

### The Intel® Xeon Phi Coprocessor

Dr-Ing. Michael Klemm Software and Services Group Intel Corporation (michael.klemm@intel.com)



## Legal Disclaimer & Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

Copyright 2014°, Intel Corporation. All rights reserved. Intel, the Intel logo, Xeon, Xeon Phi, Core, VTune, and Cilk are trademarks of Intel Corporation in the U.S. and other countries.

#### **Optimization Notice**

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804





## Intel Technologies for HPC

#### Processors

Intel® Xeon® Processor



Coprocessor Intel<sup>®</sup> Many Integrated Core



Network & Fabric



I/O & Storage



#### Software & Services





#### Transforming the Economics of HPC



#### Executing to Moore's Law

Predictable Silicon Track Record – well and alive at Intel. Enabling new devices with higher performance and functionality while controlling power, cost, and size





### **Driving Innovation and Integration**

Enabled by Leading Edge Process Technologies



#### Integrated Today



#### Coming in the Future





## **The Magic of Integration**

Moore's Law at Work & Architecture Innovations







1970s 150 MFLOPS CRAY-1 2013 1000000 MFLOPS Intel® Xeon Phi™





**#1** TOP500 June 2014 **33** PFI OPS HPI 54 PELOPS Peak **32000** Intel<sup>®</sup> Xeon<sup>®</sup> E5v2 Processors **48000** Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessors







## Intel® Many Integrated Core Architecture (Intel® MIC) & Intel® Xeon Phi™ Coprocessor

#### Intel Architecture Multicore and Manycore More cores. Wider vectors. Co-Processors.

Images do not ro	eflect actual die siz	res. Actual production	die may differ from	images.					
	Intel' Xeon' processor 64-bit	Intel Xeon processor 5100 series	Intel Xeon processor 5500 series	Intel Xeon processor 5600 series	Intel Xeon processor E5 Product Family	Intel Xeon processor code name Ivy Bridge	Intel Xeon processor code name Haswell		Intel <sup>®</sup> Xeon Phi™ Coprocess code name Knights Corrner
Core(s)	1	2	4	6	8	12	18	1	61
Threads	2	2	8	12	16	24	36		244

Intel® Xeon Phi<sup>™</sup> Coprocessor extends established CPU architecture and programming concepts to highly parallel applications





# Each Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor core is a fully functional multi-thread execution unit



Core unit based on Intel<sup>®</sup> Pentium<sup>®</sup> processor family

- Two pipelines (U and V)
  - Dual-issue on scalar instructions
- Scalar pipeline 1 clock latency
- 64-bit data path
- 4 hardware threads per core
- Thread context: GPRs, ST0-7, etc.
- "Smart" round-robin scheduling
  - Prefetch buffers 2 inst-bundles / context
  - Next ready context selected in order



# Each Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor core is a fully functional multi-thread execution unit



Instruction decoder is fully pipelined but is designed as a 2-cycle unit

- Enables significant increase to maximum core frequency, but...
  - Core cannot issue instructions from same context in adjacent cycles
  - Means minimum two threads per core to use all available compute cycles



# Each Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor core is a fully functional multi-thread vector unit



Vector unit width 512 bits!

- 32 512-bit vector registers per context
  - Each holds 16 floats or 8 doubles
  - ALUs support int32/float32 operations, float64 arithmetic, int64 logic ops
  - Ternary ops including Fused-Multiply-Add
  - Broadcast/swizzle support, float16 up-convert
  - 8 vector mask registers for per lane conditional operations
  - Most ops: 4-cycle latency 1-cycle throughput
    - Matches 4-cycle round robin of integer unit
  - Mostly IEEE 754 2008 compliant
    - Not supported: MMX<sup>™</sup> technology, Streaming SIMD Extensions (SSE), Intel<sup>®</sup> Advanced Vector Extensions (Intel<sup>®</sup> AVX)



12

## Individual cores are tied together via fully coherent caches into a bidirectional ring on the Intel® Xeon Phi<sup>™</sup> coprocessor



Optimization Notice

13

## Cache Hierarchy

Parameter	L1	L2	
Coherence	MESI	MESI	
Size	32KB + 32 KB	512 KB	
Associativity	8-way	8-way	
Line Size	64 Bytes	64 Bytes	
Banks	8	8	
Access Time	2 cycle	11 cycle	
Policy	Pseudo LRU	Pseudo LRU	
Duty Cycle	1 per clock	1 per clock	
Ports	Read or Write	Read or Write	

