



Knights Landing – An Overview for Developers

James Reinders, Intel
June 7, 2016

The image shows the cover of a book titled 'Intel Xeon Phi Knights Landing Edition'. The cover features a top-down view of a complex, multi-colored microchip. The text is white and blue, set against a dark background. The authors' names are listed below the title. The Intel MK logo is in the bottom left corner.

**INTEL® XEON PHI™
PROCESSOR
HIGH PERFORMANCE
PROGRAMMING**

KNIGHTS LANDING EDITION

Jim Jeffers | James Reinders | Avinash Sodani

MK
MICRO KNUIGHTS

Knights Landing

2nd Generation Intel® Xeon Phi™ products

ALL ABOUT PARALEL PROGRAMMING

Threading, Vectorization, Data Locality
Fortran, C, C++ (plus a little Python)
OpenMPI, MPI, Intel® TBB

New with Knights Landing:

AVX-512,

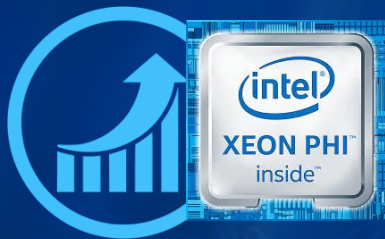
High Bandwidth Memory (MCDRAM),

Cluster Mode,

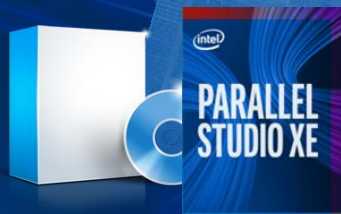
Omni-Path

XeonPhiDeveloper.com

You can buy you very own Knight Landing development system today!



Highly-Parallel Performance
to develop on



All the Software Tools & Libraries
you need



Support & Training
to help you succeed

Leading edge platform capabilities, performance to deliver multi-threaded, vectorized software for today's HPC workloads !

Why Intel® Xeon Phi™ Processors?



Imagine you could design anything

What is you wanted:

- *High performance*
- *Preserve your investment in code*
- *Reuse your knowledge of programming languages, tools, techniques*

What if the *ONLY* thing we were willing to “give up” was running non-parallel codes well?

*In other words:
what if we assumed ONLY parallel applications
would be run?*

Imagine you could design anything

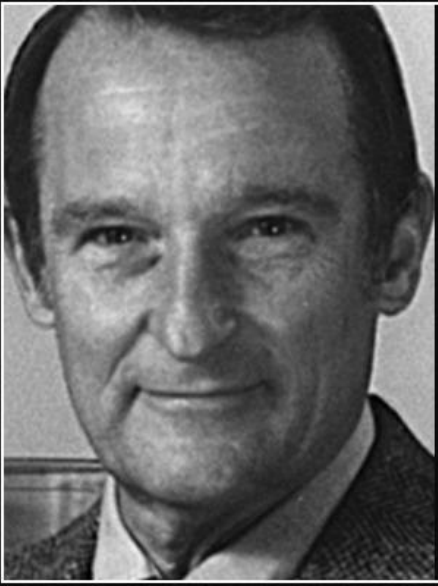
What is you wanted:

- *High performance*
- *Preserve your investment in code*
- *Reuse your knowledge of programming languages, tools, techniques*

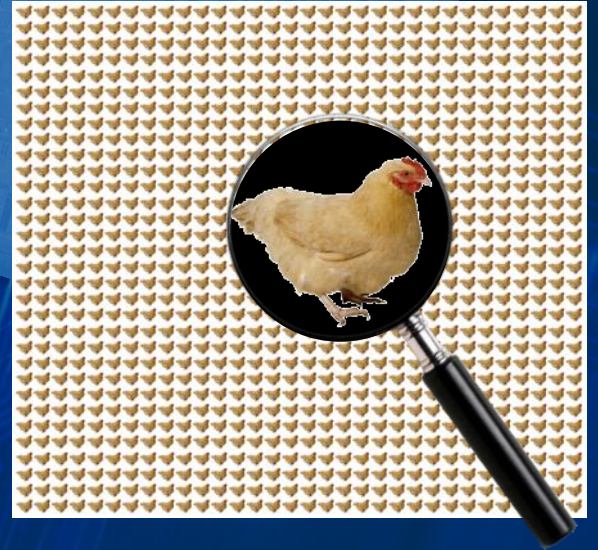
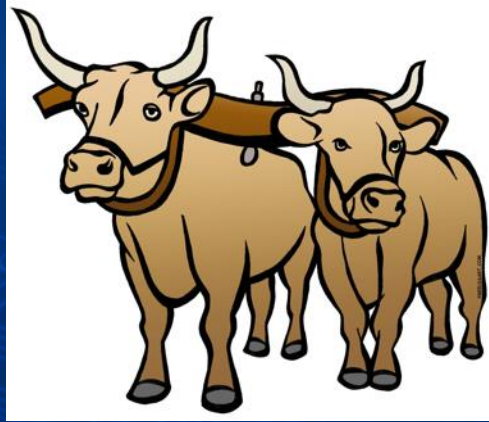
What if the ONLY thing we were willing to "give up" was running non-parallel codes well?

*In other words:
what if we assumed ONLY parallel applications
would be run?*

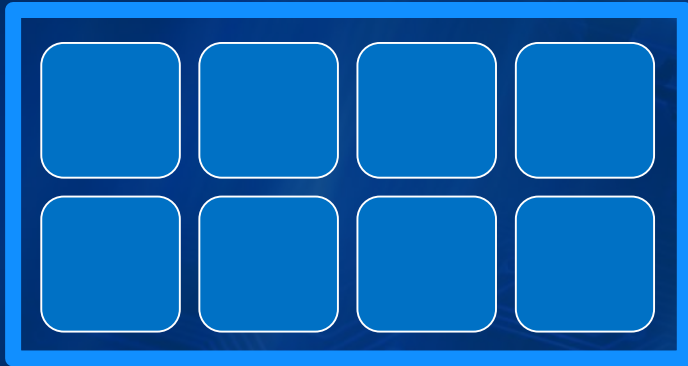




If you were plowing a field,
which would you rather use...
two strong oxen, or
1024 chickens?

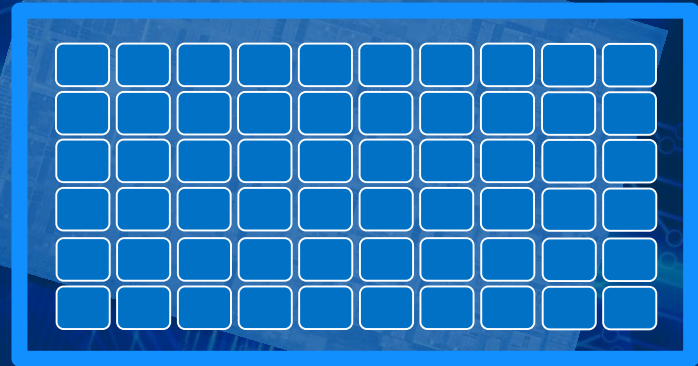


Design Question - Best for Computing?



A few powerful

vs.



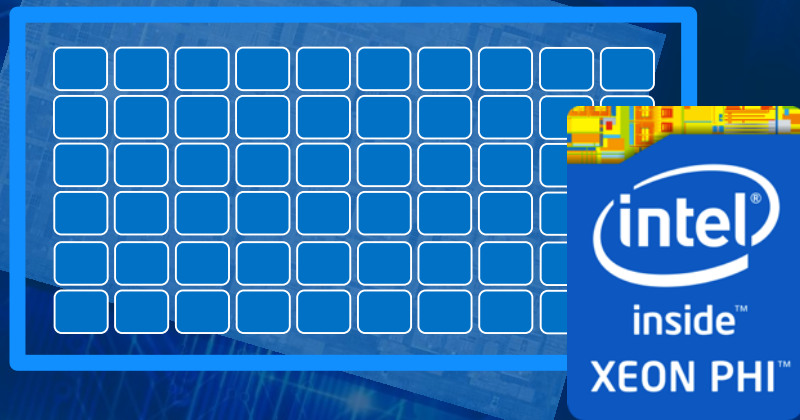
Many less powerful.

Design Question – *Our difference: Same programming models, languages, optimizations and tools*



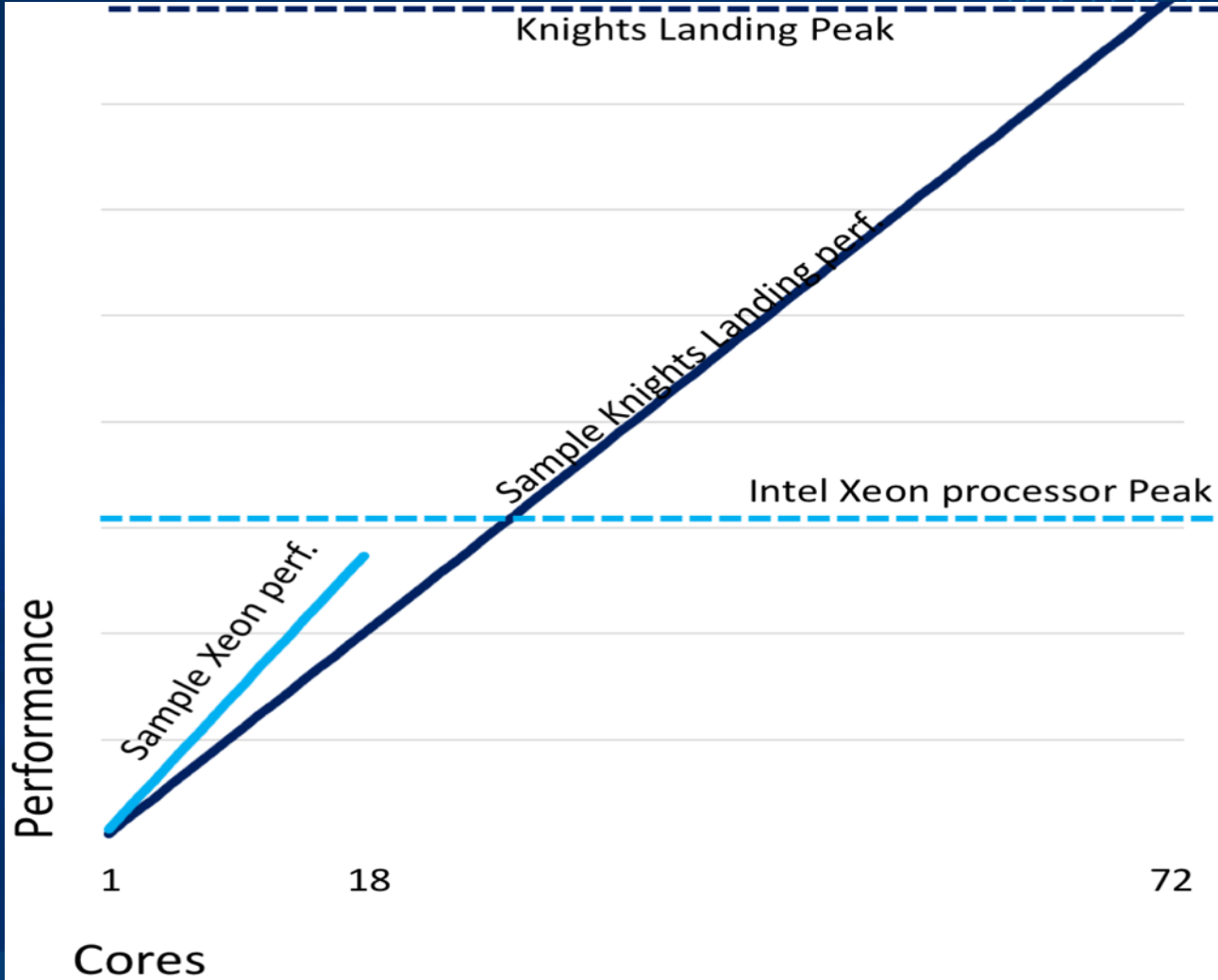
A few powerful

vs.



Many less powerful.

