be corrected properly.

The precise branch events can be found under event code C4H in Appendix A, “Performance Monitoring Events” of *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3B*.

There is one source of sampling artifact associated with precise events. It is due to the time delay between the PMU counter overflow and the arming of the PEBS hardware. During this period events cannot be detected due to the timing shadow. To illustrate the effect, consider a function call chain where a long duration function, “foo”, which calls a chain of 3 very short duration functions, “foo1” calling “foo2” which calls “foo3”, followed by a long duration function “foo4”. If the durations of foo1, foo2 and foo3 are less than the shadow period the distribution of PEBS sampled calls will be severely distorted. For example,

- if the overflow occurs on the call to foo, the PEBS mechanism is armed by the time the call to foo1 is executed and samples will be taken showing the call to foo1 from foo.
- if the overflow occurs due to the call to foo1, foo2 or foo3 however, the PEBS mechanism will not be armed until execution is in the body of foo4. Thus the calls to foo2, foo3 and foo4 cannot appear as PEBS sampled calls.

Shadowing can effect the distribution of all PEBS events. It will also effect the distribution of basic block execution counts identified by using the combination of a branch retired event (PEBS or not) and the last entry in the LBR. If there were no delay between the PMU counter overflow and the LBR freeze, the last LBR entry could be used to sample taken retired branches and from that the basic block execution counts. All the instructions between the last taken branch and the previous target are executed once.

Such a sampling could be used to generate a “software” instruction retired event with uniform sampling, which in turn can be used to identify basic block execution counts. Unfortunately the shadowing causes the branches at the end of short basic blocks to not be the last entry in the LBR, distorting the measurement. Since all the instructions in a basic block are by definition executed the same number of times.

The shadowing effect on call counts and basic block execution counts can be alleviated to a large degree by averaging over the entries in the LBR. This will be discussed in the section on LBRs.

Typically, branches account for more than 10% of all instructions in a workload, loop optimization needs to focus on those loops with high tripcounts. For counted loops, it is very common for the induction variable to be compared to the tripcount in the termination condition evaluation. This is particularly true if the induction variable is used within the body of the loop, even in the face of heavy optimization. Thus a loop sequence of unrolled operation by eight times may resemble:

```
add    rcx, 8
cmp    rcx, rax
```

```
jnge triad+0x27
```

In this case the two registers, `rax` and `rcx` are the tripcount and induction variable. If the PEBS buffer is captured for the conditional branches retired event, the average values of the two registers in the compare can be evaluated. The one with the larger average will be the tripcount. Thus the average, RMS, min and max can be evaluated and even a distribution of the recorded values.

B.2.3.4 Last Branch Record (LBR)

The LBR captures the source and target of each retired taken branch. Processors based on Intel microarchitecture code name Nehalem can track 16 pair of source/target addresses in a rotating buffer. Filtering of the branch instructions by types and privilege levels are permitted using a dedicated facility, `MSR_LBR_SELECT`. This means that the LBR mechanism can be programmed to capture branches occurring at ring0 or ring3 or both (default) privilege levels. Further the types of taken branches that are recorded can also be filtered. The list of filtering options that can be specified using `MSR_LBR_SELECT` is described in Chapter 16, “Debugging, Branch Profiles and Time-Stamp Counter” of *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3B*.

The default is to capture all branches at all privilege levels (all bits zero). Another reasonable programming would set all bits to 1 except bit 1 (capture ring 3) and bit 3 (capture near calls) and bits 6 and 7. This would leave only ring 3 calls and unconditional jumps in the LBR. Such a programming would result in the LBR having the last 16 taken calls and unconditional jumps retired and their targets in the buffer.

A PMU sampling driver could then capture this restricted “call chain” with any event, thereby providing a “call tree” context. The inclusion of the unconditional jumps will unfortunately cause problems, particularly when there are if-else structures within loops.

In the case of frequent function calls at all levels, the inclusion of returns could be added to clarify the context. However this would reduce the call chain depth that could be captured. A fairly obvious usage would be to trigger the sampling on extremely long latency loads, to enrich the sample with accesses to heavily contended locked variables, and then capture the call chain to identify the context of the lock usage.

Call Counts and Function Arguments

If the LBRs are captured for PMIs triggered by the `BR_INST_RETIRED.NEAR_CALL` event, then the call count per calling function can be determined by simply using the last entry in LBR. As the PEBS IP will equal the last target IP in the LBR, it is the entry point of the calling function. Similarly, the last source in the LBR buffer was the call site from within the calling function. If the full PEBS record is captured as well, then for functions with limited numbers of arguments on 64-bit OS’s, you can sample both the call counts and the function arguments.

LBRs and Basic Block Execution Counts

Another interesting usage is to use the BR_INST_RETIRED.ALL_BRANCHES event and the LBRs with no filter to evaluate the execution rate of basic blocks. As the LBRs capture all taken branches, all the basic blocks between a branch IP (source) and the previous target in the LBR buffer were executed one time. Thus a simple way to evaluate the basic block execution counts for a given load module is to make a map of the starting locations of every basic block. Then for each sample triggered by the PEBS collection of BR_INST_RETIRED.ALL_BRANCHES, starting from the PEBS address (a target but perhaps for a not taken branch and thus not necessarily in the LBR buffer) and walking backwards through the LBRs until finding an address not corresponding to the load module of interest, count all the basic blocks that were executed. Calling this value "number_of_basic_blocks", increment the execution counts for all of those blocks by $1/(\text{number_of_basic_blocks})$. This technique also yields the taken and not taken rates for the active branches. All branch instructions between a source IP and the previous target IP (within the same module) were not taken, while the branches listed in the LBR were taken. This is illustrated in the graphics below

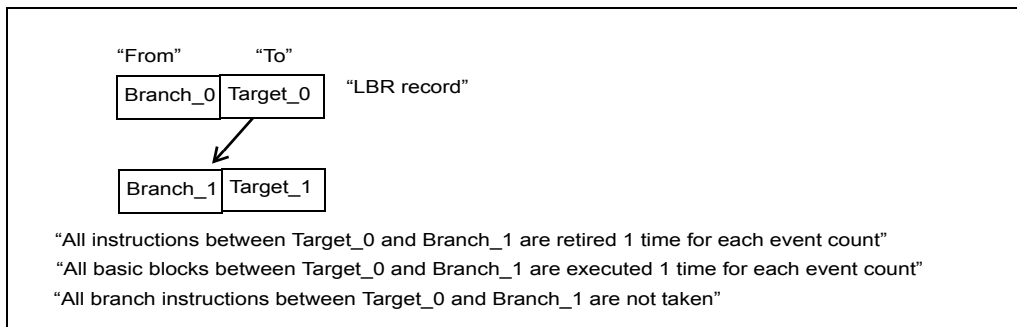


Figure B-3. LBR Records and Basic Blocks

The 16 sets LBR records can help rectify the artifact of PEBS samples aggregating disproportionately to certain instructions in the sampling process. The situation of skewed distribution of PEBS sample is illustrated below in Figure B-4.

Consider a number of basic blocks in the flow of normal execution, some basic block takes 20 cycles to execute, others taking 2 cycles, and shadowing takes 10 cycles. Each time an overflow condition occurs, the delay of PEBS being armed is at least 10 cycles. Once the PEBS is armed, PEBS record is captured on the next eventing condition. The skewed distribution of sampled instruction address using PEBS record will be skewed as shown in the middle of Figure B-4. In this conceptual example, we assume every branch is taken in these basic blocks.

In the skewed distribution of PEBS samples, the branch IP of the last basic block will be recorded 5 times as much as the least sampled branch IP address (the 2nd basic block).

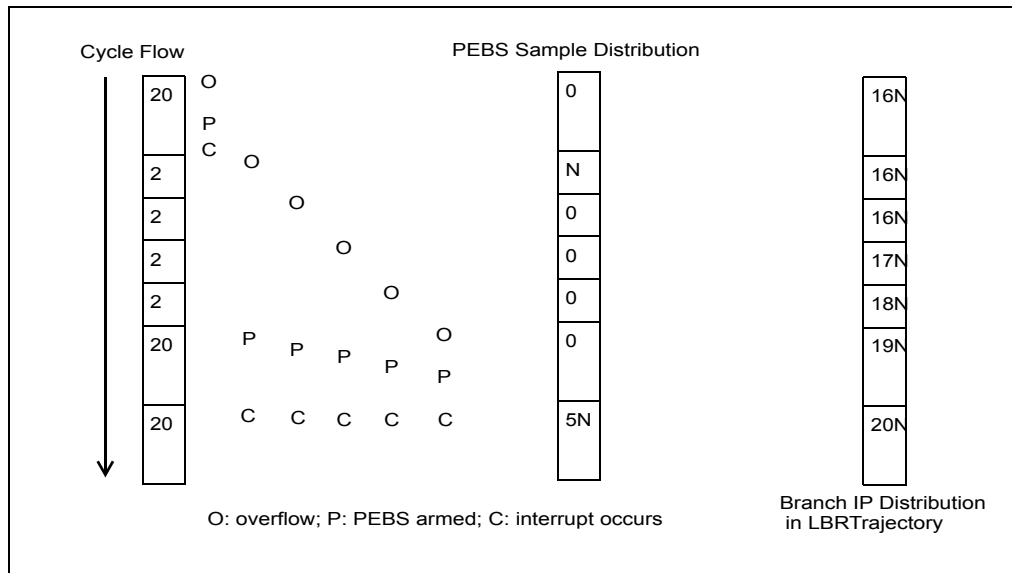


Figure B-4. Using LBR Records to Rectify Skewed Sample Distribution

This situation where some basic blocks would appear to never get samples and some have many times too many. Weighting each entry by $1/(\text{num of basic blocks in the LBR trajectory})$, in this example would result in dividing the numbers in the right most table by 16. Thus we end up with far more accurate execution counts $((1.25 \rightarrow 1.0) * N)$ in all of the basic blocks, even those that never directly caused a PEBS sample.

As on Intel® Core™2 processors there is a precise instructions retired event that can be used in a wide variety of ways. In addition there are precise events for `uops_retired`, various SSE instruction classes, FP assists. It should be noted that the FP assist events only detect x87 FP assists, not those involving SSE FP instructions. Detecting all assists will be discussed in the section on the pipeline Front End.

The instructions retired event has a few special uses. While its distribution is not uniform, the totals are correct. If the values recorded for all the instructions in a basic block are averaged, a measure of the basic block execution count can be extracted. The ratios of basic block executions can be used to estimate loop tripcounts when the counted loop technique discussed above cannot be applied.

The PEBS version (general counter) instructions retired event can further be used to profile OS execution accurately even in the face of STI/CLI semantics, because the PEBS interrupt then occurs after the critical section has completed, but the data was frozen correctly. If the `cmask` value is set to some very high value and the invert condition is applied, the result is always true, and the event will count core cycles (halted + unhalted).

Consequently both cycles and instructions retired can be accurately profiled. The UOPS_RETIRED.ANY event, which is also precise can also be used to profile Ring 0 execution and really gives a more accurate display of execution. The precise events available for this purpose are listed under event code C0H, C2H, C7H, F7H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.2.3.5 Measuring Core Memory Access Latency

Drilling down performance issues associated with locality or cache coherence issues will require using performance monitoring events. In each processor core, there is a super queue that allocates entries to buffer requests of memory access traffic due to an L2 miss to the uncore sub-system. Table B-9 lists various performance events available in the core PMU that can drill down performance issues related to L2 misses.

Table B-9. Core PMU Events to Drill Down L2 Misses

Core PMU Events	Umask	Event Code
OFFCORE_REQUESTS.DEMAND.READ_DATA ¹	01H	B0H
OFFCORE_REQUESTS.DEMAND.READ_CODE ¹	02H	B0H
OFFCORE_REQUESTS.DEMAND.RFO ¹	04H	B0H
OFFCORE_REQUESTS.ANY.READ	08H	B0H
OFFCORE_REQUESTS.ANY.RFO	10H	B0H
OFFCORE_REQUESTS.UNCACHED_MEM	20H	B0H
OFFCORE_REQUESTS.L1D.WRITEBACK	40H	B0H
OFFCORE_REQUESTS.ANY	80H	B0H

NOTES:

1. The *DEMAND* events also include any requests made by the L1D cache hardware prefetchers.

Table B-10 lists various performance events available in the core PMU that can drill down performance issues related to super queue operation.

Table B-10. Core PMU Events for Super Queue Operation

Core PMU Events	Umask	Event Code
OFFCORE_REQUESTS_BUFFER_FULL	01H	B2H
SQ_STALL	01H	F6H

Additionally, L2 misses can be drilled down further by data origin attributes and response attributes. The matrix to specify data origin and response type attributes is done by a dedicated MSR OFFCORE_RSP_0 at address 1A6H. See Table B-11 and Table B-12.

Table B-11. Core PMU Event to Drill Down OFFCore Responses

Core PMU Events	OFFCORE_RSP_0 MSR	Umask	Event Code
OFFCORE_RESPONSE	see Table B-12	01H	B7H

Table B-12. OFFCORE_RSP_0 MSR Programming

	Position	Description	Note
Request type	0	Demand Data Rd = DCU reads (includes partials, DCU Prefetch)	
	1	Demand RFO = DCU RFOs	
	2	Demand Ifetch = IFU Fetches	
	3	Writeback = L2_EVICT/DCUWB	
	4	PF Data Rd = L2 Prefetcher Reads	
	5	PF RFO= L2 Prefetcher RFO	
	6	PF Ifetch= L2 Prefetcher Instruction fetches	
	7	Other	Include non-temporal stores

Table B-12. OFFCORE_RSP_0 MSR Programming

	Position	Description	Note
Response type	8	L3_HIT_UNCORE_HIT	exclusive line
	9	L3_HIT_OTHER_CORE_HIT_SNP	clean line
	10	L3_HIT_OTHER_CORE_HITM	modified line
	11	L3_MISS_REMOTE_HIT_SCRUB	Used by multiple cores
	12	L3_MISS_REMOTE_FWD	Clean line used by one core
	13	L3_MISS_REMOTE_DRAM	
	14	L3_MISS_LOCAL_DRAM	
	15	Non-DRAM	Non-DRAM requests

Although Table B-12 allows 2^{16} combinations of setting in MSR_OFFCORE_RSP_0 in theory, it is more useful to consider combining the subsets of 8-bit values to specify "Request type" and "Response type". The more common 8-bit mask values are listed in Table B-13.

Table B-13. Common Request and Response Types for OFFCORE_RSP_0 MSR

Request Type	Mask	Response Type	Mask
ANY_DATA	xx11H	ANY_CACHE_DRAM	7FxxH
ANY_IFETCH	xx44H	ANY_DRAM	60xxH
ANY_REQUEST	xxFFH	ANY_L3_MISS	F8xxH
ANY_RFO	xx22H	ANY_LOCATION	FFxxH
CORE_WB	xx08H	IO	80xxH
DATA_IFETCH	xx77H	L3_HIT_NO_OTHER_CORE	01xxH
DATA_IN	xx33H	L3_OTHER_CORE_HIT	02xxH
DEMAND_DATA	xx03H	L3_OTHER_CORE_HITM	04xxH
DEMAND_DATA_RD	xx01H	LOCAL_CACHE	07xxH
DEMAND_IFETCH	xx04H	LOCAL_CACHE_DRAM	47xxH

Table B-13. Common Request and Response Types for OFFCORE_RSP_0 MSR

Request Type	Mask	Response Type	Mask
DEMAND_RFO	xx02H	LOCAL_DRAM	40xxH
OTHER ¹	xx80H	REMOTE_CACHE	18xxH
PF_DATA	xx30H	REMOTE_CACHE_DRAM	38xxH
PF_DATA_RD	xx10H	REMOTE_CACHE_HIT	10xxH
PF_IFETCH	xx40H	REMOTE_CACHE_HITM	08xxH
PF_RFO	xx20H	REMOTE-DRAM	20xxH
PREFETCH	xx70H		

NOTES:

1. The PMU may report incorrect counts with setting MSR_OFFCORE_RSP_0 to the value of 4080H. Non-temporal stores to the local DRAM is not reported in the count.

B.2.3.6 Measuring Per-Core Bandwidth

Measuring the bandwidth of all memory traffic for an individual core is complicated, the core PMU and uncore PMU do provide capability to measure the important components of per-core bandwidth.

At the microarchitectural level, there is the buffering of L3 for writebacks/evictions from L2 (similarly to some degree with the non temporal writes). The eviction of modified lines from the L2 causes a write of the line back to the L3. The line in L3 is only written to memory when it is evicted from the L3 some time later (if at all). And L3 is part of the uncore sub-system, not part of the core.

The writebacks to memory due to eviction of modified lines from L3 cannot be associated with an individual core in the uncore PMU logic. The net result of this is that the total write bandwidth for all the cores can be measured with events in the uncore PMU. The read bandwidth and the non-temporal write bandwidth can be measured on a per core basis. In a system populated with two physical processor, the NUMA nature of memory bandwidth implies the measurement for those 2 components has to be divided into bandwidths for the core on a per-socket basis.

The per-socket read bandwidth can be measured with the events:

- OFFCORE_RESPONSE_0.DATA_IFETCH.L3_MISS_LOCAL_DRAM
- OFFCORE_RESPONSE_0.DATA_IFETCH.L3_MISS_REMOTE_DRAM

The total read bandwidth for all sockets can be measured with the event:

- OFFCORE_RESPONSE_0.DATA_IFETCH.ANY_DRAM

The per-socket non-temporal store bandwidth can be measured with the events:

OFFCORE_RESPONSE_0.OTHER.L3_MISS_LOCAL_CACHE_DRAM
 OFFCORE_RESPONSE_0.OTHER.L3_MISS_REMOTE_DRAM

The total non-temporal store bandwidth can be measured with the event:

OFFCORE_RESPONSE_0.OTHER.ANY.CACHE_DRAM

The use of "CACHE_DRAM" encoding is to work around the defect in the footnote of Table B-13. Note that none of the above includes the bandwidth associated with writebacks of modified cacheable lines.

B.2.3.7 Miscellaneous L1 and L2 Events for Cache Misses

In addition to the OFFCORE_RESPONSE_0 event and the precise events that will be discussed later, there are several other events that can be used as well. There are additional events that can be used to supplement the offcore_response_0 events, because the offcore_response_0 event code is supported on counter 0 only.

L2 misses can also be counted with the architecturally defined event LONGEST_LAT_CACHE_ACCESS, however as this event also includes requests due to the L1D and L2 hardware prefetchers, its utility may be limited. Some of the L2 access events can be used for both drilling down L2 accesses and L2 misses by type, in addition to the OFFCORE_REQUESTS events discussed earlier. The L2_RQSTS and L2_DATA_RQSTS events can be used to discern assorted access types. In all of the L2 access events the designation PREFETCH only refers to the L2 hardware prefetch. The designation DEMAND includes loads and requests due to the L1D hardware prefetchers.

The L2_LINES_IN and L2_LINES_OUT events have been arranged slightly differently than the equivalent events on Intel® Core™2 processors. The L2_LINES_OUT event can now be used to decompose the evicted lines by clean and dirty (i.e. a Writeback) and whether they were evicted by an L1D request or an L2 HW prefetch.

The event L2_TRANSACTIONS counts all interactions with the L2.

Writes and locked writes are counted with a combined event, L2_WRITE.

The details of the numerous derivatives of L2_RQSTS, L2_DATA_RQSTS, L2_LINES_IN, L2_LINES_OUT, L2_TRANSACTIONS, L2_WRITE, can be found under event codes 24H, 26H, F1H, F2H, F0H, and 27H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.2.3.8 TLB Misses

The next largest set of memory access delays are associated with the TLBs when linear-to-physical address translation is mapped with a finite number of entries in the TLBs. A miss in the first level TLBs results in a very small penalty that can usually be hidden by the OOO execution and compiler's scheduling. A miss in the shared TLB

results in the Page Walker being invoked and this penalty can be noticeable in the execution.

The (non-PEBS) TLB miss events break down into three sets:

- DTLB misses and its derivatives are programmed with event code 49H,
- Load DTLB misses and its derivatives are programmed with event code 08H,
- ITLB misses and its derivatives are programmed with event code 85H.

Store DTLB misses can be evaluated from the difference of the DTLB misses and the Load DTLB misses. Each then has a set of sub events programmed with the umask value. The Umask details of the numerous derivatives of the above events are listed in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.2.3.9 L1 Data Cache

There are PMU events that can be used to analyze L1 data cache operations. These events can only be counted with the first 2 of the 4 general counters, i.e. IA32_PMC0 and IA32_PMC1. Most of the L1D events are self explanatory.

The total number of references to the L1D can be counted with L1D_ALL_REF, either just cacheable references or all. The cacheable references can be divided into loads and stores with L1D_CACHE_LOAD and L1D_CACHE.STORE. These events are further subdivided by MESI states through their Umask values, with the I state references indicating the cache misses.

The evictions of modified lines in the L1D result in writebacks to the L2. These are counted with the L1D_WB_L2 events. The umask values break these down by the MESI state of the version of the line in the L2.

The locked references can be counted also with the L1D_CACHE_LOCK events. Again these are broken down by MES states for the lines in L1D.

The total number of lines brought into L1D, the number that arrived in an M state and the number of modified lines that get evicted due to receiving a snoop are counted with the L1D event and its Umask variations.

The L1D events are listed under event codes 28H, 40H, 41H, 42H, 43H, 48H, 4EH, 51H, 52H, 53H, 80H, and 83H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

There are few cases of loads not being able to forward from active store buffers. The predominant situations have to do with larger loads overlapping smaller stores. There is not event that detects when this occurs. There is also a "false store forwarding" case where the addresses only match in the lower 12 address bits. This is sometimes referred to as 4K aliasing. This can be detected with the event "PARTIAL_ADDRESS_ALIAS" which has event code 07H and Umask 01H.

B.2.4 Front End Monitoring Events

Branch misprediction effects can sometimes be reduced through code changes and enhanced inlining. Most other front end performance limitations have to be dealt with by the code generation. The analysis of such issues is mostly of use by compiler developers.

B.2.4.1 Branch Mispredictions

In addition to branch retired events that was discussed in conjunction with PEBS in Section B.2.3.3. These are enhanced by use of the LBR to identify the branch location to go along with the target location captured in the PEBS buffer. Aside from those usage, many other PMU events (event code E6, E5, E0, 68, 69) associated with branch predictions are more relevant to hardware design than performance tuning.

Branch mispredictions are not in and of themselves an indication of a performance bottleneck. They have to be associated with dispatch stalls and the instruction starvation condition, `UOPS_ISSUED:C1:I1 – RESOURCE_STALLS.ANY`. Such stalls are likely to be associated with icache misses and ITLB misses. The precise ITLB miss event can be useful for such issues. The icache and ITLB miss events are listed under event code 80H, 81H, 82H, 85H, AEH.

B.2.4.2 Front End Code Generation Metrics

The remaining front end events are mostly of use in identifying when details of the code generation interact poorly with the instructions decoding and uop issue to the OOO engine. Examples are length changing prefix issues associated with the use of 16 bit immediates, rob read port stalls, instruction alignment interfering with the loop detection and instruction decoding bandwidth limitations. The activity of the LSD is monitored using `CMASK` values on a signal monitoring activity. Some of these events are listed under event code 17H, 18H, 1EH, 1FH, 87H, A6H, A8H, D0H, D2H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

Some instructions (FSIN, FCOS, and other transcendental instructions) are decoded with the assistance of MS-ROM. Frequent occurrences of instructions that required assistance of MS-ROM to decode complex uop flows are opportunity to improve instruction selection to reduce such occurrences. The `UOPS_DECODED.MS` event can be used to identify code regions that could benefit from better instruction selection.

Other situations that can trigger this event are due to FP assists, like performing a numeric operation on denormalized FP values or QNaNs. In such cases the penalty is essentially the uops required for the assist plus the pipeline clearing required to ensure the correct state.

Consequently this situation has a very clear signature consisting of `MACHINE_CLEAR.CYCLES` and uops being inserted by the microcode sequencer, `UOPS_DECODED.MS`. The execution penalty being the sum of these two contributions. The event codes for these are listed under D1H and C3H.

B.2.5 Uncore Performance Monitoring Events

The uncore sub-system includes the L3, IMC and Intel QPI units in the diagram shown in Figure B-1. Within the uncore sub-system, the uncore PMU consists of eight general-purpose counters and one fixed counter. The fixed counter in uncore monitors the unhalted clock cycles in the uncore clock domain, which runs at a different frequency than the core.

The uncore cannot by itself generate a PMI interrupt. While the core PMU can raise PMI at a per-logical-processor specificity, the uncore PMU can cause PMI at a per-core specificity using the interrupt hardware in the processor core. When an uncore counter overflows, a bit pattern is used to specify which cores should be signaled to raise a PMI. The uncore PMU is unaware of the core, Processor ID or Thread ID that caused the event that overflowed a counter. Consequently the most reasonable approach for sampling on uncore events is to raise a PMI on all the logical processors in the package.

There are a wide variety of events that monitor queue occupancies and inserts. There are others that count cacheline transfers, dram paging policy statistics, snoop types and responses, and so on. The uncore is the only place the total bandwidth to memory can be measured. This will be discussed explicitly after all the uncore components and their events are described.

B.2.5.1 Global Queue Occupancy

Each processor core has a super queue that buffers requests of memory access traffic due to an L2 miss. The uncore has a global queue (GQ) to service transaction requests from the processor cores and buffers data traffic that arrive from L3, IMC, or Intel QPI links.

Within the GQ, there are 3 "trackers" in the GQ for three types of transactions:

- on-package read requests, its tracker queue has 32 entries.
- on-package writeback requests, its tracker queue has 16 entries
- requests that arrive from a "peer", its tracker queue has 12 entries.

A "peer" refers to any requests coming from the Intel® QuickPath Interconnect.

The occupancies, inserts, cycles full and cycles not empty for all three trackers can be monitored. Further as load requests go through a series of stages the occupancy and inserts associated with the stages can also be monitored, enabling a "cycle accounting" breakdown of the uncore memory accesses due to loads.

When a uncore counter is first programmed to monitor a queue occupancy, for any of the uncore queues, the queue must first be emptied. This is accomplished by the driver of the monitoring software tool issuing a bus lock. This only needs to be done when the counter is first programmed. From that point on the counter will correctly reflect the state of the queue, so it can be repeatedly sampled for example without another bus lock being issued.

The uncore events that monitor GQ allocation (UNC_GQ_ALLOC) and GQ tracker occupancy (UNC_GQ_TRACKER_OCCUP) are listed under the event code 03H and 02H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*. The selection between the three trackers is specified from the Umask value. The mnemonic of these derivative events use the notation: "RT" signifying the read tracker, "WT", the write tracker and "PPT" the peer probe tracker.

Latency can be measured by the average duration of the queue occupancy, if the occupancy stops as soon as the data has been delivered. Thus the ratio of $\text{UNC_GQ_TRACKER_OCCUP.X/UNC_GQ_ALLOC.X}$ measures an average duration of queue occupancy, where 'X' represents a specific Umask value. The total occupancy period of the read tracker as measured by

$$\text{Total Read Period} = \text{UNC_GQ_TRACKER_OCCUP.RT/UNC_GQ_ALLOC.RT}$$

Is longer than the data delivery latency due to it including time for extra bookkeeping and cleanup. The measurement

$$\text{LLC response Latency} = \frac{\text{UNC_GQ_TRACKER_OCCUP.RT_TO_LLC_RESP}}{\text{UNC_GQ_ALLOC.RT_TO_LLC_RESP}}$$

is essentially a constant. It does not include the total time to snoop and retrieve a modified line from another core for example, just the time to scan the L3 and see if the line is or is not present in this socket.

An overall latency for an L3 hit is the weighted average of three terms:

- the latency of a simple hit, where the line has only been used by the core making the request,
- the latencies for accessing clean lines by multiple cores,
- the latencies for accessing dirty lines that have been accessed by multiple cores.

These three components of the L3 hit for loads can be decomposed using the derivative events of OFFCORE_RESPONSE:

- OFFCORE_RESPONSE_0.DEMAND_DATA.L3_HIT_NO_OTHER_CORE,
- OFFCORE_RESPONSE_0.DEMAND_DATA.L3_HIT_OTHER_CORE_HIT,
- OFFCORE_RESPONSE_0.DEMAND_DATA.L3_HIT_OTHER_CORE_HITM.

The event OFFCORE_RESPONSE_0.DEMAND_DATA.LOCAL_CACHE should be used as the denominator to obtain latencies. The individual latencies could have to be measured with microbenchmarks, but the use of the precise latency event will be far more effective as any bandwidth loading effects will be included.

The L3 miss component is the weighted average over three terms:

- the latencies of L3 hits in a cache on another socket (this is described in the previous paragraph) ,
- the latencies to local DRAM,
- the latencies to remote DRAM.

USING PERFORMANCE MONITORING EVENTS

The local dram access and the remote socket access can be decomposed with more uncore events. This will be discussed a bit later in this paper.

$$\text{Miss to fill latency} = \text{UNC_GQ_TRACKER_OCCUP_RT_LLC_MISS} / \text{UNC_GQ_ALLOC_RT_LLC_MISS}$$

The uncore GQ events using Umask value associated with *RTID* mnemonic allow the monitoring of a sub component of the Miss to fill latency associated with the communications between the GQ and the QHL.

There are uncore PMU events which monitor cycles when the three trackers are not empty (≥ 1 entry) or full. These events are listed under the event code 00H and 01H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*

Because the uncore PMU generally does not differentiate which processor core causes a particular eventing condition, the technique of dividing the latencies by the average queue occupancy in order to determine a penalty does not work for the uncore. Overlapping entries from different cores do not result in overlapping penalties and thus a reduction in stalled cycles. Each core suffers the full latency independently.

To evaluate the correction on a per core basis one needs the number of cycles there is an entry from the core in question. A *NOT_EMPTY_CORE_N* type event would be needed. There is no such event. Consequently, in the cycle decomposition one must use the full latency for the estimate of the penalty. As has been stated before it is best to use the PEBS latency event as the data sources are also collected with the latency for the individual sample.

The individual components of the read tracker, discussed above, can also be monitored as busy or full by setting the cmask value to 1 or 32 and applying it to the assorted RT occupancy events.

Table B-14. Uncore PMU Events for Occupancy Cycles

Uncore PMU Events	Cmask	Umask	Event Code
UNC_GQ_TRACKER_OCCUP_RT_L3_MISS_FULL	32	02H	02H
UNC_GQ_TRACKER_OCCUP_RT_TO_L3_RESP_FULLL	32	04H	02H
UNC_GQ_TRACKER_OCCUP_RT_TO_RTID_ACCQ_UIRED_FULL	32	08H	02H
UNC_GQ_TRACKER_OCCUP_RT_L3_MISS_BUSY	1	02H	02H
UNC_GQ_TRACKER_OCCUP_RT_TO_L3_RESP_BUSY	1	04H	02H

Table B-14. Uncore PMU Events for Occupancy Cycles

Uncore PMU Events	Cmask	Umask	Event Code
UNC_GQ_TRACKER_OCCUP.RT_TO_RTID_ACCQ UIRED_BUSY	1	08H	02H

B.2.5.2 Global Queue Port Events

The GQ data buffer traffic controls the flow of data to and from different sub-systems via separate ports:

- Core traffic: two ports handles data traffic, each port dedicated to a pair of processor cores.
- L3 traffic: one port service L3 data traffic
- Intel QPI traffic: one service traffic to QPI logic
- IMC traffic: one service data traffic to integrated memory controller.

The ports for L3 and core traffic transfer a fixed number of bits per cycle. However the Intel® QuickPath Interconnect protocols can result in either 8 or 16 bytes being transferred on the read Intel QPI and IMC ports. Consequently these events cannot be used to measure total data transfers and bandwidths.

The uncore PMU events that can distinguish traffic flow are listed under the event code 04H and 05H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.2.5.3 Global Queue Snoop Events

Cacheline requests from the cores or from a remote package or the I/O Hub are handled by the GQ. When the uncore receives a cacheline request from one of the cores, the GQ first checks the L3 to see if the line is on the package. Because the L3 is inclusive, this answer can be quickly ascertained. If the line is in the L3 and was owned by the requesting core, data can be returned to the core from the L3 directly. If the line is being used by multiple cores, the GQ will snoop the other cores to see if there is a modified copy. If so the L3 is updated and the line is sent to the requesting core.