#### There is no L3 cache!





#### Intel<sup>®</sup> Xeon Phi<sup>™</sup> Product Family based on Intel<sup>®</sup> Many Integrated Core (MIC) Architecture



\*Per Intel's announced products or planning process for future products





#### Next Intel<sup>®</sup> Xeon Phi<sup>™</sup> Processor Codename: Knights Landing



## Designed using Intel's cutting-edge **14nm process**

#### Not bound by "offloading" bottlenecks Standalone CPU or PCIe Coprocessor

Leadership compute & memory bandwidth Integrated On-Package Memory

All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.





#### Next Intel<sup>®</sup> Xeon Phi<sup>™</sup> Processor Codename: Knights Landing



#### Compute: Energy-efficient IA cores<sup>2</sup>

- Microarchitecture enhanced for HPC<sup>3</sup>
- 3x Single Thread Performance vs Knights Corner<sup>4</sup>
- Intel Xeon Processor Binary Compatible<sup>5</sup>

#### **On-Package Memory:**

- Jointly Developed with Micron Technology

All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.







### Programming for Intel Architectures

## Highly Parallel Applications



Efficient vectorization, threading, and parallel execution drives higher performance for suitable scalable applications





## Parallel Programming for Intel® Architecture

NODES	Use Intel <sup>®</sup> MPI, Co-Array Fortran		
CORES	Use threads directly (pthreads) or via OpenMP*, C++11 Use tasking, Intel® TBB / Cilk™ Plus		
VECTORS	Intrinsics, auto-vectorization, vector-libraries Language extensions for vector programming (SIMD)		
BLOCKING	Use caches to hide memory latency Organize memory access for data reuse		
DATA LAYOUTStructure of arrays facilitates vector loads / stores, unit stride Align data for vector accesses			
Parallel programming to utilize the hardware resources,			

in an abstracted and portable way





## Heterogeneous Programming





Copyright° 2014, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

#### Native Programming for Intel Xeon Phi





#### Flexible Execution Models for Heterogeneous Platforms SOURCE CODE SERIAL AND MODERATELLY HIGHLY PARALLEL CODE Compilers, Libraries, PARALLEL CODE **Runtime Systems** MAIN() MAIN() MAIN() MAIN() MAIN() **XEON XEON XEON XEON**® **XEON**® **XEON**® **XEON**® **PHI**<sup>™</sup> **PHI**<sup>™</sup> **PHI**<sup>™</sup> RESULTS RESULTS RESULTS RESULTS RESULTS **Multicore Only Multicore Hosted with Symmetric** Manycore Only Manycore Offload (Native) Optimization Notice

Copyright° 2014, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.



### Case Study: NWChem CCSD(T)

## Finding Offload Candidates

#### Requirements for offload candidates

- Compute-intensive code regions (kernels)
- Highly parallel
- Compute scaling stronger than data transfer, e.g., compute O(n<sup>3</sup>) vs. data size O(n<sup>2</sup>)

Finding offload candidates

- Create a benchmark to trigger application code of interest
- Find hotspots in the application
- Determine input and output data
- Determine data sizes transferred





## Offload Analysis Methodolgy







## Example: NWChem Hotspots

#### NWChem hotspot profile

Basic Hotspots Hotspots by CPU Usage viewpoint (change) 2							
Analysis Target 🙏 Analysis Type 📟 Collection Log 🕅 Summary Communication							
Grouping: Function / Call Stack							
Function / Call Stack CPU Time by Utili  ★ D Overhead approxime Total of 38% of							
▷comex_make_progress	52.0%	0.250s nw	vchem spent ir	$n sd_t dX_Y$			
▶sd_t_d2_2	4.2%	0s nw	vchem				
▶sd_t_d2_8	4.2%	0s nw	vchem	0x16a272!			
▶sd_t_d2_9	4.2%	0s nv	Ver 11	0x16a2e3			
▶sd_t_d2_5	4.2%	nw ever	vchem	0x16a1a4(			
♦sd_t_d2_3	4.1%	Os nw	vchem	0x16a3aaa			
♦sd_t_d2_6	4.1%	0s nw	vchem	0x16a345			
▶sd_t_d2_1	3.1%	0s nw	vchem	0x16a40b			
▶sd_t_d2_4	3.0%	0s nw	vchem	0x16a1424			
▶sd_t_d2_7	1.4%		vchem	0x16a211(			
▶_mq_test	1.2%	0.030s nw	vchem	0x29565a			
Selected 18 row(s):	38.0%	0s		~			
< III >							