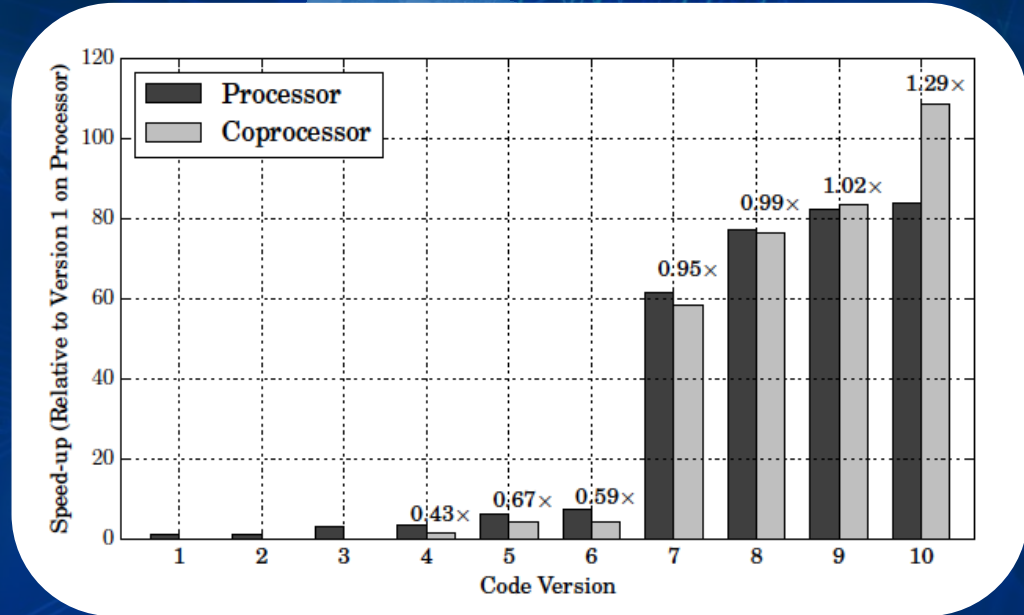
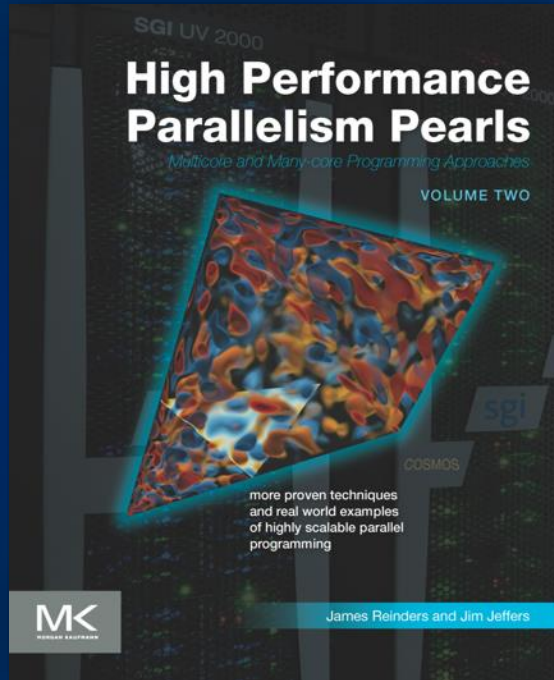
Same programming models, languages, optimizations and tools.



vision

span from *few cores* to *many cores*
with consistent models,
languages, tools, and techniques

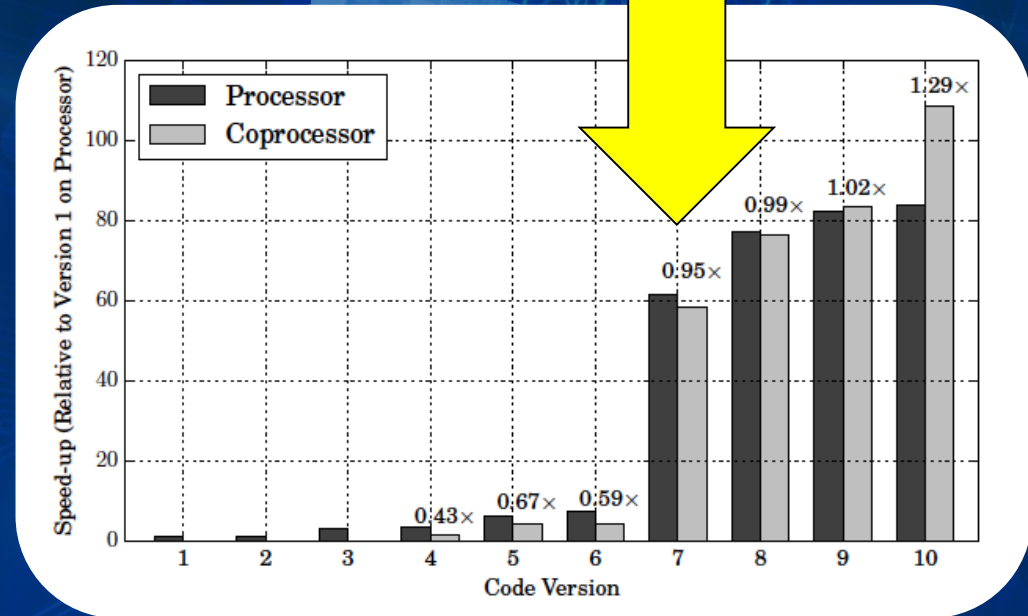
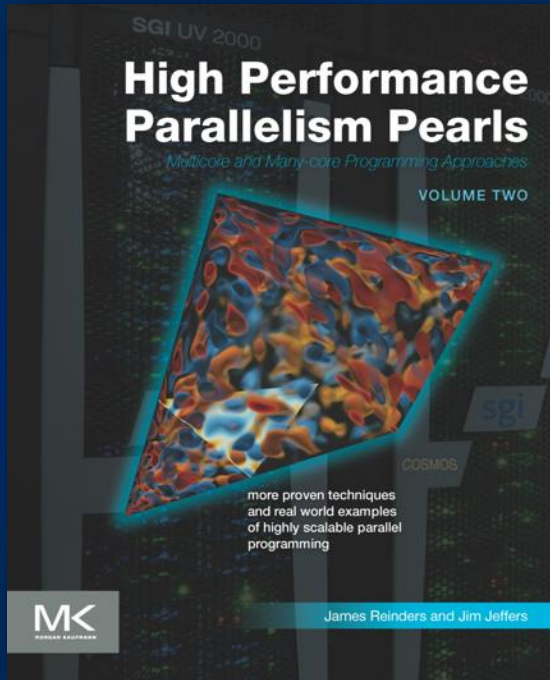
#Moderncode: COSMOS



Book Cover Background: Photo of the COSMOS@DiRAC SGI UV2000 based Supercomputer manufactured by SGI, Inc and operated by the Stephen Hawking Centre for Theoretical Cosmology, University of Cambridge. Photo courtesy of Philip Mynott. Book Cover Foreground: 3D visualization of statistical fluctuations in the Cosmic Microwave Background, the remnant of the first measurable light after the Big Bang. CMB data is from the Planck satellite and is the topic of Chapter 10 providing insights into new physics and how the universe evolved. Visualization rendered with Intel's OSPRay ray tracing open source software by Gregory P. Johnson and Timothy Rowley, Intel Corporation.

#Moderncode: COSMOS

What?



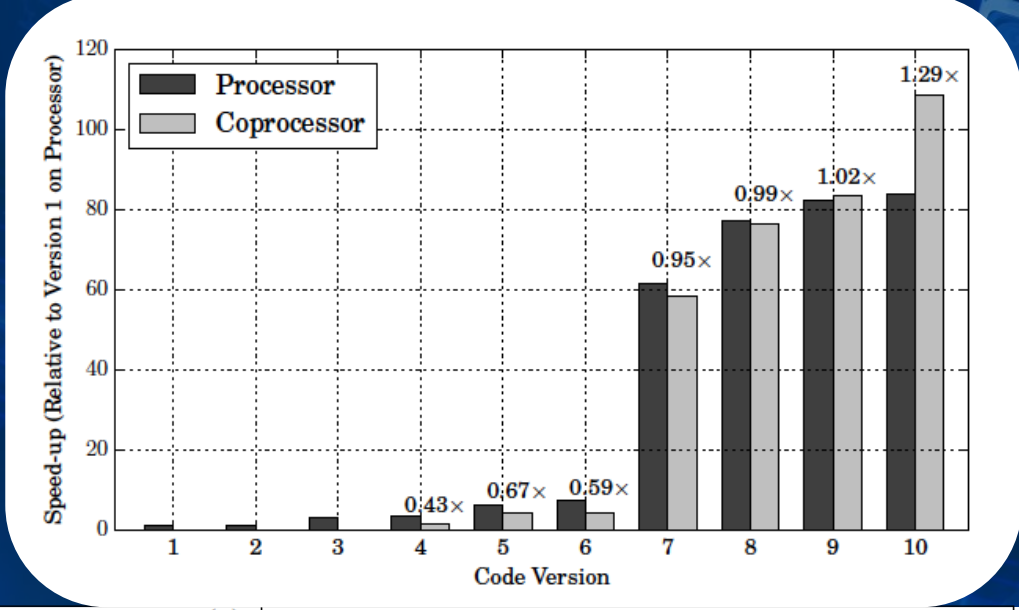
Book Cover Background: Photo of the COSMOS@DiRAC SGI UV2000 based Supercomputer manufactured by SGI, Inc and operated by the Stephen Hawking Centre for Theoretical Cosmology, University of Cambridge. Photo courtesy of Philip Mynott. Book Cover Foreground: 3D visualization of statistical fluctuations in the Cosmic Microwave Background, the remnant of the first measurable light after the Big Bang. CMB data is from the Planck satellite and is the topic of Chapter 10 providing insights into new physics and how the universe evolved. Visualization rendered with Intel's OSPRay ray tracing open source software by Gregory P. Johnson and Timothy Rowley, Intel Corporation.

High Performance Parallelism Pearls

Cosmic Microwave Background Analysis: Nested Parallelism in Practice

Volume 2, Chapter 10

We find that using a simple trapezium rule integrator combined with hand-selected sampling points (to improve accuracy in areas of interest) provides sufficient numerical accuracy to obtain a physically meaningful result, and the reduced space and time requirements of this simplified method give a speed-up of O(10x).



Version	Processor (s)	Coprocessor (s)	Comment
1	2887.0	-	Original code.
2	2610.0	-	Loop simplification.
3	882.0	-	Intel® MKL integration routines and function inlining.
4	865.9	1991.6	Flattened loops and introduced OpenMP threads.
5	450.6	667.9	Loop reordering and manual nested threading.
6	385.6	655.0	Blocked version of the loop (for cache).
7	46.9	49.5	Numerical integration routine (Trapezium Rule).
8	37.4	37.7	Reduction with DGEMM.
9	35.1	34.5	Data alignment (for vectorization).
10	34.3	26.6	Tuning of software prefetching distances.

The image shows the cover of a book titled 'Intel Xeon Phi Knights Landing Edition'. The cover features a close-up, colorful image of a microchip. The text on the cover is white and blue. At the top, it says 'INTEL XEON PHI PROCESSOR HIGH PERFORMANCE PROGRAMMING'. Below that, in a blue bar, it says 'KNIGHTS LANDING EDITION'. At the bottom left, there is a logo for 'MK' (Morgan Kaufmann) and the authors' names: 'Jim Jeffers | James Reinders | Avinash Sodani'.

INTEL® XEON PHI™
PROCESSOR
HIGH PERFORMANCE
PROGRAMMING

KNIGHTS LANDING EDITION

Jim Jeffers | James Reinders | Avinash Sodani

MK
MORGAN KAUFMANN

Knights Landing

2nd Generation Intel® Xeon Phi™

ALL ABOUT PARALEL PROGRAMMING

Threading, Vectorization, Data Locality
Fortran, C, C++ (plus a little Python)
OpenMPI, MPI, TBB

New with Knights Landing:

AVX-512,

High Bandwidth Memory (MCDRAM),

Cluster Mode,

Omni-Path

INTEL® XEON PHI™ PROCESSOR HIGH PERFORMANCE PROGRAMMING

KNIGHTS LANDING EDITION

Jim Jeffers | James Reinders | Avinash Sodani

MK
MORGAN KAUFMANN

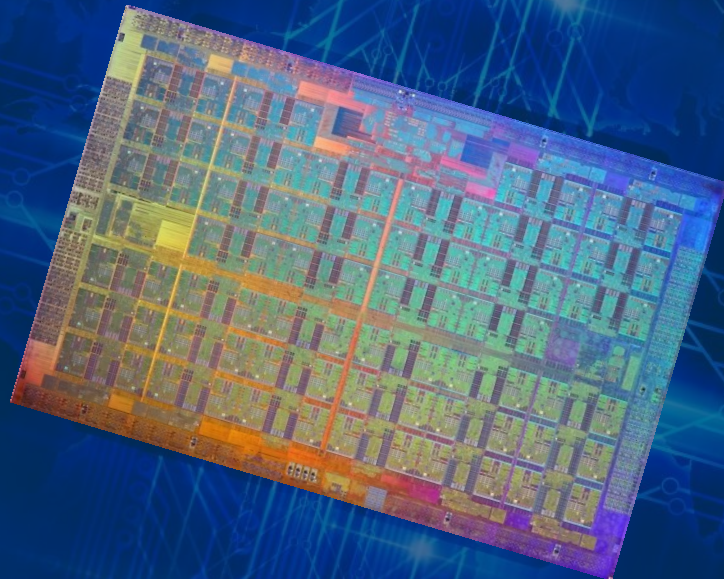


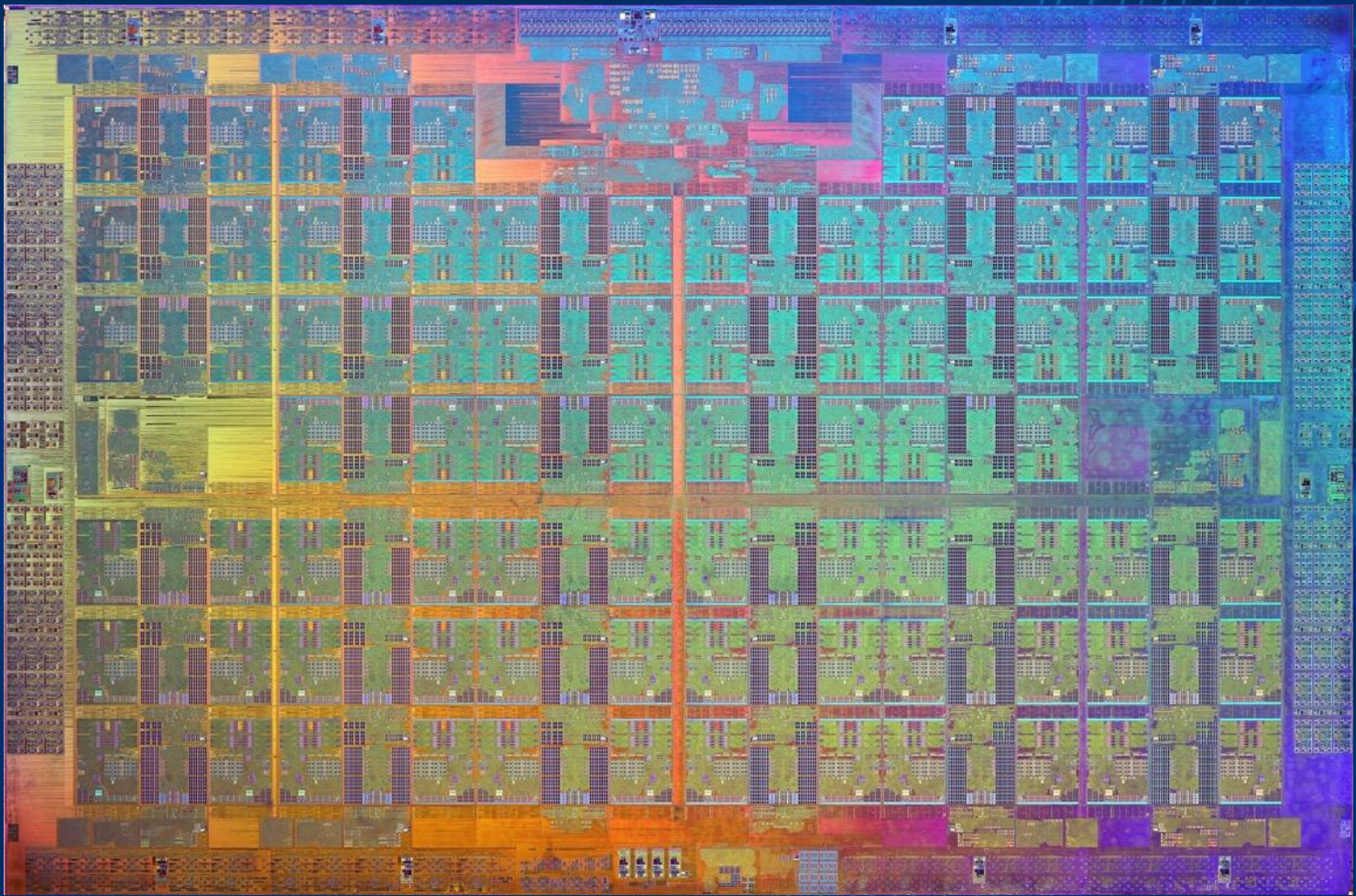
INTEL® XEON PHI™ PROCESSOR HIGH PERFORMANCE PROGRAMMING

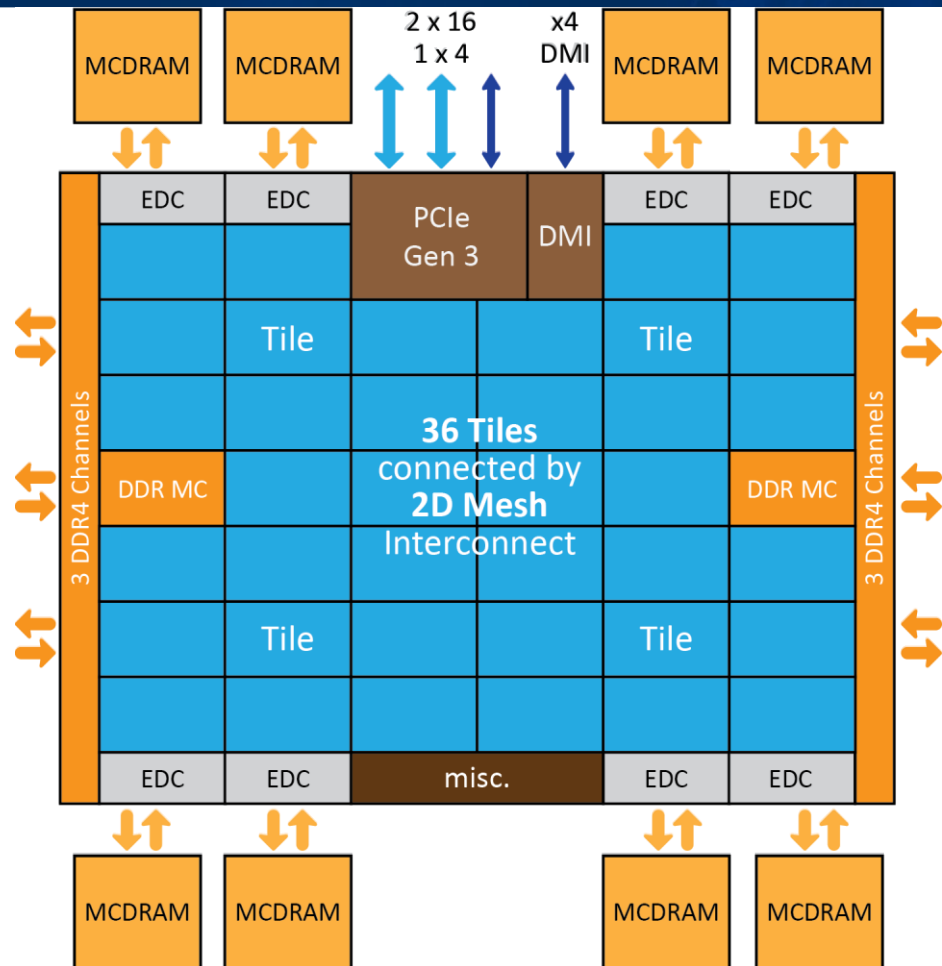
KNIGHTS LANDING EDITION

Jim Jeffers | James Reinders | Avinash Sodani

MK
MORGAN KAUFMANN



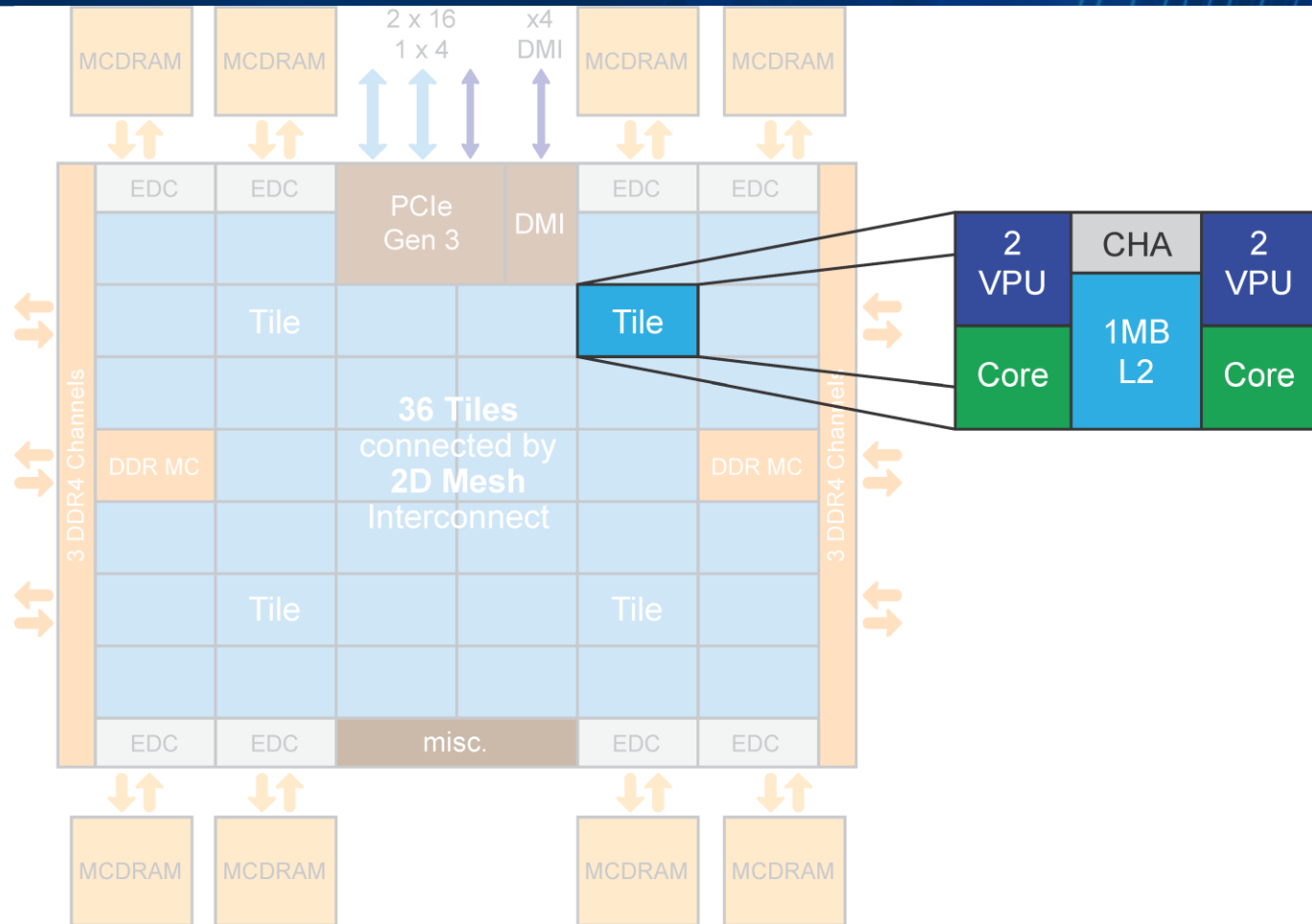