In the event of an L3 miss, the GQ must send out requests to the local memory controller (or over the Intel QPI links) for the line. A request through the Intel QPI to a remote L3 (or remote DRAM) must be made if data exists in a remote L3 or does not exist in local DRAM. As each physical package has its own local integrated memory controller the GQ must identify the "home" location of the requested cache-line from the physical address. If the address identifies home as being on the local package then the GQ makes a simultaneous request to the local memory controller. If home is identified as belonging to the remote package, the request sent over the Intel QPI will also access the remote IMC.

The GQ handles the snoop responses for the cacheline requests that come in from the Intel® QuickPath Interconnect. These snoop traffic correspond to the queue entries in the peer probe tracker.

The snoop responses are divided into requests for locally homed data and remotely homed data. If the line is in a modified state and the GQ is responding to a read request, the line also must be written back to memory. This would be a wasted effort for a response to a RFO as the line will just be modified again, so no Writeback is done for RFOs.

The snoop responses of local home events that can be monitored by an uncore PMU are listed under event code 06H in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*. The snoop responses of remotely home events are listed under event code 07H.

Some related events count the MESI transitions in response to snoops from other caching agents (processors or IOH). Some of these rely on programming MSR so they can only be measured one at a time, as there is only one MSR. The Intel performance tools will schedule this correctly by restricting these events to a single general uncore counter.

B.2.5.4 L3 Events

Although the number of L3 hits and misses can be determined from the GQ tracker allocation events, Several uncore PMU event is simpler to use. They are listed under event code 08H and 09H in the uncore event list of Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

The MESI states breakdown of lines allocated and victimized can also be monitored with LINES_IN, LINES_OUT events in the uncore using event code 0AH and 0BH. Details are listed in Appendix A, "Performance Monitoring Events" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.2.6 Intel QuickPath Interconnect Home Logic (QHL)

When a data misses L3 and causing the GQ of the uncore to send out a transaction request, the Intel QPI fabric will fulfill the request either from the local DRAM controller or from a remote DRAM controller in another physical package. The GQ must identify the "home" location of the requested cacheline from the physical address. If the address identifies home as being on the local package then the GQ makes a simultaneous request to the local memory controller, the Integrated memory controller (IMC). If home is identified as belonging to the remote package, the request is sent to the Intel QPI first and then to access the remote IMC.

The Intel QPI logic and IMC are distinct units in the uncore sub-system. The Intel QPI logic distinguish the local IMC relative to remote IMC using the concept of "caching agent" and "home agent". Specifically, the Intel QPI protocol considers each socket as having a "caching agent": and a "home agent":

- Caching Agent is the GQ and L3 in the uncore (or an IOH if present),
- Home Agent is the IMC.

An L3 miss results in simultaneous queries for the line from all the Caching Agents and the Home agent (wherever it is).

QHL requests can be superseded when another source can supply the required line more quickly. L3 misses to locally homed lines, due to on package requests, are simultaneously directed to the QHL and Intel QPI. If a remote caching agent supplies the line first then the request to the QHL is sent a signal that the transaction is complete. If the remote caching agent returns a modified line in response to a read request then the data in dram must be updated with a writeback of the new version of the line.

There is a similar flow of control signals when the Intel QPI simultaneously sends a snoop request for a locally homed line to both the GQ and the QHL. If the L3 has the line, the QHL must be signaled that the transaction was completed by the L3/GQ. If the line in L3 (or the cores) was modified and the snoop request from the remote package was for a load, then a writeback must be completed by the QHL and the QHL forwards the line to the Intel QPI to complete the transaction.

Uncore PMU provides events for monitoring these cacheline access and writeback traffic in the uncore by using the QHL opcode matching capability. The opcode matching facility is described in Chapter 30, "Performance Monitoring" of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*. The uncore PMU event that uses the opcode matching capability is listed under event code 35H. Several of the more useful settings to program QHL opcode matching is shown in Table B-15

Table B-15. Common QHL Opcode Matching Facility Programming

Load Latency Precise Events	MSR 0x396	Umask	Event Code
UNC_ADDR_OPCODE_MATCH.IOH.NONE	0	1H	35H
UNC_ADDR_OPCODE_MATCH.IOH.RSPFWDI	40001900_00000000	1H	35H
UNC_ADDR_OPCODE_MATCH.IOH.RSPFWDS	40001A00_00000000	1H	35H
UNC_ADDR_OPCODE_MATCH.IOH.RSPIWB	40001D00_00000000	1H	35H
UNC_ADDR_OPCODE_MATCH.REMOTE.NON E	0	2H	35H
UNC_ADDR_OPCODE_MATCH.REMOTE.RSP FWDI	40001900_00000000	2H	35H

Table B-15. Common QHL Opcode Matching Facility Programming

Load Latency Precise Events	MSR 0x396	Uma sk	Event Code
UNC_ADDR_OPCODE_MATCH.REMOTE.RSP FWDS	40001A00_00000000	2H	35H
UNC_ADDR_OPCODE_MATCH.REMOTE.RSPI WB	40001D00_00000000	2H	35H
UNC_ADDR_OPCODE_MATCH.LOCAL.NONE	0	4H	35H
UNC_ADDR_OPCODE_MATCH.LOCAL.RSPF WDI	40001900_00000000	1H	35H
UNC_ADDR_OPCODE_MATCH.LOCAL.RSPF WDS	40001A00_00000000	1H	35H
UNC_ADDR_OPCODE_MATCH.LOCAL.RSPIW B	40001D00_00000000	1H	35H

These predefined opcode match encodings can be used to monitor HITM accesses. It is the only event that allows profiling the code requesting HITM transfers.

The diagrams Figure B-5 through Figure B-12 show a series of Intel QPI protocol exchanges associated with Data Reads and Reads for Ownership (RFO), after an L3 miss, under a variety of combinations of the local home of the cacheline, and the MESI state in the remote cache. Of particular note are the cases where the data comes from the remote QHL even when the data was in the remote L3. These are the Read Data with the remote L3 having the line in an M state.

Whether the line is locally or remotely “homed” it has to be written back to dram before the originating GQ receives the line, so it always appears to come from a QHL. The RFO does not do this. However, when responding to a remote RFO (SnpInvOwn) and the line is in an S or F state, the cacheline gets invalidated and the line is sent from the QHL. The point is that the data source might not always be so obvious.

B.2.7 Measuring Bandwidth From the Uncore

Read bandwidth can be measured on a per core basis using events like `OFFCORE_RESPONSE_0.DATA_IN.LOCAL_DRAM` and `OFFCORE_RESPONSE_0.DATA_IN.REMOTE_DRAM`. The total bandwidth includes writes and these cannot be monitored from the core as they are mostly caused by evictions of modified lines in the L3. Thus a line used and modified by one core can end up being written back to dram when it is evicted due to a read on another core doing some completely unrelated task. Modified cached lines and writebacks of uncached lines (e.g. written with non temporal streaming stores) are handled differently in the uncore and their writebacks increment various events in different ways.

All full lines written to DRAM are counted by the `UNC_IMC_WRITES.FULL.*` events. This includes the writebacks of modified cached lines and the writes of uncached lines, for example generated by non-temporal SSE stores. The uncached line writebacks from a remote socket will be counted by `UNC_QHL_REQUESTS.REMOTE_WRITES`. The uncached writebacks from the local cores are not counted by `UNC_QHL_REQUESTS.LOCAL_WRITES`, as this event only counts writebacks of locally cached lines.

The `UNC_IMC_NORMAL_READS.*` events only count the reads. The `UNC_QHL_REQUESTS.LOCAL_READS` and the `UNC_QHL_REQUESTS.REMOTE_READS` count the reads and the “InvtoE” transactions, which are issued for the uncacheable writes, eg USWC/UC writes. This allows the evaluation of the uncacheable writes, by computing the difference of `UNC_QHL_REQUESTS.LOCAL_READS +`

`UNC_QHL_REQUESTS.REMOTE_READS - UNC_IMC_NORMAL_READS.ANY`.

These uncore PMU events that are useful for bandwidth evaluation are listed under event code 20H, 2CH, 2FH in Appendix A, “Performance Monitoring Events” of *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3B*.

B.3 PERFORMANCE TUNING TECHNIQUES FOR INTEL® MICROARCHITECTURE CODE NAME SANDY BRIDGE

This section covers various performance tuning techniques using performance monitoring events. Some techniques can be adapted in general to other microarchitectures, most of the performance events are specific to Intel microarchitecture code name Sandy Bridge.

B.3.1 Correlating Performance Bottleneck to Source Location

Performance analysis tools often sample events to identify hot spots of instruction pointer addresses to help programmers identify source locations of potential performance bottlenecks.

The sampling technique requires a service routine to respond to the performance monitoring interrupt (PMI) generated from an overflow condition of the performance counter. There is a finite delay between the performance monitoring event detection of the eventing condition relative to the capture of the instruction pointer address. This is known as “skid”. In other words, the event skid is the distance between the instruction or instructions that caused the issue and the instruction where the event is tagged. There are a few things to note in general:

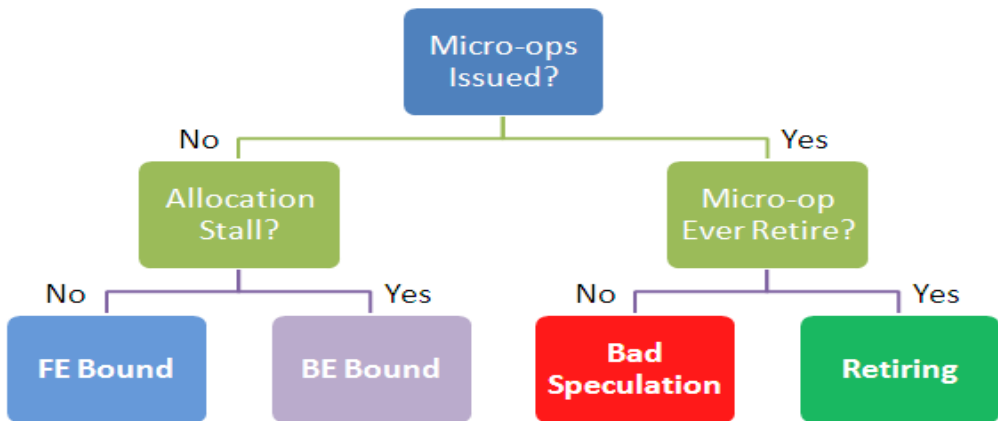
- Precise events have a defined event skid of 1 IP to the next IP retired. In the case when the offending instruction is a branch, the event is tagged with the branch target, which can be separated from the branch instruction. Thus sampling with precise events are likely to have less noise in pin-pointing source locations of bottlenecks.
- Using a performance event with eventing condition that carries a larger performance impact generally has a shorter skid and vice versa. The following examples illustrate this rule:
 - A store forward block issue can cause a delay of more than 10 cycles. Sampling a store forward block event almost always tags to the next couple of IPs after the blocked load.
 - On the other hand, sampling loads that forwarded successfully will have larger skids, and less helpful for performance tuning.
- The closer the eventing condition is to the retirement of the instruction, the shorter the skid. The events in the front end of the pipeline tend to tag to IPs further from the responsible IP than events that are taken at execution or retirement.
- Cycles counted with the event `CPU_CLK_UNHALTED.THREAD` often tag in greater counts on the IP after larger bottlenecks in the pipeline. If cycles are accumulated on an IP this is probably due to a bottleneck on the instruction at the previous IP.
- It is very difficult to determine the source of issues with a low cost that occur in the front end. Front-end events can also skid to IPs that precede the actual instructions that are causing the issue.

B.3.2 Locating Stalls in the Microarchitecture Pipeline

To determine what portion of the pipeline is stalled, you need to look at a buffer that queues the micro-ops supplied by the front end and feeds the out-of-order back end (see Section 2.1.1). This buffer is called the micro-op queue. From the micro-op queue viewpoint, there may be four different types of stalls:

- Front end stalls - The front end is delivering less than four micro-ops per cycle when the back end of the pipeline is requesting micro-ops. When these stalls happen, the rename/allocate part of the OOO engine will starved. Thus, execution is said to be front-end bound.
- Back end stalls – No micro-ops are being delivered from the micro-op queue due to lack of required resources for accepting more micro-ops in the back end of the pipeline. When these stalls happen, execution is said to be back-end bound.
- Bad speculation - The pipeline performs speculative execution of instructions that never successfully retire. The most common case is a branch misprediction where the pipeline predicts a branch target in order to keep the pipeline full instead of waiting for the branch to execute. If the processor prediction is incorrect it has to flush the pipeline without retiring the speculated instructions.
- Retiring – The micro-op queue delivers micro-ops that eventually retire. In the common case, the micro-ops originate from the program code. One exception is with assists where the microcode sequencer generates micro-ops to deal with issues in the pipeline.

The following figure illustrates how the execution opportunities are logically divided:



It is possible to estimate the amount of execution slots spent in each category using the following formulas in conjunction with core PMU performance events in Intel microarchitecture code name Sandy Bridge:

%FE_Bound =

$$100 * (IDQ_UOPS_NOT_DELIVERED.CORE / N);$$

%Bad_Speculation =

$$100 * ((UOPS_ISSUED.ANY - UOPS_RETIRED.RETIRE_SLOTS) / N);$$

%Retiring = $100 * (UOPS_RETIRED.RETIRE_SLOTS / N);$

%BE_Bound = $100 * (1 - (FE_Bound + Retiring + Bad_Speculation));$

N represents total execution slots opportunities. Execution opportunities are the number of cycles multiplied by four.

- $N = 4 * CPU_CLK_UNHALTED.THREAD$

The following sections explain the source for penalty cycles in three categories: back end stalls, front end stalls and bad speculation. They use formulas that can be applied to process, module, function, and instruction granularity.

B.3.3 Back End Stalls

Back end stalls have two main sources: memory sub-system stalls and execution stalls. As a first step to understanding the source of back end stalls, use the resource stall event.

Before putting micro-ops into the scheduler, the rename stage has to have certain resources allocated. When an application encounters a significant bottleneck at the back end of the pipeline, it runs out of these resources as the pipeline backs up. The RESOURCE_STALLS event tracks stall cycles when a resource could not be allocated. The event breaks up each resource into a separate sub-event so you can track which resource is not available for allocation. Counting these events can help identifying the reason for issues in the back end of the pipeline.

The resource stall ratios described below can be accomplished at process, module, function and even instruction granularities with the cycles, counted by CPU_CLK_UNHALTED.THREAD, representing the penalty tagged at the same granularity.

Usages of Specific Events

RESOURCE_STALLS.ANY - Counts stall cycles that the rename stage is unable to put micro-ops into the scheduler, due to lack of resources that have to be allocated at this stage. The event skid tends to be low since it is close to the retirement of the blocking instruction. This event accounts for all stalls counted by other RESOURCE_STALL sub events and also includes the sub-events of RESOURCE_STALLS2. If this ratio is high, count the included sub-events to get a better isolation of the reason for the stall.

$$\%RESOURCE_STALLS.COST = 100 * RESOURCE_STALLS.ANY / CPU_CLK_UNHALTED.THREAD;$$

RESOURCE_STALLS.SB - Occurs when a store micro-op is ready for allocation and all store buffer entries are in use, usually due to long latency stores in progress. Typically this event tags to the IP after the store instruction that is stalled at allocation.

$$\%RESOURCE_STALLS.SB.COST = 100 * RESOURCE_STALLS.SB / CPU_CLK_UNHALTED.THREAD;$$

RESOURCE_STALLS.LB - Counts cycles in which a load micro-op is ready for allocation and all load buffer entries are taken, usually due to long latency loads in progress. In many cases the queue to the scheduler becomes full by micro-ops that depend on the long latency loads, before the load buffer gets full.

```
%RESOURCE.STALLS.LB.COST =
    100 * RESOURCE_STALLS.LB / CPU_CLK_UNHALTED.THREAD;
```

In the above cases the event RESOURCE_STALLS.RS will often count in parallel. The best methodology to further investigate loss in data locality is the high cache line replacement study described in Section B.3.4.2, concentrating on L1 DCache replacements first

RESOURCE_STALLS.RS - Scheduler slots are typically the first resource that runs out when the pipeline is backed up. However, this can be due to almost any bottleneck in the back end, including long latency loads and instructions backed up at the execute stage. Thus it is recommended to investigate other resource stalls, before digging into the stalls tagged to lack of scheduler entries. The skid tends to be low on this event.

```
%RESOURCE.STALLS.RS.COST =
    100 * RESOURCE_STALLS.RS / CPU_CLK_UNHALTED.THREAD;
```

RESOURCE_STALLS.ROB - Counts cycles when allocation stalls because all the reorder buffer (ROB) entries are taken. This event occurs less frequently than the RESOURCE_STALLS.RS and typically indicates that the pipeline is being backed up by a micro-op that is holding all younger micro-ops from retiring because they have to retire in order.

```
%RESOURCE.STALLS.ROB.COST =
    100 * RESOURCE_STALLS.ROB / CPU_CLK_UNHALTED.THREAD;
```

RESOURCE_STALLS2.BOB_FULL - Counts when allocation is stalled due to a branch micro-op that is ready for allocation, but the number of branches in progress in the processor has reached the limit.

```
%RESOURCE.STALLS.BOB.COST =
    100 * RESOURCE_STALLS2.BOB / CPU_CLK_UNHALTED.THREAD;
```

B.3.4 Memory Sub-System Stalls

The following subsections discuss using specific performance monitoring events in Intel microarchitecture code name Sandy Bridge to identify stalls in the memory sub-systems.

B.3.4.1 Accounting for Load Latency

The breakdown of load operation locality can be accomplished at any granularity including process, module, function and instruction. When you find that a load instruction is a bottleneck, investigate it further with the precise load breakdown. If this does not explain the bottleneck, check for other issues which can impact loads.

You can use these events to estimate the costs of the load causing a bottleneck, and to obtain a percentage breakdown of memory hierarchy level. Not all tools provide support for precise event sampling. If the precise version (event name ends with a suffix PS) of these event is not supported in a given tool, you can use the non-precise

version. The precise load breakdown events listed below have an advantage over the non-precise equivalents because the load latency events rely on a statistical sampling technique in the time domain. Sampling the non-precise events can lead to skewed data.

The precise load events tag the event to the next instruction retired (IP+1). See the load latency at each hierarchy level in section Table 2-6.

Required events

`MEM_LOAD_UOPS_RETIRED.L1_HIT_PS` - Counts demand loads that hit the first level of the data cache, the L1 DCache. Demand loads are non speculative load micro-ops.

`MEM_LOAD_UOPS_RETIRED.L2_HIT_PS` - Counts demand loads that hit the 2nd level cache, the L2.

`MEM_LOAD_UOPS_RETIRED.L3_HIT_PS` - Counts demand loads that hit the 3rd level shared cache, the LLC.

`MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_MISS` - Counts demand loads that hit the 3rd level shared cache and are assumed to be present also in a cache of another core but the cache line was already evicted from there.

`MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_HIT_PS` - Counts demand loads that hit a cache line in a cache of another core and the cache line has not been modified.

`MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_HITM_PS` - Counts demand loads that hit a cache line in the cache of another core and the cache line has been written to by that other core. This event is important for many performance bottlenecks that can occur in multi-threaded applications, such as lock contention and false sharing.

`MEM_LOAD_MISC_RETIRED.LLC_MISS_PS` - Counts demand loads that missed the LLC. This means that the load is satisfied from memory.

`MEM_LOAD_UOPS_RETIRED.HIT_LFB_PS` - Counts demand loads that hit in the line fill buffer (LFB). A LFB entry is allocated every time a miss occurs in the L1 DCache. When a load hits at this location it means that a previous load, store or hardware prefetch has already missed in the L1 DCache and the data fetch is in progress. Therefore the cost of a hit in the LFB varies. This event may count cache-line split loads that miss in the L1 DCache but do not miss the LLC.

On 32-byte Intel AVX loads, all loads that miss in the L1 DCache show up as hits in the L1 DCache or hits in the LFB. They never show hits on any other level of memory hierarchy. Most loads arise from the line fill buffer (LFB) when Intel AVX loads miss in the L1 DCache.

Precise Load Breakdown

The percentage breakdown of each load source can be tagged at any granularity including a single IP, function, module, or process. This is particularly useful at a single instruction to determine the breakdown of where the load was found in the cache hierarchy. The following formula shows how to calculate the percentage of

time a load was satisfied by the LLC. Similar formulas can be used for all other hierarchy levels.

```
%LocL3.HIT =
    100 * MEM_LOAD_UOPS_RETIRED.L3_HIT_PS / $$SumOf_PRECISE_LOADS;
```

```
$$SumOf_PRECISE_LOADS =
    MEM_LOAD_UOPS_RETIRED.HIT_LFB_PS + MEM_LOAD_UOPS_RETIRED.L1_HIT_PS +
    MEM_LOAD_UOPS_RETIRED.L2_HIT_PS + MEM_LOAD_UOPS_RETIRED.L3_HIT_PS +
    MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_MISS +
    MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_HIT_PS +
    MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_HITM_PS +
    MEM_LOAD_UOPS_MISC_RETIRED.LLC_MISS_PS;
```

Estimated Load Penalty

The formulas below help estimating to what degree loads from a certain memory hierarchy are responsible for a slowdown. The CPU_CLK_UNHALTED.THREAD programmable event represents the penalty in cycles tagged at the same granularity. At the instruction level, the cycle cost of an expensive load tends to only skid one IP, similar to the precise event. The calculations below apply to any granularity process, module, function or instruction, since the events are precise. Anything representing 10%, or higher, of the total clocks in a granularity of interest should be investigated.

If the code has highly dependent loads you can use the MEM_LOAD_UOPS_RETIRED.L1_HIT_PS event to determine if the loads are hit by the five cycle latency of the L1 DCache.

Estimated cost of L2 latency

```
%L2.COST =
    12 * MEM_LOAD_UOPS_RETIRED.L2_HIT_PS / CPU_CLK_UNHALTED.THREAD;
```

Estimated cost of L3 hits

```
%L3.COST =
    26 * MEM_LOAD_UOPS_RETIRED.L3_HIT_PS / CPU_CLK_UNHALTED.THREAD;
```

Estimated cost of hits in the cache of other cores

```
%HIT.COST =
    43 * MEM_LOAD_UOPS_LLC_HIT_RETIRE.XSNP_HIT_PS /
    CPU_CLK_UNHALTED.THREAD;
```

Estimated cost of memory latency

```
%MEMORY.COST =
    200 * MEM_LOAD_UOPS_MISC_RETIRED.LLC_MISS_PS /
    CPU_CLK_UNHALTED.THREAD;
```

Actual memory latency can vary greatly depending on memory parameters. The amount of concurrent memory traffic often reduces the effect cost of a given memory

hierarchy. Typically, the estimates above may be on the pessimistic side (like pointer-chasing situations).

Often, cache misses will manifest as delaying and bunching on the retirement of instructions. The precise loads breakdown can provide estimates of the distribution of hierarchy levels where the load is satisfied.

Given a significant impact from a particular cache level, the first step is to find where heavy cache line replacements are occurring in the code. This could coincide with your hot portions of code detected by the memory hierarchy breakdown, but often does not. For instance, regular traversal of a large data structure can unintentionally clear out levels of cache.

If hits of non modified or modified data in another core have high estimated cost and are hot at locations in the code, it can be due to locking, sharing or false sharing issues between threads.

If load latency in memory hierarchy levels further from the L1 DCache does not justify the amount of cycles spent on a load, try one of the following:

- Eliminate unneeded load operations such as spilling general purpose registers to XMM registers rather than memory.
- Continue searching for issues impacting load instructions described in Section B.3.4.4.

B.3.4.2 Cache-line Replacement Analysis

When an application has many cache misses, it is a good idea to determine where cache lines are being replaced at the highest frequency. The instructions responsible for high amount of cache replacements are not always where the application is spending the majority of its time, since replacements can be driven by the hardware prefetchers and store operations.

Typically traversing large arrays or data structures can cause heavy cache line replacements. Do not be surprised if a frequently executed cold function such as `setjmp`, for exception handling, is causing the majority of L1 DCache replacements.

Required events

`L1D.REPLACEMENT` - Replacements in the 1st level data cache.

`L2_LINES_IN.ANY` - Cache lines being brought into the L2 cache.

`OFFCORE_RESPONSE.DATA_IN_SOCKET.LLC_MISS_LOCAL.DRAM_0` - Cache lines being brought into the LLC.

Usages of events

Identifying the replacements that potentially cause performance loss can be done at process, module, and function level. Do it in two steps:

- Use the precise load breakdown to identify the memory hierarchy level at which loads are satisfied and cause the highest penalty.
- Identify, using the formulas below, which portion of code causes the majority of the replacements in the level below the one that satisfies these high penalty loads.

For example, if there is high penalty due to loads hitting the LLC, check the code which is causing replacements in the L2 and the L1. In the formulas below, the nominators are the replacements accounted for a module or function. The sum of the replacements in the denominators is the sum of all replacements in a cache level for the entire process. This enables you to identify the portion of code that causes the majority of the replacements.

L1D Cache Replacements

```
%L1D.REPLACEMENT =
    L1D.REPLACEMENT / SumOverAllProcesses(L1D.REPLACEMENT );
```

L2 Cache Replacements

```
%L2.REPLACEMENT =
    L2_LINES_IN.ANY / SumOverAllProcesses(L2_LINES_IN.ANY );
```

L3 Cache Replacements

```
%L3.REPLACEMENT =
    OFFCORE_RESPONSE.DATA_IN_SOCKET.LLC_MISS_LOCAL.DRAM_0/ SumOverAllProcesses(OFFCORE_RESPONSE.DATA_IN_SOCKET.LLC_MISS_LOCAL.DRAM_0 );
```

B.3.4.3 Lock Contention Analysis

The amount of contention on locks is critical in scalability analysis of multi-threaded applications. A typical ring3 lock almost always results in the execution of an atomic instruction. An atomic instruction is either an XCHG instruction involving a memory address or one of the following instructions with memory destination and lock prefix: ADD, ADC, AND, BTC, BTR, BTS, CMPXCHG, CMPXCH8B, DEC, INC, NEG, NOT, OR, SBB, SUB, XOR or XADD. Precise events enable you to get an idea of the contention on any lock. Many locking APIs start by an atomic instruction in ring3 locks and back off a contended lock by jumping into ring0. This means many locking APIs can be very costly in low contention scenarios. To estimate the amount of contention on a locked instruction, you can measure the number of times the cache line containing the memory destination of an atomic instruction is found modified in another core.

Required events

MEM_UOP_RETIRED.LOCK_LOADS_PS - Counts the number of atomic instructions which are retired with a precise skid of IP+1.

MEM_LOAD_UOPS_LLC_HIT_RETIRED.XSNP_HITM_PS - Counts the occurrences that the load hits a modified cache line in another core. This event is important for

many performance bottlenecks that can occur in multi-core systems, such as lock contention, and false sharing.

Usages of events

The lock contention factor gives the percentage of locked operations executed that contend with another core and therefore have a high penalty. Usually a lock contention factor over 5% is worth investigating on a hot lock. A heavily contended lock may impact the performance of multiple threads.

```
%LOCK.CONTENTION =  
    100 * MEM_LOAD_UOPS_LLC_HIT_RETIRED.XSNP_HITM_PS /  
    MEM_UOP_RETIRED.LOCK_LOAD_PS;
```

B.3.4.4 Other Memory Access Issues

Store Forwarding Blocked

When store forwarding is not possible the dependent loads are blocked. The average penalty for store forward block is 13 cycles. Since many cases of store forwarding blocked were fixed in prior architectures, the most common case in code today involves storing to a smaller memory space than an ensuing larger load.

Required events

LD_BLOCKS.STORE_FORWARD - Counts the number of times a store forward opportunity is blocked due to the inability of the architecture to forward a small store to a larger load and some rare alignment cases.

Usages of Events

Use the following formula to estimate the cost of the store forward block. The event LD_BLOCKS.STORE_FORWARD tends to be tagged to the next IP after the attempted load, so it is recommended to look at this issue at the instruction level. However it is possible to inspect the ratio at any granularity: process, module, function or IP.

```
%STORE.FORWARD.BLOCK.COST =  
    100 * LD_BLOCKS.STORE_FORWARD * 13 / CPU_CLK_UNHALTED.THREAD;
```

After you find a load that is blocked from store forwarding, you need to find the location of the store. Typically, about 60% of all store forwarded blocked issue are caused by stores in the last 10 instructions executed prior to the load.

The most common case where we see store forward blocked is a small store that is unable to forward to a larger load. For example the code below generated writes to a byte pointer address and then reads from a four byte (dword) memory space:

```
and    byte ptr [ebx],7f  
and    dword ptr [ebx],ecx
```

To fix a store forward block it's usually best to fix the store operation and not the load.

Cache Line Splits

Starting from the Intel microarchitecture code name Nehalem, the L1 DCache has split registers which enable it to handle loads and stores that span two cache lines. This puts the cost of split loads at about five cycles, as long as split registers are available, instead of the 20 cycles required in earlier microarchitectures. Handling of split stores handling is usually hidden, but if there are many of them they can stall allocation due to a full store buffer, or they can consume split registers that may be needed for handling split loads. You can still get solid quantifiable gains from eliminating cache line splits.

Required events

MEM_UOP_RETIRED.SPLIT_LOADS_PS - Counts the number of demand loads that span two cache lines. The event is precise.

MEM_UOP_RETIRED.SPLIT_STORES_PS - Counts the number of stores that span two cache lines. The event is precise.

Usages of events

Finding split loads is fairly easy because they usually tag the majority of their cost to the next IP which is executed. The ratio below can be used at any granularity: process, module, function, and IP after split.

%SPLIT.LOAD.COST =

$$100 * \text{MEM_UOP_RETIRED.SPLIT_STORES_PS} * 5 / \text{CPU_CLK_UNHALTED.THREAD};$$

Split store penalty is more difficult to find using an estimated cost, because in typical cases stores do not push out the retirement of instructions. To detect significant amount of split stores divide their number by the total number of stores retired at that IP.

SPLIT.STORE.RATIO =

$$\text{MEM_UOP_RETIRED.SPLIT_STORES_PS} / \text{MEM_UOP_RETIRED.ANY_STORES_PS};$$

4k Aliasing

A 4k aliasing conflict between loads and stores causes a reissue on the load. Five cycles is used as an estimate in the model below.

Required Events

LD_BLOCKS_PARTIAL.ADDRESS_ALIAS - Counts the number of loads that have partial address match with preceding stores, causing the load to be reissued.

Usages of events

%4KALIAS.COST =

$$100 * \text{LD_BLOCK_PARTIAL.ADDRESS_ALIAS} * 5 / \text{CPU_CLK_UNHALTED.THREAD};$$

Load and Store Address Translation

USING PERFORMANCE MONITORING EVENTS

There are two levels of translation look-aside buffer (TLB) for linear to physical address translation. A miss in the DTLB, the first level TLB, that hits in the STLB, the second level TLB, incurs a seven cycle penalty.

Missing in the STLB requires the processor to walk through page table entries that contain the address translation. These walks have variable cost depending on the location of the page table entries. The walk duration is a fairly good estimate of the cost of STLB misses.

Required events

DTLB_LOAD_MISSES.STLB_HIT - Counts loads that miss the DTLB and hit in the STLB. This event has a low skid and hence can be used at the IP level.

DTLB_LOAD_MISSES.WALK_DURATION - Duration of a page walks in cycles following STLB misses. Event skid is typically one instruction, enabling you to detect the issue at instruction, function, module or process granularities.

MEM_UOP_RETIRED.STLB_MISS_LOADS_PS - Precise event for loads which have their translation miss the STLB. The event counts only the first load from a page that initiates the page walk.

Usage of events

Cost of STLB hits on loads:

```
%STLB.HIT.COST =  
    100 * DTLB_LOAD_MISSES.STLB_HIT * 7 / CPU_CLK_UNHALTED.THREAD;
```

Cost of page walks:

```
%STLB.LOAD.MISS.WALK.COST =  
    100 * DTLB_LOAD_MISSES.WALK_DURATION / CPU_CLK_UNHALTED.THREAD;
```

Use the precise STLB miss event at the IP level to determine exactly which instruction and source line suffers from frequent STLB misses.

```
%STLB.LOAD.MISS =  
    100 * MEM_UOP_RETIRED.STLB_MISS_LOADS_PS /  
    MEM_UOP_RETIRED.ANY_LOADS_PS;
```

Large walk durations, of hundreds of cycles, are an indication that the page tables have been thrown out of the LLC. To determine the average cost of a page walk use the following ratio:

```
STLB.LOAD.MISS.AVGCOST =  
    DTLB_LOAD_MISSES.WALK_DURATION /  
    DTLB_LOAD_MISSES.WALK_COMPLETED;
```

To a lesser extent than loads, STLB misses on stores can be a bottleneck. If the store itself is a large bottleneck, cycles will tag to the next IP after the store.

```
%STLB.STORE.MISS =  
    100 * MEM_UOP_RETIRED.STLB_MISS_STORES_PS /  
    MEM_UOP_RETIRED.ANY_STORES_PS;
```

Reducing DTLB/STLB misses increases data locality. One may consider using an commercial-grade memory allocators to improve data locality. Compilers which offer profile guided optimizations may reorder global variables to increase data locality, if the compiler can operate on the whole module. For issues with a large amount of time spent in page walks, server and HPC applications may be able to use large pages for data.