## Example: NWChem Hotspots

Optimi Not

#### Call-tree analysis shows relationship of hotspots

Basic Hotspots Hotspots by CPU Usage viewpoint ( <u>change</u> This function is the								
🕢 🖶 Analysis Target 🛝 Analysis Type 📟 Collection Log 📓 Summary 🗞 Bottom-up 🛛 COMMON anchor fo								
Grouping: Function / Ca	Grouping: Function / Call Stack all hotspots.							
Function / Call Stack	Sunction / Call Stack CPU Time by Utili							
Tunction / can stack	🔲 Idle 📕 Poor 📙	Function Stack	CPUe: Tota	l by Utilization <del>、</del> 🔊	^			
▷comex_mak_progress	52.0%	Tunction Stuck	Phare Poor Ok	Ideal 📘 Over	🛛 Idle			
⊽sd_t_d2_2	4.270	▽ccsd_t_doubles_l_2	1333.239s		0			
	4 2%	♦get_nash_block_i	335.529s					
▼sd_t_d2_8	4.2%	▽get_hash_block	289.569s					
	4.2%		289.569s		Ξ			
⊽sd_t_d2_9	4.2%	Dga_get_	289.559s					
	4.2%	Þutil_wallsec	0.010s		0			
Þsd_t_d2_5	4.2%	sd_t_d2_2	92.330s		92			
Þsd_t_d2_3	4.1%	sd_t_d2_8	91.870s 📒		91			
Þsd_t_d2_6	4.1%	sd_t_d2_9	91.570s 📒		91			
Þsd_t_d2_1	3.1%	sd_t_d2_5	90.490s		90			
Selected 1 row(s):	-	sd_t_d2_3	89.930s		89			
≤ m ) < sd_t_d2_6 88.020s								
tion	Selected 1 row(s): 1333.239s							
Copyriaht° 2014, Inte	Corporation, All rig	<	< III		>			



## Example: Loop Analysis

All kernels expose the same structure

7 perfectly nested loops

Trip count per loop is equal to "tile size" (20-30)

Naïve per-kernel solution is obvious

<pre>subroutine sd_t_d1_1(h3d,h2d,h1d,p6d,p5d,p4d,</pre>	
implicit pope	
integen h2d h2d h1d n6d nEd n4d h7d	offload
integer h $3$ h $2$ h $1$ nf nf nf h $7$	
double procision triplosy(b2d b2d b1d p6d pEd p	44)
double precision $\pm 2 \operatorname{sub}(h7d, p1d, p5d, h1d)$	40)
double precision v2sub( $h^2d$ , $p^2d$ , $p^2d$ , $h^2d$ )	
do $p=1$ $p=1$	
do n5=1.n5d multi-thro	eading
do n6=1.n6d	
do h1=1.h1d	
do $h2=1,h2d$	
do h3=1,h3d	
do h7=1,h7d	
<pre>triplesx(h3,h2,h1,p6,p5,p4)=</pre>	
1 triplesx(h3,h2,h1,p6,p5,p4)	
1 - t2sub(h7,p4,p5,h1)*v2sub(h3,h2,p6,h7)	SIMD
enddo	SITE
enddo	
end	





#### Issues w/ Naïve Offload Solution





## Optimization of Data Transfers

#### Use call-tree analysis to find common anchor for hotspots

- Hoist data transfers up as high as possible
- Make offload regions as large as possible

<pre>cdir\$ offload_transfer target(mic) nocopy(triplesx:length(triplesx_1) ALLOC) cdir\$ offload_transfer target(mic) nocopy(t2sub:length(t2sub_1) ALLOC) cdir\$ offload_transfer target(mic) nocopy(v2sub:length(v2sub_1) ALLOC) cdir\$ offload target(mic) nocopy(triplesx:length(0) REUSE)</pre>	data env.
cdir\$ offload target(mic) in(triplesx:length(0),REUSE)       offload         1       in(t2sub:length(2sub_1),REUSE)         3       in(v2sub:length(v2sub),REUSE)         2       in(h3d,h1d,p6d,p5d,p4d,h7d)	
endif	
<pre>c sd_t_d1_2 until sd_t_d1_9     enddo</pre>	
<pre>c Similar structure for sd_t_d2_1 until sd_t_d2_9</pre>	
<pre>cdir\$ offload_transfer target(mic) out(triplesx:length(triplesx_1) REUSE)</pre>	