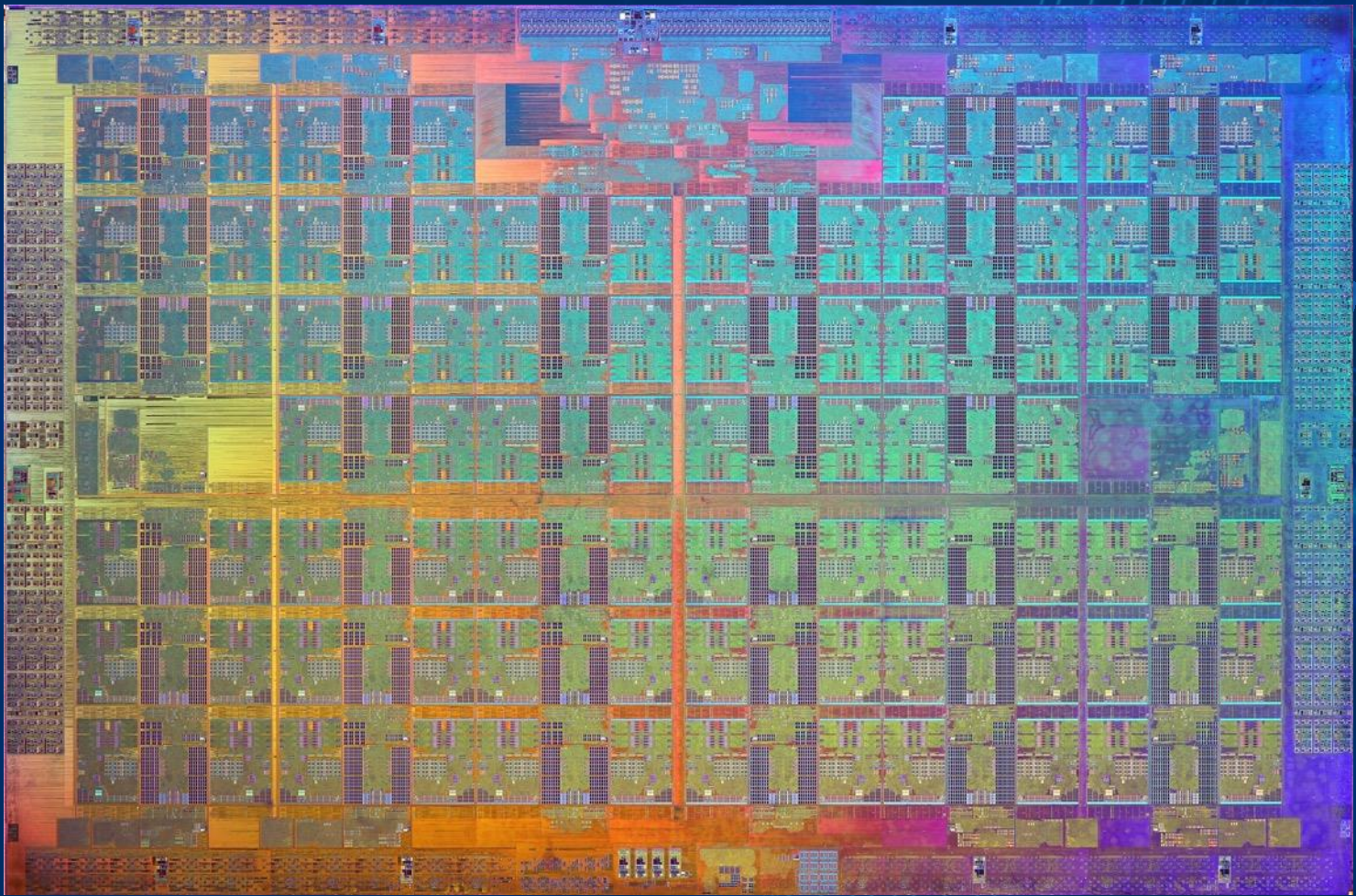


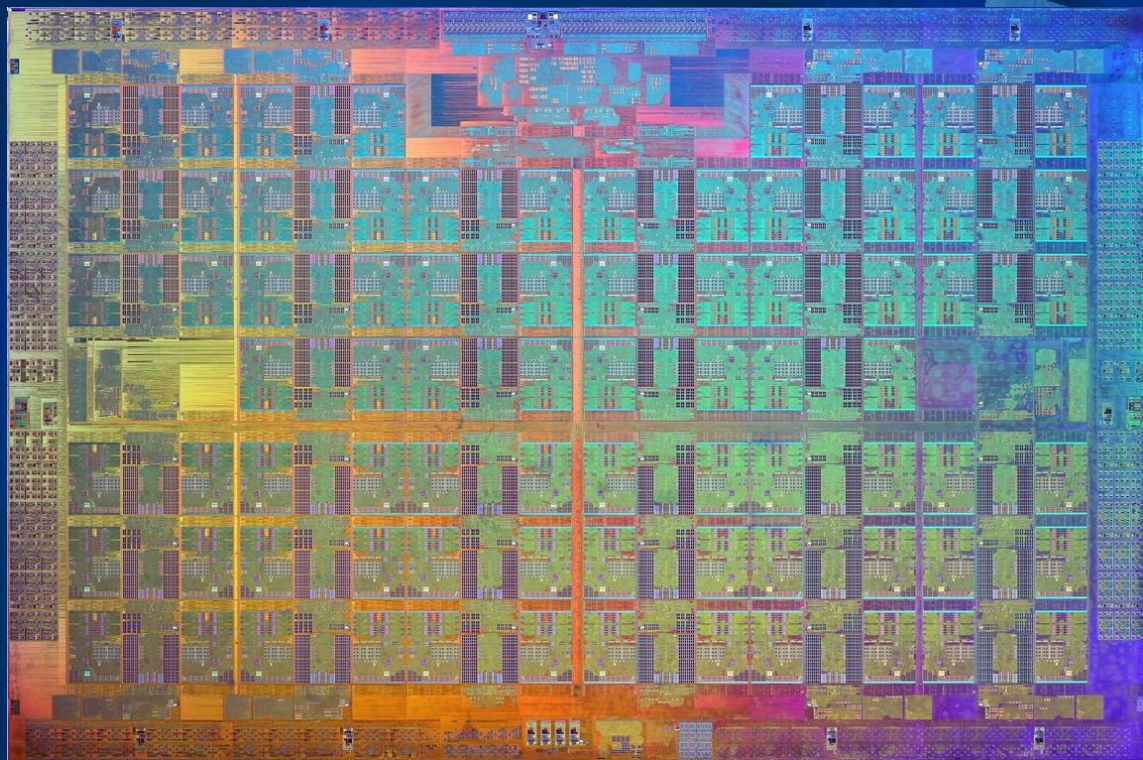
2nd Generation Intel® Xeon Phi™ Products

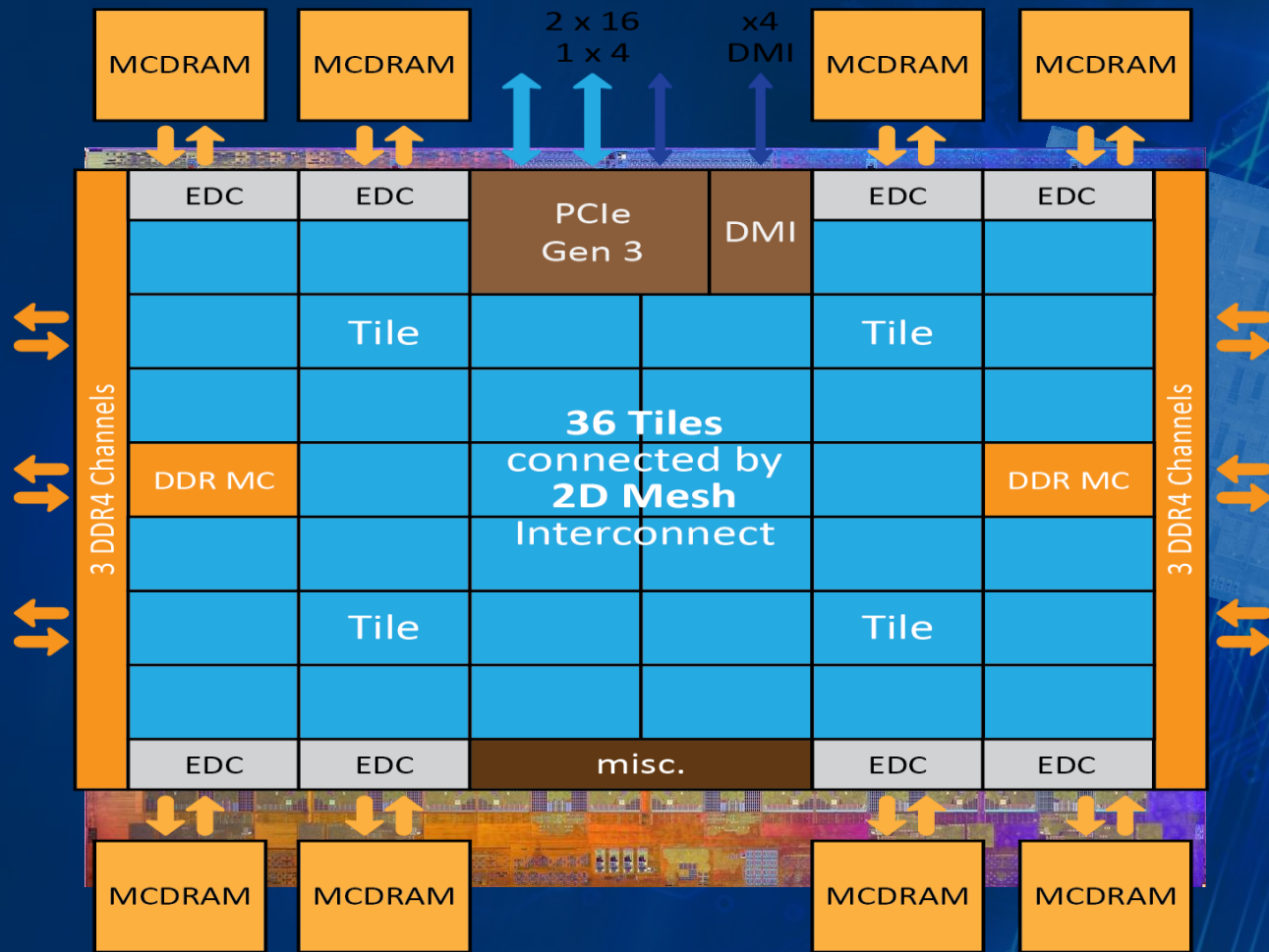
Codename:

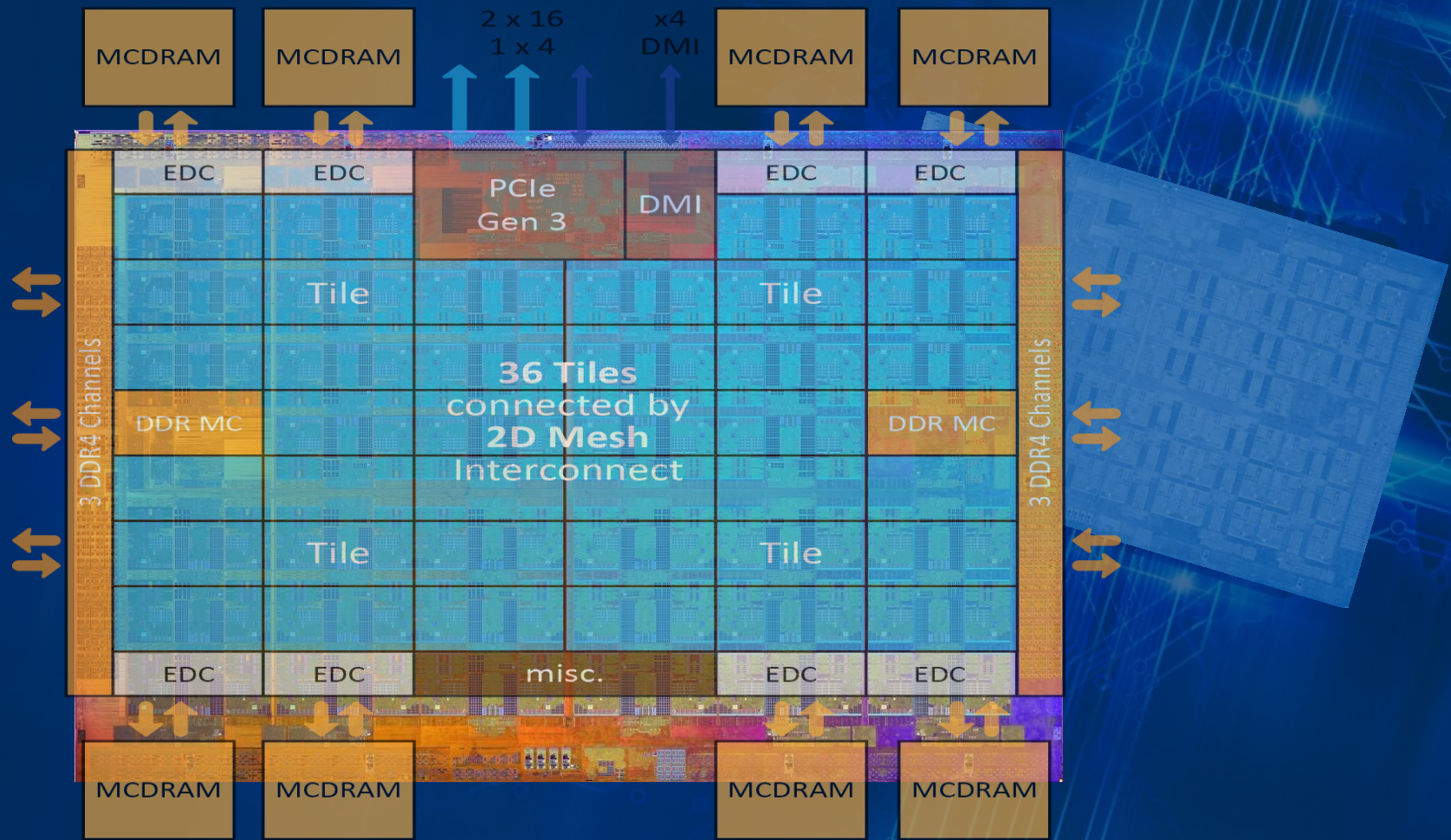
Knights Landing











Knights Landing: Next Intel® Xeon Phi™ Processor

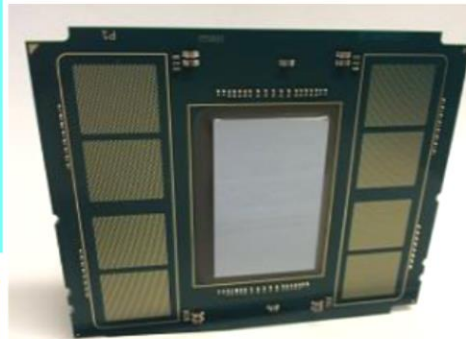
Intel® Many-Core Processor targeted for HPC and Supercomputing

First **self-boot** Intel® Xeon Phi™ processor that is **binary compatible** with main line IA. Boots standard OS.