B.3.5 Execution Stalls

The following subsections discusses using specific performance monitoring events in Intel microarchitecture code name Sandy Bridge to identify stalls in the out-of-order engine.

B.3.5.1 Longer Instruction Latencies

Some microarchitectural changes manifested in longer latency for some legacy instructions in existing code. It is possible to detect some of these situations:

- Three-operand slow LEA instructions (see Section 3.5.1.3)
- Flags merge micro-op - These merges are primarily required by "shift cl" instructions (see Section 3.5.2.6).

Required events:

PARTIAL_RAT_STALLS.ANY_CYCLES - Counts the number of cycles that a longer latency micro-op or the PAUSE micro-op was allocated. It is possible to count the cycles for each of the scenarios separately.

Usually, the sub-events tend to have a skid as high as a 10 instructions because the eventing condition is detected early in the pipeline.

Event usage:

To use this event effectively without being distracted by the event skid, you can use it to locate performance issue at the process, module and function granularities, but not at the instruction level. To identify issue at the instruction IP granularity, one can perform static analysis on the functions identified by this event. To estimate the contribution of these events to the code latency, divide them by the cycles at the same granularity. To estimate the overall impact, start with the total cycles due to these issues and if significant continue to search for the exact reason using the sub events.

Total cycles spent in the specified scenarios:

```
%PARTIAL.STALL.COST =
    100 * PARTIAL_RAT_STALLS.ANY_CYCLES / CPU_CLK_UNHALTED.THREAD;
```

Flags Merge micro-op ratio:

```
%FLAGS.MERGE.UOP =
    100 * PARTIAL_RAT_STALLS.FLAGS_MERGE_UOP_CYCLES /
```

```
CPU_CLK_UNHALTED.THREAD;
```

Slow LEA instructions allocated:

```
%SLOW.LEA.WINDOW =  
    100 * PARTIAL_RAT_STALLS.FLAGS_MERGE_UOP_CYCLES /  
    CPU_CLK_UNHALTED.THREAD;
```

B.3.5.2 Assists

Assists usually involve the microcode sequencer that helps handle the assist. Determining the number of cycles where microcode is generated from the microcode sequencer is often a good methodology to determine the total cost of the assist. If the overall cost of assists are high, a breakdown of assists into specific types will be useful.

Estimating the total cost of assists using microcode sequencer cycles:

```
%ASSISTS.COST =  
    100 * IDQ.MS_CYCLES / CPU_CLK_UNHALTED.THREAD;
```

Floating point assists:

Denormal inputs for X87 instructions require an FP assist, potentially costing hundreds of cycles.

```
%FP.ASSISTS =  
    100 * FP_ASSISTS.ANY / INST_RETIRED.ANY;
```

Transitions between Intel SSE and Intel AVX

The transitions between SSE and AVX code are explained in detail in Section 11.3.1. The typical cost is about 75 cycles.

```
%AVX2SSE.TRANSITION.COST =  
    100 * OTHER_ASSISTS.AVX_TO_SSE / CPU_CLK_UNHALTED.THREAD;
```

```
%SSE2AVX.TRANSITION.COST =  
    100 * OTHER_ASSISTS.SSE_TO_AVX / CPU_CLK_UNHALTED.THREAD;
```

AVX Store Assists

32-byte AVX store instructions that span two pages require an assist that costs roughly 150 cycles. A large amount of microcode tagged to the IP after a 32-byte AVX store is a good sign that an assist has occurred.

```
%AVX.STORE.ASSIST.COST =  
    100 * OTHER_ASSISTS.AVX_STORE / CPU_CLK_UNHALTED.THREAD;
```

B.3.6 Bad Speculation

This section discusses mispredicted branch instructions resulting in a pipeline flush.

B.3.6.1 Branch Mispredicts

The largest challenge with mispredicted branches is finding the branch which caused them. Branch mispredictions incur penalty of about 20 cycles. The cost varies based upon the misprediction, and whether the correct path is found in the Decoded ICache or in the legacy decode pipeline.

Required Events

BR_MISP_RETIRED.ALL_BRANCHES_PS is a precise event that counts branches that incorrectly predicted the branch target. Since this is a precise event that skids to the next instruction, it tags to the first instruction in the correct path after the branch misprediction. This study can be performed at the process, module, function or instruction granularity.

Usages of Events

Use the following ratio to estimate the cost of mispredicted branches:

%BR.MISP.COST =

$$20 * BR_MISP_RETIRED.ALL_BRANCHES_PS / CPU_CLK_UNHALTED.THREAD;$$

B.3.7 Front End Stalls

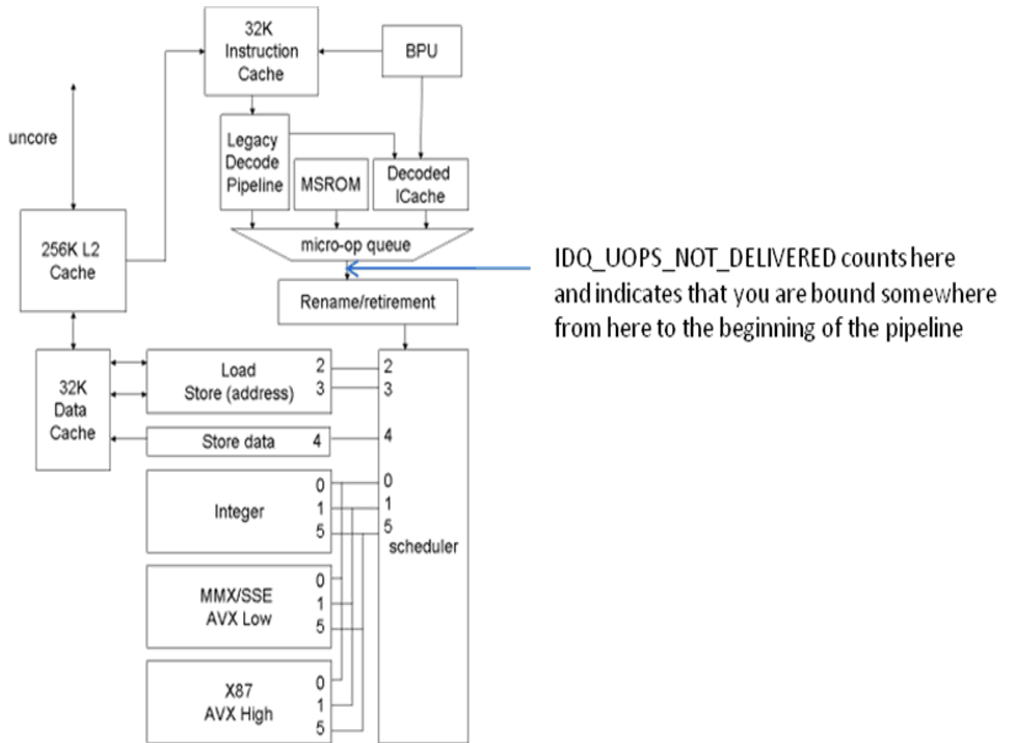
Stalls in the front end should not be investigated unless the analysis in Section B.3.2 showed at least 30% of a granularity being bound in the front end. This section explains the main issues that can cause delays in the front end of the pipeline. Events detected in the front end have unpredictable skid. Therefore do not try and associate the penalty at the IP level. Stay at the function, module, and process level for these events.

B.3.7.1 Understanding the Micro-op Delivery Rate

Usages of Counters

The event IDQ_UOPS_NOT_DELIVERED counts when the maximum of four micro-ops are not delivered to the rename stage, while it is requesting micro-ops. When the pipeline is backed up the rename stage does not request any further micro-ops from the front end. The diagram above shows how this event tracks micro-ops between the micro-op queue and the rename stage.

USING PERFORMANCE MONITORING EVENTS



You can use the IDQ_UOPS_NOT_DELIVERED event to breakdown the distribution of cycles when 0, 1, 2, 3 micro-ops are delivered from the front end.

Percentage of cycles the front end is effective, or execution is back end bound:

%FE.DELIVERING =

$$100 * (\text{CPU_CLK_UNHALTED.THREAD} - \text{IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_3_UOP_DELIV.CORE}) / \text{CPU_CLK_UNHALTED.THREAD};$$

Percentage of cycles the front end is delivering three micro-ops per cycle:

%FE.DELIVER.3UOPS =

$$100 * (\text{IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_3_UOP_DELIV.CORE} - \text{IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_2_UOP_DELIV.CORE}) / \text{CPU_CLK_UNHALTED.THREAD};$$

Percentage of cycles the front end is delivering two micro-ops per cycle:

%FE.DELIVER.2UOPS =

$$100 * (\text{IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_2_UOP_DELIV.CORE} - \text{IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_1_UOP_DELIV.CORE}) /$$

```
CPU_CLK_UNHALTED.THREAD;
```

Percentage of cycles the front end is delivering one micro-ops per cycle:

```
%FE.DELIVER.1UOPS =
    100 * ( IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_1_UOP_DELIV.CORE -
            IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_0_UOP_DELIV.CORE ) /
    CPU_CLK_UNHALTED.THREAD;
```

Percentage of cycles the front end is delivering zero micro-ops per cycle:

```
%FE.DELIVER.1UOPS =
    100 * ( IDQ_UOPS_NOT_DELIVERED.CYCLES_LE_0_UOP_DELIV.CORE ) /
    CPU_CLK_UNHALTED.THREAD;
```

Average Micro-ops Delivered per Cycle: This ratio assumes that the front end could potentially deliver four micro-ops per cycle when bound in the back end.

```
AVG.uops.per.cycle =
    (4 * (%FE.DELIVERING) + 3 * (%FE.DELIVER.3UOPS) + 2 * (%FE.DELIVER.2UOPS) +
     (%FE.DELIVER.1UOPS)) / 100
```

Seeing the distribution of the micro-ops being delivered in a cycle is a hint at the front end bottlenecks that might be occurring. Issues such as LCPs and penalties from switching from the decoded ICache to the legacy decode pipeline tend to result in zero micro-ops being delivered for several cycles. Fetch bandwidth issues and decoder stalls result in less than four micro-ops delivered per cycle.

B.3.7.2 Understanding the Sources of the Micro-op Queue

The micro-op queue can get micro-ops from the following sources:

- Decoded ICache
- Legacy decode pipeline
- Microcode sequencer (MS)

A typical distribution is approximately 80% of the micro-ops coming from the Decoded ICache, 15% coming from legacy decode pipeline and 5% coming from the microcode sequencer. Excessive micro-ops coming from the legacy decode pipeline can be a warning sign that the Decoded ICache is not working effectively. A large portion of micro-ops coming from the microcode sequencer may be benign, such as complex instructions, or string operations, but can also be due to code assists handling undesired situations like Intel SSE to Intel AVX code transitions.

Description of Counters Required

IDQ.DSB_UOPS - Micro-ops delivered to the micro-op queue from the Decoded ICache.

IDQ.MITE_UOPS - Micro-ops delivered to the micro-op queue from the legacy decode pipeline.

IDQ.MS_UOPS -Micro-ops delivered from the microcode sequencer.

Usage of Counters

Percentage of micro-ops coming from Decoded ICache:

$$\%UOPS.DSB = \frac{IDQ.DSB_UOPS}{ALL_IDQ_UOPS};$$

Percentage of micro-ops coming from legacy decoder pipeline:

$$\%UOPS.MITE = \frac{IDQ.MITE_UOPS}{ALL_IDQ_UOPS};$$

Percentage of micro-ops coming from micro-sequencer:

$$\%UOPS.MS = \frac{IDQ.MS_UOPS}{ALL_IDQ_UOPS};$$

$$ALL_IDQ_UOPS = (IDQ.DSB_UOPS + IDQ.MITE_UOPS + IDQ.MS_UOPS);$$

If your application is not bound in the front end then whether micro-ops are coming from the legacy decode pipeline or Decoded ICache is of lesser importance. Excessive micro-ops coming from the microcode sequencer are worth investigating further to see if assists might be a problem.

Cases to investigate are listed below:

- ($\%FE_BOUND > 30\%$) and ($\%UOPS.DSB < 70\%$)
We use a threshold of 30% to define a "front end bound" case. This threshold may be applicable to many situations, but may also vary somewhat across different workloads.
 - Investigate why micro-ops are not coming from the Decoded ICache
 - Investigate issues which can impact the legacy decode pipeline
- ($\%FE_BOUND > 30\%$) and ($\%UOP_DSB > 70\%$)
 - Investigate switches from Decoded ICache to legacy decode pipeline since it may be switching to run portions of code that are too small to be effective
 - Look at the amount of bad speculation, since branch mispredictions still impact FE performance
 - Determine the average number of micro-ops being delivered per 32-byte chunk hit. If there are many taken branches from one 32-byte chunk into another, it impacts the micro-ops being delivered per cycle.
 - Micro-op delivery from the Decoded ICache may be an issue which is not covered
- ($\%FE_BOUND < 20\%$) and ($\%UOPS_MS > 25\%$)
We use a threshold of 20% to define a "front end not bound" case. This threshold may be applicable to many situations, but may also vary somewhat across different workloads.

The following steps can help determine why micro-ops came from the microcode, in order of most common to least common.

- Long latency instructions - Any instruction over four micro-ops starts the microcode sequencer. Some instructions such as transcendentals can generate many micro-ops from the microcode.
- String operations - string operations can produce a large amount of microcode. In some cases there are assists which can occur due to string operations such as REP MOVSB with trip count greater than 3, which costs 70+ cycles.
- Assists - See Section B.3.5.2.

B.3.7.3 The Decoded ICache

The Decoded ICache has many advantages over the legacy decode pipeline. It eliminates many bottlenecks of the legacy decode pipeline such as instructions decoded into more than one micro-op and length changing prefix (LCP) stalls.

A switch to the legacy decode pipeline from the Decoded ICache only occurs when a lookup in the Decoded ICache fails and usually costs anywhere from zero to three cycles in the front end of the pipeline.

Required events

The Decoded ICache events all have large skids and the exact instruction where they are tagged is usually not the source of the problem so only look for this issue at the process, module and function granularities.

DSB2MITE_SWITCHES.PENALTY_CYCLES - Counts the cycles attributed to the switch from the Decoded ICache to the legacy decode pipeline, excluding cycles when the micro-op queue cannot accept micro-ops because it is back end bound.

DSB2MITE_SWITCHES.COUNT - Counts the number of switches between the Decoded ICache and the legacy decode pipeline.

DSB_FILL.ANY_CANCEL - Counts when fills to the Decoded ICache are canceled.

DSB_FILL.EXCEED_WAY_LIMIT - Counts when a fill is canceled because the allocated lines for Decoded ICache has exceeded three for the 32-byte chunk.

Usage of Events

Since these studies involve front end events, do not try to tag the event to a specific instruction.

Determining cost of switches from the Decoded ICache to the legacy decode pipeline.

```
%DSB2MITE_SWITCHES.COST =
    100 * DSB2MITE_SWITCHES.PENALTY_CYCLES / CPU_CLK_UNHALTED.THREAD;
```

Determining the average cost per Decoded ICache switch to the legacy front end:

```
AVG.DSB2MITE_SWITCHES.COST =
    DSB2MITE_SWITCHES.PENALTY_CYCLES / DSB2MITE_SWITCHES.COUNT;
```

Determining causes of misses in the decoded ICache

There are no partial hits in the Decoded ICache. If any micro-op that is part of that lookup on the 32-byte chunk is missing, a Decoded ICache miss occurs on all micro-ops for that transaction.

There are three primary reasons for missing micro-ops in the Decoded ICache:

- Portions of a 32-byte chunk of code were not able to fit within three ways of the Decoded ICache.
- A frequently run portion of your code section is too large for the Decoded ICache. This case is more common on server applications since client applications tend to have a smaller set of code which is "hot".
- The Decoded ICache is getting flushed for example when an ITLB entry is evicted.

To determine if a portion of the 32-byte code is unable to fit into three lines within the Decoded ICache use the DSB_FILL.EXCEED_WAY_LIMIT event at the process, module or function granularities

```
%DSB.EXCEED.WAY.LIMIT =  
    100 * DSB_FILL.EXCEED_WAY_LIMIT / DSB_FILL.ANY_CANCEL;
```

B.3.7.4 Issues in the Legacy Decode Pipeline

If a large percentage of the micro-ops going to the micro-op queue are being delivered from the legacy decode pipeline, you should check to see if there are bottlenecks impacting that stage. The most common bottlenecks in the legacy decode pipeline are:

- Fetch not providing enough instructions
This happens when hot code is poorly aligned. For example if the hot code being fetched to be run is on the 15th byte, then only one byte is fetched.
- Length changing prefix stalls in the instruction length decoder
Instructions that are decoded into two to four micro-ops may introduce a bubble in the decoder throughput. If the instruction queue, preceding the decoders, becomes full, this indicates that these instructions may cause a penalty.

```
%ILD.STALL.COST =  
    100 * ILD_STALL.LCP * 3 / CPU_CLK_UNHALTED.THREAD;;
```

- Determine the percentage of cycles which the instruction queue is full.

```
%IQ.FULL =  
    100 * ILD_STALL.IQ_FULL / CPU_CLK_UNHALTED.THREAD;;
```

- Determine if there are instructions which are decoded into two to four micro-ops.

```
%DECODERO =  
    100 * UOPS_DECODED.DECO_UOPS / IDQ.MITE_UOPS;;
```

B.3.7.5 Instruction Cache

Applications with large hot code sections tend to run into many issues with the instruction cache and ITLB. This is more typical in server applications.

Required events

ICACHE.MISSES - Counts the number of instruction byte fetches that miss the ICache

Usage of events

To determine whether ICache misses are causing the issue, compare them to the instructions retired event count, using the same granularity (process, model, or function). Anything over 1% of instructions retired can be a significant issue.

ICACHE.PER.INST.RET =

$$\text{ICACHE.MISSES} / \text{INST_RETIRED.ANY};$$

If ICache misses are causing a significant problem, try to reduce the size of your hot code section, using the profile guided optimizations. Most compilers have options for text reordering which helps reduce the number of pages and, to a lesser extent, the number of pages your application is covering.

If the application makes significant use of macros, try to either convert them to functions, or use intelligent linking to eliminate repeated code.

B.4 USING PERFORMANCE EVENTS OF INTEL® CORE™ SOLO AND INTEL® CORE™ DUO PROCESSORS

There are performance events specific to the microarchitecture of Intel Core Solo and Intel Core Duo processors. See also: Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

B.4.1 Understanding the Results in a Performance Counter

Each performance event detects a well-defined microarchitectural condition occurring in the core while the core is active. A core is active when:

- It's running code (excluding the halt instruction).
- It's being snooped by the other core or a logical processor on the platform. This can also happen when the core is halted.

Some microarchitectural conditions are applicable to a sub-system shared by more than one core and some performance events provide an event mask (or unit mask) that allows qualification at the physical processor boundary or at bus agent boundary.

Some events allow qualifications that permit the counting of microarchitectural conditions associated with a particular core versus counts from all cores in a physical

processor (see L2 and bus related events in Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*).

When a multi-threaded workload does not use all cores continuously, a performance counter counting a core-specific condition may progress to some extent on the halted core and stop progressing or a unit mask may be qualified to continue counting occurrences of the condition attributed to either processor core. Typically, one can adjust the highest two bits (bits 15:14 of the IA32_PERFVTSELx MSR) in the unit mask field to distinguish such asymmetry (See Chapter 18, "Debugging and Performance Monitoring," of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*).

There are three cycle-counting events which will not progress on a halted core, even if the halted core is being snooped. These are: Unhalted core cycles, Unhalted reference cycles, and Unhalted bus cycles. All three events are detected for the unit selected by event 3CH.

Some events detect microarchitectural conditions but are limited in their ability to identify the originating core or physical processor. For example, bus_drdy_clocks may be programmed with a unit mask of 20H to include all agents on a bus. In this case, the performance counter in each core will report nearly identical values. Performance tools interpreting counts must take into account that it is only necessary to equate bus activity with the event count from one core (and not use not the sum from each core).

The above is also applicable when the core-specificity sub field (bits 15:14 of IA32_PERFVTSELx MSR) within an event mask is programmed with 11B. The result of reported by performance counter on each core will be nearly identical.

B.4.2 Ratio Interpretation

Ratios of two events are useful for analyzing various characteristics of a workload. It may be possible to acquire such ratios at multiple granularities, for example: (1) per-application thread, (2) per logical processor, (3) per core, and (4) per physical processor.

The first ratio is most useful from a software development perspective, but requires multi-threaded applications to manage processor affinity explicitly for each application thread. The other options provide insights on hardware utilization.

In general, collect measurements (for all events in a ratio) in the same run. This should be done because:

- If measuring ratios for a multi-threaded workload, getting results for all events in the same run enables you to understand which event counter values belongs to each thread.
- Some events, such as writebacks, may have non-deterministic behavior for different runs. In such a case, only measurements collected in the same run yield meaningful ratio values.

B.4.3 Notes on Selected Events

This section provides event-specific notes for interpreting performance events listed in Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.

- **L2_Reject_Cycles, event number 30H** — This event counts the cycles during which the L2 cache rejected new access requests.
- **L2_No_Request_Cycles, event number 32H** — This event counts cycles during which no requests from the L1 or prefetches to the L2 cache were issued.
- **Unhalted_Core_Cycles, event number 3C, unit mask 00H** — This event counts the smallest unit of time recognized by an active core.

In many operating systems, the idle task is implemented using HLT instruction. In such operating systems, clock ticks for the idle task are not counted. A transition due to Enhanced Intel SpeedStep Technology may change the operating frequency of a core. Therefore, using this event to initiate time-based sampling can create artifacts.

- **Unhalted_Ref_Cycles, event number 3C, unit mask 01H** — This event guarantees a uniform interval for each cycle being counted. Specifically, counts increment at bus clock cycles while the core is active. The cycles can be converted to core clock domain by multiplying the bus ratio which sets the core clock frequency.
- **Serial_Execution_Cycles, event number 3C, unit mask 02H** — This event counts the bus cycles during which the core is actively executing code (non-halted) while the other core in the physical processor is halted.
- **L1_Pref_Req, event number 4FH, unit mask 00H** — This event counts the number of times the Data Cache Unit (DCU) requests to prefetch a data cache line from the L2 cache. Requests can be rejected when the L2 cache is busy. Rejected requests are re-submitted.
- **DCU_Snoop_to_Share, event number 78H, unit mask 01H** — This event counts the number of times the DCU is snooped for a cache line needed by the other core. The cache line is missing in the L1 instruction cache or data cache of the other core; or it is set for read-only, when the other core wants to write to it. These snoops are done through the DCU store port. Frequent DCU snoops may conflict with stores to the DCU, and this may increase store latency and impact performance.
- **Bus_Not_In_Use, event number 7DH, unit mask 00H** — This event counts the number of bus cycles for which the core does not have a transaction waiting for completion on the bus.
- **Bus_Snoops, event number 77H, unit mask 00H** — This event counts the number of CLEAN, HIT, or HITM responses to external snoops detected on the bus.

In a single-processor system, CLEAN and HIT responses are not likely to happen. In a multiprocessor system this event indicates an L2 miss in one processor that did not find the missed data on other processors.

In a single-processor system, an HITM response indicates that an L1 miss (instruction or data) found the missed cache line in the other core in a modified state. In a multiprocessor system, this event also indicates that an L1 miss (instruction or data) found the missed cache line in another core in a modified state.

B.5 DRILL-DOWN TECHNIQUES FOR PERFORMANCE ANALYSIS

Software performance intertwines code and microarchitectural characteristics of the processor. Performance monitoring events provide insights to these interactions. Each microarchitecture often provides a large set of performance events that target different sub-systems within the microarchitecture. Having a methodical approach to select key performance events will likely improve a programmer's understanding of the performance bottlenecks and improve the efficiency of code-tuning effort.

Recent generations of Intel 64 and IA-32 processors feature microarchitectures using an out-of-order execution engine. They are also accompanied by an in-order front end and retirement logic that enforces program order. Superscalar hardware, buffering and speculative execution often complicates the interpretation of performance events and software-visible performance bottlenecks.

This section discusses a methodology of using performance events to drill down on likely areas of performance bottleneck. By narrowed down to a small set of performance events, the programmer can take advantage of Intel VTune Performance Analyzer to correlate performance bottlenecks with source code locations and apply coding recommendations discussed in Chapter 3 through Chapter 8. Although the general principles of our method can be applied to different microarchitectures, this section will use performance events available in processors based on Intel Core microarchitecture for simplicity.

Performance tuning usually centers around reducing the time it takes to complete a well-defined workload. Performance events can be used to measure the elapsed time between the start and end of a workload. Thus, reducing elapsed time of completing a workload is equivalent to reducing measured processor cycles.

The drill-down methodology can be summarized as four phases of performance event measurements to help characterize interactions of the code with key pipe stages or sub-systems of the microarchitecture. The relation of the performance event drill-down methodology to the software tuning feedback loop is illustrated in Figure B-13.

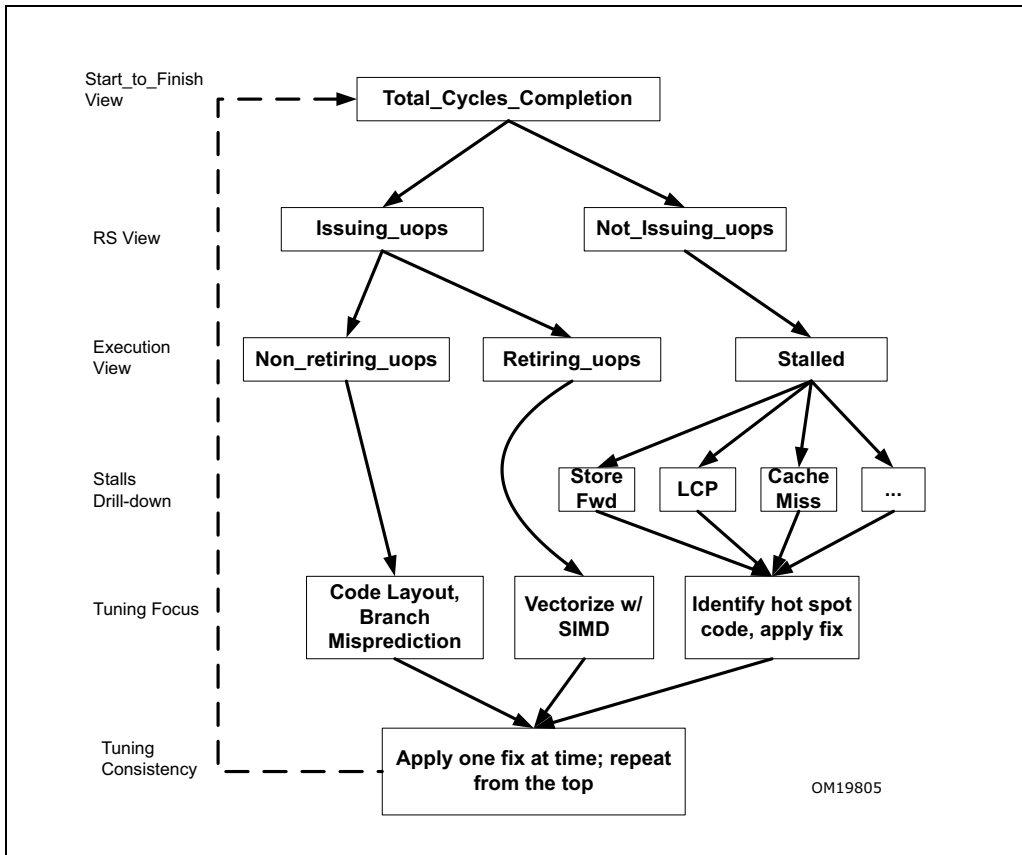


Figure B-13. Performance Events Drill-Down and Software Tuning Feedback Loop

Typically, the logic in performance monitoring hardware measures microarchitectural conditions that varies across different counting domains, ranging from cycles, micro-ops, address references, instances, etc. The drill-down methodology attempts to provide an intuitive, cycle-based view across different phases by making suitable approximations that are described below:

- **Total cycle measurement** — This is the start to finish view of total number of cycle to complete the workload of interest. In typical performance tuning situations, the metric Total_cycles can be measured by the event CPU_CLK_UNHALTED.CORE. See Appendix A, "Performance Monitoring Events," of *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*.
- **Cycle composition at issue port** — The reservation station (RS) dispatches micro-ops for execution so that the program can make forward progress. Hence the metric Total_cycles can be decomposed as consisting of two exclusive components: Cycles_not_issuing_uops representing cycles that the RS is not

issuing micro-ops for execution, and `Cycles_issuing_uops` cycles that the RS is issuing micro-ops for execution. The latter component includes μ ops in the architected code path or in the speculative code path.

- **Cycle composition of OOO execution** — The out-of-order engine provides multiple execution units that can execute μ ops in parallel. If one execution unit stalls, it does not necessarily imply the program execution is stalled. Our methodology attempts to construct a cycle-composition view that approximates the progress of program execution. The three relevant metrics are: `Cycles_stalled`, `Cycles_not_retiring_uops`, and `Cycles_retiring_uops`.
- **Execution stall analysis** — From the cycle compositions of overall program execution, the programmer can narrow down the selection of performance events to further pin-point unproductive interaction between the workload and a micro-architectural sub-system.

When cycles lost to a stalled microarchitectural sub-system, or to unproductive speculative execution are identified, the programmer can use VTune Analyzer to correlate each significant performance impact to source code location. If the performance impact of stalls or misprediction is insignificant, VTune can also identify the source locations of hot functions, so the programmer can evaluate the benefits of vectorization on those hot functions.

B.5.1 Cycle Composition at Issue Port

Recent processor microarchitectures employ out-of-order engines that execute streams of μ ops natively, while decoding program instructions into μ ops in its front end. The metric `Total_cycles` alone, is opaque with respect to decomposing cycles that are productive or non-productive for program execution. To establish a consistent cycle-based decomposition, we construct two metrics that can be measured using performance events available in processors based on Intel Core microarchitecture. These are:

- **Cycles_not_issuing_uops** — This can be measured by the event `RS_UOPS_DISPATCHED`, setting the INV bit and specifying a counter mask (CMASK) value of 1 in the target performance event select (IA32_PERFVSELx) MSR (See Chapter 18 of the *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3B*). In VTune Analyzer, the special values for CMASK and INV is already configured for the VTune event name `RS_UOPS_DISPATCHED.CYCLES_NONE`.
- **Cycles_issuing_uops** — This can be measured using the event `RS_UOPS_DISPATCHED`, clear the INV bit and specifying a counter mask (CMASK) value of 1 in the target performance event select MSR

Note the cycle decomposition view here is approximate in nature; it does not distinguish specificities, such as whether the RS is full or empty, transient situations of RS being empty but some in-flight uops is getting retired.

B.5.2 Cycle Composition of OOO Execution

In an OOO engine, speculative execution is an important part of making forward progress of the program. But speculative execution of μ ops in the shadow of mispredicted code path represent un-productive work that consumes execution resources and execution bandwidth.

`Cycles_not_issuing_uops`, by definition, represents the cycles that the OOO engine is stalled (`Cycles_stalled`). As an approximation, this can be interpreted as the cycles that the program is not making forward progress.