#### Kernel Optimizations





- Loop ordering not optimal for SIMD execution
  - Too low trip count for inner loop
- Index analysis shows that loops can be reordered
  - Swap h7 and h3
  - Swap h7 and h2 again







- Loop ordering not optimal for SIMD execution
  - Too low trip count for inner loop
- Index analysis shows that loops can be reordered
  - Swap h7 and h3
  - Swap h7 and h2 again
- Loops h2 and h3 can be collapsed

<pre>subroutine sd_t_d1_1(h3d,h2d,h1d,p6d, 1 h7d,triplesx,t2s implicit none integer h3d,h2d,h1d,p6d,p5d,p4d,h7d integer h3,h2,h1,p6,p5,p4,h7 double precision triplesx(h3d,h2d,h1d double precision t2sub(h7d,p4d,p5d,h1 double precision v2sub(h3d,h2d,p6d,h7</pre>	,p5d,p4d, sub,v2sub) d,p6d,p5d,p4d) Ld) 7d)
<pre>!\$omp parallel do collapse(3)</pre>	multi-threading
do p5=1,p5d	
do p6=1,p6d	
do h1=1,h1d	
do h7=1,h7d	
do h2=1,h2d	
do h3=1,h3d	
<pre>triplesx(h3,h2,h1,p6,p5,p4)=triples</pre>	sx(h3,h2,h1,p6,p5,p4)
<pre>1 - t2sub(h7,p4,p5,h1)*v2sub(h3,h2,</pre>	,p6,h7)
enddo	
<pre>!\$omp end parallel do</pre>	
end	



- Loop ordering not optimal for SIMD execution
  - Too low trip count for inner loop
- Index analysis shows that loops can be reordered
  - Swap h7 and h3
  - Swap h7 and h2 again
- Loops h2 and h3 can be collapsed

<pre>subroutine sd_t_d1_(h3d,h2d,h1d,p 1 h7d,triplesx, implicit none integer h3d,h2d,h1d,p6d,p5d,p4d,h7 integer h3,h2,h1,p6,p5,p4,h7 double precision triplesx(h3d,h2d, double precision t2sub(h7d,p4d,p5d double precision v2sub(h3d,h2d,p6d</pre>	06d,p5d,p4d, t2sub,v2sub) 7d h1d,p6d,p5d,p4d) l,h1d) l,h7d)
<pre>!\$omp parallel do collapse(3)</pre>	multi-threading
do p5=1,p5d	
do p6=1,p6d	
do h1=1,h1d	
do h7=1,h7d	
do h2=1,h2d	
do h3=1,h3d	
<pre>triplesx(h3,h2,h1,p6,p5,p4)=trip</pre>	lesx( <mark>h3,h2</mark> ,h1,p6,p5,p4)
<pre>1 - t2sub(h7,p4,p5,h1)*v2sub(h3,</pre>	h2,p6,h7)
enddo	
!\$omp end parallel do	
end	



- Loop ordering not optimal for SIMD execution
  - Too low trip count for inner loop
- Index analysis shows that loops can be reordered
  - Swap h7 and h3
  - Swap h7 and h2 again
- Loops h2 and h3 can be collapsed

<pre>subroutine sd_t_d1_1(h3d,h2d,h1d,p6d,p5 1</pre>	d,p4d,	
implicit none	, 2300)	
integen h3d h2d h1d n6d n5d n4d h7d		
integer h3 h2 h1 n6 n5 n4 h7		
double precision triplesy(h3d*h2d h1d n	6d n5d n4d)	
double precision triplesx(ind n2d, ind, p	οα,ροα,ρ+α)	
double precision v2sub( $h/d$ , $p+d$ , $p-d$ , $h/d$ )		
Isomn narallel do collanse(3)		
do n4=1 n4d	multi-thr	eading
do $p=1, p=0$		
do $p=1, p=0$		
do = 1, pod		
$d_0 h_{1-1}^{-1} h_{7d}^{-1}$		
do h2h3=1, h2d*h3d		
triplesx(h2h3,h1,n6,n5,n4)=triplesx(h	2h3.h1.n6.n5.n4)	
$1 - t^2 sub(h^2, n^4, n^5, h^1) * v^2 sub(h^2h^3, n^6, h^2h^2)$	h7)	CTMD
enddo	,	SIMD
enddo		
!\$omp end parallel do		
end		