Significant improvement in scalar and vector performance

Integration of **Memory on package**: innovative memory architecture for high bandwidth and high capacity

Integration of **Fabric on package**



Three products

KNL Self-Boot

KNL Self-Boot w/ Fabric

KNL Card

(Baseline)

(Fabric Integrated)

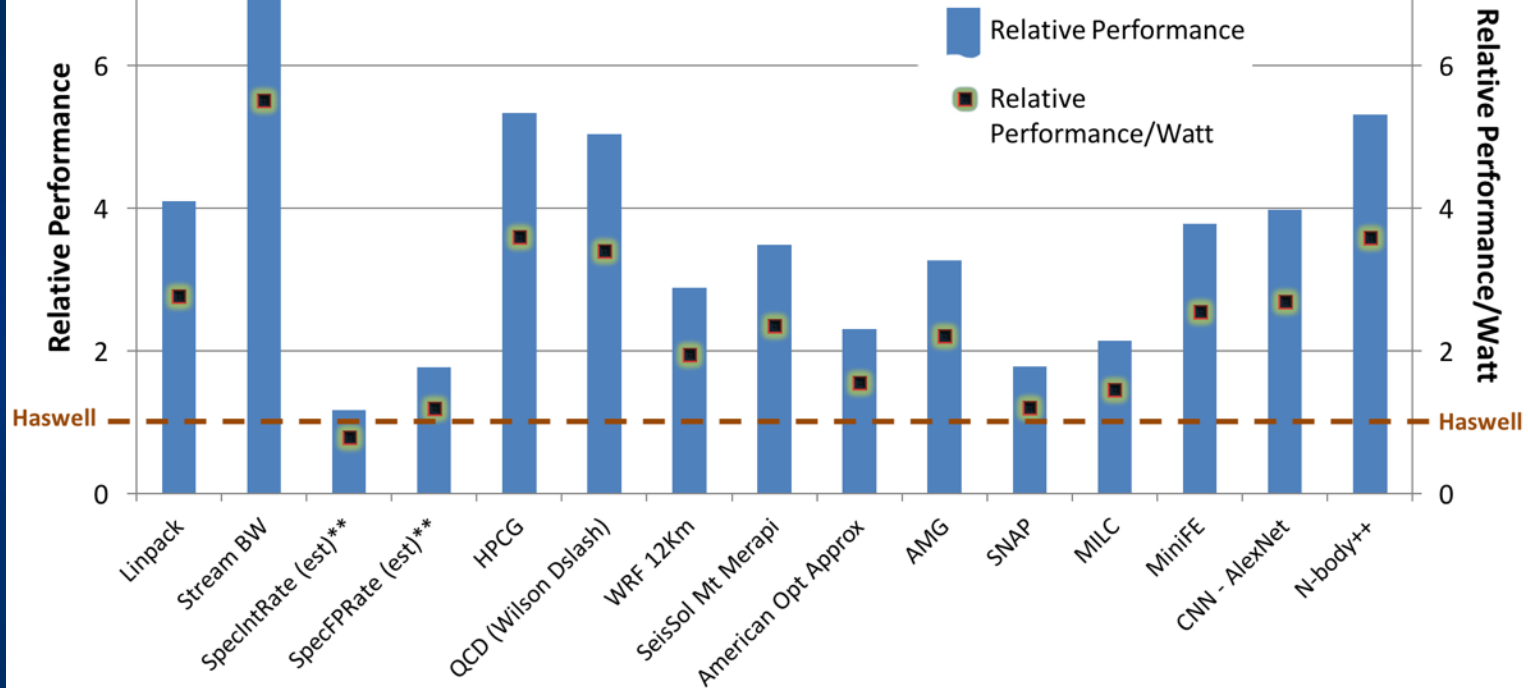
(PCIe-Card)

Potential future options subject to change without notice.

All timeframes, features, products and dates are preliminary forecasts and subject to change without further notification.

action turning signals of digital signals and data processing
and data processing signals and data processing
and data processing signals and data processing
and data processing signals and data processing

pre-production Knights Landing (A0) processor vs. a Haswell (Intel® Xeon® E5-2697-V3 processor)



Measurements on a pre-production Knights Landing (A0) processor. Results subject to change on production parts.

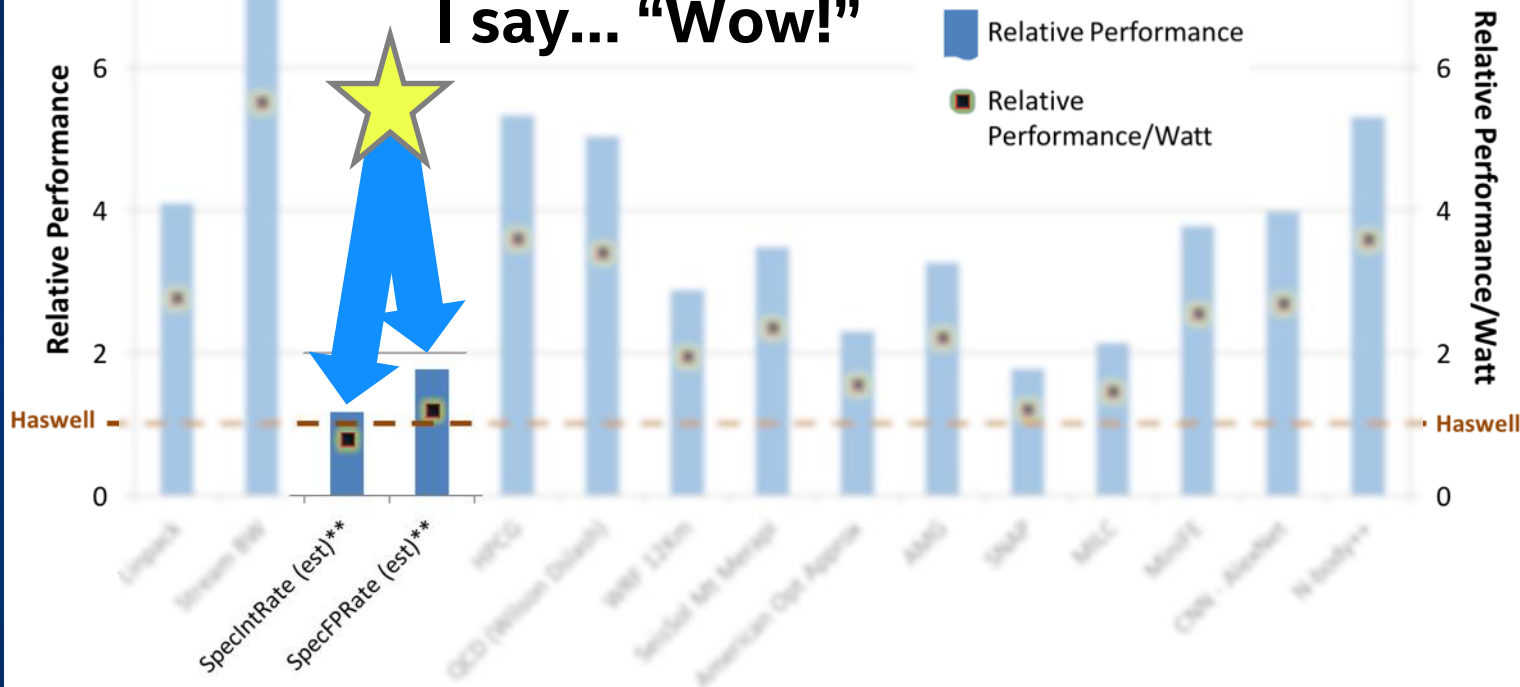
** Early A0 processor runs, estimated SPEC (not fully compliant nor submitted runs)

++ Code and measurement done by colfaxresearch.com

Data courtesy of Intel Corporation

pre-production Knights Landing (A0) processor vs. a Haswell (Intel® Xeon® ES-2697-V3 processor)

I say... "Wow!"



Measurements on a pre-production Knights Landing (A0) processor. Results subject to change on production parts.

** Early A0 processor runs, estimated SPEC (not fully compliant nor submitted runs)

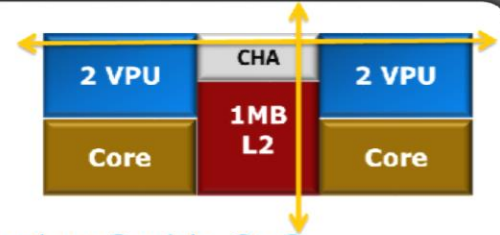
** Code and measurement done by colfaxresearch.com

Data courtesy of Intel Corporation

AVX-512, High Bandwidth Memory, Cluster Mode, Omni-Path

512-bit
vectors
via
AVX-512
means:
High
performance
and
Binary
compatibility

KNL Tile: 2 Cores, each with 2 VPU
1M L2 shared between two Cores



Core: Changed from Knights Corner (KNC) to KNL. Based on 2-wide OoO Silvermont™ Microarchitecture, but with many changes for HPC.

4 thread/core. Deeper OoO. Better RAS. Higher bandwidth. Larger TLBs.

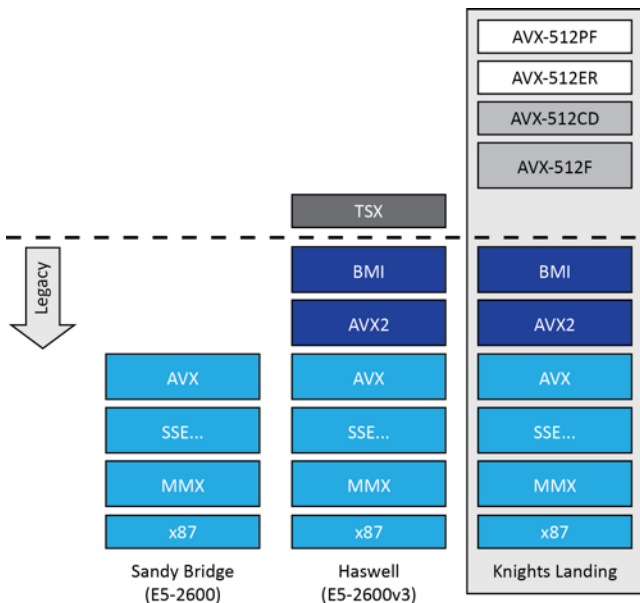
2 VPU: 2x AVX512 units. 32SP/16DP per unit. X87, SSE, AVX1, AVX2 and EMU

L2: 1MB 16-way. 1 Line Read and ½ Line Write per cycle. Coherent across all Tiles

CHA: Caching/Home Agent. Distributed Tag Directory to keep L2s coherent. MESIF protocol. 2D-Mesh connections for Tile

AVX-512, High Bandwidth Memory, Cluster Mode, Omni-Path

Support of vectors under 512-bits in size means: **Binary compatibility**



KNL implements all legacy instructions

- Legacy binary runs w/o recompilation
- KNC binary requires recompilation

KNL introduces AVX-512 Extensions

- 512-bit FP/Integer Vectors
- 32 registers, & 8 mask registers
- Gather/Scatter

Conflict Detection: Improves Vectorization

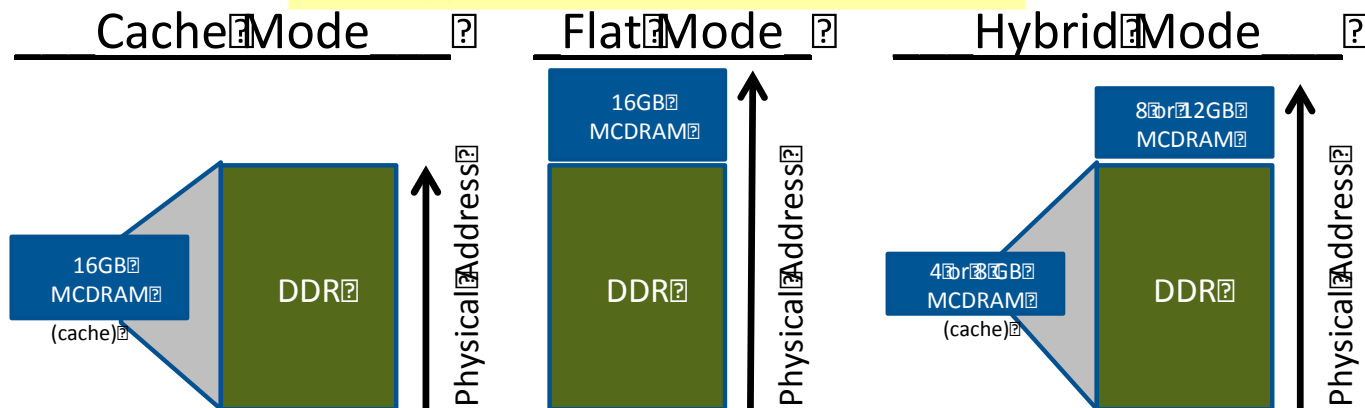
Prefetch: Gather and Scatter Prefetch

Exponential and Reciprocal Instructions

16MB
on package
DRAM
means:
Performance
with
options

Memory Modes

Three Modes. Selected at boot



- SW-Transparent, Mem-side cache
- Direct mapped. 64B lines.
- Tags part of line
- Covers whole DDR range

- MCDRAM as regular memory
- SW-Managed
- Same address space

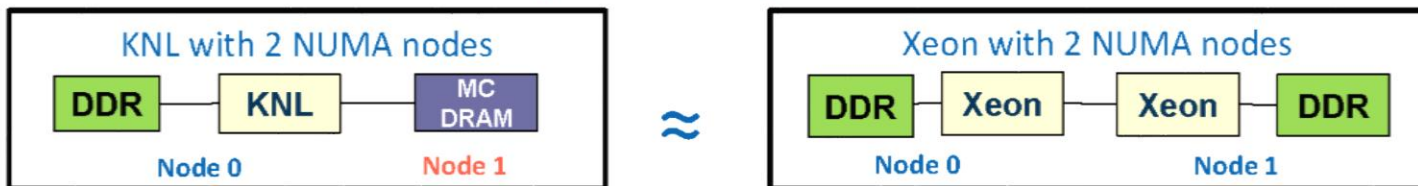
- Part cache, Part memory
- 25% or 50% cache
- Benefits of both

Standard
“NUMA”
node
recognition
by BIOS, OS,
and
applications.