The μ ops that are issued for execution do not necessarily end in retirement. Those μ ops that do not reach retirement do not help forward progress of program execution. Hence, a further approximation is made in the formalism of decomposition of `Cycles_issuing_uops` into:

- **Cycles_non_retiring_uops** — Although there isn't a direct event to measure the cycles associated with non-retiring μ ops, we will derive this metric from available performance events, and several assumptions:
 - A constant issue rate of μ ops flowing through the issue port. Thus, we define: `uops_rate` = "`Dispatch_uops/Cycles_issuing_uops`, where `Dispatch_uops` can be measured with `RS_UOPS_DISPATCHED`, clearing the `INV` bit and the `CMASK`.
 - We approximate the number of non-productive, non-retiring μ ops by [`non_productive_uops = Dispatch_uops - executed_retired_uops`], where `executed_retired_uops` represent productive μ ops contributing towards forward progress that consumed execution bandwidth.
 - The `executed_retired_uops` can be approximated by the sum of two contributions: `num_retired_uops` (measured by the event `UOPS_RETIRED.ANY`) and `num_fused_uops` (measured by the event `UOPS_RETIRED.FUSED`).

Thus, `Cycles_non_retiring_uops = non_productive_uops / uops_rate`.

- **Cycles_retiring_uops** — This can be derived from `Cycles_retiring_uops = num_retired_uops / uops_rate`.

The cycle-decomposition methodology here does not distinguish situations where productive uops and non-productive μ ops may be dispatched in the same cycle into the OOO engine. This approximation may be reasonable because heuristically high contribution of non-retiring uops likely correlates to situations of congestions in the OOO engine and subsequently cause the program to stall.

Evaluations of these three components: `Cycles_non_retiring_uops`, `Cycles_stalled`, `Cycles_retiring_uops`, relative to the `Total_cycles`, can help steer tuning effort in the following directions:

- If the contribution from `Cycles_non_retiring_uops` is high, focusing on code layout and reducing branch mispredictions will be important.
- If both the contributions from `Cycles_non_retiring_uops` and `Cycles_stalled` are insignificant, the focus for performance tuning should be directed to vectorization or other techniques to improve retirement throughput of hot functions.

- If the contributions from `Cycles_stalled` is high, additional drill-down may be necessary to locate bottlenecks that lies deeper in the microarchitecture pipeline.

B.5.3 Drill-Down on Performance Stalls

In some situations, it may be useful to evaluate cycles lost to stalls associated with various stress points in the microarchitecture and sum up the contributions from each candidate stress points. This approach implies a very gross simplification and introduce complications that may be difficult to reconcile with the superscalar nature and buffering in an OOO engine.

Due to the variations of counting domains associated with different performance events, cycle-based estimation of performance impact at each stress point may carry different degree of errors due to over-estimation of exposures or under-estimations.

Over-estimation is likely to occur when overall performance impact for a given cause is estimated by multiplying the per-instance-cost to an event count that measures the number of occurrences of that microarchitectural condition. Consequently, the sum of multiple contributions of lost cycles due to different stress points may exceed the more accurate metric `Cycles_stalled`.

However an approach that sums up lost cycles associated with individual stress point may still be beneficial as an iterative indicator to measure the effectiveness of code tuning loop effort when tuning code to fix the performance impact of each stress point. The remaining of this sub-section will discuss a few common causes of performance bottlenecks that can be counted by performance events and fixed by following coding recommendations described in this manual.

The following items discuss several common stress points of the microarchitecture:

- **L2 Miss Impact** — An L2 load miss may expose the full latency of memory subsystem. The latency of accessing system memory varies with different chipset, generally on the order of more than a hundred cycles. Server chipset tend to exhibit longer latency than desktop chipsets. The number L2 cache miss references can be measured by `MEM_LOAD_RETIRED.L2_LINE_MISS`.

An estimation of overall L2 miss impact by multiplying system memory latency with the number of L2 misses ignores the OOO engine's ability to handle multiple outstanding load misses. Multiplication of latency and number of L2 misses imply each L2 miss occur serially.

To improve the accuracy of estimating L2 miss impact, an alternative technique should also be considered, using the event `BUS_REQUEST_OUTSTANDING` with a `CMASK` value of 1. This alternative technique effectively measures the cycles that the OOO engine is waiting for data from the outstanding bus read requests. It can overcome the over-estimation of multiplying memory latency with the number of L2 misses.

- **L2 Hit Impact** — Memory accesses from L2 will incur the cost of L2 latency (See Table 2-3). The number cache line references of L2 hit can be measured by the

difference between two events: MEM_LOAD_RETIRED.L1D_LINE_MISS - MEM_LOAD_RETIRED.L2_LINE_MISS.

An estimation of overall L2 hit impact by multiplying the L2 hit latency with the number of L2 hit references ignores the OOO engine's ability to handle multiple outstanding load misses.

- **L1 DTLB Miss Impact** — The cost of a DTLB lookup miss is about 10 cycles. The event MEM_LOAD_RETIRED.DTLB_MISS measures the number of load micro-ops that experienced a DTLB miss.
- **LCP Impact** — The overall impact of LCP stalls can be directly measured by the event ILD_STALLS. The event ILD_STALLS measures the number of times the slow decoder was triggered, the cost of each instance is 6 cycles
- **Store forwarding stall Impact** — When a store forwarding situation does not meet address or size requirements imposed by hardware, a stall occurs. The delay varies for different store forwarding stall situations. Consequently, there are several performance events that provide fine-grain specificity to detect different store-forwarding stall conditions. These include:
 - A load blocked by preceding store to unknown address: This situation can be measure by the event Load_Blocks.Sta. The per-instance cost is about 5 cycles.
 - Load partially overlaps with proceeding store or 4-KByte aliased address between a load and a proceeding store: these two situations can be measured by the event Load_Blocks.Overlap_store.
 - A load spanning across cache line boundary: This can be measured by Load_Blocks.Until_Retire. The per-instance cost is about 20 cycles.

B.6 EVENT RATIOS FOR INTEL CORE MICROARCHITECTURE

Appendix B.6 provides examples of using performance events to quickly diagnose performance bottlenecks. This section provides additional information on using performance events to evaluate metrics that can help in wide range of performance analysis, workload characterization, and performance tuning. Note that many performance event names in the Intel Core microarchitecture carry the format of XXXX.YYY. this notation derives from the general convention that XXXX typically corresponds to a unique event select code in the performance event select register (IA32_PERFVSELx), while YYY corresponds to a unique sub-event mask that uniquely defines a specific microarchitectural condition (See Chapter 18 and Appendix A of the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3B*).

B.6.1 Clocks Per Instructions Retired Ratio (CPI)

1. Clocks Per Instruction Retired Ratio (CPI): $\text{CPU_CLK_UNHALTED.CORE} / \text{INST_RETIRED.ANY}$.

The Intel Core microarchitecture is capable of reaching CPI as low as 0.25 in ideal situations. But most of the code has higher CPI. The greater value of CPI for a given workload indicate it has more opportunity for code tuning to improve performance. The CPI is an overall metric, it does not provide specificity of what microarchitectural sub-system may be contributing to a high CPI value.

The following subsections defines a list of event ratios that are useful to characterize interactions with the front end, execution, and memory.

B.6.2 Front-end Ratios

2. RS Full Ratio: $\text{RESOURCE_STALLS.RS_FULL} / \text{CPU_CLK_UNHALTED.CORE} * 100$
3. ROB Full Ratio: $\text{RESOURCE_STALLS.ROB_FULL} / \text{CPU_CLK_UNHALTED.CORE} * 100$
4. Load or Store Buffer Full Ratio: $\text{RESOURCE_STALLS.LD_ST} / \text{CPU_CLK_UNHALTED.CORE} * 100$

When there is a low value for the ROB Full Ratio, RS Full Ratio, and Load Store Buffer Full Ratio, and high CPI it is likely that the front end cannot provide instructions and micro-ops at a rate high enough to fill the buffers in the out-of-order engine, and therefore it is starved waiting for micro-ops to execute. In this case check further for other front end performance issues.

B.6.2.1 Code Locality

5. Instruction Fetch Stall: $\text{CYCLES_L1I_MEM_STALLED} / \text{CPU_CLK_UNHALTED.CORE} * 100$

The Instruction Fetch Stall ratio is the percentage of cycles during which the Instruction Fetch Unit (IFU) cannot provide cache lines for decoding due to cache and Instruction TLB (ITLB) misses. A high value for this ratio indicates potential opportunities to improve performance by reducing the working set size of code pages and instructions being executed, hence improving code locality.

6. ITLB Miss Rate: $\text{ITLB_MISS_RETIRED} / \text{INST_RETIRED.ANY}$

A high ITLB Miss Rate indicates that the executed code is spread over too many pages and cause many Instructions TLB misses. Retired ITLB misses cause the pipeline to naturally drain, while the miss stalls fetching of more instructions.

7. L1 Instruction Cache Miss Rate: $\text{L1I_MISSES} / \text{INST_RETIRED.ANY}$

A high value for L1 Instruction Cache Miss Rate indicates that the code working set is bigger than the L1 instruction cache. Reducing the code working set may improve performance.

8. L2 Instruction Cache Line Miss Rate: $L2_IFETCH.SELF.I_STATE / INST_RETIRED.ANY$

L2 Instruction Cache Line Miss Rate higher than zero indicates instruction cache line misses from the L2 cache may have a noticeable performance impact of program performance.

B.6.2.2 Branching and Front-end

9. BACLEAR Performance Impact: $7 * BACLEARS / CPU_CLK_UNHALTED.CORE$

A high value for BACLEAR Performance Impact ratio usually indicates that the code has many branches such that they cannot be consumed by the Branch Prediction Unit.

10. Taken Branch Bubble: $(BR_TKN_BUBBLE_1 + BR_TKN_BUBBLE_2) / CPU_CLK_UNHALTED.CORE$

A high value for Taken Branch Bubble ratio indicates that the code contains many taken branches coming one after the other and cause bubbles in the front-end. This may affect performance only if it is not covered by execution latencies and stalls later in the pipe.

B.6.2.3 Stack Pointer Tracker

11. ESP Synchronization: $ESP.SYNCH / ESP.ADDITIONS$

The ESP Synchronization ratio calculates the ratio of ESP explicit use (for example by load or store instruction) and implicit uses (for example by PUSH or POP instruction). The expected ratio value is 0.2 or lower. If the ratio is higher, consider rearranging your code to avoid ESP synchronization events.

B.6.2.4 Macro-fusion

12. Macro-Fusion: $UOPS_RETIRED.MACRO_FUSION / INST_RETIRED.ANY$

The Macro-Fusion ratio calculates how many of the retired instructions were fused to a single micro-op. You may find this ratio is high for a 32-bit binary executable but significantly lower for the equivalent 64-bit binary, and the 64-bit binary performs slower than the 32-bit binary. A possible reason is the 32-bit binary benefited from macro-fusion significantly.

B.6.2.5 Length Changing Prefix (LCP) Stalls

13. LCP Delays Detected: $ILD_STALL / CPU_CLK_UNHALTED.CORE$

A high value of the LCP Delays Detected ratio indicates that many Length Changing Prefix (LCP) delays occur in the measured code.

B.6.2.6 Self Modifying Code Detection

14. Self Modifying Code Clear Performance Impact: $\text{MACHINE_NUKES.SMC} * 150 / \text{CPU_CLK_UNHALTED.CORE} * 100$

A program that writes into code sections and shortly afterwards executes the generated code may incur severe penalties. Self Modifying Code Performance Impact estimates the percentage of cycles that the program spends on self-modifying code penalties.

B.6.3 Branch Prediction Ratios

Appendix B.7.2.2, discusses branching that impacts the front-end performance. This section describes event ratios that are commonly used to characterize branch mispredictions.

B.6.3.1 Branch Mispredictions

15. Branch Misprediction Performance Impact: $\text{RESOURCE_STALLS.BR_MISS_CLEAR} / \text{CPU_CLK_UNHALTED.CORE} * 100$

With the Branch Misprediction Performance Impact, you can tell the percentage of cycles that the processor spends in recovering from branch mispredictions.

16. Branch Misprediction per Micro-Op Retired:
 $\text{BR_INST_RETIRED.MISPRED} / \text{UOPS_RETIRED.ANY}$

The ratio Branch Misprediction per Micro-Op Retired indicates if the code suffers from many branch mispredictions. In this case, improving the predictability of branches can have a noticeable impact on the performance of your code.

In addition, the performance impact of each branch misprediction might be high. This happens if the code prior to the mispredicted branch has high CPI, such as cache misses, which cannot be parallelized with following code due to the branch misprediction. Reducing the CPI of this code will reduce the misprediction performance impact. See other ratios to identify these cases.

You can use the precise event `BR_INST_RETIRED.MISPRED` to detect the actual targets of the mispredicted branches. This may help you to identify the mispredicted branch.

B.6.3.2 Virtual Tables and Indirect Calls

17. Virtual Table Usage: $\text{BR_IND_CALL_EXEC} / \text{INST_RETIRED.ANY}$

A high value for the ratio Virtual Table Usage indicates that the code includes many indirect calls. The destination address of an indirect call is hard to predict.

18. Virtual Table Misuse: $\text{BR_CALL_MISSP_EXEC} / \text{BR_INST_RETIRED.MISPRED}$

A high value of Branch Misprediction Performance Impact ratio (Ratio 15) together with high Virtual Table Misuse ratio indicate that significant time is spent due to mispredicted indirect function calls.

In addition to explicit use of function pointers in C code, indirect calls are used for implementing inheritance, abstract classes, and virtual methods in C++.

B.6.3.3 Mispredicted Returns

19. Mispredicted Return Instruction Rate: $\text{BR_RET_MISSP_EXEC} / \text{BR_RET_EXEC}$

The processor has a special mechanism that tracks CALL-RETURN pairs. The processor assumes that every CALL instruction has a matching RETURN instruction. If a RETURN instruction restores a return address, which is not the one stored during the matching CALL, the code incurs a misprediction penalty.

B.6.4 Execution Ratios

This section covers event ratios that can provide insights to the interactions of microops with RS, ROB, execution units, and so forth.

B.6.4.1 Resource Stalls

A high value for the RS Full Ratio (Ratio 2) indicates that the Reservation Station (RS) often gets full with μops due to long dependency chains. The μops that get into the RS cannot execute because they wait for their operands to be computed by previous μops , or they wait for a free execution unit to be executed. This prevents exploiting the parallelism provided by the multiple execution units.

A high value for the ROB Full Ratio (Ratio 3) indicates that the reorder buffer (ROB) often gets full with μops . This usually implies on long latency operations, such as L2 cache demand misses.

B.6.4.2 ROB Read Port Stalls

20. ROB Read Port Stall Rate: $\text{RAT_STALLS.ROB_READ_PORT} / \text{CPU_CLK_UNHALTED.CORE}$

The ratio ROB Read Port Stall Rate identifies ROB read port stalls. However it should be used only if the number of resource stalls, as indicated by Resource Stall Ratio, is low.

B.6.4.3 Partial Register Stalls

21. Partial Register Stalls Ratio: $\text{RAT_STALLS.PARTIAL_CYCLES} / \text{CPU_CLK_UNHALTED.CORE} * 100$

Frequent accesses to registers that cause partial stalls increase access latency and decrease performance. Partial Register Stalls Ratio is the percentage of cycles when partial stalls occur.

B.6.4.4 Partial Flag Stalls

22. Partial Flag Stalls Ratio: $\text{RAT_STALLS.FLAGS} / \text{CPU_CLK_UNHALTED.CORE}$

Partial flag stalls have high penalty and they can be easily avoided. However, in some cases, Partial Flag Stalls Ratio might be high although there are no real flag stalls. There are a few instructions that partially modify the RFLAGS register and may cause partial flag stalls. The most popular are the shift instructions (SAR, SAL, SHR, and SHL) and the INC and DEC instructions.

B.6.4.5 Bypass Between Execution Domains

23. Delayed Bypass to FP Operation Rate: $\text{DELAYED_BYPASS.FP} / \text{CPU_CLK_UNHALTED.CORE}$

24. Delayed Bypass to SIMD Operation Rate: $\text{DELAYED_BYPASS.SIMD} / \text{CPU_CLK_UNHALTED.CORE}$

25. Delayed Bypass to Load Operation Rate: $\text{DELAYED_BYPASS.LOAD} / \text{CPU_CLK_UNHALTED.CORE}$

Domain bypass adds one cycle to instruction latency. To identify frequent domain bypasses in the code you can use the above ratios.

B.6.4.6 Floating Point Performance Ratios

26. Floating Point Instructions Ratio: $\text{X87_OPS_RETIRED.ANY} / \text{INST_RETIRED.ANY} * 100$

Significant floating-point activity indicates that specialized optimizations for floating-point algorithms may be applicable.

27. FP Assist Performance Impact: $\text{FP_ASSIST} * 80 / \text{CPU_CLK_UNHALTED.CORE} * 100$

Floating Point assist is activated for non-regular FP values like denormals and NaNs. FP assist is extremely slow compared to regular FP execution. Different assists incur different penalties. FP Assist Performance Impact estimates the overall impact.

28. Divider Busy: $\text{IDLE_DURING_DIV} / \text{CPU_CLK_UNHALTED.CORE} * 100$

A high value for the Divider Busy ratio indicates that the divider is busy and no other execution unit or load operation is in progress for many cycles. Using this ratio ignores L1 data cache misses and L2 cache misses that can be executed in parallel and hide the divider penalty.

29. Floating-Point Control Word Stall Ratio: $\text{RESOURCE_STALLS.FPCW} / \text{CPU_CLK_UNHALTED.CORE} * 100$

Frequent modifications to the Floating-Point Control Word (FPCW) might significantly decrease performance. The main reason for changing FPCW is for changing rounding mode when doing FP to integer conversions.

B.6.5 Memory Sub-System - Access Conflicts Ratios

A high value for Load or Store Buffer Full Ratio (Ratio 4) indicates that the load buffer or store buffer are frequently full, hence new micro-ops cannot enter the execution pipeline. This can reduce execution parallelism and decrease performance.

30. Load Rate: $L1D_CACHE_LD.MESI / CPU_CLK_UNHALTED.CORE$

One memory read operation can be served by a core each cycle. A high "Load Rate" indicates that execution may be bound by memory read operations.

31. Store Order Block: $STORE_BLOCK.ORDER / CPU_CLK_UNHALTED.CORE * 100$

Store Order Block ratio is the percentage of cycles that store operations, which miss the L2 cache, block committing data of later stores to the memory sub-system. This behavior can further cause the store buffer to fill up (see Ratio 4).

B.6.5.1 Loads Blocked by the L1 Data Cache

32. Loads Blocked by L1 Data Cache Rate:
 $LOAD_BLOCK.L1D / CPU_CLK_UNHALTED.CORE$

A high value for "Loads Blocked by L1 Data Cache Rate" indicates that load operations are blocked by the L1 data cache due to lack of resources, usually happening as a result of many simultaneous L1 data cache misses.

B.6.5.2 4K Aliasing and Store Forwarding Block Detection

33. Loads Blocked by Overlapping Store Rate:
 $LOAD_BLOCK.OVERLAP_STORE / CPU_CLK_UNHALTED.CORE$

4K aliasing and store forwarding block are two different scenarios in which loads are blocked by preceding stores due to different reasons. Both scenarios are detected by the same event: `LOAD_BLOCK.OVERLAP_STORE`. A high value for "Loads Blocked by Overlapping Store Rate" indicates that either 4K aliasing or store forwarding block may affect performance.

B.6.5.3 Load Block by Preceding Stores

34. Loads Blocked by Unknown Store Address Rate: $LOAD_BLOCK.STA / CPU_CLK_UNHALTED.CORE$

A high value for "Loads Blocked by Unknown Store Address Rate" indicates that loads are frequently blocked by preceding stores with unknown address and implies performance penalty.

35. Loads Blocked by Unknown Store Data Rate: `LOAD_BLOCK.STD / CPU_CLK_UNHALTED.CORE`

A high value for “Loads Blocked by Unknown Store Data Rate” indicates that loads are frequently blocked by preceding stores with unknown data and implies performance penalty.

B.6.5.4 Memory Disambiguation

The memory disambiguation feature of Intel Core microarchitecture eliminates most of the non-required load blocks by stores with unknown address. When this feature fails (possibly due to flaky load - store disambiguation cases) the event `LOAD_BLOCK.STA` will be counted and also `MEMORY_DISAMBIGUATION.RESET`.

B.6.5.5 Load Operation Address Translation

36. L0 DTLB Miss due to Loads - Performance Impact: `DTLB_MISSES.L0_MISS_LD * 2 / CPU_CLK_UNHALTED.CORE`

High number of DTLB0 misses indicates that the data set that the workload uses spans a number of pages that is bigger than the DTLB0. The high number of misses is expected to impact workload performance only if the CPI (Ratio 1) is low - around 0.8. Otherwise, it is likely that the DTLB0 miss cycles are hidden by other latencies.

B.6.6 Memory Sub-System - Cache Misses Ratios

B.6.6.1 Locating Cache Misses in the Code

Intel Core microarchitecture provides you with precise events for retired load instructions that miss the L1 data cache or the L2 cache. As precise events they provide the instruction pointer of the instruction following the one that caused the event. Therefore the instruction that comes immediately prior to the pointed instruction is the one that causes the cache miss. These events are most helpful to quickly identify on which loads to focus to fix a performance problem. The events are:

`MEM_LOAD_RETIRE.L1D_MISS`

`MEM_LOAD_RETIRE.L1D_LINE_MISS`

`MEM_LOAD_RETIRE.L2_MISS`

`MEM_LOAD_RETIRE.L2_LINE_MISS`

B.6.6.2 L1 Data Cache Misses

37. L1 Data Cache Miss Rate: `L1D_REPL / INST_RETIRED.ANY`

A high value for L1 Data Cache Miss Rate indicates that the code misses the L1 data cache too often and pays the penalty of accessing the L2 cache. See also Loads Blocked by L1 Data Cache Rate (Ratio 32).

You can count separately cache misses due to loads, stores, and locked operations using the events `L1D_CACHE_LD.I_STATE`, `L1D_CACHE_ST.I_STATE`, and `L1D_CACHE_LOCK.I_STATE`, accordingly.

B.6.6.3 L2 Cache Misses

38. L2 Cache Miss Rate: `L2_LINES_IN.SELF.ANY / INST_RETIRED.ANY`

A high L2 Cache Miss Rate indicates that the running workload has a data set larger than the L2 cache. Some of the data might be evicted without being used. Unless all the required data is brought ahead of time by the hardware prefetcher or software prefetching instructions, bringing data from memory has a significant impact on the performance.

39. L2 Cache Demand Miss Rate: `L2_LINES_IN.SELF.DEMAND / INST_RETIRED.ANY`

A high value for L2 Cache Demand Miss Rate indicates that the hardware prefetchers are not exploited to bring the data this workload needs. Data is brought from memory when needed to be used and the workload bears memory latency for each such access.

B.6.7 Memory Sub-system - Prefetching

B.6.7.1 L1 Data Prefetching

The event `L1D_PREFETCH.REQUESTS` is counted whenever the DCU attempts to prefetch cache lines from the L2 (or memory) to the DCU. If you expect the DCU prefetchers to work and to count this event, but instead you detect the event `MEM_LOAD_RETIRE.L1D_MISS`, it might be that the IP prefetcher suffers from load instruction address collision of several loads.

B.6.7.2 L2 Hardware Prefetching

With the event `L2_LD.SELF.PREFETCH.MESI` you can count the number of prefetch requests that were made to the L2 by the L2 hardware prefetchers. The actual number of cache lines prefetched to the L2 is counted by the event `L2_LD.SELF.PREFETCH.I_STATE`.

B.6.7.3 Software Prefetching

The events for software prefetching cover each level of prefetching separately.

40. Useful PrefetchNTA Ratio: $\text{SSE_PRE_MISS.NTA} / \text{SSE_PRE_EXEC.NTA} * 100$

41. Useful PrefetchT0 Ratio: $\text{SSE_PRE_MISS.L1} / \text{SSE_PRE_EXEC.L1} * 100$

42. Useful PrefetchT1 and PrefetchT2 Ratio: $\text{SSE_PRE_MISS.L2} / \text{SSE_PRE_EXEC.L2} * 100$

A low value for any of the prefetch usefulness ratios indicates that some of the SSE prefetch instructions prefetch data that is already in the caches.

43. Late PrefetchNTA Ratio: $\text{LOAD_HIT_PRE} / \text{SSE_PRE_EXEC.NTA}$

44. Late PrefetchT0 Ratio: $\text{LOAD_HIT_PRE} / \text{SSE_PRE_EXEC.L1}$

45. Late PrefetchT1 and PrefetchT2 Ratio: $\text{LOAD_HIT_PRE} / \text{SSE_PRE_EXEC.L2}$

A high value for any of the late prefetch ratios indicates that software prefetch instructions are issued too late and the load operations that use the prefetched data are waiting for the cache line to arrive.

B.6.8 Memory Sub-system - TLB Miss Ratios

46. TLB miss penalty: $\text{PAGE_WALKS.CYCLES} / \text{CPU_CLK_UNHALTED.CORE} * 100$

A high value for the TLB miss penalty ratio indicates that many cycles are spent on TLB misses. Reducing the number of TLB misses may improve performance. This ratio does not include DTLB0 miss penalties (see Ratio 37).

The following ratios help to focus on the kind of memory accesses that cause TLB misses most frequently. See "ITLB Miss Rate" (Ratio 6) for TLB misses due to instruction fetch.

47. DTLB Miss Rate: $\text{DTLB_MISSES.ANY} / \text{INST_RETIRED.ANY}$

A high value for DTLB Miss Rate indicates that the code accesses too many data pages within a short time, and causes many Data TLB misses.

48. DTLB Miss Rate due to Loads: $\text{DTLB_MISSES.MISS_LD} / \text{INST_RETIRED.ANY}$

A high value for DTLB Miss Rate due to Loads indicates that the code accesses loads data from too many pages within a short time, and causes many Data TLB misses. DTLB misses due to load operations may have a significant impact, since the DTLB miss increases the load operation latency. This ratio does not include DTLB0 miss penalties (see Ratio 37).

To precisely locate load instructions that caused DTLB misses you can use the precise event `MEM_LOAD_RETIRE.DTLB_MISS`.

49. DTLB Miss Rate due to Stores: $\text{DTLB_MISSES.MISS_ST} / \text{INST_RETIRED.ANY}$

A high value for DTLB Miss Rate due to Stores indicates that the code accesses too many data pages within a short time, and causes many Data TLB misses due to store

operations. These misses can impact performance if they do not occur in parallel to other instructions. In addition, if there are many stores in a row, some of them missing the DTLB, it may cause stalls due to full store buffer.

B.6.9 Memory Sub-system - Core Interaction

B.6.9.1 Modified Data Sharing

50. Modified Data Sharing Ratio: $\text{EXT_SNOOP.ALL_AGENTS.HITM} / \text{INST_RETIRED.ANY}$

Frequent occurrences of modified data sharing may be due to two threads using and modifying data laid in one cache line. Modified data sharing causes L2 cache misses. When it happens unintentionally (aka false sharing) it usually causes demand misses that have high penalty. When false sharing is removed code performance can dramatically improve.

51. Local Modified Data Sharing Ratio: $\text{EXT_SNOOP.THIS_AGENT.HITM} / \text{INST_RETIRED.ANY}$

Modified Data Sharing Ratio indicates the amount of total modified data sharing observed in the system. For systems with several processors you can use Local Modified Data Sharing Ratio to indicates the amount of modified data sharing between two cores in the same processor. (In systems with one processor the two ratios are similar).

B.6.9.2 Fast Synchronization Penalty

52. Locked Operations Impact: $(\text{L1D_CACHE_LOCK_DURATION} + 20 * \text{L1D_CACHE_LOCK.MESI}) / \text{CPU_CLK_UNHALTED.CORE} * 100$

Fast synchronization is frequently implemented using locked memory accesses. A high value for Locked Operations Impact indicates that locked operations used in the workload have high penalty. The latency of a locked operation depends on the location of the data: L1 data cache, L2 cache, other core cache or memory.

B.6.9.3 Simultaneous Extensive Stores and Load Misses

53. Store Block by Snoop Ratio: $(\text{STORE_BLOCK.SNOOP} / \text{CPU_CLK_UNHALTED.CORE}) * 100$

A high value for "Store Block by Snoop Ratio" indicates that store operations are frequently blocked and performance is reduced. This happens when one core executes a dense stream of stores while the other core in the processor frequently snoops it for cache lines missing in its L1 data cache.

B.6.10 Memory Sub-system - Bus Characterization

B.6.10.1 Bus Utilization

54. Bus Utilization: $\text{BUS_TRANS_ANY.ALL_AGENTS} * 2 / \text{CPU_CLK_UNHALTED.BUS} * 100$

Bus Utilization is the percentage of bus cycles used for transferring bus transactions of any type. In single processor systems most of the bus transactions carry data. In multiprocessor systems some of the bus transactions are used to coordinate cache states to keep data coherency.

55. Data Bus Utilization: $\text{BUS_DRDY_CLOCKS.ALL_AGENTS} / \text{CPU_CLK_UNHALTED.BUS} * 100$

Data Bus Utilization is the percentage of bus cycles used for transferring data among all bus agents in the system, including processors and memory. High bus utilization indicates heavy traffic between the processor(s) and memory. Memory sub-system latency can impact the performance of the program. For compute-intensive applications with high bus utilization, look for opportunities to improve data and code locality. For other types of applications (for example, copying large amounts of data from one memory area to another), try to maximize bus utilization.

56. Bus Not Ready Ratio: $\text{BUS_BNR_DRV.ALL_AGENTS} * 2 / \text{CPU_CLK_UNHALTED.BUS} * 100$

Bus Not Ready Ratio estimates the percentage of bus cycles during which new bus transactions cannot start. A high value for Bus Not Ready Ratio indicates that the bus is highly loaded. As a result of the Bus Not Ready (BNR) signal, new bus transactions might defer and their latency will have higher impact on program performance.

57. Burst Read in Bus Utilization: $\text{BUS_TRANS_BRD.SELF} * 2 / \text{CPU_CLK_UNHALTED.BUS} * 100$

A high value for Burst Read in Bus Utilization indicates that bus and memory latency of burst read operations may impact the performance of the program.

58. RFO in Bus Utilization: $\text{BUS_TRANS_RFO.SELF} * 2 / \text{CPU_CLK_UNHALTED.BUS} * 100$

A high value for RFO in Bus Utilization indicates that latency of Read For Ownership (RFO) transactions may impact the performance of the program. RFO transactions may have a higher impact on the program performance compared to other burst read operations (for example, as a result of loads that missed the L2). See also Ratio 31.

B.6.10.2 Modified Cache Lines Eviction

59. L2 Modified Lines Eviction Rate: $\text{L2_M_LINES_OUT.SELF.ANY} / \text{INST_RETIRED.ANY}$

When a new cache line is brought from memory, an existing cache line, possibly modified, is evicted from the L2 cache to make space for the new line. Frequent evic-

tions of modified lines from the L2 cache increase the latency of the L2 cache misses and consume bus bandwidth.

60. Explicit WB in Bus Utilization: $\text{BUS_TRANS_WB.SELF} * 2 / \text{CPU_CLK_UNHALTED.BUS} * 100$

Explicit Write-back in Bus Utilization considers modified cache line evictions not only from the L2 cache but also from the L1 data cache. It represents the percentage of bus cycles used for explicit write-backs from the processor to memory.

USING PERFORMANCE MONITORING EVENTS

APPENDIX C

INSTRUCTION LATENCY AND THROUGHPUT

This appendix contains tables showing the latency and throughput are associated with commonly used instructions¹. The instruction timing data varies across processors family/models. It contains the following sections:

- **Appendix C.1, “Overview”** — Provides an overview of issues related to instruction selection and scheduling.
- **Appendix C.2, “Definitions”** — Presents definitions.
- **Appendix C.3, “Latency and Throughput”** — Lists instruction throughput, latency associated with commonly-used instruction.

C.1 OVERVIEW

This appendix provides information to assembly language programmers and compiler writers. The information aids in the selection of instruction sequences (to minimize chain latency) and in the arrangement of instructions (assists in hardware processing). The performance impact of applying the information has been shown to be on the order of several percent. This is for applications not dominated by other performance factors, such as:

- cache miss latencies
- bus bandwidth
- I/O bandwidth

Instruction selection and scheduling matters when the programmer has already addressed the performance issues discussed in Chapter 2:

- observe store forwarding restrictions
- avoid cache line and memory order buffer splits
- do not inhibit branch prediction
- minimize the use of `xchg` instructions on memory locations

1. Although instruction latency may be useful in some limited situations (e.g., a tight loop with a dependency chain that exposes instruction latency), software optimization on super-scalar, out-of-order microarchitecture, in general, will benefit much more on increasing the effective throughput of the larger-scale code path. Coding techniques that rely on instruction latency alone to influence the scheduling of instruction is likely to be sub-optimal as such coding technique is likely to interfere with the out-of-order machine or restrict the amount of instruction-level parallelism.