about 50% speed-up



#### Kernel Optimizations, Multi-versioning

subroutine sd_	t_d1_1(h3d,h2d,h1d,p6d,p5d,p4d,	с	continued from left column	
<pre>1 h7d,triplesx,t2sub,v2sub)</pre>			else	
implicit none		<pre>!\$omp parallel do collapse(3)</pre>		
integer h3d,h2	d,h1d,p6d,p5d,p4d,h7d		do p4=1,p4d	
<pre>integer h3,h2,</pre>	h1,p6,p5,p4,h7		do p5=1,p5d	
integer rmndr			do p6=1,p6d	
double precisi	on triplesx(h3d*h2d,h1d,p6d,p5d,p4d)		do h1=1,h1d	
double precisi	on t2sub(h7d,p4d,p5d,h1d)		do h7=1,h7d	
double precisi	on v2sub(h3d*h2d,p6d,h7d)		do h2h3=1,h2d*h3d	
rmndr = mod(h3	d,8) + mod(h2d,8) + mod(h1d,8) +		<pre>triplesx(h2h3,h1,p6,p5,p4)=triplesx(h2h3,h1,p6,p5,p4)</pre>	
1 mod(p6	(d,8) + mod(p5d,8) + mod(p4d,8) +		1 - t2sub(h7,p4,p5,h1)*v2sub(h2h3,p6,h7)	
2 mod(h7	d,8)		enddo	
if (rmndr.eq.0	) then		enddo	
<pre>!\$omp parallel do co</pre>	llapse(3)	enddo enddo		
do p4=1,p4d				
do p5=1,p5d		enddo		
do p6=1,p6d		enddo enddo		
do h1=1,h1d				
do h7=1,h7d		<pre>!\$omp end parallel do</pre>		
<pre>!dec\$ vector aligned</pre>			endif	
do h2h3=1,h2	ld*h3d	end		
triplesx(h	2h3,h1,p6,p5,p4)=triplesx(h2h3,h1,p6,p5,p4)			
1 - t2sub(h7	,p4,p5,h1)*v2sub(h2h3,p6,h7)			
enddo			about 15% speed-up	
<pre>!\$omp end parallel do</pre>			about 1370 speed-up	



Copyright° 2014, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

## Device Partitioning

#### Host executes several MPI ranks

- Utilize coprocessor from several host processes concurrently
- Utilize host CPUs for increased performance

#### Partition coprocessors through OpenMP\* runtime

- Less threading overhead, better overall system utilization
- Rank 0: OFFLOAD\_DEVICES=0 KMP\_PLACE\_THREADS=30c,4t,0o
- Rank 1: OFFLOAD\_DEVICES=0 KMP\_PLACE\_THREADS=30c,4t,300
- Rank 4: OFFLOAD\_DEVICES=1 KMP\_PLACE\_THREADS=30c,4t,0o
- Rank 5: OFFLOAD\_DEVICES=1 KMP\_PLACE\_THREADS=30c,4t,30o







#### Performance Results



Performance tests are measured using specific computer systems, components, software, operations, and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. System configuration: Atipa Visione vf442 server with two Intel Xeon E5-2670 8-core processors at 2.6 GHz (128 GB DDR3 with 1333 MHz, Scientific Linux release 6.5) and Intel C600 IOH, two Intel Xeon Phi coprocessors 5110P (GDDR5 with 3.6 GT/sec, driver v3.1.2-1, flash image/micro OS 2.1.02.0390, Intel Composer XE 14.0.1.106). Benchmark perturbative triples correction to the CCSD(T) correlation energy of the 1,3,4,5-tetrasilylimidazol-2-ylidene molecule (formula  $Si_4C_3N_2H_{12}$ ) in its triplet state.





## Software

°2014, Intel Corporation. All rights reserved. Intel, the Intel logo, Intel Inside, Intel Xeon, and Intel Xeon Phi are trademarks of Intel Corporation in the U.S. and/or other countries. \*Other names and brands may be claimed as the property of others.