Systems will
eventually
have this as
a “norm.”

Flat MCDRAM: SW Architecture

MCDRAM exposed as a separate NUMA node



Memory allocated in DDR by default → Keeps non-critical data out of MCDRAM.

Apps explicitly allocate critical data in MCDRAM. Using two methods:

- “Fast Malloc” functions in High BW library (<https://github.com/memkind>)
 - Built on top to existing *libnuma* API
- “FASTMEM” Compiler Annotation for Intel Fortran

Flat MCDRAM with existing NUMA support in Legacy OS

Flat MCDRAM SW Usage: Code Snippets

C/C++
“high
bandwidth”
malloc
or new

Fortran
“high
bandwidth”
allocatables

C/C++ ([*https://github.com/memkind](https://github.com/memkind))

Allocate into DDR

```
float *fv;  
fv = (float *)malloc(sizeof(float)*100);
```



Allocate into MCDRAM

```
float *fv;  
fv = (float *)hbw_malloc(sizeof(float) * 100);
```

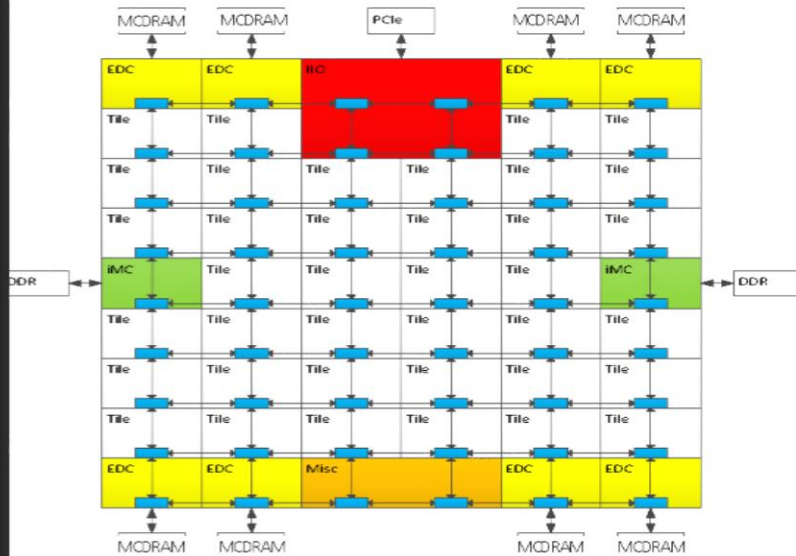
Intel Fortran

Allocate into MCDRAM

```
c Declare arrays to be dynamic  
REAL, ALLOCATABLE :: A(:)  
  
!DEC$ ATTRIBUTES, FASTMEM :: A  
  
NSIZE=1024  
c allocate array 'A' from MCDRAM  
c  
ALLOCATE (A(1:NSIZE))
```

Mesh interconnect means:
Higher Performance with options

KNL Mesh Interconnect



Mesh of Rings

- Every row and column is a (half) ring
- YX routing: Go in Y → Turn → Go in X
- Messages arbitrate at injection and on turn

Cache Coherent Interconnect

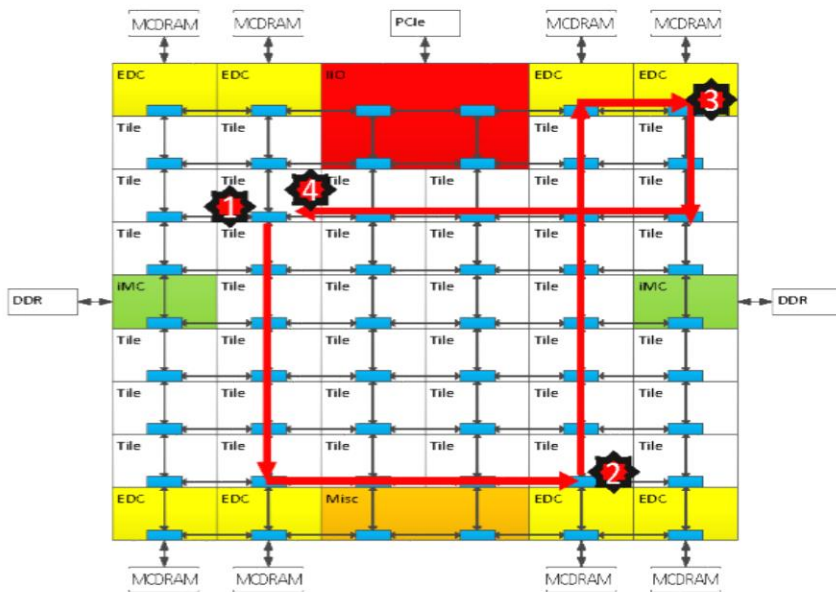
- MESIF protocol (F = Forward)
- Distributed directory to filter snoops

Three Cluster Modes

- (1) All-to-All (2) Quadrant (3) Sub-NUMA Clustering

All-to-all
to
use all
cores
together
uniformly
for
instance –
OpenMP
across
all of KNL

Cluster Mode: All-to-All



Address uniformly hashed across all distributed directories

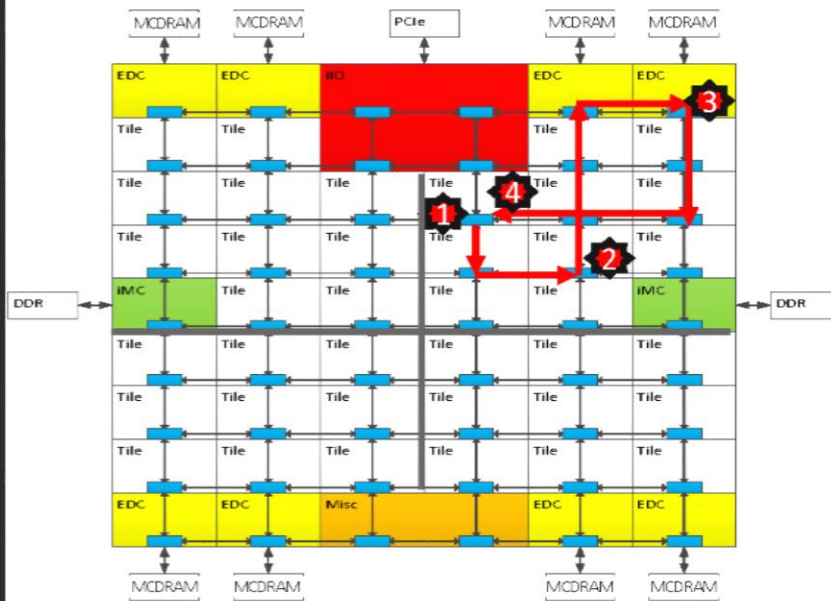
No affinity between Tile, Directory and Memory

Most general mode. Lower performance than other modes.

Typical Read L2 miss

1. L2 miss encountered
2. Send request to the distributed directory
3. Miss in the directory. Forward to memory
4. Memory sends the data to the requestor

Cluster Mode: Sub-NUMA Clustering (SNC)



Each Quadrant (Cluster) exposed as a separate NUMA domain to OS.

Looks analogous to 4-Socket Xeon

Affinity between Tile, Directory and Memory

Local communication. Lowest latency of all modes.

SW needs to NUMA optimize to get benefit.

1) L2 miss, 2) Directory access, 3) Memory access, 4) Data return

1) L2 miss, 2) Directory access, 3) Memory access, 4) Data return

Integrated Fabric
means
Lower Cost,
power;
Higher Density;
Plus the
Advancements
of Omni-Path
in latency,
bandwidth &
scaling

KNL with Omni-Path™

Omni-Path™ Fabric integrated *on package*

First product with integrated fabric

Connected to KNL die via 2 x16 PCIe* ports

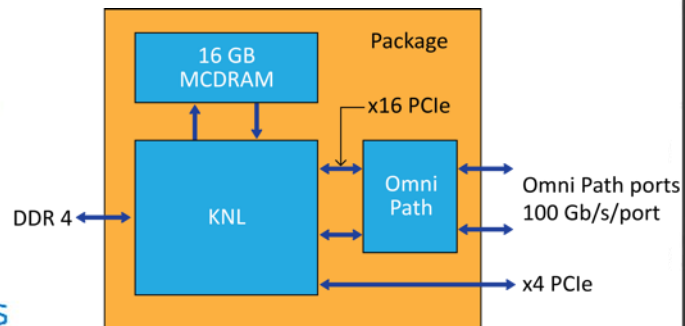
Output: 2 Omni-Path ports

- 25 GB/s/port (bi-dir)

Benefits

- Lower cost, latency and power
- Higher density and bandwidth
- Higher scalability

*On package connect with PCIe semantics, with MCP optimizations for physical layer



Intel® Parallel Studio XE 2017 Beta

Submitted by RAVI (Intel) on March 28, 2016 [Translate](#)

[f Share](#) [Tweet](#) [g+ Share](#)

Contents

- [How to enroll in the Beta program](#)
- [What's New in the 2017 Beta](#)
- [Frequently Asked Questions](#)
- [Beta duration and schedule](#)
- [Support](#)
- [Beta webinars](#)
- [Beta Release Notes](#)
- [Known issues](#)
- [Next steps](#)

How to enroll in the Beta program

Complete the pre-beta survey at [registration link](#)

[http://
software.intel.com/
articles/
intel-parallel-studio-xe-2017-beta](http://software.intel.com/articles/intel-parallel-studio-xe-2017-beta)

BETA for “2017” Product – NOW

Vectorization advisor

Many factors impact achieving good vectorization for our applications. The Vectorization Advisor directly analyzes an application and provides feedback on the extent of current vectorization and on possible steps to achieve more effective vectorization. Vectorization Advisor works with any compiler although some features in the Intel® compilers will increase the effectiveness of advice from the Vectorization Advisor tool. It is like having an expert sitting next to us who never tires of digging into an application to analyze what is really happening.

The Vectorization Advisor is one of the two major *workflows* (feature sets) available in the Intel® Advisor “2016” and later versions. The Intel Advisor also includes a thread prototyping feature set which can be useful for analysis of scaling for threads. In this chapter, we focus on using the Vectorization Advisor to help us maximize our vectorization performance.

What is new with Knights Landing in this chapter?

AVX-512 and the Vectorization Advisor within the Intel® Advisor tool.

How close is my application to maximum performance? Insight into this is helped by a “roofline model” analysis, in the Advisor Roofline Report section.

Intel® Advisor



Development > Tools > Resources >

Intel® Advisor



Vectorization Optimization and Thread Prototyping

- Vectorize & thread code or performance “dies”
- Easy workflow + data + tips = faster code faster
- Prioritize, Prototype & Predict performance gain

I will talk about some NEW “2017” features – which help Intel Xeon processors tuning and Intel Xeon Phi processor tuning BOTH – Of Course!

Memory Access Pattern Report

MEMORY ACCESS PATTERN REPORT

An initial **Survey** analysis of hot loops often identifies inefficient memory access patterns as a main bottleneck. Memory access patterns issues are the toughest and most frequent performance problem in code not yet modernized for vector SIMD parallelism.

Applying *straightforward* SIMD and threading optimizations often does not provide desirable speedups because some parts of applications (including vectorized hot loops) become *memory bound*. Memory-access-patterns-bound code is just one sub-type of a larger memory-bound class of problems, along with *memory-bandwidth-bound* and partially overlapping with *memory-latency-bound* sub-types.

Can recommend:

- AoS to SoA
- AoSoA
- Use of SDLT
- Use of MCDRAM

The screenshot shows a software interface with three tabs: Source, Assembly, and Details. The Details tab is active and shows a table with columns for Source, Stride, and Loop instance footprint. A row is highlighted in blue, showing the source code line `b[:n-1:2] = a[:n/2]+c[:n-1:2]` with a stride of `[1] [32]` and a loop instance footprint of 188B. To the right of the table, there are details for a 'Gather (irregular)' operation, including operand size (32 bits), operand type (float32, int32), vector length (8), and memory access footprint (188B). Below this, 'Gather/scatter details' are shown, including a pattern of 'Constant (non-unit)', instruction access details, horizontal stride (8 bytes), vertical stride (128 bytes), a constant mask, and 100% active elements in the mask.