While several items on the above list involve selecting the right instruction, this appendix focuses on the following issues. These are listed in priority order, though which item contributes most to performance varies by application:

- Maximize the flow of μ ops into the execution core. Instructions which consist of more than four μ ops require additional steps from microcode ROM. Instructions with longer μ op flows incur a delay in the front end and reduce the supply of μ ops to the execution core.

In Pentium 4 and Intel Xeon processors, transfers to microcode ROM often reduce how efficiently μ ops can be packed into the trace cache. Where possible, it is advisable to select instructions with four or fewer μ ops. For example, a 32-bit integer multiply with a memory operand fits in the trace cache without going to microcode, while a 16-bit integer multiply to memory does not.

- Avoid resource conflicts. Interleaving instructions so that they don't compete for the same port or execution unit can increase throughput. For example, alternate PADDQ and PMULUDQ (each has a throughput of one issue per two clock cycles). When interleaved, they can achieve an effective throughput of one instruction per cycle because they use the same port but different execution units. Selecting instructions with fast throughput also helps to preserve issue port bandwidth, hide latency and allows for higher software performance.
- Minimize the latency of dependency chains that are on the critical path. For example, an operation to shift left by two bits executes faster when encoded as two adds than when it is encoded as a shift. If latency is not an issue, the shift results in a denser byte encoding.

In addition to the general and specific rules, coding guidelines and the instruction data provided in this manual, you can take advantage of the software performance analysis and tuning toolset available at <http://developer.intel.com/software/products/index.htm>. The tools include the Intel VTune Performance Analyzer, with its performance-monitoring capabilities.

C.2 DEFINITIONS

The data is listed in several tables. The tables contain the following:

- **Instruction Name** — The assembly mnemonic of each instruction.
- **Latency** — The number of clock cycles that are required for the execution core to complete the execution of all of the μ ops that form an instruction.
- **Throughput** — The number of clock cycles required to wait before the issue ports are free to accept the same instruction again. For many instructions, the throughput of an instruction can be significantly less than its latency.

C.3 LATENCY AND THROUGHPUT

This section presents the latency and throughput information for commonly-used instructions including: MMX technology, Streaming SIMD Extensions, subsequent generations of SIMD instruction extensions, and most of the frequently used general-purpose integer and x87 floating-point instructions.

Due to the complexity of dynamic execution and out-of-order nature of the execution core, the instruction latency data may not be sufficient to accurately predict realistic performance of actual code sequences based on adding instruction latency data.

- Instruction latency data is useful when tuning a dependency chain. However, dependency chains limit the out-of-order core's ability to execute micro-ops in parallel. Instruction throughput data are useful when tuning parallel code unencumbered by dependency chains.
- Numeric data in the tables is:
 - approximate and subject to change in future implementations of the microarchitecture.
 - not meant to be used as reference for instruction-level performance benchmarks. Comparison of instruction-level performance of microprocessors that are based on different microarchitectures is a complex subject and requires information that is beyond the scope of this manual.

Comparisons of latency and throughput data between different microarchitectures can be misleading.

Appendix C.3.1 provides latency and throughput data for the register-to-register instruction type. Appendix C.3.3 discusses how to adjust latency and throughput specifications for the register-to-memory and memory-to-register instructions.

In some cases, the latency or throughput figures given are just one half of a clock. This occurs only for the double-speed ALUs.

C.3.1 Latency and Throughput with Register Operands

Instruction latency and throughput data are presented in Table C-4 through Table C-15. Tables include AESNI, SSE4.2, SSE4.1, Supplemental Streaming SIMD Extension 3, Streaming SIMD Extension 3, Streaming SIMD Extension 2, Streaming SIMD Extension, MMX technology and most common Intel 64 and IA-32 instructions. Instruction latency and throughput for different processor microarchitectures are in separate columns.

Processor instruction timing data is implementation specific; it can vary between model encodings within the same family encoding (e.g. model = 3 vs model < 2). Separate sets of instruction latency and throughput are shown in the columns for CPUID signature 0xF2n and 0xF3n. The column represented by 0xF3n also applies to Intel processors with CPUID signature 0xF4n and 0xF6n. The notation 0xF2n represents the hex value of the lower 12 bits of the EAX register reported by CPUID instruction with input value of EAX = 1; 'F' indicates the family encoding value is 15,

INSTRUCTION LATENCY AND THROUGHPUT

'2' indicates the model encoding is 2, 'n' indicates it applies to any value in the stepping encoding.

Intel Core Solo and Intel Core Duo processors are represented by 06_0EH. Processors based on 65 nm Intel Core microarchitecture are represented by 06_0FH. Processors based on Enhanced Intel Core microarchitecture are represented by 06_17H and 06_1DH. CPUID family/Model signatures of processors based on Intel microarchitecture code name Nehalem are represented by 06_1AH, 06_1EH, 06_1FH, and 06_2EH. Intel microarchitecture code name Westmere are represented by 06_25H, 06_2CH and 06_2FH. Intel microarchitecture code name Sandy Bridge are represented by 06_2AH.

Instruction latency varies by microarchitectures. Table C-1 lists SIMD extensions support in recent microarchitectures. Each microarchitecture may be associated with more than one signature values given by CPUID's "display_family" and "display_model". Not all instruction set extensions are enabled in all processors associated with a particular family/model designation. To determine whether a given instruction set extension is supported, software must use the appropriate CPUID feature flag as described in *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2A*.

Table C-1. SIMD Instruction Extensions Support by Microarchitectures (CPUID Signature)

SIMD Instruction Extensions	DisplayFamily_DisplayModel						
	06_2AH	06_25H 06_2CH 06_2FH	06_1AH 06_1EH 06_1FH 06_2EH	06_17H 06_1DH	06_0FH	06_0EH	0F_06H
AVX	Yes	No	No	No	No	No	No
AESNI, PCLMULQDQ	Yes	Yes	No	No	No	No	No
SSE4.2, POPCNT	Yes	Yes	Yes	No	No	No	No
SSE4.1	Yes	Yes	Yes	Yes	No	No	No
SSSE3	Yes	Yes	Yes	Yes	Yes	No	No
SSE3	Yes	Yes	Yes	Yes	Yes	Yes	Yes
SSE2	Yes	Yes	Yes	Yes	Yes	Yes	Yes
SSE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
MMX	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Table C-2. 256-bit AVX Instructions

Instruction	Latency ¹	Throughput
DisplayFamily_DisplayModel	06_2AH	06_2AH
VADDPD/PS ymm1, ymm2, ymm3	3	1
VADDSUBPD/PS ymm1, ymm2, ymm3	3	1
VANDNPD/PS ymm1, ymm2, ymm3	1	1
VANDPD/PS ymm1, ymm2, ymm3	1	1
VBLENDPD/PS ymm1, ymm2, ymm3, imm	1	0.5
VBLENDVPD/PS ymm1, ymm2, ymm3, ymm	2	1
VCMPDP/PS ymm1, ymm2, ymm3	3	1
VCVTDQ2PD ymm1, ymm2	4	1
VCVTDQ2PS ymm1, ymm2	3	1
VCVT(T)PD2DQ ymm1, ymm2	4	1
VCVTPD2PS ymm1, ymm2	4	1
VCVT(T)PS2DQ ymm1, ymm2	3	1
VCVTPS2PD ymm1, ymm2	2	1
VDIVPD ymm1, ymm2, ymm3	45	44
VDIVPS ymm1, ymm2, ymm3	29	28
VDPPS ymm1, ymm2, ymm3	12	2
VEXTRACTF128 ymm1, ymm2, imm	1	1
VMAXPD/PS ymm1, ymm2, ymm3	3	1
VMINPD/PS ymm1, ymm2, ymm3	3	1
VMOVAPD/PS ymm1, ymm2	1	1
VMOVDUP ymm1, ymm2	1	1
VMOVDQA/U ymm1, ymm2	1	0.5
VMOVMSKPD/PS ymm1, ymm2	2	1
VMOVSHDUP ymm1, ymm2	1	1
VMOVSLDUP ymm1, ymm2	1	1
VMOVUPD/PS ymm1, ymm2	1	1
VMULPD/PS ymm1, ymm2, ymm3	5	1
VORPD/PS ymm1, ymm2, ymm3	1	1
VPERM2F128 ymm1, ymm2, ymm3, imm	1	1

Table C-2. 256-bit AVX Instructions (Contd.)

Instruction	Latency ¹	Throughput
DisplayFamily_DisplayModel	06_2AH	06_2AH
VPERMILPD/PS ymm1, ymm2, ymm3	1	1
VRCPPS ymm1, ymm2	7	1
VROUNDPD/PS ymm1, ymm2, imm	3	1
VRSQRTPS ymm1, ymm2	7	1
VSHUFPD/PS ymm1, ymm2, ymm3, imm	1	1
VSQRTPD ymm1, ymm2	45	44
VSQRTPS ymm1, ymm2	29	28
VSUBPD/PS ymm1, ymm2, imm	3	1
VTESTPD/PS ymm1, ymm2	1	1
VUNPCKHPD/PS ymm1, ymm2, ymm3	1	1
VUNPCKLPD/PS ymm1, ymm2, ymm3	1	1
VXORPD/PS ymm1, ymm2, ymm3	1	1
VZEROUPPER	0	1

Latency of VEX.128 encoded AVX instructions should refer to corresponding legacy 128-bit instructions.

Table C-3. AESNI and PCLMULQDQ Instructions

Instruction	Latency ¹		Throughput	
	06_2AH	06_25H, 06_2CH	06_2AH	06_25H, 06_2CH
AESDEC/AESDECLAST xmm1, xmm2	8	6	1	2
AESENC/AESENCLAST xmm1, xmm2	8	6	1	2
AESIMC xmm1, xmm2	12	6	2	2
AESKEYGENASSIST xmm1, xmm2, imm	10	6	2	2
PCLMULQDQ xmm1, xmm2, imm	14	14	8	8

Table C-4. SSE4.2 Instructions

Instruction	Latency ¹		Throughput	
	06_2AH	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH	06_2AH	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH
DisplayFamily_DisplayModel				
CRC32 r32, r32	3	3	1	1
PCMPSTRM xmm1, xmm2, imm	11	7	3	2
PCMPSTRM xmm1, xmm2, imm	11	8	3	2
PCMPISTRM xmm1, xmm2, imm	11	7	3	2
PCMPISTRM xmm1, xmm2, imm	11	8	3	2
PCMPGTQ xmm1, xmm2	5	3	1	1
POPCNT r32, r32	3	3	1	1
POPCNT r64, r64	3	3	1	1

Table C-5. SSE4.1 Instructions

Instruction	Latency ¹			Throughput		
	06_2AH,	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH	06_17H 06_1DH	06_2AH	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH	06_17H 06_1DH
DisplayFamily_DisplayModel						
BLENDDP/S xmm1, xmm2, imm	1	1	1	0.5	1	1
BLENDVDP/S xmm1, xmm2	2	2	2	1	2	2
DPPD xmm1, xmm2	9	9	9	1	2	2
DPPS xmm1, xmm2	12	11	11	2	2	2
EXTRACTPS xmm1, xmm2, imm	3	2	5	1	1	1
INSERTPS xmm1, xmm2, imm	1	1	1	1	1	1
MPSADBW xmm1, xmm2, imm	5	4	4	1	1	2
PACKUSDW xmm1, xmm2	1	1	1	0.5	0.5	1
PBLENVB xmm1, xmm2	2	1	2	1	1	2

Table C-5. SSE4.1 Instructions (Contd.)

Instruction	Latency ¹			Throughput		
	06_2AH,	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH	06_17H 06_1DH	06_2AH	06_25H, 06_2CH, 06_1AH, 06_1EH, 06_1FH, 06_2EH	06_17H 06_1DH
DisplayFamily_DisplayModel						
PBLENDW xmm1, xmm2, imm	1	1	1	0.5	1	1
PCMPEQQ xmm1, xmm2	1	1	1	0.5	0.5	1
PEXTRB/W/D reg, xmm1, imm	3	3	5	1	1	1
PHMINPOSUW xmm1, xmm2	5	3		1	1	
PINSRB/W/D xmm1, reg, imm	2	2	4	1	1	1
PMAXSB/SD xmm1, xmm2	1	1	1	0.5	1	1
PMAXUW/UD xmm1, xmm2	1	1	1	0.5	1	1
PMINSB/SD xmm1, xmm2	1	1	1	0.5	1	1
PMINUW/UD xmm1, xmm2	1	1	1	0.5	1	1
PMOVSXBD/BW/BQ xmm1, xmm2	1	1	1	0.5	1	1
PMOVXWD/WQ/DQ xmm1, xmm2	1	1	1	0.5	1	1
PMOVZXBBD/BW/BQ xmm1, xmm2	1	1	1	0.5	1	1
PMOVZXWD/WQ/DQ xmm1, xmm2	1	1	1	0.5	1	1
PMULDQ xmm1, xmm2	3	3	3	1	1	1
PMULLD xmm1, xmm2	3	6	5	1	2	2
PTEST xmm1, xmm2	2	2	2	1	1	1
ROUNDPD/PS xmm1, xmm2, imm	3	3	1	1	1	1
ROUNDSD/SS xmm1, xmm2, imm	3	3	1	1	1	1

Table C-6. Supplemental Streaming SIMD Extension 3 Instructions

Instruction	Latency ¹				Throughput			
	06_2A	06_25, 06_2C, 06_1A, 06_1E, 06_1F, 06_2E	06_17 06_1D	06_0F	06_2A	06_25, 06_2C, 06_1A, 06_1E, 06_1F, 06_2E	06_17 06_1D	06_0F
PALIGNR xmm1, xmm2, imm	1	1	1	2	0.5	1	1	1
PHADDD xmm1, xmm2	3	3	3	5	1.5	1.5	2	3
PHADDW xmm1, xmm2	3	3	3	6	1.5	1.5	2	4
PHADDSW xmm1, xmm2	3	3	3	6	1.5	1.5	2	4
PHSUBD xmm1, xmm2	3	3	3	5	1.5	1.5	2	3
PHSUBW xmm1, xmm2	3	3	3	6	1.5	1.5	2	4
PHSUBSW xmm1, xmm2	3	3	3	6	1.5	1.5	2	4
PMADDUBSW xmm1, xmm2	3	3	3	3	1	1	1	1
PMULHRSW xmm1, xmm2	3	3	3	3	1	1	1	1
PSHUFB xmm1, xmm2	1	1	1	3	0.5	0.5	1	2
PSIGNB/D/W xmm1, xmm2	1	1	1	1	0.5	0.5	0.5	0.5
PABSB/D/W xmm1, xmm2	1	1	1	1	0.5	0.5	0.5	0.5

Table C-7. Streaming SIMD Extension 3 SIMD Floating-point Instructions

Instruction	Latency ¹	Throughput	Execution Unit
CPUID	0F_03H	0F_03H	0F_03H
ADDSUBPD/ADDSUBPS	5	2	FP_ADD
HADDPD/HADDPs	13	4	FP_ADD,FP_MISC
HSUBPD/HSUBPS	13	4	FP_ADD,FP_MISC
MOVDDUP xmm1, xmm2	4	2	FP_MOVE
MOVSHDUP xmm1, xmm2	6	2	FP_MOVE
MOVSLDUP xmm1, xmm2	6	2	FP_MOVE

See Appendix C.3.2, "Table Footnotes"

Table C-7a. Streaming SIMD Extension 3 SIMD Floating-point Instructions

Instruction	Latency ¹				Throughput			
	06_2A	06_25, 06_2C, 06_1A, 06_1E, 06_1F, 06_2E	06_17 06_1D	06_0F	06_2A	06_25, 06_2C, 06_1A, 06_1E, 06_1F, 06_2E	06_17 06_1D	06_0F
ADDSUBPD/ADDSUBPS	3	3	3	3	1	1	1	1
HADDPD xmm1, xmm2	5	5	6	5	2	2	1	2
HADDPs xmm1, xmm2	5	5	7	9	2	2	2	4
HSUBPD xmm1, xmm2	5	5	6	5	2	2	1	2
HSUBPS xmm1, xmm2	5	5	7	9	2	2	2	4
MOVDDUP xmm1, xmm2	1	1	1	1	1	1	1	1
MOVSHDUP xmm1, xmm2	1			2	1			1
MOVSLDUP xmm1, xmm2	1				1			

Table C-8. Streaming SIMD Extension 2 128-bit Integer Instructions

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID					OF_2H
CVTDQ2PS ³ xmm, xmm	5	5	2	2	FP_ADD
CVTPS2DQ ³ xmm, xmm	5	5	2	2	FP_ADD
CVTTPS2DQ ³ xmm, xmm	5	5	2	2	FP_ADD
MOVD xmm, r32	6	6	2	2	MMX_MISC,MMX_S HFT
MOVD r32, xmm	10	10	1	1	FP_MOVE, FP_MISC
MOVDQA xmm, xmm	6	6	1	1	FP_MOVE
MOVDQU xmm, xmm	6	6	1	1	FP_MOVE
MOVDQ2Q mm, xmm	8	8	2	2	FP_MOVE, MMX_ALU
MOVQ2DQ xmm, mm	8	8	2	2	FP_MOVE, MMX_SHFT
MOVQ xmm, xmm	2	2	2	2	MMX_SHFT
PACKSSWB/PACKSSDW/ PACKUSWB xmm, xmm	4	4	2	2	MMX_SHFT
PADDB/PADDW/PADDD xmm, xmm	2	2	2	2	MMX_ALU
PADDSB/PADDSW/ PADDUSB/PADDUSW xmm, xmm	2	2	2	2	MMX_ALU
PADDQ mm, mm	2	2	1	1	FP_MISC
PSUBQ mm, mm	2	2	1	1	FP_MISC
PADDQ/ PSUBQ ³ xmm, xmm	6	6	2	2	FP_MISC
PAND xmm, xmm	2	2	2	2	MMX_ALU
PANDN xmm, xmm	2	2	2	2	MMX_ALU
PAVGB/PAVGW xmm, xmm	2	2	2	2	MMX_ALU
PCMPEQB/PCMPEQD/ PCMPEQW xmm, xmm	2	2	2	2	MMX_ALU
PCMPGTB/PCMPGTD/PCMPGTW xmm, xmm	2	2	2	2	MMX_ALU
PEXTRW r32, xmm, imm8	7	7	2	2	MMX_SHFT, FP_MISC

Table C-8. Streaming SIMD Extension 2 128-bit Integer Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID					OF_2H
PINSRW xmm, r32, imm8	4	4	2	2	MMX_SHFT,MMX_MISC
PMADDWD xmm, xmm	9	8	2	2	FP_MUL
PMAX xmm, xmm	2	2	2	2	MMX_ALU
PMIN xmm, xmm	2	2	2	2	MMX_ALU
PMOVMASK ³ r32, xmm	7	7	2	2	FP_MISC
PMULHUW/PMULHW/ PMULLW ³ xmm, xmm	9	8	2	2	FP_MUL
PMULUDQ mm, mm	9	8		1	FP_MUL
PMULUDQ xmm, xmm	9	8	2	2	FP_MUL
POR xmm, xmm	2	2	2	2	MMX_ALU
PSADBW xmm, xmm	4	4	2	2	MMX_ALU
PSHUFD xmm, xmm, imm8	4	4	2	2	MMX_SHFT
PSHUFW xmm, xmm, imm8	2	2	2	2	MMX_SHFT
PSHUFLW xmm, xmm, imm8	2	2	2	2	MMX_SHFT
PSLLDQ xmm, imm8	4	4	2	2	MMX_SHFT
PSLLW/PSLLD/PSLLQ xmm, xmm/imm8	2	2	2	2	MMX_SHFT
PSRAW/PSRAD xmm, xmm/imm8	2	2	2	2	MMX_SHFT
PSRLDQ xmm, imm8	4	4	2	2	MMX_SHFT
PSRLW/PSRLD/PSRLQ xmm, xmm/imm8	2	2	2	2	MMX_SHFT
PSUBB/PSUBW/PSUBD xmm, xmm	2	2	2	2	MMX_ALU
PSUBSB/PSUBSW/PSUBUSB/PSUBU SW xmm, xmm	2	2	2	2	MMX_ALU
PUNPCKHBW/PUNPCKHWD/PUNPC KHDQ xmm, xmm	4	4	2	2	MMX_SHFT
PUNPCKHQDQ xmm, xmm	4	4	2	2	MMX_SHFT
PUNPCKLBW/PUNPCKLWD/PUNPCK LDQ xmm, xmm	2	2	2	2	MMX_SHFT
PUNPCKLQDQ ³ xmm, xmm	4	4	1	1	FP_MISC
PXOR xmm, xmm	2	2	2	2	MMX_ALU

See Appendix C.3.2, "Table Footnotes"

Table C-8a. Streaming SIMD Extension 2 128-bit Integer Instructions

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_OF H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_OF H
CPUID								
CVTPS2DQ xmm, xmm	3	3	3	3	1	1	1	1
CVTTPS2DQ xmm, xmm	3	3	3	3	1	1	1	1
MASKMOVDQU xmm, xmm				8				2
MOVD xmm, r32	1	1	1	1	1	0.33	0.33	0.5
MOVD xmm, r64	1	1		1	1	0.33		0.5
MOVD r32, xmm	2	1	1	1	1	0.33	0.33	0.33
MOVD r64, xmm		1		1		0.33		0.33
MOVDQA xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVDQU xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.5
MOVDQ2Q mm, xmm	2	1		1	1	0.33	0.33	
MOVQ2DQ xmm, mm	1	1		1	0.33	0.33	0.33	1
MOVQ xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
PACKSSWB/PACKSS DW/ PACKUSWB xmm, xmm	1	1	1	2	0.5	0.5	1	2
PADDB/PADDW/PA DDD xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PADDSB/PADDSW/ PADDUSB/PADDUS W xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PADDQ/ PSUBQ ³ xmm, xmm	1	1	2	2	0.5	1	1	1
PAND xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33

Table C-8a. Streaming SIMD Extension 2 128-bit Integer Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_0F H
CPUID								
PANDN xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
PAVGB/PAVGW xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.5
PCMPEQB/PCMPEQ D/ PCMPEQW xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PCMPGTB/PCMPGT D/PCMPGTW xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PEXTRW r32, xmm, imm8	3	3		2	1	1	1	1
PINSRW xmm, r32, imm8	2	2		3	1	1	1	1
PMADDWD xmm, xmm	3	3	3	3	1	1	1	1
PMADDUBSW xmm, xmm	3	3			1	1		
PMAX xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.5
PMIN xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.5
PMOVBMSKB ³ r32, xmm	2	1			1	1	1	1
PMULHUW/PMULH W/ PMULLW xmm, xmm	3	3	3	3	1	1	1	1
PMULUDQ mm, mm	3	3	3	3	1	1	1	1
PMULUDQ xmm, xmm	3	3	3	3	1	1	1	1
POR xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
PSADBW xmm, xmm	3	3	3	3	1	1	1	1
PSHUFD xmm, xmm, imm8	1	1	1	2	0.5	0.5	1	1
PSHUFW xmm, xmm, imm8	1	1	1	1	0.5	0.5	1	1

Table C-8a. Streaming SIMD Extension 2 128-bit Integer Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H, 06_1DH	06_0F H
CPUID								
PSHUFLW xmm, xmm, imm8	1	1	1	1	0.5	0.5	1	1
PSLLDQ xmm, imm8	1	1	1	3	0.5	0.5	1	2
PSLLW/PSLLD/PSLL Q xmm, imm8	1	1	1	2	1	1	1	1
PSLL/PSRL xmm, xmm	2	2			1	1		
PSRAW/PSRAD xmm, imm8	1	1	1	2	1	1	1	1
PSRAW/PSRAD xmm, xmm	2	2			1	1		
PSRLDQ xmm, imm8	1	1	1	2	0.5	0.5	1	1
PSRLW/PSRLD/PSR LQ xmm, imm8	1	1	1	2	1	1	1	1
PSUBB/PSUBW/PSU BD xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PSUBSB/PSUBSW/P SUBUSB/PSUBUSW xmm, xmm	1	1	1	1	0.5	0.5	0.5	0.33
PUNPCKHBW/PUNP CKHWD/PUNPCKHD Q xmm, xmm	1	1	1	2	0.5	0.5	1	2
PUNPCKHQDQ xmm, xmm	1	1	1	1	0.5	0.5	1	1
PUNPCKLBW/PUNP CKLWD/PUNPCKLD Q xmm, xmm	1	1	1	2	0.5	0.5	1	2
PUNPCKLQDQ xmm, xmm	1	1	1	1	0.5	0.5	1	1
PXOR xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33

**Table C-9. Streaming SIMD Extension 2 Double-precision
Floating-point Instructions**

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID					
ADDPD xmm, xmm	5	4	2	2	FP_ADD
ADDSD xmm, xmm	5	4	2	2	FP_ADD
ANDNPD ³ xmm, xmm	4	4	2	2	MMX_ALU
ANDPD ³ xmm, xmm	4	4	2	2	MMX_ALU
CMPPD xmm, xmm, imm8	5	4	2	2	FP_ADD
CMPSD xmm, xmm, imm8	5	4	2	2	FP_ADD
COMISD xmm, xmm	7	6	2	2	FP_ADD, FP_MISC
CVTDQ2PD xmm, xmm	8	8	3	3	FP_ADD, MMX_SHFT
CVTPD2PI mm, xmm	12	11	3	3	FP_ADD, MMX_SHFT, MMX_ALU
CVTPD2DQ xmm, xmm	10	9	2	2	FP_ADD, MMX_SHFT
CVTPD2PS ³ xmm, xmm	11	10	2	2	FP_ADD, MMX_SHFT
CVTPI2PD xmm, mm	12	11	2	4	FP_ADD, MMX_SHFT, MMX_ALU
CVTPS2PD ³ xmm, xmm	3	2		2	FP_ADD, MMX_SHFT, MMX_ALU
CVTSD2SI r32, xmm	9	8	2	2	FP_ADD, FP_MISC
CVTSD2SS ³ xmm, xmm	17	16	2	4	FP_ADD, MMX_SHFT
CVTSI2SD ³ xmm, r32	16	15	2	3	FP_ADD, MMX_SHFT, MMX_MISC
CVTSS2SD ³ xmm, xmm	9	8	2	2	
CVTTPD2PI mm, xmm	12	11	3	3	FP_ADD, MMX_SHFT, MMX_ALU

Table C-9. Streaming SIMD Extension 2 Double-precision Floating-point Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	OF_2H
CVTTPD2DQ xmm, xmm	10	9	2	2	FP_ADD, MMX_SHFT
CVTTSD2SI r32, xmm	8	8	2	2	FP_ADD, FP_MISC
DIVPD xmm, xmm	70	69	70	69	FP_DIV
DIVSD xmm, xmm	39	38	39	38	FP_DIV
MAXPD xmm, xmm	5	4	2	2	FP_ADD
MAXSD xmm, xmm	5	4	2	2	FP_ADD
MINPD xmm, xmm	5	4	2	2	FP_ADD
MINSD xmm, xmm	5	4	2	2	FP_ADD
MOVAPD xmm, xmm	6	6	1	1	FP_MOVE
MOVMSKPD r32, xmm	6	6	2	2	FP_MISC
MOVSD xmm, xmm	6	6	2	2	MMX_SHFT
MOVUPD xmm, xmm	6	6	1	1	FP_MOVE
MULPD xmm, xmm	7	6	2	2	FP_MUL
MULSD xmm, xmm	7	6	2	2	FP_MUL
ORPD ³ xmm, xmm	4	4	2	2	MMX_ALU
SHUFPD ³ xmm, xmm, imm8	6	6	2	2	MMX_SHFT
SQRTPD xmm, xmm	70	69	70	69	FP_DIV
SQRTSD xmm, xmm	39	38	39	38	FP_DIV
SUBPD xmm, xmm	5	4	2	2	FP_ADD
SUBSD xmm, xmm	5	4	2	2	FP_ADD
UCOMISD xmm, xmm	7	6	2	2	FP_ADD, FP_MISC
UNPCKHPD xmm, xmm	6	6	2	2	MMX_SHFT
UNPCKLPD ³ xmm, xmm	4	4	2	2	MMX_SHFT
XORPD ³ xmm, xmm	4	4	2	2	MMX_ALU

See Appendix C.3.2, "Table Footnotes"

**Table C-9a. Streaming SIMD Extension 2 Double-precision
Floating-point Instructions**

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17/ 1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17/ 1DH	06_0F H
CPUID								
ADDPD xmm, xmm	3	3	3	3	1	1	1	1
ADDSD xmm, xmm	3	3	3	3	1	1	1	1
ANDNPD xmm, xmm	1	1	1	1	1	0.33	0.33	1
ANDPD xmm, xmm	1	1	1	1	1	0.33	0.33	1
CMPPD xmm, xmm, imm8	3	3	3	3	1	1	1	1
CMPSD xmm, xmm, imm8	3	3	3	3	1	1	1	1
COMISD xmm, xmm	2	1	1	1	1	1	1	1
CVTDQ2PD xmm, xmm	4	4	4		1	1	1	1
CVTDQ2PS xmm, xmm	3	3	3	4	1	1	1	1
CVTPD2PI mm, xmm	4	9	7		1	1	1	1
CVTPD2DQ xmm, xmm	4	4	4	4	1	1	1	1
CVTPD2PS xmm, xmm	4	4	4	4	1	1	1	1
CVTPI2PD xmm, mm	4	4		4				1
CVT[T]PS2DQ xmm, xmm	3	3		3				1
CVTPS2PD xmm, xmm	2	2	2	2	1	1	2	2
CVTSD2SI r32, xmm	5	3		3	1	1	1	1
CVT[T]SD2SI r64, xmm	5	3		3	1	1		1
CVTSD2SS xmm, xmm	4	4	4	4	1	1	1	1
CVTSI2SD xmm, r32	4	3			1	1	3	1
CVTSI2SD xmm, r64	4	3		4	1	1		1
CVTSS2SD xmm, xmm	1	1	2	2	1	1	2	2
CVTTPD2PI mm, xmm	4	4			1	1		1

Table C-9a. Streaming SIMD Extension 2 Double-precision Floating-point Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17/ 1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17/ 1DH	06_0F H
CVTTPD2DQ xmm, xmm	4	4	4	4	1	1	1	1
CVTTSD2SI r32, xmm	5	3		3	1	1	1	1
DIVPD xmm, xmm ¹	22	<24	<32	<35	22	<20	<26	<30
DIVSD xmm, xmm	22	<24	<32	<35	22	<20	<26	<30
MAXPD xmm, xmm	3	3	3	3	1	1	1	1
MAXSD xmm, xmm	3	3	3	3	1	1	1	1
MINPD xmm, xmm	3	3	3	3	1	1	1	1
MINSD xmm, xmm	3	3	3	3	1	1	1	1
MOVAPD xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVMSKPD r32, xmm	2	1		1	1	1		1
MOVMSKPD r64, xmm	2	1		1	1	1		1
MOVSD xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVUPD xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.5
MULPD xmm, xmm	5	5	5	5	1	1	1	1
MULSD xmm, xmm	5	5	5	5	1	1	1	1
ORPD xmm, xmm	1	1	1	1	0.33	0.33	0.33	1
SHUFPD xmm, xmm, imm8	1	1	1	1	1	1	1	1
SQRTPD xmm, xmm ²	22	<34	<31	<60	22	<30	<25	<57
SQRTSD xmm, xmm	22	<34	<31	<60	22	<30	<25	<57
SUBPD xmm, xmm	3	3	3	3	1	1	1	1
SUBSD xmm, xmm	3	3	3	3	1	1	1	1
UCOMISD xmm, xmm	2	1		1	1	1		1
UNPCKHPD xmm, xmm	1	1	1	1	1	1	1	1
UNPCKLPD xmm, xmm	1	1	1	1	1	1	1	1
XORPD ³ xmm, xmm	1	1	1	1	0.33	0.33	0.33	1

NOTES:

1. The latency and throughput of DIVPD/DIVSD can vary with input values. For certain values, hardware can complete quickly, throughput may be as low as ~ 6 cycles. Similarly, latency for certain input values may be as low as less than 10 cycles.
2. The latency throughput of SQRTPD/SQRTSD can vary with input value. For certain values, hardware can complete quickly, throughput may be as low as ~ 6 cycles. Similarly, latency for certain input values may be as low as less than 10 cycles.

Table C-10. Streaming SIMD Extension Single-precision Floating-point Instructions

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID					
ADDPS xmm, xmm	5	4	2	2	FP_ADD
ADDSS xmm, xmm	5	4	2	2	FP_ADD
ANDNPS ³ xmm, xmm	4	4	2	2	MMX_ALU
ANDPS ³ xmm, xmm	4	4	2	2	MMX_ALU
CMPPS xmm, xmm	5	4	2	2	FP_ADD
CMPSS xmm, xmm	5	4	2	2	FP_ADD
COMISS xmm, xmm	7	6	2	2	FP_ADD,FP_MISC
CVTPI2PS xmm, mm	12	11	2	4	MMX_ALU,FP_ADD,MMX_SHFT
CVTPS2PI mm, xmm	8	7	2	2	FP_ADD,MMX_ALU
CVTSI2SS ³ xmm, r32	12	11	2	2	FP_ADD,MMX_SHFT, MMX_MISC
CVTSS2SI r32, xmm	9	8	2	2	FP_ADD,FP_MISC
CVTTPS2PI mm, xmm	8	7	2	2	FP_ADD,MMX_ALU
CVTTSS2SI r32, xmm	9	8	2	2	FP_ADD,FP_MISC
DIVPS xmm, xmm	40	39	40	39	FP_DIV
DIVSS xmm, xmm	32	23	32	23	FP_DIV
MAXPS xmm, xmm	5	4	2	2	FP_ADD
MAXSS xmm, xmm	5	4	2	2	FP_ADD

Table C-10. Streaming SIMD Extension Single-precision Floating-point Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	OF_2H
CPUID					
MINPS xmm, xmm	5	4	2	2	FP_ADD
MINSS xmm, xmm	5	4	2	2	FP_ADD
MOVAPS xmm, xmm	6	6	1	1	FP_MOVE
MOVHLPS ³ xmm, xmm	6	6	2	2	MMX_SHFT
MOVLHPS ³ xmm, xmm	4	4	2	2	MMX_SHFT
MOVMSKPS r32, xmm	6	6	2	2	FP_MISC
MOVSS xmm, xmm	4	4	2	2	MMX_SHFT
MOVUPS xmm, xmm	6	6	1	1	FP_MOVE
MULPS xmm, xmm	7	6	2	2	FP_MUL
MULSS xmm, xmm	7	6	2	2	FP_MUL
ORPS ³ xmm, xmm	4	4	2	2	MMX_ALU
RCPSPS ³ xmm, xmm	6	6	4	4	MMX_MISC
RCPSS ³ xmm, xmm	6	6	2	2	MMX_MISC, MMX_SHFT
RSQRTPS ³ xmm, xmm	6	6	4	4	MMX_MISC
RSQRTSS ³ xmm, xmm	6	6	4	4	MMX_MISC, MMX_SHFT
SHUFPS ³ xmm, xmm, imm8	6	6	2	2	MMX_SHFT
SQRTPS xmm, xmm	40	39	40	39	FP_DIV
SQRTSS xmm, xmm	32	23	32	23	FP_DIV
SUBPS xmm, xmm	5	4	2	2	FP_ADD
SUBSS xmm, xmm	5	4	2	2	FP_ADD
UCOMISS xmm, xmm	7	6	2	2	FP_ADD, FP_MISC
UNPCKHPS ³ xmm, xmm	6	6	2	2	MMX_SHFT
UNPCKLPS ³ xmm, xmm	4	4	2	2	MMX_SHFT
XORPS ³ xmm, xmm	4	4	2	2	MMX_ALU
FXRSTOR		150			
FXSAVE		100			

See Appendix C.3.2

Table C-10a. Streaming SIMD Extension Single-precision Floating-point Instructions

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H 06_1DH	06_0F H
CPUID								
ADDPS xmm, xmm	3	3	3	3	1	1	1	1
ADDSS xmm, xmm	3	3	3	3	1	1	1	1
ANDNPS xmm, xmm	1	1	1	1	1	0.33	0.33	1
ANDPS xmm, xmm	1	1	1	1	1	0.33	0.33	1
CMPPS xmm, xmm	3	3	3	3	1	1	1	1
CMPSS xmm, xmm	3	3	3	3	1	1	1	1
COMISS xmm, xmm	2	1	1	1	1	1	1	1
CVTPI2PS xmm, mm	3	3		3	1	1		1
CVTPS2PI mm, xmm	3	3		3	1	1		1
CVTSI2SS xmm, r32	5	3	6	4	1	1	3	1
CVTSS2SI r32, xmm	5	3	6	3	1	1	1	1
CVT[T]SS2SI r64, xmm	5	3		4	1	1		1
CVTTPS2PI mm, xmm	3	3		3	1	1		1
CVTTSS2SI r32, xmm	5	3	6	3	1	1	1	1
DIVPS xmm, xmm ¹	14	<16	<21	<21	14	<12	<14	<16
DIVSS xmm, xmm	14	<16	<21	<21	14	<12	<14	<16
MAXPS xmm, xmm	3	3	3	3	1	1	1	1
MAXSS xmm, xmm	3	3	3	3	1	1	1	1
MINPS xmm, xmm	3	3	3	3	1	1	1	1
MINSS xmm, xmm	3	3	3	3	1	1	1	1
MOVAPS xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVHPS xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVLPS xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
MOVMSKPS r32, xmm	2	1		1	1	1		1
MOVMSKPS r64, xmm	2	1		1	1	1		1

Table C-10a. Streaming SIMD Extension Single-precision Floating-point Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2E	06_17H 06_1DH	06_0F H
CPUID								
MOVSS xmm, xmm	1	1	1	1	01	0.33	0.33	0.33
MOVUPS xmm, xmm	1	1	1	1	1	0.33	0.33	0.5
MULPS xmm, xmm	5	4	4	4	1	1	1	1
MULSS xmm, xmm	5	4	4	4	1	1	1	1
ORPS xmm, xmm	1	1	1	1	1	0.33	0.33	0.33
RCPPS xmm, xmm	5	3	3	3	1	2	2	1
RCPSS xmm, xmm	5	3	3	3	1	3	3	1
RSQRTPS xmm, xmm	5	3	3	3	1	2	2	2
RSQRTSS xmm, xmm	5	3	3	3	1	3	3	2
SHUFPS xmm, xmm, imm8	1	1	1	2	1	1	1	1
SQRTPS xmm, xmm ²	14	<20	<21	<32	14	<16	<14	<27
SQRTSS xmm, xmm	14	<20	<21	<32	14	<16	<14	<27
SUBPS xmm, xmm	3	3	3	3	1	1	1	1
SUBSS xmm, xmm	3	3	3	3	1	1	1	1
UCOMISS xmm, xmm	2			1	1			1
UNPCKHPS xmm, xmm	1	1	1	2	1	1	1	1
UNPCKLPS xmm, xmm	1	1	1	2	1	1	1	1
XORPS xmm, xmm	1	1	1	1	0.33	0.33	0.33	0.33
FXRSTOR								
FXSAVE								

NOTES:

1. The latency and throughput of DIVPS/DIVSS can vary with input values. For certain values, hardware can complete quickly, throughput may be as low as ~ 6 cycles. Similarly, latency for certain input values may be as low as less than 10 cycles.
2. The latency and throughput of SQRTPS/SQRTSS can vary with input values. For certain values, hardware can complete quickly, throughput may be as low as ~ 6 cycles. Similarly, latency for certain input values may be as low as less than 10 cycles

Table C-11. Streaming SIMD Extension 64-bit Integer Instructions

Instruction	Latency ¹		Throughput		Execution Unit
	OF_3 H	OF_2 H	OF_3 H	OF_2 H	
CPUID					OF_2H
PAVGB/PAVGW mm, mm	2	2	1	1	MMX_ALU
PEXTRW r32, mm, imm8	7	7	2	2	MMX_SHFT, FP_MISC
PINSRW mm, r32, imm8	4	4	1	1	MMX_SHFT, MMX_MISC
PMAX mm, mm	2	2	1	1	MMX_ALU
PMIN mm, mm	2	2	1	1	MMX_ALU
PMOVMSKB ³ r32, mm	7	7	2	2	FP_MISC
PMULHUW ³ mm, mm	9	8	1	1	FP_MUL
PSADBW mm, mm	4	4	1	1	MMX_ALU
PSHUFW mm, mm, imm8	2	2	1	1	MMX_SHFT

See Appendix C.3.2, "Table Footnotes"

Table C-11a. Streaming SIMD Extension 64-bit Integer Instructions

Instruction	Latency ¹				Throughput			
	06_17 H	06_0F H	06_0E H	06_0D H	06_17 H	06_0F H	06_0EH	06_0D H
CPUID								
MASKMOVQ mm, mm		3				1		
PAVGB/PAVGW mm, mm	1	1	1	1	0.5	0.5	0.5	0.5
PEXTRW r32, mm, imm8		2*	2	2	1	1	1	1
PINSRW mm, r32, imm8		1	1	1	1	1	1	1
PMAX mm, mm	1	1	1	1	0.5	0.5	0.5	0.5
PMIN mm, mm	1	1	1	1	0.5	0.5	0.5	0.5
PMOVMSKB r32, mm			1	1	1	1	1	1
PMULHUW mm, mm	3	3	3	3	1	1	1	1

Table C-11a. Streaming SIMD Extension 64-bit Integer Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_17 H	06_0F H	06_0E H	06_0D H	06_17 H	06_0F H	06_0EH	06_0D H
PSADBW mm, mm	3	3	5	5	1	1	2	2
PSHUFW mm, mm, imm8	1	1	1	1	1	1	1	1
See Appendix C.3.2, "Table Footnotes"								

Table C-12. MMX Technology 64-bit Instructions

Instruction	Latency ¹		Throughput		Execution Unit ²
	0F_3 H	0F_2H	0F_3 H	0F_2 H	
MOVD mm, r32	2	2	1	1	MMX_ALU
MOVD ³ r32, mm	5	5	1	1	FP_MISC
MOVQ mm, mm	6	6	1	1	FP_MOV
PACKSSWB/PACKSSDW/PACKUSWB mm, mm	2	2	1	1	MMX_SHFT
PADDB/PADDW/PADDD mm, mm	2	2	1	1	MMX_ALU
PADDSB/PADDSW/PADDUSB/PADDUSW mm, mm	2	2	1	1	MMX_ALU
PAND mm, mm	2	2	1	1	MMX_ALU
PANDN mm, mm	2	2	1	1	MMX_ALU
PCMPEQB/PCMPEQD/PCMPEQW mm, mm	2	2	1	1	MMX_ALU
PCMPGTB/PCMPGTD/PCMPGTW mm, mm	2	2	1	1	MMX_ALU
PMADDWD ³ mm, mm	9	8	1	1	FP_MUL
PMULHW/PMULLW ³ mm, mm	9	8	1	1	FP_MUL
POR mm, mm	2	2	1	1	MMX_ALU
PSLLQ/PSLLW/PSLLD mm, mm/imm8	2	2	1	1	MMX_SHFT
PSRAW/PSRAD mm, mm/imm8	2	2	1	1	MMX_SHFT

Table C-12. MMX Technology 64-bit Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	0F_3H	0F_2H	0F_3H	0F_2H	
PSRLQ/PSRLW/PSRLD mm, mm/imm8	2	2	1	1	MMX_SHFT
PSUBB/PSUBW/PSUBD mm, mm	2	2	1	1	MMX_ALU
PSUBSB/PSUBSW/PSUBU SB/PSUBUSW mm, mm	2	2	1	1	MMX_ALU
PUNPCKHBW/PUNPCKHW D/PUNPCKHDQ mm, mm	2	2	1	1	MMX_SHFT
PUNPCKLBW/PUNPCKLWD /PUNPCKLDQ mm, mm	2	2	1	1	MMX_SHFT
PXOR mm, mm	2	2	1	1	MMX_ALU
EMMS ¹		12		12	

See Appendix C.3.2, "Table Footnotes"

Table C-13. MMX Technology 64-bit Instructions

Instruction	Latency ¹				Throughput			
	06_17H	06_0FH	06_0EH	06_0DH	06_17H	06_0FH	06_0EH	06_0DH
MOVD mm, r32		1	1	1		0.5	0.5	0.5
MOVD r32, mm		1	1	1		0.33	0.5	0.5
MOVQ mm, mm	1	1	1	1	0.33	0.5	0.5	0.5
PACKSSWB/PACKSSDW /PACKUSWB mm, mm	1	1	1	1	1	1	1	1
PADDB/PADDW/PADD D mm, mm	1	1	1	1	0.5	0.33	1	1
PADDSB/PADDSW /PADDUSB/PADDUSW mm, mm	1	1	1	1	0.5	0.33	1	1
PAND mm, mm	1	1	1	1	0.33	0.33	0.5	0.5
PANDN mm, mm	1	1	1	1	0.33	0.33	0.5	0.5
PCMPEQB/PCMPEQD PCMPEQW mm, mm	0.5	1	1	1	0.5	0.33	0.5	0.5

Table C-13. MMX Technology 64-bit Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_17 H	06_0F H	06_0E H	06_0D H	06_17 H	06_0F H	06_0EH	06_0D H
PCMPGTB/PCMPGTD/ PCMPGTW mm, mm	0.5	1	1	1	0.5	0.33	0.5	0.5
PMADDWD mm, mm	3	3	3	3	1	1	1	1
PMULHW/PMULLW ³ mm, mm	3	3	3	3	1	1	1	1
POR mm, mm	1	1	1	1	0.33	0.33	0.5	0.5
PSLLQ/PSLLW/ PSLLD mm, mm/imm8	1	1	1	1	1	1	1	1
PSRAW/PSRAD mm, mm/imm8	1	1	1	1	1	1	1	1
PSRLQ/PSRLW/PSRLD mm, mm/imm8	1	1	1	1	1	1	1	1
PSUBB/PSUBW/PSUBD mm, mm	0.5	1	1	1	0.5	0.33	0.5	0.5
PSUBSB/PSUBSW/PSU BUSB/PSUBUSW mm, mm	0.5	1	1	1	0.5	0.33	0.5	0.5
PUNPCKHBW/PUNPCK HWD/PUNPCKHDQ mm, mm	1	1	1	1	1	1	1	1
PUNPCKLBW/PUNPCKL WD/PUNPCKLDQ mm, mm	1	1	1	1	1	1	1	1
PXOR mm, mm	0.33	1	1	1	0.33	0.33	0.5	0.5
EMMS ¹			6	6	6		5	5

See Appendix C.3.2, "Table Footnotes"

Table C-14. x87 Floating-point Instructions

Instruction	Latency ¹		Throughput		Execution Unit ²
	0F_3H	0F_2H	0F_3H	0F_2H	
FABS	3	2	1	1	FP_MISC
FADD	6	5	1	1	FP_ADD
FSUB	6	5	1	1	FP_ADD

Table C-14. x87 Floating-point Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
FMUL	8	7	2	2	FP_MUL
FCOM	3	2	1	1	FP_MISC
FCHS	3	2	1	1	FP_MISC
FDIV Single Precision	30	23	30	23	FP_DIV
FDIV Double Precision	40	38	40	38	FP_DIV
FDIV Extended Precision	44	43	44	43	FP_DIV
FSQRT SP	30	23	30	23	FP_DIV
FSQRT DP	40	38	40	38	FP_DIV
FSQRT EP	44	43	44	43	FP_DIV
F2XM1 ⁴	100-200	90-150		60	
FCOS ⁴	180-280	190-240		130	
FPATAN ⁴	220-300	150-300		140	
FPTAN ⁴	240-300	225-250		170	
FSIN ⁴	160-200	160-180		130	
FSINCOS ⁴	170-250	160-220		140	
FYL2X ⁴	100-250	140-190		85	
FYL2XP1 ⁴		140-190		85	
FSCALE ⁴		60		7	
FRNDINT ⁴		30		11	
FXCH ⁵		0		1	FP_MOVE
FLDZ ⁶		0			
FINCSTP/FDECSTP ⁶		0			

See Appendix C.3.2, "Table Footnotes"

Table C-14a. x87 Floating-point Instructions

Instruction	Latency ¹				Throughput			
	06_17 H	06_OF H	06_OEH	06_OD H	06_17 H	06_OFH	06_OEH	06_OD H
CPUID								
FABS	1	1	1	1		1	1	1
FADD	3	3	3	3	1		1	1
FSUB	3	3	3	3	1	1	1	1
FMUL	5	5	5	5	2	2	2	2
FCOM		1	1	1		1	1	1
FCHS	1				0			
FDIV Single Precision	6	32			5	32		
FDIV Double Precision	6	32			5	32		
FDIV Extended Precision								
FSQRT	6	58	58	58		58	58	58
F2XM1 ⁴	45		69	69			67	67
FCOS ⁴	97		119	119			117	117
FPATAN ⁴			147	147			147	147
FPTAN ⁴			123	123			83	83
FSIN ⁴	82		119	119			116	116
FSINCOS ⁴			119	119			85	85
FYL2X ⁴			96	96			92	92
FYL2XP1 ⁴			98	98			93	93
FSCALE ⁴			17	17			15	15
FRNDINT ⁴	21		21	21			20	20
FXCH ⁵	1							
FLDZ ⁶		1	1	1		1	1	1
FINCSTP/ FDECSTP ⁶			1	1			1	1

See Appendix C.3.2, "Table Footnotes"

Table C-15. General Purpose Instructions

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID	OF_3H	OF_2H	OF_3H	OF_2H	OF_2H
ADC/SBB reg, reg	8	8	3	3	
ADC/SBB reg, imm	8	6	2	2	ALU
ADD/SUB	1	0.5	0.5	0.5	ALU
AND/OR/XOR	1	0.5	0.5	0.5	ALU
BSF/BSR	16	8	2	4	
BSWAP	1	7	0.5	1	ALU
BTC/BTR/BTS	8-9		1		
CLI				26	
CMP/TEST	1	0.5	0.5	0.5	ALU
DEC/INC	1	1	0.5	0.5	ALU
IMUL r32	10	14	1	3	FP_MUL
IMUL imm32		14	1	3	FP_MUL
IMUL		15-18		5	
IDIV	66-80	56-70	30	23	
IN/OUT ¹		<225		40	
Jcc ⁷		Not Applicable		0.5	ALU
LOOP		8		1.5	ALU
MOV	1	0.5	0.5	0.5	ALU
MOVSB/MOVS _W	1	0.5	0.5	0.5	ALU
MOVZB/MOVZ _W	1	0.5	0.5	0.5	ALU
NEG/NOT/NOP	1	0.5	0.5	0.5	ALU
POP r32		1.5		1	MEM_LOAD, ALU
PUSH		1.5		1	MEM_STORE, ALU
RCL/RCR reg, 1 ⁸	6	4	1	1	
ROL/ROR	1	4	0.5	1	
RET		8		1	MEM_LOAD, ALU
SAHF	1	0.5	0.5	0.5	ALU

Table C-15. General Purpose Instructions (Contd.)

Instruction	Latency ¹		Throughput		Execution Unit ²
	OF_3H	OF_2H	OF_3H	OF_2H	
CPUID					OF_2H
SAL/SAR/SHL/SHR	1	4	0.5	1	
SCAS		4		1.5	ALU, MEM_LOAD
SETcc		5		1.5	ALU
STI				36	
STOSB		5		2	ALU, MEM_STORE
XCHG	1.5	1.5	1	1	ALU
CALL		5		1	ALU, MEM_STORE
MUL	10	14-18	1	5	
DIV	66-80	56-70	30	23	

See Appendix C.3.2, "Table Footnotes"

Table C-15a. General Purpose Instructions

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0FH
ADC/SBB reg, reg	2	2	2	2	1	1.5	1.5	1
ADC/SBB reg, imm	2	2	2	2	1	1.5	1.5	1.5
ADD/SUB	1	1	1	1	0.33	0.33	0.33	0.33
AND/OR/XOR	1	1	1	1	0.33	0.33	0.33	0.33
BSF/BSR	3	3	1	2	1	1	1	1
BSWAP	2	3	4	2	1	1	1	0.5
BT	1	1	1	1	0.5	0.5	1	0.33
BTC/BTR/BTS	1	1	1	1	0.5	0.5	1	0.33
CBW	1	1	1	1	0.33	0.33		0.33
CDQ	1				0.5			
CDQE	1	1			0.33	0.33		
CLC	1	1		1	1	0.33	0.33	0.33

Table C-15a. General Purpose Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0FH
CMC	1	1		1	0.33	0.33	0.33	0.33
CLI	5			9	2			9
CLFLUSH		4						
CMOVcc	2	2	2	2	0.66	1	1	0.5
CMOVBE/NBE/A/NA	3				1			
CMP/TEST	1	1	1	1	0.33	0.33	0.33	0.33
CPUID (EAX = 0)						~200	~200	~190
CMPXCHG		5						
DEC/INC	1	1	1	1	0.33	0.33	0.33	0.33
IMUL r32	3	3	3	3	1	1	1	0.5
IMUL imm32	3	3	3	3	1	1	1	0.5
IDIV r32	14-25 ⁹	11-21 ⁹	13-23 ⁹	17-41 ¹⁰	11-17 ⁹	5-13 ⁹	5-14 ⁹	12-36 ¹⁰
MOVSB/MOVS _W	1			1	0.33			0.33
MOVZB/MOVZ _W	1			1	0.33			0.33
DIV r64 ¹¹	80~90	70~80		~90				
MUL r64 ¹²	3, 4	3, 10	6, 7	6, 7				
NEG/NOT	1			1	0.33			0.33
NOP	0			1	0.25			0.33
PUSH				3				1
RCL/RCR reg, 1	2	2	2	2	1	1	1.5	1.5
RCL/RCR	6	12	10	10	6	12	9	9
RDTSC					~28	~31	~31	~65
ROL/ROR reg, im	1	1	1	1	1	0.5	1	1
ROL/ROR reg, cl	2	1	1	1	1.5	0.5	1	1
SAHF	1	1	1	1	0.33	0.33	0.33	0.33
SAL/SAR/SHL/ SHR reg, reg, imm	1	1	1	1	0.5	0.33	0.33	0.33

Table C-15a. General Purpose Instructions (Contd.)

Instruction	Latency ¹				Throughput			
	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0F H	06_2A	06_25/ 2C/1A/ 1E/1F/ 2EH	06_17H 06_1DH	06_0FH
SAL/SAR/SHL/ SHR reg, reg, cl	3	1	1	1	1.5	0.33	0.33	0.33
SETcc	1	1	1	1	0.33	0.33	0.33	0.33
SHLD r32, r32, cl	2	4	2	1	1.5	1	1	1
SHLD r64, r64, cl	2	4	2		1.5	1	1	
SHLD r32, r32, imm	1	4	2	2	0.5	1	1	1
SHLD r64, r64, imm	1	4	2			1		
XADD	3	3			1			
XCHG reg, reg	1.5		2.5	3	1	1	1	1

See Appendix C.3.2, "Table Footnotes"

C.3.2 Table Footnotes

The following footnotes refer to all tables in this appendix.

1. Latency information for many instructions that are complex (> 4 μ ops) are estimates based on conservative (worst-case) estimates. Actual performance of these instructions by the out-of-order core execution unit can range from somewhat faster to significantly faster than the latency data shown in these tables.
2. The names of execution units apply to processor implementations of the Intel NetBurst microarchitecture with a CPUID signature of family 15, model encoding = 0, 1, 2. They include: ALU, FP_EXECUTE, FPMOVE, MEM_LOAD, MEM_STORE. See Figure 2-22 for execution units and ports in the out-of-order core. Note the following:
 - The FP_EXECUTE unit is actually a cluster of execution units, roughly consisting of seven separate execution units.
 - The FP_ADD unit handles x87 and SIMD floating-point add and subtract operation.
 - The FP_MUL unit handles x87 and SIMD floating-point multiply operation.
 - The FP_DIV unit handles x87 and SIMD floating-point divide square-root operations.

INSTRUCTION LATENCY AND THROUGHPUT

- The MMX_SHFT unit handles shift and rotate operations.
 - The MMX_ALU unit handles SIMD integer ALU operations.
 - The MMX_MISC unit handles reciprocal MMX computations and some integer operations.
 - The FP_MISC designates other execution units in port 1 that are separated from the six units listed above.
3. It may be possible to construct repetitive calls to some Intel 64 and IA-32 instructions in code sequences to achieve latency that is one or two clock cycles faster than the more realistic number listed in this table.
 4. Latency and Throughput of transcendental instructions can vary substantially in a dynamic execution environment. Only an approximate value or a range of values are given for these instructions.
 5. The FXCH instruction has 0 latency in code sequences. However, it is limited to an issue rate of one instruction per clock cycle.
 6. The load constant instructions, FINCSTP, and FDECSTP have 0 latency in code sequences.
 7. Selection of conditional jump instructions should be based on the recommendation of section Section 3.4.1, "Branch Prediction Optimization," to improve the predictability of branches. When branches are predicted successfully, the latency of jcc is effectively zero.
 8. RCL/RCR with shift count of 1 are optimized. Using RCL/RCR with shift count other than 1 will be executed more slowly. This applies to the Pentium 4 and Intel Xeon processors.
 9. The latency and throughput of IDIV in Enhanced Intel Core microarchitecture varies with operand sizes and with the number of significant digits of the quotient of the division. If the quotient is zero, the minimum latency can be 13 cycles, and the minimum throughput can be 5 cycles. Latency and throughput of IDIV increases with the number of significant digit of the quotient. The latency and throughput of IDIV with 64-bit operand are significantly slower than those with 32-bit operand. Latency of DIV is similar to IDIV. Generally, the latency of DIV may be one cycle less.
 10. The latency and throughput of IDIV in Intel Core microarchitecture varies with the number of significant digits of the quotient of the division. Latency and throughput of IDIV may increases with the number of significant digit of the quotient. The latency and throughput of IDIV with 64-bit operand are significantly slower than those with 32-bit operand.
 11. The latency of DIV r64 varies. If the dividend (RDX:RAX) is no more than 64-bits, the latency is similar to DIV r32. If the dividend is more than 64-bits, the latency is in the neighborhood of 70 cycles or more.

12. The latency of MUL r64 into 128-bit result has two sets of numbers, the read-to-use latency of the low 64-bit result (RAX) is smaller. The latency of the high 64-bit of the 128 bit result (RDX) is larger.

C.3.3 Instructions with Memory Operands

The latency of an Instruction with memory operand can vary greatly due to a number of factors, including data locality in the memory/cache hierarchy and characteristics that are unique to each microarchitecture. Generally, software can approach tuning for locality and instruction selection independently. Thus Table C-4 through Table C-15 can be used for the purpose of instruction selection. Latency and throughput of data movement in the cache/memory hierarchy can be dealt with independent of instruction latency and throughput. Latency data for the cache hierarchy can be found in Chapter 2.

INSTRUCTION LATENCY AND THROUGHPUT

APPENDIX D

STACK ALIGNMENT

This appendix details on the alignment of the stacks of data for Streaming SIMD Extensions and Streaming SIMD Extensions 2.

D.4 STACK FRAMES

This section describes the stack alignment conventions for both ESP-based (normal), and EBP-based (debug) stack frames. A stack frame is a contiguous block of memory allocated to a function for its local memory needs. It contains space for the function's parameters, return address, local variables, register spills, parameters needing to be passed to other functions that a stack frame may call, and possibly others. It is typically delineated in memory by a stack frame pointer (ESP) that points to the base of the frame for the function and from which all data are referenced via appropriate offsets. The convention on Intel 64 and IA-32 is to use the ESP register as the stack frame pointer for normal optimized code, and to use EBP in place of ESP when debug information must be kept. Debuggers use the EBP register to find the information about the function via the stack frame.

It is important to ensure that the stack frame is aligned to a 16-byte boundary upon function entry to keep local `__m128` data, parameters, and XMM register spill locations aligned throughout a function invocation. The Intel C++ Compiler for Win32* Systems supports conventions presented here help to prevent memory references from incurring penalties due to misaligned data by keeping them aligned to 16-byte boundaries. In addition, this scheme supports improved alignment for `__m64` and double type data by enforcing that these 64-bit data items are at least eight-byte aligned (they will now be 16-byte aligned).

For variables allocated in the stack frame, the compiler cannot guarantee the base of the variable is aligned unless it also ensures that the stack frame itself is 16-byte aligned. Previous software conventions, as implemented in most compilers, only ensure that individual stack frames are 4-byte aligned. Therefore, a function called from a Microsoft-compiled function, for example, can only assume that the frame pointer it used is 4-byte aligned.

Earlier versions of the Intel C++ Compiler for Win32 Systems have attempted to provide 8-byte aligned stack frames by dynamically adjusting the stack frame pointer in the prologue of main and preserving 8-byte alignment of the functions it compiles. This technique is limited in its applicability for the following reasons:

- The main function must be compiled by the Intel C++ Compiler.
- There may be no functions in the call tree compiled by some other compiler (as might be the case for routines registered as callbacks).
- Support is not provided for proper alignment of parameters.

STACK ALIGNMENT

The solution to this problem is to have the function's entry point assume only 4-byte alignment. If the function has a need for 8-byte or 16-byte alignment, then code can be inserted to dynamically align the stack appropriately, resulting in one of the stack frames shown in Figure D-1.

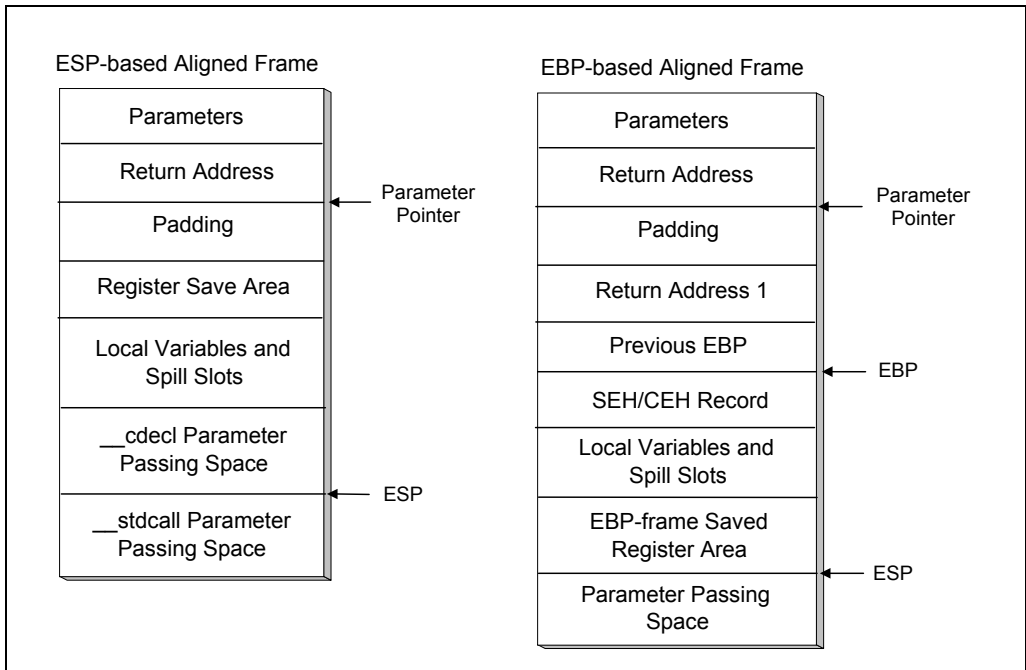


Figure D-1. Stack Frames Based on Alignment Type

As an optimization, an alternate entry point can be created that can be called when proper stack alignment is provided by the caller. Using call graph profiling of the VTune analyzer, calls to the normal (unaligned) entry point can be optimized into calls to the (alternate) aligned entry point when the stack can be proven to be properly aligned. Furthermore, a function alignment requirement attribute can be modified throughout the call graph so as to cause the least number of calls to unaligned entry points.

As an example of this, suppose function F has only a stack alignment requirement of 4, but it calls function G at many call sites, and in a loop. If G's alignment requirement is 16, then by promoting F's alignment requirement to 16, and making all calls to G go to its aligned entry point, the compiler can minimize the number of times that control passes through the unaligned entry points. Example D-1 and Example D-2 in the following sections illustrate this technique. Note the entry points `foo` and `foo.aligned`; the latter is the alternate aligned entry point.

D.4.1 Aligned ESP-Based Stack Frames

This section discusses data and parameter alignment and the `declspec(align)` extended attribute, which can be used to request alignment in C and C++ code. In creating ESP-based stack frames, the compiler adds padding between the return address and the register save area as shown in Example 4-16. This frame can be used only when debug information is not requested, there is no need for exception handling support, inlined assembly is not used, and there are no calls to `alloca` within the function.