Source	Stride	Loop instance footprint
<code>subroutine size(ntimes,ld,n,ctime,dtime,a,b,c,d,e,aa,bb,cc)</code>		
<code>induction variables</code>		
<code>coupled induction variables</code>		
<code>integer ntimes,ld,n,i,nl,j,k</code>		
<code>real a(n),b(n),c(n),d(n),e(n),aa(ld,n),bb(ld,n),cc(ld,n)</code>		
<code>real t1,t2,chksum,ctime,dtime,csld</code>		
<code>call init(ld,n,a,b,c,d,e,aa,bb,cc,'s128')</code>		
<code>call fortttime(t1)</code>		
<code>do n1= 1,2*ntimes</code>		
<code>a[:n/2] = b[:n-1:2]-d[:n/2]</code>		
<code>b[:n-1:2] = a[:n/2]+c[:n-1:2]</code>	<code>[1] [32]</code>	188B
<code>call dummy(ld,n,a,b,c,d,e,aa,bb,cc,1.)</code>		
<code>enddo</code>		
<code>call fortttime(t2)</code>		
<code>t2= t2-t1-ctime-(dtime*float(2*ntimes))</code>		
<code>chksum= csld(n,a)+csld(n,b)</code>		
<code>call check(chksum,2*ntimes*(n/2),n,t2,'s128')</code>		
<code>return</code>		
<code>end</code>		

Gather/Scatter Report



Mask Utilization and FLOPS Profiler

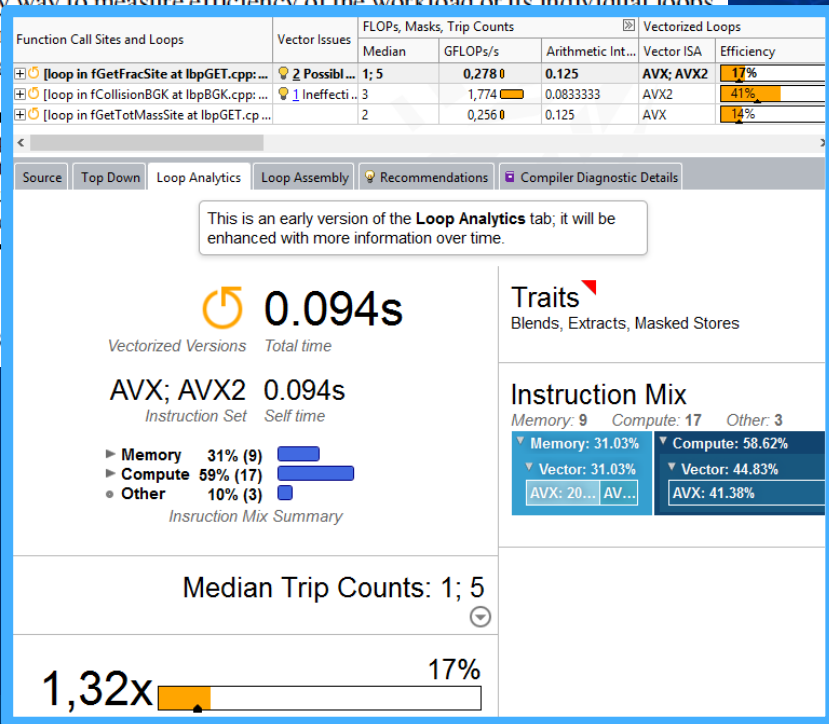
MASK UTILIZATION AND FLOPS PROFILER

Counting FLOPs on Knights Landing is not ~~directly~~ supported by the hardware because there is no accounting for the values in mask registers when AVX-512 instructions are counted. ~~Certain~~ capabilities of the Advisor ~~tools~~ can make up for this lack of ~~direct~~ hardware support.

FLOP/s is a key way to measure efficiency of the workload or its individual loops or kernels. Measure the performance of target hardware.

In this book, and especially for floating point operations. Elsewhere, it is not used to measure FLOP/s. We need to

While masked vectorization, they also



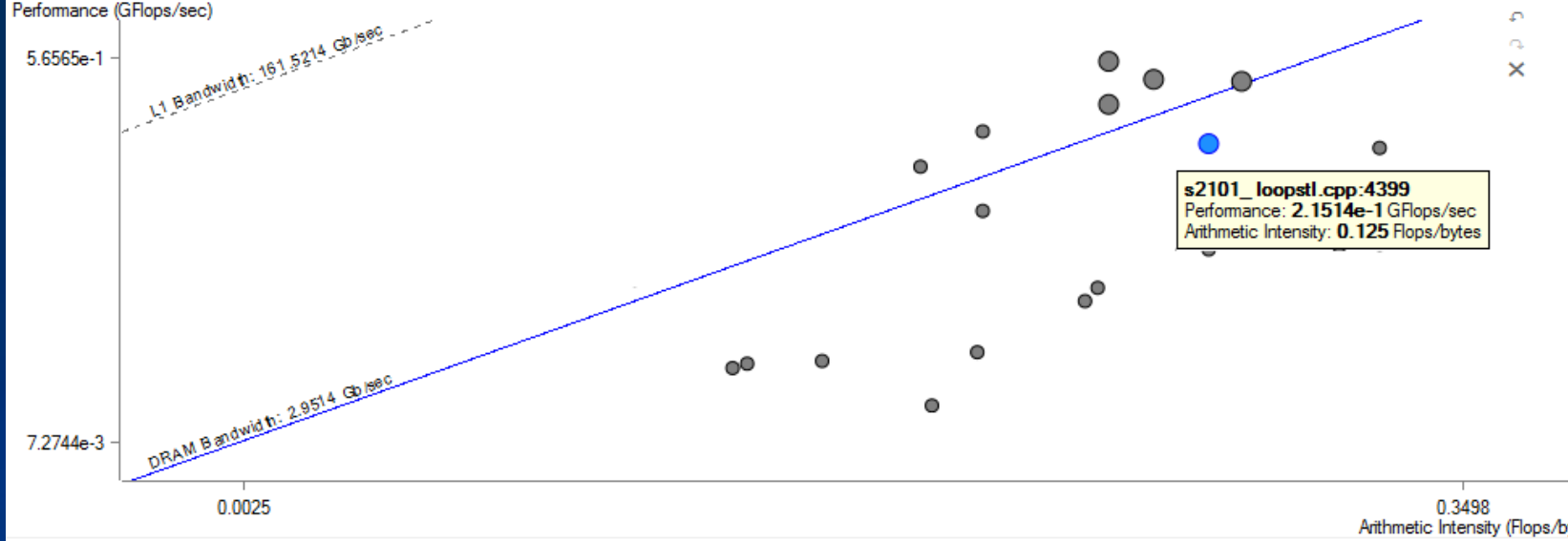
Mask-aware:

- FLOPs Report
- Vector Efficiency
- Memory Access Pattern, (coming soon):
- Roofline Analysis Graph

Vectorization efficiency and FLOP/s in Survey Report and Loop Analytics.



Roofline Report in Intel® Advisor 2017 (soon)



Source | Top Down | Loop Assembly | Recommendations | Compiler Diagnostic Details

File: loopstl.cpp:4399 s2101_

Line	Source	Total Time	%	Loop Time	%	Traits
4399	for (i__ = 1; i__ <= i__2; ++i__)	0,031s		1,859s		
4400	aa[i__ + i__ * aa_dim1] += bb[i__ + i__ * bb_dim1] * cc[i__ + i__	1,828s				FMA
4401	* cc_dim1];					

...supplements AI-based analysis with a dynamic FLOP/s profile and peak FLOPs and memory sub-system throughput levels providing enlightening “bounds and bottlenecks” analysis for complex workloads.

Intel® Distribution for Python*

- Faster NumPy/SciPy performance powered by Intel® Math Kernel Library and Intel® Threading Building Blocks and Intel® Data Analytics Acceleration Library
- ~3X speedups on single thread
- Easy installation
- Python 2.7 & 3.5
- Windows & Linux & OS X



CALLING ALL PYTHON DEVELOPERS

Join the Intel® Distribution for Python* 2017 Beta program with new features to get faster performance from your Python applications.

[Access Beta](#)

Intel® Distribution for Python* 2017 Beta

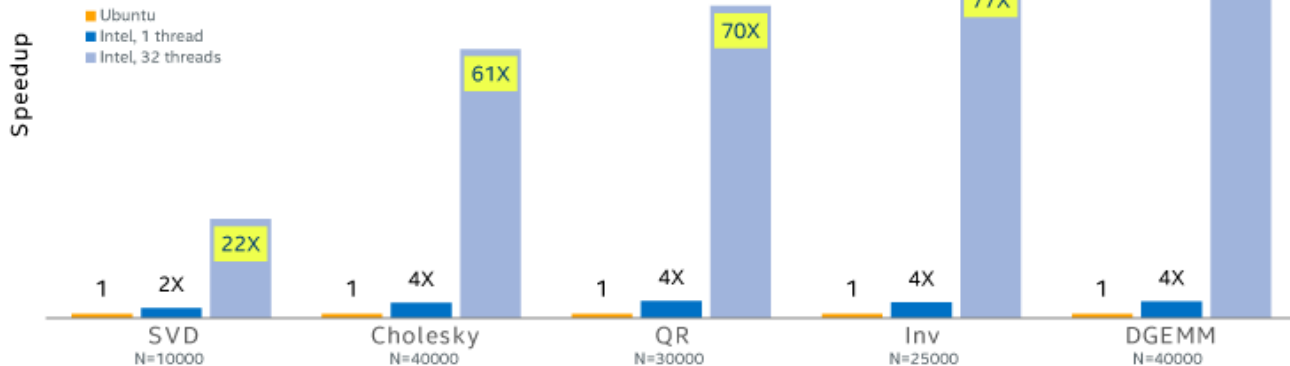
Faster performance from Python packages powered by Intel® Math Kernel Library (Intel® MKL)

<http://bit.ly/intel-python>

Some Performance Benchmarks

Intel® Distribution for Python*

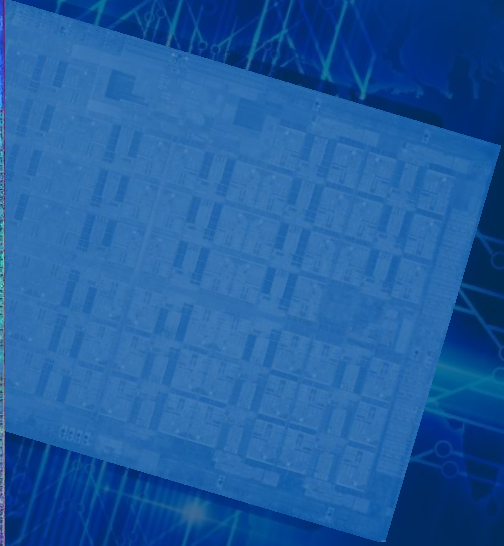
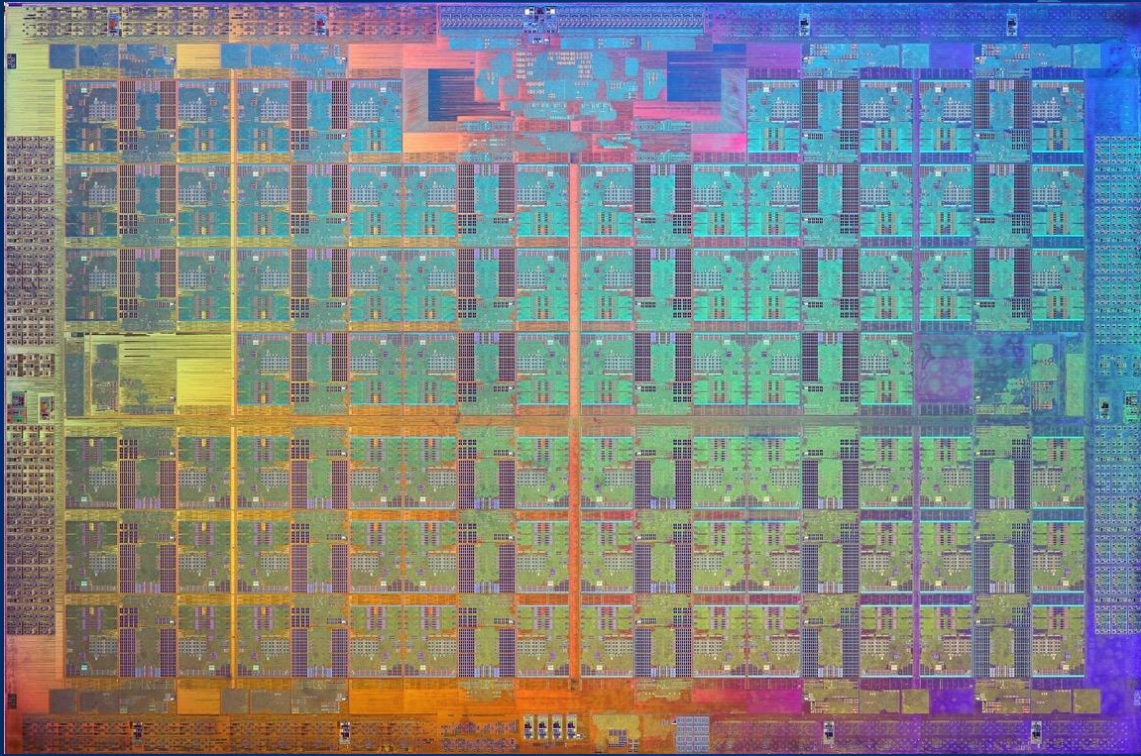
Python* Performance Boost on Select Numerical Functions Intel® Distribution for Python (Technical Preview) vs. Ubuntu Python