If the above conditions are not met, an aligned EDP-based frame must be used. When using this type of frame, the sum of the sizes of the return address, saved registers, local variables, register spill slots, and parameter space must be a multiple of 16 bytes. This causes the base of the parameter space to be 16-byte aligned. In addition, any space reserved for passing parameters for `stdcall` functions also must be a multiple of 16 bytes. This means that the caller needs to clean up some of the stack space when the size of the parameters pushed for a call to a `stdcall` function is not a multiple of 16. If the caller does not do this, the stack pointer is not restored to its pre-call value.

In Example D-1, we have 12 bytes on the stack after the point of alignment from the caller: the return pointer, EBX and EDX. Thus, we need to add four more to the stack pointer to achieve alignment. Assuming 16 bytes of stack space are needed for local variables, the compiler adds $16 + 4 = 20$ bytes to ESP, making ESP aligned to a 0 mod 16 address.

Example D-1. Aligned esp-Based Stack Frame

```
void _cdecl foo (int k)
{
  int j;
  foo: // See Note A below
    push  ebx
    mov   ebx, esp

    sub   esp, 0x00000008
    and   esp, 0xfffffff0
    add   esp, 0x00000008
    jmp   common
foo.aligned:
  push  ebx
  mov   ebx, esp
```

Example D-1. Aligned esp-Based Stack Frame (Contd.)

```

common:                // See Note B below
    push    edx
    sub     esp, 20
    j = k;
    mov     edx, [ebx + 8]
    mov     [esp + 16], edx

foo(5);
    mov     [esp], 5
    call   foo.aligned
return j;
    mov     eax, [esp + 16]
    add     esp, 20
    pop     edx
    mov     esp, ebx
    pop     ebx
    ret

// NOTES:
// (A) Aligned entry points assume that parameter block beginnings are aligned. This places the
// stack pointer at a 12 mod 16 boundary, as the return pointer has been pushed. Thus, the
// unaligned entry point must force the stack pointer to this boundary
// (B) The code at the common label assumes the stack is at an 8 mod 16 boundary, and adds
// sufficient space to the stack so that the stack pointer is aligned to a 0 mod 16 boundary.

```

D.4.2 Aligned EDP-Based Stack Frames

In EDP-based frames, padding is also inserted immediately before the return address. However, this frame is slightly unusual in that the return address may actually reside in two different places in the stack. This occurs whenever padding must be added and exception handling is in effect for the function. Example D-2 shows the code generated for this type of frame. The stack location of the return address is aligned 12 mod 16. This means that the value of EDP always satisfies the condition $(EDP \& 0x0f) == 0x08$. In this case, the sum of the sizes of the return address, the previous EDP, the exception handling record, the local variables, and the spill area must be a multiple of 16 bytes. In addition, the parameter passing space must be a multiple of 16 bytes. For a call to a stdcall function, it is necessary for the caller to

reserve some stack space if the size of the parameter block being pushed is not a multiple of 16.

Example D-2. Aligned ebp-based Stack Frames

```

void _stdcall foo (int k)
{
  int j;
  foo:
    push   ebx
    mov    ebx, esp
    sub    esp, 0x00000008
    and    esp, 0xfffffff0
    add    esp, 0x00000008    // esp is (8 mod 16) after add
    jmp   common
  foo.aligned:
    push   ebx                // esp is (8 mod 16) after push
    mov    ebx, esp
  common:
    push   ebp                // this slot will be used for
                                // duplicate return pt
    push   ebp                // esp is (0 mod 16) after push
                                // (rtn,ebx,ebp,ebp)
    mov    ebp, [ebx + 4]     // fetch return pointer and store
                                // relative to ebp
                                // (rtn,ebx,rtn,ebp)
    mov    ebp, esp          // ebp is (0 mod 16)
    sub    esp, 28           // esp is (4 mod 16)
                                // see Note A below
    push   edx                // esp is (0 mod 16) after push
                                // goal is to make esp and ebp
                                // (0 mod 16) here

    j = k;
    mov    edx, [ebx + 8]     // k is (0 mod 16) if caller
                                // aligned its stack
    mov    [ebp - 16], edx    // j is (0 mod 16)
  foo(5);
    add    esp, -4           // normal call sequence to
                                // unaligned entry
    mov    [esp],5
    call  foo                // for stdcall, callee
                                // cleans up stack

```

Example D-2. Aligned ebp-based Stack Frames (Contd.)

```

foo.aligned(5);
    add    esp,-16                // aligned entry, this should
                                // be a multiple of 16

    mov   [esp],5
    call  foo.aligned
    add   esp,12                 // see Note B below

return j;
    mov   eax,[ebp-16]
    pop   edx
    mov   esp,ebp
    pop   ebp
    mov   esp,ebx
    pop   ebx
ret 4
}

// NOTES:
// (A) Here we allow for local variables. However, this value should be adjusted so that, after
// pushing the saved registers, esp is 0 mod 16.
// (B) Just prior to the call, esp is 0 mod 16. To maintain alignment, esp should be adjusted by 16.
// When a callee uses the stdcall calling sequence, the stack pointer is restored by the callee. The
// final addition of 12 compensates for the fact that only 4 bytes were passed, rather than
// 16, and thus the caller must account for the remaining adjustment.

```

D.4.3 Stack Frame Optimizations

The Intel C++ Compiler provides certain optimizations that may improve the way aligned frames are set up and used. These optimizations are as follows:

- If a procedure is defined to leave the stack frame 16-byte-aligned and it calls another procedure that requires 16-byte alignment, then the callee's aligned entry point is called, bypassing all of the unnecessary aligning code.
- If a static function requires 16-byte alignment, and it can be proven to be called only by other functions that require 16-byte alignment, then that function will not have any alignment code in it. That is, the compiler will not use EBX to point to the argument block and it will not have alternate entry points, because this function will never be entered with an unaligned frame.

D.5 INLINED ASSEMBLY AND EBX

When using aligned frames, the EBX register generally should not be modified in inlined assembly blocks since EBX is used to keep track of the argument block. Programmers may modify EBX only if they do not need to access the arguments and provided they save EBX and restore it before the end of the function (since ESP is restored relative to EBX in the function's epilog).

NOTE

Do not use the EBX register in inline assembly functions that use dynamic stack alignment for `double`, `__m64`, and `__m128` local variables unless you save and restore EBX each time you use it. The Intel C++ Compiler uses the EBX register to control alignment of variables of these types, so the use of EBX, without preserving it, will cause unexpected program execution.

STACK ALIGNMENT

APPENDIX E

SUMMARY OF RULES AND SUGGESTIONS

This appendix summarizes the rules and suggestions specified in this manual. Please be reminded that coding recommendations are ranked in importance according to these two criteria:

- Local impact (referred to earlier as “impact”) – the difference that a recommendation makes to performance for a given instance.
- Generality – how frequently such instances occur across all application domains.

Again, understand that this ranking is intentionally very approximate, and can vary depending on coding style, application domain, and other factors. Throughout the chapter you observed references to these criteria using the high, medium and low priorities for each recommendation. In places where there was no priority assigned, the local impact or generality has been determined not to be applicable.

E.1 ASSEMBLY/COMPILER CODING RULES

Assembler/Compiler Coding Rule 1. (MH impact, M generality) Arrange code to make basic blocks contiguous and eliminate unnecessary branches. 3-7

Assembler/Compiler Coding Rule 2. (M impact, ML generality) Use the SETCC and CMOV instructions to eliminate unpredictable conditional branches where possible. Do not do this for predictable branches. Do not use these instructions to eliminate all unpredictable conditional branches (because using these instructions will incur execution overhead due to the requirement for executing both paths of a conditional branch). In addition, converting a conditional branch to SETCC or CMOV trades off control flow dependence for data dependence and restricts the capability of the out-of-order engine. When tuning, note that all Intel 64 and IA-32 processors usually have very high branch prediction rates. Consistently mispredicted branches are generally rare. Use these instructions only if the increase in computation time is less than the expected cost of a mispredicted branch. 3-7

Assembler/Compiler Coding Rule 3. (M impact, H generality) Arrange code to be consistent with the static branch prediction algorithm: make the fall-through code following a conditional branch be the likely target for a branch with a forward target, and make the fall-through code following a conditional branch be the unlikely target for a branch with a backward target. 3-10

Assembler/Compiler Coding Rule 4. (MH impact, MH generality) Near calls must be matched with near returns, and far calls must be matched with far

- returns. Pushing the return address on the stack and jumping to the routine to be called is not recommended since it creates a mismatch in calls and returns. 3-12
- Assembler/Compiler Coding Rule 5. (MH impact, MH generality)** Selectively inline a function if doing so decreases code size or if the function is small and the call site is frequently executed. 3-12
- Assembler/Compiler Coding Rule 6. (H impact, H generality)** Do not inline a function if doing so increases the working set size beyond what will fit in the trace cache. 3-12
- Assembler/Compiler Coding Rule 7. (ML impact, ML generality)** If there are more than 16 nested calls and returns in rapid succession; consider transforming the program with inline to reduce the call depth. 3-12
- Assembler/Compiler Coding Rule 8. (ML impact, ML generality)** Favor inlining small functions that contain branches with poor prediction rates. If a branch misprediction results in a RETURN being prematurely predicted as taken, a performance penalty may be incurred.) 3-12
- Assembler/Compiler Coding Rule 9. (L impact, L generality)** If the last statement in a function is a call to another function, consider converting the call to a jump. This will save the call/return overhead as well as an entry in the return stack buffer. 3-12
- Assembler/Compiler Coding Rule 10. (M impact, L generality)** Do not put more than four branches in a 16-byte chunk. 3-12
- Assembler/Compiler Coding Rule 11. (M impact, L generality)** Do not put more than two end loop branches in a 16-byte chunk. 3-12
- Assembler/Compiler Coding Rule 12. (M impact, H generality)** All branch targets should be 16-byte aligned. 3-13
- Assembler/Compiler Coding Rule 13. (M impact, H generality)** If the body of a conditional is not likely to be executed, it should be placed in another part of the program. If it is highly unlikely to be executed and code locality is an issue, it should be placed on a different code page. 3-13
- Assembler/Compiler Coding Rule 14. (M impact, L generality)** When indirect branches are present, try to put the most likely target of an indirect branch immediately following the indirect branch. Alternatively, if indirect branches are common but they cannot be predicted by branch prediction hardware, then follow the indirect branch with a UD2 instruction, which will stop the processor from decoding down the fall-through path. 3-13
- Assembler/Compiler Coding Rule 15. (H impact, M generality)** Unroll small loops until the overhead of the branch and induction variable accounts (generally) for less than 10% of the execution time of the loop. 3-16
- Assembler/Compiler Coding Rule 16. (H impact, M generality)** Avoid unrolling loops excessively; this may thrash the trace cache or instruction cache. 3-16
- Assembler/Compiler Coding Rule 17. (M impact, M generality)** Unroll loops that are frequently executed and have a predictable number of iterations to reduce the number of iterations to 16 or fewer. Do this unless it increases code size so that the working set no longer fits in the trace or instruction cache. If the

- loop body contains more than one conditional branch, then unroll so that the number of iterations is 16/(# conditional branches)..... 3-16
- Assembler/Compiler Coding Rule 18. (ML impact, M generality)** For improving fetch/decode throughput, Give preference to memory flavor of an instruction over the register-only flavor of the same instruction, if such instruction can benefit from micro-fusion. 3-17
- Assembler/Compiler Coding Rule 19. (M impact, ML generality)** Employ macro-fusion where possible using instruction pairs that support macro-fusion. Prefer TEST over CMP if possible. Use unsigned variables and unsigned jumps when possible. Try to logically verify that a variable is non-negative at the time of comparison. Avoid CMP or TEST of MEM-IMM flavor when possible. However, do not add other instructions to avoid using the MEM-IMM flavor. 3-19
- Assembler/Compiler Coding Rule 20. (M impact, ML generality)** Software can enable macro fusion when it can be logically determined that a variable is non-negative at the time of comparison; use TEST appropriately to enable macro-fusion when comparing a variable with 0. 3-21
- Assembler/Compiler Coding Rule 21. (MH impact, MH generality)** Favor generating code using imm8 or imm32 values instead of imm16 values..... 3-23
- Assembler/Compiler Coding Rule 22. (M impact, ML generality)** Ensure instructions using 0xF7 opcode byte does not start at offset 14 of a fetch line; and avoid using these instruction to operate on 16-bit data, upcast short data to 32 bits. 3-24
- Assembler/Compiler Coding Rule 23. (MH impact, MH generality)** Break up a loop long sequence of instructions into loops of shorter instruction blocks of no more than the size of LSD. 3-25
- Assembler/Compiler Coding Rule 24. (MH impact, M generality)** Avoid unrolling loops containing LCP stalls, if the unrolled block exceeds the size of LSD. 3-25
- Assembler/Compiler Coding Rule 25. (M impact, M generality)** Avoid putting explicit references to ESP in a sequence of stack operations (POP, PUSH, CALL, RET). 3-28
- Assembler/Compiler Coding Rule 26. (ML impact, L generality)** Use simple instructions that are less than eight bytes in length. 3-28
- Assembler/Compiler Coding Rule 27. (M impact, MH generality)** Avoid using prefixes to change the size of immediate and displacement. 3-28
- Assembler/Compiler Coding Rule 28. (M impact, H generality)** Favor single-micro-operation instructions. Also favor instruction with shorter latencies. .. 3-29
- Assembler/Compiler Coding Rule 29. (M impact, L generality)** Avoid prefixes, especially multiple non-0F-prefixed opcodes. 3-29
- Assembler/Compiler Coding Rule 30. (M impact, L generality)** Do not use many segment registers. 3-29
- Assembler/Compiler Coding Rule 31. (ML impact, M generality)** Avoid using complex instructions (for example, enter, leave, or loop) that have more than

- four μ ops and require multiple cycles to decode. Use sequences of simple instructions instead. 3-29
- Assembler/Compiler Coding Rule 32. (M impact, H generality)** INC and DEC instructions should be replaced with ADD or SUB instructions, because ADD and SUB overwrite all flags, whereas INC and DEC do not, therefore creating false dependencies on earlier instructions that set the flags. 3-30
- Assembler/Compiler Coding Rule 33. (ML impact, L generality)** If an LEA instruction using the scaled index is on the critical path, a sequence with ADDs may be better. If code density and bandwidth out of the trace cache are the critical factor, then use the LEA instruction. 3-32
- Assembler/Compiler Coding Rule 34. (ML impact, L generality)** Avoid ROTATE by register or ROTATE by immediate instructions. If possible, replace with a ROTATE by 1 instruction. 3-33
- Assembler/Compiler Coding Rule 35. (M impact, ML generality)** Use dependency-breaking-idiom instructions to set a register to 0, or to break a false dependence chain resulting from re-use of registers. In contexts where the condition codes must be preserved, move 0 into the register instead. This requires more code space than using XOR and SUB, but avoids setting the condition codes. 3-35
- Assembler/Compiler Coding Rule 36. (M impact, MH generality)** Break dependences on portions of registers between instructions by operating on 32-bit registers instead of partial registers. For moves, this can be accomplished with 32-bit moves or by using MOVZX. 3-35
- Assembler/Compiler Coding Rule 37. (M impact, M generality)** Try to use zero extension or operate on 32-bit operands instead of using moves with sign extension. 3-36
- Assembler/Compiler Coding Rule 38. (ML impact, L generality)** Avoid placing instructions that use 32-bit immediates which cannot be encoded as sign-extended 16-bit immediates near each other. Try to schedule μ ops that have no immediate immediately before or after μ ops with 32-bit immediates. 3-36
- Assembler/Compiler Coding Rule 39. (ML impact, M generality)** Use the TEST instruction instead of AND when the result of the logical AND is not used. This saves μ ops in execution. Use a TEST of a register with itself instead of a CMP of the register to zero, this saves the need to encode the zero and saves encoding space. Avoid comparing a constant to a memory operand. It is preferable to load the memory operand and compare the constant to a register. 3-37
- Assembler/Compiler Coding Rule 40. (ML impact, M generality)** Eliminate unnecessary compare with zero instructions by using the appropriate conditional jump instruction when the flags are already set by a preceding arithmetic instruction. If necessary, use a TEST instruction instead of a compare. Be certain

that any code transformations made do not introduce problems with overflow.3-37

- Assembler/Compiler Coding Rule 41. (H impact, MH generality)** For small loops, placing loop invariants in memory is better than spilling loop-carried dependencies. 3-39
- Assembler/Compiler Coding Rule 42. (M impact, ML generality)** Avoid introducing dependences with partial floating point register writes, e.g. from the MOVSD XMMREG1, XMMREG2 instruction. Use the MOVAPD XMMREG1, XMMREG2 instruction instead. 3-45
- Assembler/Compiler Coding Rule 43. (ML impact, L generality)** Instead of using MOVUPD XMMREG1, MEM for a unaligned 128-bit load, use MOVSD XMMREG1, MEM; MOVSD XMMREG2, MEM+8; UNPCKLPD XMMREG1, XMMREG2. If the additional register is not available, then use MOVSD XMMREG1, MEM; MOVHPD XMMREG1, MEM+8..... 3-46
- Assembler/Compiler Coding Rule 44. (M impact, ML generality)** Instead of using MOVUPD MEM, XMMREG1 for a store, use MOVSD MEM, XMMREG1; UNPCKHPD XMMREG1, XMMREG1; MOVSD MEM+8, XMMREG1 instead..... 3-46
- Assembler/Compiler Coding Rule 45. (H impact, H generality)** Align data on natural operand size address boundaries. If the data will be accessed with vector instruction loads and stores, align the data on 16-byte boundaries. 3-60
- Assembler/Compiler Coding Rule 46. (H impact, M generality)** Pass parameters in registers instead of on the stack where possible. Passing arguments on the stack requires a store followed by a reload. While this sequence is optimized in hardware by providing the value to the load directly from the memory order buffer without the need to access the data cache if permitted by store-forwarding restrictions, floating point values incur a significant latency in forwarding. Passing floating point arguments in (preferably XMM) registers should save this long latency operation. 3-63
- Assembler/Compiler Coding Rule 47. (H impact, M generality)** A load that forwards from a store must have the same address start point and therefore the same alignment as the store data. 3-64
- Assembler/Compiler Coding Rule 48. (H impact, M generality)** The data of a load which is forwarded from a store must be completely contained within the store data. 3-64
- Assembler/Compiler Coding Rule 49. (H impact, ML generality)** If it is necessary to extract a non-aligned portion of stored data, read out the smallest aligned portion that completely contains the data and shift/mask the data as necessary. This is better than incurring the penalties of a failed store-forward.3-65
- Assembler/Compiler Coding Rule 50. (MH impact, ML generality)** Avoid several small loads after large stores to the same area of memory by using a single large read and register copies as needed..... 3-65
- Assembler/Compiler Coding Rule 51. (H impact, MH generality)** Where it is possible to do so without incurring other penalties, prioritize the allocation of

- variables to registers, as in register allocation and for parameter passing, to minimize the likelihood and impact of store-forwarding problems. Try not to store-forward data generated from a long latency instruction - for example, MUL or DIV. Avoid store-forwarding data for variables with the shortest store-load distance. Avoid store-forwarding data for variables with many and/or long dependence chains, and especially avoid including a store forward on a loop-carried dependence chain.3-68
- Assembler/Compiler Coding Rule 52. (M impact, MH generality)** Calculate store addresses as early as possible to avoid having stores block loads.3-69
- Assembler/Compiler Coding Rule 53. (H impact, M generality)** Try to arrange data structures such that they permit sequential access.3-71
- Assembler/Compiler Coding Rule 54. (H impact, M generality)** If 64-bit data is ever passed as a parameter or allocated on the stack, make sure that the stack is aligned to an 8-byte boundary.3-72
- Assembler/Compiler Coding Rule 55. (H impact, M generality)** Avoid having a store followed by a non-dependent load with addresses that differ by a multiple of 4 KBytes. Also, lay out data or order computation to avoid having cache lines that have linear addresses that are a multiple of 64 KBytes apart in the same working set. Avoid having more than 4 cache lines that are some multiple of 2 KBytes apart in the same first-level cache working set, and avoid having more than 8 cache lines that are some multiple of 4 KBytes apart in the same first-level cache working set.3-74
- Assembler/Compiler Coding Rule 56. (M impact, L generality)** If (hopefully read-only) data must occur on the same page as code, avoid placing it immediately after an indirect jump. For example, follow an indirect jump with its mostly likely target, and place the data after an unconditional branch.3-75
- Assembler/Compiler Coding Rule 57. (H impact, L generality)** Always put code and data on separate pages. Avoid self-modifying code wherever possible. If code is to be modified, try to do it all at once and make sure the code that performs the modifications and the code being modified are on separate 4-KByte pages or on separate aligned 1-KByte subpages.3-75
- Assembler/Compiler Coding Rule 58. (H impact, L generality)** If an inner loop writes to more than four arrays (four distinct cache lines), apply loop fission to break up the body of the loop such that only four arrays are being written to in each iteration of each of the resulting loops.3-77
- Assembler/Compiler Coding Rule 59. (H impact, M generality)** Minimize changes to bits 8-12 of the floating point control word. Changes for more than two values (each value being a combination of the following bits: precision, rounding and infinity control, and the rest of bits in FCW) leads to delays that are on the order of the pipeline depth.....3-94
- Assembler/Compiler Coding Rule 60. (H impact, L generality)** Minimize the number of changes to the rounding mode. Do not use changes in the rounding

- mode to implement the floor and ceiling functions if this involves a total of more than two values of the set of rounding, precision, and infinity bits. 3-96
- Assembler/Compiler Coding Rule 61. (H impact, L generality)** Minimize the number of changes to the precision mode. 3-97
- Assembler/Compiler Coding Rule 62. (M impact, M generality)** Use FXCH only where necessary to increase the effective name space. 3-97
- Assembler/Compiler Coding Rule 63. (M impact, M generality)** Use Streaming SIMD Extensions 2 or Streaming SIMD Extensions unless you need an x87 feature. Most SSE2 arithmetic operations have shorter latency than their X87 counterpart and they eliminate the overhead associated with the management of the X87 register stack. 3-98
- Assembler/Compiler Coding Rule 64. (M impact, L generality)** Try to use 32-bit operands rather than 16-bit operands for FILD. However, do not do so at the expense of introducing a store-forwarding problem by writing the two halves of the 32-bit memory operand separately. 3-99
- Assembler/Compiler Coding Rule 65. (H impact, M generality)** Use the 32-bit versions of instructions in 64-bit mode to reduce code size unless the 64-bit version is necessary to access 64-bit data or additional registers. 9-2
- Assembler/Compiler Coding Rule 66. (M impact, MH generality)** When they are needed to reduce register pressure, use the 8 extra general purpose registers for integer code and 8 extra XMM registers for floating-point or SIMD code. ... 9-2
- Assembler/Compiler Coding Rule 67. (ML impact, M generality)** Prefer 64-bit by 64-bit integer multiplies that produce 64-bit results over multiplies that produce 128-bit results. 9-2
- Assembler/Compiler Coding Rule 68. (ML impact, M generality)** Stagger accessing the high 64-bit result of a 128-bit multiply after accessing the low 64-bit results. 9-2
- Assembler/Compiler Coding Rule 69. (ML impact, M generality)** Use the 64-bit versions of multiply for 32-bit integer multiplies that require a 64 bit result. 9-7
- Assembler/Compiler Coding Rule 70. (ML impact, M generality)** Use the 64-bit versions of add for 64-bit adds. 9-8
- Assembler/Compiler Coding Rule 71. (L impact, L generality)** If software prefetch instructions are necessary, use the prefetch instructions provided by SSE. 9-8
- Assembler/Compiler Coding Rule 72. (H impact, H generality)** Whenever 256-bit AVX code and 128-bit SSE code might execute together, use the VZERoupper instruction whenever a transition from "Modified/Unsaved" state is expected. 11-9
- Assembler/Compiler Coding Rule 73. (H impact, H generality)** Add VZERoupper instruction after 256-bit AVX instructions are executed and before

any function call that might execute SSE code. Add VZEROUPPER at the end of any function that uses 256-bit AVX instructions. 11-10

Assembler/Compiler Coding Rule 74. (H impact, M generality) Align data to 32-byte boundary when possible. Prefer store alignment over load alignment.... 11-23

Assembler/Compiler Coding Rule 75. (M impact, H generality) Align data to 32-byte boundary when possible. Prefer store alignment over load alignment.... 11-25

Assembler/Compiler Coding Rule 76. (M impact, M generality) Use Blend instructions in lieu of shuffle instruction in AVX whenever possible. 11-36

E.2 USER/SOURCE CODING RULES

User/Source Coding Rule 1. (M impact, L generality) If an indirect branch has two or more common taken targets and at least one of those targets is correlated with branch history leading up to the branch, then convert the indirect branch to a tree where one or more indirect branches are preceded by conditional branches to those targets. Apply this “peeling” procedure to the common target of an indirect branch that correlates to branch history 3-14

User/Source Coding Rule 2. (H impact, M generality) Use the smallest possible floating-point or SIMD data type, to enable more parallelism with the use of a (longer) SIMD vector. For example, use single precision instead of double precision where possible. 3-47

User/Source Coding Rule 3. (M impact, ML generality) Arrange the nesting of loops so that the innermost nesting level is free of inter-iteration dependencies. Especially avoid the case where the store of data in an earlier iteration happens

- lexically after the load of that data in a future iteration, something which is called a lexically backward dependence. 3-47
- User/Source Coding Rule 4. (M impact, ML generality)** Avoid the use of conditional branches inside loops and consider using SSE instructions to eliminate branches 3-47
- User/Source Coding Rule 5. (M impact, ML generality)** Keep induction (loop) variable expressions simple 3-47
- User/Source Coding Rule 6. (H impact, M generality)** Pad data structures defined in the source code so that every data element is aligned to a natural operand size address boundary 3-69
- User/Source Coding Rule 7. (M impact, L generality)** Beware of false sharing within a cache line (64 bytes) and within a sector of 128 bytes on processors based on Intel NetBurst microarchitecture 3-71
- User/Source Coding Rule 8. (H impact, ML generality)** Consider using a special memory allocation library with address offset capability to avoid aliasing. .. 3-74
- User/Source Coding Rule 9. (M impact, M generality)** When padding variable declarations to avoid aliasing, the greatest benefit comes from avoiding aliasing on second-level cache lines, suggesting an offset of 128 bytes or more 3-74
- User/Source Coding Rule 10. (H impact, H generality)** Optimization techniques such as blocking, loop interchange, loop skewing, and packing are best done by the compiler. Optimize data structures either to fit in one-half of the first-level cache or in the second-level cache; turn on loop optimizations in the compiler to enhance locality for nested loops 3-79
- User/Source Coding Rule 11. (M impact, ML generality)** If there is a blend of reads and writes on the bus, changing the code to separate these bus transactions into read phases and write phases can help performance 3-79
- User/Source Coding Rule 12. (H impact, H generality)** To achieve effective amortization of bus latency, software should favor data access patterns that result in higher concentrations of cache miss patterns, with cache miss strides that are significantly smaller than half the hardware prefetch trigger threshold . 3-79
- User/Source Coding Rule 13. (M impact, M generality)** Enable the compiler's use of SSE, SSE2 or SSE3 instructions with appropriate switches 3-90
- User/Source Coding Rule 14. (H impact, ML generality)** Make sure your application stays in range to avoid denormal values, underflows. 3-91
- User/Source Coding Rule 15. (M impact, ML generality)** Do not use double precision unless necessary. Set the precision control (PC) field in the x87 FPU control word to "Single Precision". This allows single precision (32-bit) computation to complete faster on some operations (for example, divides due to early out). However, be careful of introducing more than a total of two values for the floating point control word, or there will be a large performance penalty. See Section 3.8.3 3-91
- User/Source Coding Rule 16. (H impact, ML generality)** Use fast float-to-int routines, FISTTP, or SSE2 instructions. If coding these routines, use the FISTTP

- instruction if SSE3 is available, or the CVTTSS2SI and CVTTSD2SI instructions if coding with Streaming SIMD Extensions 2.3-91
- User/Source Coding Rule 17. (M impact, ML generality)** Removing data dependence enables the out-of-order engine to extract more ILP from the code. When summing up the elements of an array, use partial sums instead of a single accumulator.3-91
- User/Source Coding Rule 18. (M impact, ML generality)** Usually, math libraries take advantage of the transcendental instructions (for example, FSIN) when evaluating elementary functions. If there is no critical need to evaluate the transcendental functions using the extended precision of 80 bits, applications should consider an alternate, software-based approach, such as a look-up-table-based algorithm using interpolation techniques. It is possible to improve transcendental performance with these techniques by choosing the desired numeric precision and the size of the look-up table, and by taking advantage of the parallelism of the SSE and the SSE2 instructions.3-91
- User/Source Coding Rule 19. (H impact, ML generality)** Denormalized floating-point constants should be avoided as much as possible3-92
- User/Source Coding Rule 20. (M impact, H generality)** Insert the PAUSE instruction in fast spin loops and keep the number of loop repetitions to a minimum to improve overall system performance.8-17
- User/Source Coding Rule 21. (M impact, L generality)** Replace a spin lock that may be acquired by multiple threads with pipelined locks such that no more than two threads have write accesses to one lock. If only one thread needs to write to a variable shared by two threads, there is no need to use a lock.8-18
- User/Source Coding Rule 22. (H impact, M generality)** Use a thread-blocking API in a long idle loop to free up the processor8-19
- User/Source Coding Rule 23. (H impact, M generality)** Beware of false sharing within a cache line (64 bytes on Intel Pentium 4, Intel Xeon, Pentium M, Intel Core Duo processors), and within a sector (128 bytes on Pentium 4 and Intel Xeon processors)8-21
- User/Source Coding Rule 24. (M impact, ML generality)** Place each synchronization variable alone, separated by 128 bytes or in a separate cache line.8-22
- User/Source Coding Rule 25. (H impact, L generality)** Do not place any spin lock variable to span a cache line boundary8-22
- User/Source Coding Rule 26. (M impact, H generality)** Improve data and code locality to conserve bus command bandwidth.8-24
- User/Source Coding Rule 27. (M impact, L generality)** Avoid excessive use of software prefetch instructions and allow automatic hardware prefetcher to work.

- Excessive use of software prefetches can significantly and unnecessarily increase bus utilization if used inappropriately. 8-25
- User/Source Coding Rule 28. (M impact, M generality)** Consider using overlapping multiple back-to-back memory reads to improve effective cache miss latencies. 8-26
- User/Source Coding Rule 29. (M impact, M generality)** Consider adjusting the sequencing of memory references such that the distribution of distances of successive cache misses of the last level cache peaks towards 64 bytes. ... 8-26
- User/Source Coding Rule 30. (M impact, M generality)** Use full write transactions to achieve higher data throughput. 8-26
- User/Source Coding Rule 31. (H impact, H generality)** Use cache blocking to improve locality of data access. Target one quarter to one half of the cache size when targeting Intel processors supporting HT Technology or target a block size that allow all the logical processors serviced by a cache to share that cache simultaneously. 8-27
- User/Source Coding Rule 32. (H impact, M generality)** Minimize the sharing of data between threads that execute on different bus agents sharing a common bus. The situation of a platform consisting of multiple bus domains should also minimize data sharing across bus domains 8-28
- User/Source Coding Rule 33. (H impact, H generality)** Minimize data access patterns that are offset by multiples of 64 KBytes in each thread. 8-30
- User/Source Coding Rule 34. (M impact, L generality)** Avoid excessive loop unrolling to ensure the LSD is operating efficiently. 8-30

E.3 TUNING SUGGESTIONS

- Tuning Suggestion 1.** In rare cases, a performance problem may be caused by executing data on a code page as instructions. This is very likely to happen when execution is following an indirect branch that is not resident in the trace cache. If this is clearly causing a performance problem, try moving the data elsewhere, or inserting an illegal opcode or a pause instruction immediately after the indirect branch. Note that the latter two alternatives may degrade performance in some circumstances. 3-75
- Tuning Suggestion 2.** If a load is found to miss frequently, either insert a prefetch before it or (if issue bandwidth is a concern) move the load up to execute earlier. 3-83
- Tuning Suggestion 3.** Optimize single threaded code to maximize execution throughput first. 8-35
- Tuning Suggestion 4.** Employ efficient threading model, leverage available tools (such as Intel Threading Building Block, Intel Thread Checker, Intel Thread

Profiler) to achieve optimal processor scaling with respect to the number of physical processors or processor cores.8-35

E.4 SSE4.2 CODING RULES

SSE4.2 Coding Rule 1. (H impact, H generality) Loop-carry dependency that depends on the ECX result of PCMPESTRI/PCMPESTRM/PCMPISTRI/PCMPISTRM for address adjustment must be minimized. Isolate code paths that expect ECX result will be 16 (bytes) or 8 (words), replace these values of ECX with constants in address adjustment expressions to take advantage of memory disambiguation hardware. 10-12

E.5 ASSEMBLY/COMPILER CODING RULES FOR THE INTEL® ATOM™ PROCESSOR

Assembly/Compiler Coding Rule 1. (MH impact, ML generality) For Intel Atom processors, minimize the presence of complex instructions requiring MSROM to take advantage the optimal decode bandwidth provided by the two decode units. 12-4

Assembly/Compiler Coding Rule 2. (M impact, H generality) For Intel Atom processors, keeping the instruction working set footprint small will help the front end to take advantage the optimal decode bandwidth provided by the two decode units. 12-4

Assembly/Compiler Coding Rule 3. (MH impact, ML generality) For Intel Atom processors, avoiding back-to-back X87 instructions will help the front end to take advantage the optimal decode bandwidth provided by the two decode units. 12-4

Assembly/Compiler Coding Rule 4. (M impact, H generality) For Intel Atom processors, place a MOV instruction between a flag producer instruction and a flag consumer instruction that would have incurred a two-cycle delay. This will prevent partial flag dependency. 12-7

Assembly/Compiler Coding Rule 5. (MH impact, H generality) For Intel Atom processors, LEA should be used for address manipulation; but software should avoid the following situations which creates dependencies from ALU to AGU: an ALU instruction (instead of LEA) for address manipulation or ESP updates; a LEA for ternary addition or non-destructive writes which do not feed address

- generation. Alternatively, hoist producer instruction more than 3 cycles above the consumer instruction that uses the AGU. 12-8
- Assembly/Compiler Coding Rule 6. (M impact, M generality)** For Intel Atom processors, sequence an independent FP or integer multiply after an integer multiply instruction to take advantage of pipelined IMUL execution. 12-9
- Assembly/Compiler Coding Rule 7. (M impact, M generality)** For Intel Atom processors, hoist the producer instruction for the implicit register count of an integer shift instruction before the shift instruction by at least two cycles. . 12-9
- Assembly/Compiler Coding Rule 8. (M impact, MH generality)** For Intel Atom processors, LEA, simple loads and POP are slower if the input is smaller than 4 bytes. 12-9
- Assembly/Compiler Coding Rule 9. (MH impact, H generality)** For Intel Atom processors, prefer SIMD instructions operating on XMM register over X87 instructions using FP stack. Use Packed single-precision instructions where possible. Replace packed double-precision instruction with scalar double-precision instructions. 12-11
- Assembly/Compiler Coding Rule 10. (M impact, ML generality)** For Intel Atom processors, library software performing sophisticated math operations like transcendental functions should use SIMD instructions operating on XMM register instead of native X87 instructions. 12-11
- Assembly/Compiler Coding Rule 11. (M impact, M generality)** For Intel Atom processors, enable DAZ and FTZ whenever possible. 12-11
- Assembly/Compiler Coding Rule 12. (H impact, L generality)** For Intel Atom processors, use divide instruction only when it is absolutely necessary, and pay attention to use the smallest data size operand. 12-12
- Assembly/Compiler Coding Rule 13. (MH impact, M generality)** For Intel Atom processors, prefer a sequence MOVAPS+PALIGN over MOVUPS. Similarly, MOVDQA+PALIGNR is preferred over MOVDQU. 12-12
- Assembly/Compiler Coding Rule 14. (MH impact, H generality)** For Intel Atom processors, ensure data are aligned in memory to its natural size. For example, 4-byte data should be aligned to 4-byte boundary, etc. Additionally, smaller access (less than 4 bytes) within a chunk may experience delay if they touch different bytes. 12-13
- Assembly/Compiler Coding Rule 15. (H impact, ML generality)** For Intel Atom processors, use segments with base set to 0 whenever possible; avoid non-zero

segment base address that is not aligned to cache line boundary at all cost. 12-14

Assembly/Compiler Coding Rule 16. (H impact, L generality) For Intel Atom processors, when using non-zero segment bases, Use DS, FS, GS; string operation should use implicit ES. 12-14

Assembly/Compiler Coding Rule 17. (M impact, ML generality) For Intel Atom processors, favor using ES, DS, SS over FS, GS with zero segment base. 12-14

Assembly/Compiler Coding Rule 18. (MH impact, M generality) For Intel Atom processors, "bool" and "char" value should be passed onto and read off the stack as 32-bit data. 12-15

Assembly/Compiler Coding Rule 19. (MH impact, M generality) For Intel Atom processors, favor register form of PUSH/POP and avoid using LEAVE; Use LEA to adjust ESP instead of ADD/SUB. 12-15

Numerics

- 64-bit mode
 - arithmetic, 9-3
 - coding guidelines, 9-1
 - compiler settings, A-2
 - CVTSI2SD instruction, 9-4
 - CVTSI2SS instruction, 9-4
 - default operand size, 9-1
 - introduction, 2-60
 - legacy instructions, 9-1
 - multiplication notes, 9-2
 - register usage, 9-2, 9-3
 - REX prefix, 9-1
 - sign-extension, 9-2
 - software prefetch, 9-5

A

- absolute difference of signed numbers, 5-21
- absolute difference of unsigned numbers, 5-20
- absolute value, 5-21
- active power, 11-1
- ADDSUBPD instruction, 6-14
- ADDSUBPS instruction, 6-14, 6-16
- algorithm to avoid changing the rounding mode, 3-82
- alignment
 - arrays, 3-56
 - code, 3-12
 - stack, 3-59
 - structures, 3-56
- Amdahl's law, 8-2
- AoS format, 4-21
- application performance tools, A-1
- arrays
 - aligning, 3-56
- assembler/compiler coding rules, E-1
- automatic vectorization, 4-13, 4-14

B

- battery life
 - guidelines for extending, 11-5
 - mobile optimization, 11-1
 - OS APIs, 11-6
 - quality trade-offs, 11-5
- branch prediction
 - choosing types, 3-13
 - code examples, 3-8
 - eliminating branches, 3-7
 - optimizing, 3-6
 - unrolling loops, 3-15

C

- C4-state, 11-4
- cache management
 - blocking techniques, 7-23
 - cache level, 7-5
 - CLFLUSH instruction, 7-12
 - coding guidelines, 7-1
 - compiler choices, 7-2
 - compiler intrinsics, 7-2
 - CPUID instruction, 3-5, 7-38
 - function leaf, 3-5
 - optimizing, 7-1
 - simple memory copy, 7-33
 - smart cache, 2-50
 - video decoder, 7-32
 - video encoder, 7-32
 - See also: optimizing cache utilization
- call graph profiling, A-12
- CD/DVD, 11-7
- changing the rounding mode, 3-82
- classes (C/C++), 4-12
- CLFLUSH instruction, 7-12
- clipping to an arbitrary signed range, 5-25
- clipping to an arbitrary unsigned range, 5-27
- coding techniques, 4-8, 8-23
 - 64-bit guidelines, 9-1
 - absolute difference of signed numbers, 5-21
 - absolute difference of unsigned numbers, 5-20
 - absolute value, 5-21
 - clipping to an arbitrary signed range, 5-25
 - clipping to an arbitrary unsigned range, 5-27
 - conserving power, 11-7
 - data in segment, 3-63
 - generating constants, 5-19
 - interleaved pack with saturation, 5-8
 - interleaved pack without saturation, 5-10
 - latency and throughput, C-1
 - methodologies, 4-9
 - non-interleaved unpack, 5-10
 - optimization options, A-2
 - rules, 3-5, E-1
 - signed unpack, 5-7
 - simplified clip to arbitrary signed range, 5-26
 - sleep transitions, 11-7
 - suggestions, 3-5, E-1
 - summary of rules, E-1
 - tuning hints, 3-5, E-1
 - unsigned unpack, 5-6
 - See also: floating-point code
- coherent requests, 7-9
- command-line options
 - floating-point arithmetic precision, A-6
 - inline expansion of library functions, A-6

INDEX

- rounding control, A-6
- vectorizer switch, A-5
- comparing register values, 3-28, 3-30
- compatibility mode, 9-1
- compatibility model, 2-60
- compiler intrinsics
 - _mm_load, 7-2, 7-32
 - _mm_prefetch, 7-2, 7-32
 - _mm_stream, 7-2, 7-32
- compilers
 - branch prediction support, 3-16
 - documentation, 1-4
 - general recommendations, 3-2
 - plug-ins, A-2
 - supported alignment options, 4-17
 - See also: Intel C++ Compiler & Intel Fortran Compiler
- computation
 - intensive code, 4-7
- converting 64-bit to 128-bit SIMD integers, 5-43
- converting code to MMX technology, 4-5
- CPUID instruction
 - AP-485, 1-4
 - cache parameters, 7-38
 - function leaf, 7-38
 - function leaf 4, 3-5
 - Intel compilers, 3-4
 - MMX support, 4-2
 - SSE support, 4-2
 - SSE2 support, 4-3
 - SSE3 support, 4-3
 - SSSE3 support, 4-4
 - strategy for use, 3-4
- C-states, 11-1, 11-3
- CVTSI2SD instruction, 9-4
- CVTSI2SS instruction, 9-4
- CVTTSP2PI instruction, 6-13
- CVTTSS2SI instruction, 6-13

D

- data
 - access pattern of array, 3-58
 - aligning arrays, 3-56
 - aligning structures, 3-56
 - alignment, 4-14
 - arrangement, 6-3
 - code segment and, 3-63
 - deswizzling, 6-9
 - prefetching, 2-51
 - swizzling, 6-6
 - swizzling using intrinsics, 6-7
- declspec(align), D-3
- deeper sleep, 11-4
- denormals-are-zero (DAZ), 6-13
- deterministic cache parameters
 - cache sharing, 7-38, 7-40
 - multicore, 7-40

- overview, 7-38
- prefetch stride, 7-40
- domain decomposition, 8-5
- Dual-core Intel Xeon processors, 2-1
- Dynamic execution, 2-21

E

- EDP-based stack frames, D-4
- eliminating branches, 3-9
- EMMS instruction, 5-2, 5-3
 - guidelines for using, 5-3
- Enhanced Intel SpeedStep Technology
 - description of, 11-8
 - multicore processors, 11-11
 - usage scenario, 11-2
- ESP-based stack frames, D-3
- extract word instruction, 5-12

F

- fencing operations, 7-7
 - LFENCE instruction, 7-11
 - MFENCE instruction, 7-12
- FIST instruction, 3-82
- FLDCW instruction, 3-82
- floating-point code
 - arithmetic precision options, A-6
 - data arrangement, 6-3
 - data deswizzling, 6-9
 - data swizzling using intrinsics, 6-7
 - guidelines for optimizing, 3-77
 - horizontal ADD, 6-11
 - improving parallelism, 3-84
 - memory access stall information, 3-53
 - operations with integer operands, 3-86
 - operations, integer operands, 3-86
 - optimizing, 3-77
 - planning considerations, 6-1
 - rules and suggestions, 6-1
 - scalar code, 6-2
 - transcendental functions, 3-86
 - unrolling loops, 3-15
 - vertical versus horizontal computation, 6-3
 - See also: coding techniques
- flush-to-zero (FTZ), 6-13
- front end
 - branching ratios, B-49
 - characterizing mispredictions, B-50
 - key practices, 8-13
 - loop unrolling, 8-13, 8-30
 - optimization, 3-6
 - Pentium M processor, 3-24
 - trace cache, 8-13
- functional decomposition, 8-5
- FXCH instruction, 3-85, 6-2

G

generating constants, 5-19
 GetActivePwrScheme, 11-6
 GetSystemPowerStatus, 11-6

H

HADDPD instruction, 6-14
 HADDPS instruction, 6-14, 6-18
 hardware multithreading
 support for, 3-5
 hardware prefetch
 cache blocking techniques, 7-27
 description of, 7-3
 latency reduction, 7-14
 memory optimization, 7-13
 operation, 7-13
 horizontal computations, 6-10
 hotspots
 definition of, 4-7
 identifying, 4-7
 VTune analyzer, 4-7
 HSUBPD instruction, 6-14
 HSUBPS instruction, 6-14, 6-18
 Hyper-Threading Technology
 avoid excessive software prefetches, 8-25
 bus optimization, 8-12
 cache blocking technique, 8-27
 conserve bus command bandwidth, 8-23
 eliminate 64-K-aliased data accesses, 8-29
 excessive loop unrolling, 8-30
 front-end optimization, 8-30
 full write transactions, 8-26
 functional decomposition, 8-5
 improve effective latency of cache misses, 8-25
 memory optimization, 8-26
 minimize data sharing between physical
 processors, 8-27
 multitasking environment, 8-3
 optimization, 8-1
 optimization guidelines, 8-11
 optimization with spin-locks, 8-18
 overview, 2-52
 parallel programming models, 8-5
 pipeline, 2-55
 placement of shared synchronization variable,
 8-21
 prevent false-sharing of data, 8-21
 processor resources, 2-53
 shared execution resources, 8-35
 shared-memory optimization, 8-27
 synchronization for longer periods, 8-18
 synchronization for short periods, 8-16
 system bus optimization, 8-23
 thread sync practices, 8-12
 thread synchronization, 8-14
 tools for creating multithreaded applications, 8-10

I

IA-32e mode, 2-60
 IA32_PERFVSELx MSR, B-48
 increasing bandwidth
 memory fills, 5-39
 video fills, 5-39
 indirect branch, 3-13
 inline assembly, 5-4
 inline expansion library functions option, A-6
 inlined-asm, 4-10
 insert word instruction, 5-13
 instruction latency/throughput
 overview, C-1
 instruction scheduling, 3-63
 Intel 64 and IA-32 processors, 2-1
 Intel 64 architecture
 and IA-32 processors, 2-60
 features of, 2-60
 IA-32e mode, 2-60
 Intel Advanced Digital Media Boost, 2-3
 Intel Advanced Memory Access, 2-13
 Intel Advanced Smart Cache, 2-2, 2-18
 Intel Core Duo processors, 2-1, 2-50
 128-bit integers, 5-44
 data prefetching, 2-51
 front end, 2-51
 microarchitecture, 2-50
 packed FP performance, 6-18
 performance events, B-39
 prefetch mechanism, 7-3
 processor perspectives, 3-3
 shared cache, 2-58
 SIMD support, 4-1
 special programming models, 8-6
 static prediction, 3-9
 Intel Core microarchitecture, 2-1, 2-2
 advanced smart cache, 2-18
 branch prediction unit, 2-6
 event ratios, B-47
 execution core, 2-9
 execution units, 2-10
 issue ports, 2-10
 front end, 2-4
 instruction decode, 2-8
 instruction fetch unit, 2-6
 instruction queue, 2-7
 advanced memory access, 2-13
 micro-fusion, 2-9
 pipeline overview, 2-3
 special programming models, 8-6
 stack pointer tracker, 2-8
 static prediction, 3-11
 Intel Core Solo processors, 2-1
 128-bit SIMD integers, 5-44
 data prefetching, 2-51
 front end, 2-51
 microarchitecture, 2-50
 performance events, B-39

INDEX

- prefetch mechanism, 7-3
 - processor perspectives, 3-3
 - SIMD support, 4-1
 - static prediction, 3-9
 - Intel Core2 Duo processors, 2-1
 - processor perspectives, 3-3
 - Intel C++ Compiler, 3-1
 - 64-bit mode settings, A-2
 - branch prediction support, 3-16
 - description, A-1
 - IA-32 settings, A-2
 - multithreading support, A-5
 - OpenMP, A-5
 - optimization settings, A-2
 - related information, 1-3
 - stack frame support, D-1
 - Intel Debugger
 - description, A-1
 - Intel developer link, 1-4
 - Intel Enhanced Deeper Sleep
 - C-state numbers, 11-3
 - enabling, 11-10
 - multiple-cores, 11-13
 - Intel Fortran Compiler
 - description, A-1
 - multithreading support, A-5
 - OpenMP, A-5
 - optimization settings, A-2
 - related information, 1-3
 - Intel Integrated Performance Primitives
 - for Linux, A-14
 - for Windows, A-14
 - Intel Math Kernel Library for Linux, A-13
 - Intel Math Kernel Library for Windows, A-13
 - Intel Mobile Platform SDK, 11-6
 - Intel NetBurst microarchitecture, 2-1
 - core, 2-37, 2-40
 - design goals, 2-34
 - front end, 2-36
 - introduction, 2-21, 2-33
 - out-of-order core, 2-40
 - pipeline, 2-35, 2-38
 - prefetch characteristics, 7-3
 - processor perspectives, 3-3
 - retirement, 2-37
 - trace cache, 3-11
 - Intel Pentium D processors, 2-1, 2-56
 - Intel Pentium M processors, 2-1
 - core, 2-50
 - data prefetching, 2-49
 - front end, 2-48
 - microarchitecture, 2-47
 - retirement, 2-50
 - Intel Performance Libraries, A-13
 - benefits, A-14, A-18
 - optimizations, A-14
 - Intel performance libraries
 - description, A-1
 - Intel Performance Tools, 3-1, A-1
 - Intel Smart Cache, 2-50
 - Intel Smart Memory Access, 2-2
 - Intel software network link, 1-4
 - Intel Thread Checker, 8-11
 - example output, A-15, A-16, A-17
 - Intel Thread Profiler
 - Intel Threading Tools, 8-11
 - Intel Threading Tools, A-15, A-17
 - Intel VTune Performance Analyzer
 - call graph, A-12
 - code coach, 4-7
 - coverage, 3-2
 - description, A-1
 - related information, 1-4
 - Intel Wide Dynamic Execution, 2-2, 2-3, 2-21
 - interleaved pack with saturation, 5-8
 - interleaved pack without saturation, 5-10
 - interprocedural optimization, A-6
 - introduction
 - chapter summaries, 1-2
 - optimization features, 2-1
 - processors covered, 1-1
 - references, 1-3
 - IPO. See interprocedural optimization
- ## L
- large load stalls, 3-54
 - latency, 7-4, 7-16
 - legacy mode, 9-1
 - LFENCE instruction, 7-11
 - links to web data, 1-3
 - load instructions and prefetch, 7-6
 - loading-storing to-from same DRAM page, 5-40
 - loop
 - blocking, 4-24
 - unrolling, 7-21, A-5
- ## M
- MASKMOVDQU instruction, 7-7
 - memory bank conflicts, 7-2
 - memory optimizations
 - loading-storing to-from same DRAM page, 5-40
 - overview, 5-35
 - partial memory accesses, 5-36, 5-40
 - performance, 4-19
 - reference instructions, 3-27
 - using aligned stores, 5-40
 - using prefetch, 7-13
 - MFENCE instruction, 7-12
 - micro-op fusion, 2-51
 - misaligned data access, 4-14
 - misalignment in the FIR filter, 4-16
 - mobile computing
 - ACPI standard, 11-1, 11-3
 - active power, 11-1

- battery life, 11-1, 11-5, 11-6
 - C4-state, 11-4
 - CD/DVD, WLAN, WiFi, 11-7
 - C-states, 11-1, 11-3
 - deep sleep transitions, 11-7
 - deeper sleep, 11-4, 11-10
 - Intel Mobil Platform SDK, 11-6
 - OS APIs, 11-6
 - OS changes processor frequency, 11-2
 - OS synchronization APIs, 11-6
 - overview, 11-1, 12-1
 - performance options, 11-5
 - platform optimizations, 11-7
 - P-states, 11-1
 - Speedstep technology, 11-8
 - spin-loops, 11-6
 - state transitions, 11-2
 - static power, 11-1
 - WM_POWERBROADCAST message, 11-8
 - MOVAPD instruction, 6-3
 - MOVAPS instruction, 6-3
 - MOVDDUP instruction, 6-14
 - move byte mask to integer, 5-15
 - MOVHLPS instruction, 6-11
 - MOVLHPS instruction, 6-11
 - MOVNTDQ instruction, 7-7
 - MOVNTI instruction, 7-7
 - MOVNTPD instruction, 7-7
 - MOVNTPS instruction, 7-7
 - MOVNTQ instruction, 7-7
 - MOVQ Instruction, 5-39
 - MOVSHDUP instruction, 6-14, 6-16
 - MOVSLDUP instruction, 6-14, 6-16
 - MOVUPD instruction, 6-3
 - MOVUPS instruction, 6-3
 - multicore processors
 - architecture, 2-1
 - C-state considerations, 11-12
 - energy considerations, 11-10
 - features of, 2-56
 - functional example, 2-56
 - pipeline and core, 2-58
 - SpeedStep technology, 11-11
 - thread migration, 11-11
 - multiprocessor systems
 - dual-core processors, 8-1
 - HT Technology, 8-1
 - optimization techniques, 8-1
 - See also: multithreading & Hyper-Threading Technology
 - multithreading
 - Amdahl's law, 8-2
 - application tools, 8-10
 - bus optimization, 8-12
 - compiler support, A-5
 - dual-core technology, 3-5
 - environment description, 8-1
 - guidelines, 8-11
 - hardware support, 3-5
 - HT technology, 3-5
 - Intel Core microarchitecture, 8-6
 - parallel & sequential tasks, 8-2
 - programming models, 8-4
 - shared execution resources, 8-35
 - specialized models, 8-6
 - thread sync practices, 8-12
 - See Hyper-Threading Technology
- ## N
- Newton-Raphson iteration, 6-1
 - non-coherent requests, 7-9
 - non-interleaved unpack, 5-10
 - non-temporal stores, 7-8, 7-31
 - NOP, 3-31
- ## O
- OpenMP compiler directives, 8-10, A-5
 - optimization
 - branch prediction, 3-6
 - branch type selection, 3-13
 - eliminating branches, 3-7
 - features, 2-1
 - general techniques, 3-1
 - spin-wait and idle loops, 3-9
 - static prediction, 3-9
 - unrolling loops, 3-15
 - optimizing cache utilization
 - cache management, 7-32
 - examples, 7-11
 - non-temporal store instructions, 7-7, 7-10
 - prefetch and load, 7-6
 - prefetch instructions, 7-5
 - prefetching, 7-5
 - SFENCE instruction, 7-11, 7-12
 - streaming, non-temporal stores, 7-7
 - See also: cache management
 - OS APIs, 11-6
- ## P
- pack instructions, 5-8
 - packed average byte or word), 5-29
 - packed multiply high unsigned, 5-28
 - packed shuffle word, 5-16
 - packed signed integer word maximum, 5-28
 - packed sum of absolute differences, 5-28, 5-29
 - parallelism, 4-8, 8-5
 - partial memory accesses, 5-36
 - PAUSE instruction, 3-9, 8-12
 - PAVGB instruction, 5-29
 - PAVGW instruction, 5-29
 - PeekMessage(), 11-6
 - Pentium 4 processors
 - inner loop iterations, 3-15
 - static prediction, 3-9

INDEX

- Pentium M processors
 - prefetch mechanisms, 7-3
 - processor perspectives, 3-3
 - static prediction, 3-9
 - Pentium Processor Extreme Edition, 2-1, 2-56
 - performance models
 - Amdahl's law, 8-2
 - multithreading, 8-2
 - parallelism, 8-1
 - usage, 8-1
 - performance monitoring events
 - analysis techniques, B-42
 - Bus_Not_In_Use, B-42
 - Bus_Snoops, B-42
 - DCU_Snoop_to_Share, B-42
 - drill-down techniques, B-42
 - event ratios, B-47
 - Intel Core Duo processors, B-39
 - Intel Core Solo processors, B-39
 - Intel Netburst architecture, B-1
 - Intel Xeon processors, B-1
 - L1_Pref_Req, B-41
 - L2_No_Request_Cycles, B-41
 - L2_Reject_Cycles, B-41
 - Pentium 4 processor, B-1
 - performance counter, B-40
 - ratio interpretation, B-40
 - See also: clock ticks
 - Serial_Execution_Cycles, B-41
 - Unhalted_Core_Cycles, B-41
 - Unhalted_Ref_Cycles, B-41
 - performance tools, 3-1
 - PEXTRW instruction, 5-12
 - PGO. See profile-guided optimization
 - PINSRW instruction, 5-13
 - PMINSW instruction, 5-28
 - PMINUB instruction, 5-28
 - PMOVMSKB instruction, 5-15
 - PMULHUW instruction, 5-28
 - predictable memory access patterns, 7-5
 - prefetch
 - 64-bit mode, 9-5
 - coding guidelines, 7-2
 - compiler intrinsics, 7-2
 - concatenation, 7-20
 - deterministic cache parameters, 7-38
 - hardware mechanism, 7-3
 - characteristics, 7-13
 - latency, 7-14
 - how instructions designed, 7-5
 - innermost loops, 7-5
 - instruction considerations
 - cache block techniques, 7-23
 - checklist, 7-17
 - concatenation, 7-19
 - hint mechanism, 7-4
 - minimizing number, 7-20
 - scheduling distance, 7-18
 - single-pass execution, 7-2, 7-28
 - spread with computations, 7-22
 - strip-mining, 7-25
 - summary of, 7-4
 - instruction variants, 7-5
 - latency hiding/reduction, 7-16
 - load Instructions, 7-6
 - memory access patterns, 7-5
 - memory optimization with, 7-13
 - minimizing number of, 7-20
 - scheduling distance, 7-2, 7-18
 - software data, 7-4
 - spreading, 7-23
 - when introduced, 7-1
 - PREFETCHNT0 instruction, 7-6
 - PREFETCHNT1 instruction, 7-6
 - PREFETCHNT2 instruction, 7-6
 - PREFETCHNTA instruction, 7-6, 7-25
 - usage guideline, 7-2
 - PREFETCHT0 instruction, 7-25
 - usage guideline, 7-2
 - producer-consumer model, 8-6
 - profile-guided optimization, A-6
 - PSADBW instruction, 5-28
 - PSHUF instruction, 5-16
 - P-states, 11-1
- ## Q
- Qparallel, 8-10
- ## R
- ratios, B-47
 - branching and front end, B-49
 - references, 1-3
 - releases of, 2-63
 - rounding control option, A-6
 - rules, E-1
- ## S
- sampling
 - event-based, A-12
 - scheduling distance (PSD), 7-18
 - Self-modifying code, 3-63
 - SFENCE Instruction, 7-11
 - SHUFPS instruction, 6-3
 - signed unpack, 5-7
 - SIMD
 - auto-vectorization, 4-13
 - cache instructions, 7-1
 - classes, 4-12
 - coding techniques, 4-8
 - data alignment for MMX, 4-17
 - data and stack alignment, 4-14
 - data alignment for 128-bits, 4-17
 - example computation, 2-60
 - history, 2-60

- identifying hotspots, 4-7
- instruction selection, 4-26
- loop blocking, 4-24
- memory utilization, 4-19
- microarchitecture differences, 4-28
- MMX technology support, 4-2
- padding to align data, 4-15
- parallelism, 4-8
- SSE support, 4-2
- SSE2 support, 4-3
- SSE3 support, 4-3
- SSSE3 support, 4-4
- stack alignment for 128-bits, 4-16
- strip-mining, 4-23
- using arrays, 4-15
- vectorization, 4-8
- VTune capabilities, 4-7
- SIMD floating-point instructions
 - data arrangement, 6-3
 - data deswizzling, 6-9
 - data swizzling, 6-6
 - different microarchitectures, 6-14
 - general rules, 6-1
 - horizontal ADD, 6-10
 - Intel Core Duo processors, 6-18
 - Intel Core Solo processors, 6-18
 - planning considerations, 6-1
 - reciprocal instructions, 6-1
 - scalar code, 6-2
 - SSE3 complex math, 6-15
 - SSE3 FP programming, 6-14
 - using
 - ADDSUBPS, 6-16
 - CVTTPS2PI, 6-13
 - CVTTSS2SI, 6-13
 - FXCH, 6-2
 - HADDPS, 6-18
 - HSUBPS, 6-18
 - MOVAPD, 6-3
 - MOVAPS, 6-3
 - MOVHLPS, 6-11
 - MOVLHPS, 6-11
 - MOVSHDUP, 6-16
 - MOVSLDUP, 6-16
 - MOVUPD, 6-3
 - MOVUPS, 6-3
 - SHUFPS, 6-3
 - vertical vs horizontal computation, 6-3
 - with x87 FP instructions, 6-2
- SIMD technology, 2-63
- SIMD-integer instructions
 - 64-bits to 128-bits, 5-43
 - data alignment, 5-4
 - data movement techniques, 5-6
 - extract word, 5-12
 - integer intensive, 5-1
 - memory optimizations, 5-35
 - move byte mask to integer, 5-15
 - optimization by architecture, 5-44
 - packed average byte or word, 5-29
 - packed multiply high unsigned, 5-28
 - packed shuffle word, 5-16
 - packed signed integer word maximum, 5-28
 - packed sum of absolute differences, 5-28
 - rules, 5-2
 - signed unpack, 5-7
 - unsigned unpack, 5-6
 - using
 - EMMS, 5-2
 - MOVQDQ, 5-39
 - MOVQ2DQ, 5-19
 - PABSW, 5-21
 - PACKSSDW, 5-8
 - PADDQ, 5-30
 - PALIGNR, 5-5
 - PAVGB, 5-29
 - PAVGW, 5-29
 - PEXTRW, 5-12
 - PINSRW, 5-13
 - PMADDWD, 5-30
 - PMAXSW, 5-28
 - PMAXUB, 5-28
 - PMINSW, 5-28
 - PMINUB, 5-28
 - PMOVMASKB, 5-15
 - PMULHUW, 5-28
 - PMULHW, 5-28
 - PMULUDQ, 5-28
 - PSADBW, 5-28
 - PSHUF, 5-16
 - PSHUFB, 5-22, 5-24
 - PSHUFLW, 5-17
 - PSLLDQ, 5-31
 - PSRLDQ, 5-31
 - PSUBQ, 5-30
 - PUNPCHQDQ, 5-18
 - PUNPCKLQDQ, 5-18
 - simplified 3D geometry pipeline, 7-16
 - simplified clipping to an arbitrary signed range, 5-27
 - single vs multi-pass execution, 7-28
 - sleep transitions, 11-7
 - smart cache, 2-50
 - SoA format, 4-21
 - software write-combining, 7-31
 - spin-loops, 11-6
 - optimization, 3-9
 - PAUSE instruction, 3-9
 - related information, 1-3
- SSE, 2-63
- SSE2, 2-63
- SSE3, 2-64
- SSSE3, 2-64, 2-65
- stack
 - aligned EDP-based frames, D-4
 - aligned ESP-based frames, D-3
 - alignment 128-bit SIMD, 4-16

INDEX

- alignment stack, 3-59
 - dynamic alignment, 3-59
 - frame optimizations, D-6
 - inlined assembly & EBX, D-7
 - Intel C++ Compiler support for, D-1
 - overview, D-1
- state transitions, 11-2
- static branch prediction algorithm, 3-10
- static power, 11-1
- static prediction, 3-9
- streaming stores, 7-7
 - coherent requests, 7-9
 - improving performance, 7-7
 - non-coherent requests, 7-9
- strip-mining, 4-23, 4-24, 7-25, 7-26
 - prefetch considerations, 7-27
- structures
 - aligning, 3-56
- suggestions, E-1
- summary of coding rules, E-1
- system bus optimization, 8-23

T

- time-based sampling, A-12
- time-consuming innermost loops, 7-5
- TLB. See transaction lookaside buffer
- transaction lookaside buffer, 7-33
- transcendental functions, 3-86

U

- unpack instructions, 5-10
- unrolling loops
 - benefits of, 3-15
 - code examples, 3-16
 - limitation of, 3-15
- unsigned unpack, 5-6
- using MMX code for copy, shuffling, 6-10

V

- vector class library, 4-13
- vectorized code
 - auto generation, A-7
 - automatic vectorization, 4-13
 - high-level examples, A-7
 - parallelism, 4-8
 - SIMD architecture, 4-8
 - switch options, A-5
- vertical vs horizontal computation, 6-3

W

- WaitForSingleObject(), 11-6
- WaitMessage(), 11-6
- weakly ordered stores, 7-7
- WiFi, 11-7
- WLAN, 11-7

- workload characterization
 - retirement throughput, A-12
- write-combining
 - buffer, 7-31
 - memory, 7-31
 - semantics, 7-8

X

- XCHG EAX,EAX, support for, 3-31

Z

- , A-1