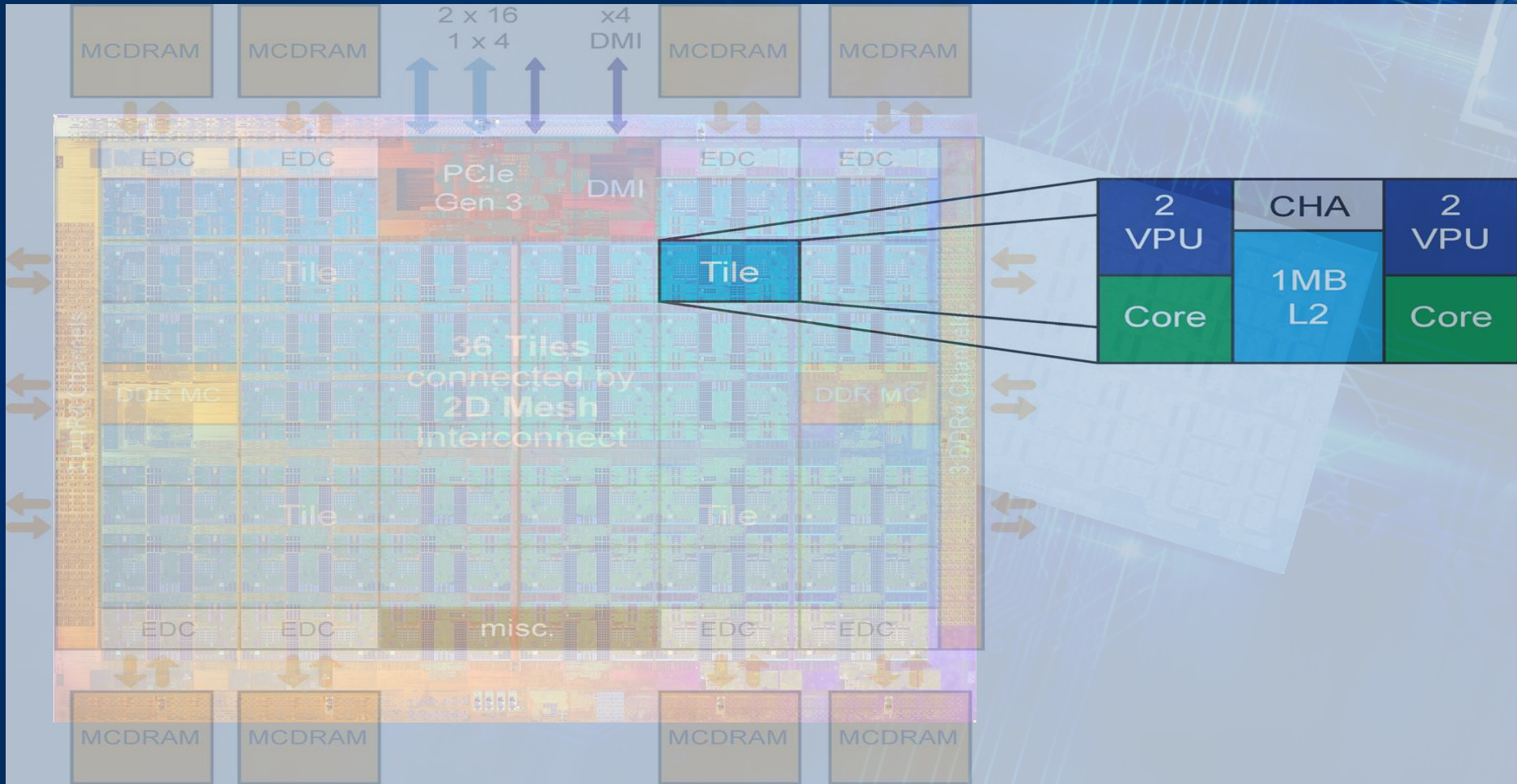


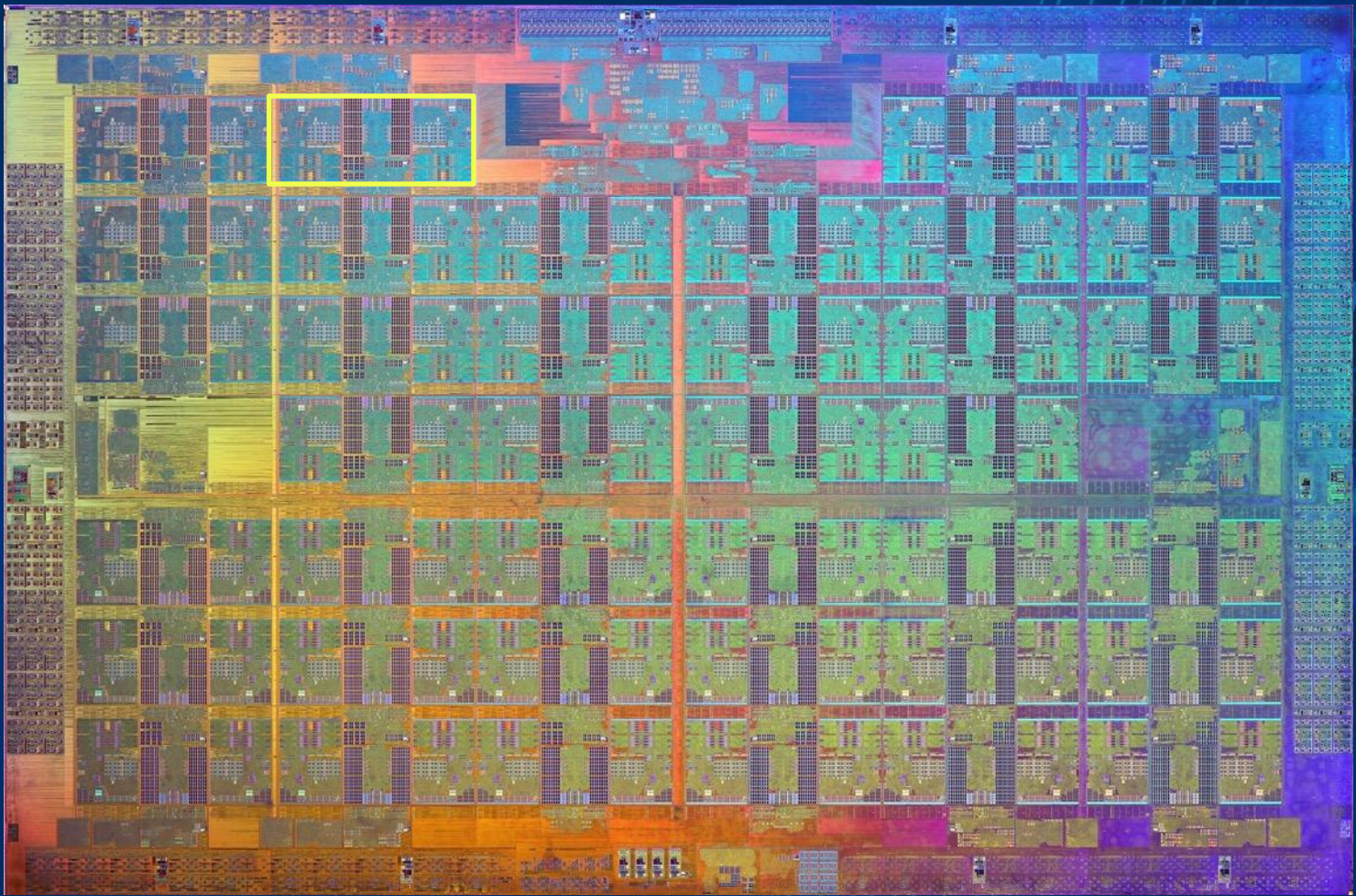
Configuration info: - Versions: Intel® Distribution for Python 2.7.10 Technical Preview 1 (Aug 03, 2015), Ubuntu® built Python*: Python 2.7.10, NumPy 1.9.2 built with gcc 4.8.4; Hardware: Intel® Xeon® CPU E5-2698 v3 @ 2.30GHz [2 sockets, 16 cores each, HT=OFF], 64 GB of RAM, 8 DIMMS of 8GB@2133MHz; Operating System: Ubuntu 14.04 LTS.

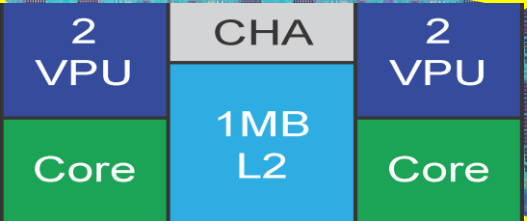
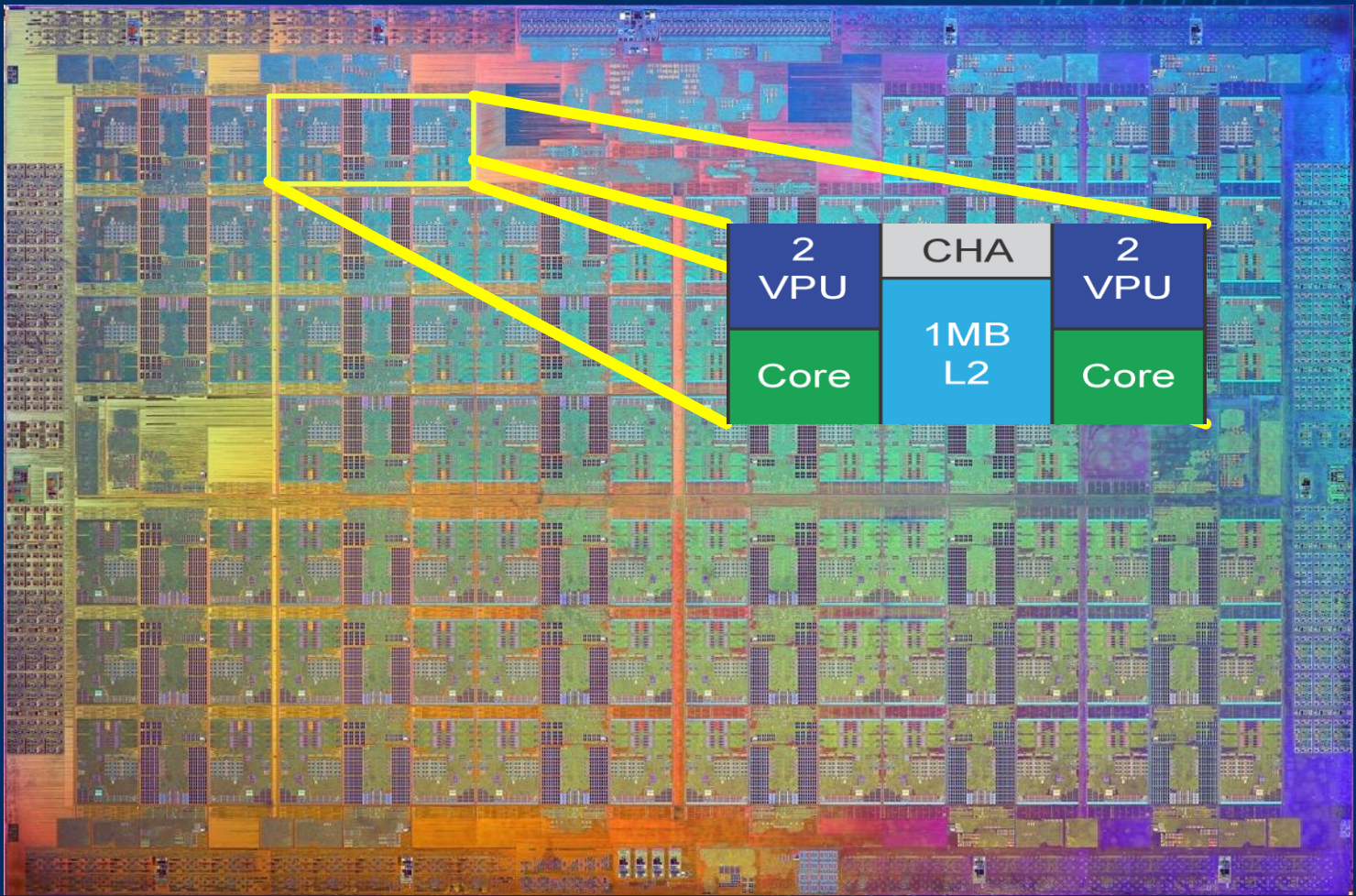
Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may 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. * Other brands and names are the property of their respective owners. Benchmark Source: Intel Corporation

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.











#ModernCode

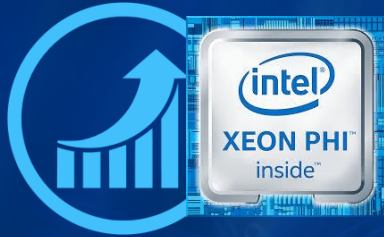
Optimizations for Intel® Xeon® and Intel® Xeon Phi™ products share the same:



- ✓ Languages
- ✓ Directives
- ✓ Libraries
- ✓ Tools

XeonPhiDeveloper.com

You can buy you very own Knight Landing development system today!



Highly-Parallel Performance
to develop on



All the Software Tools & Libraries
you need



Support & Training
to help you succeed

Leading edge platform capabilities, performance to deliver multi-threaded, vectorized software for today's HPC workloads !

Q&A

Thank you!

