# Knights Landing Intel® Xeon Phi™ CPU: Path to Parallelism with General Purpose Programming

Avinash Sodani

Knights Landing Chief Architect

Senior Principal Engineer, Intel Corp.

# Legal

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# Roadmap

Why parallelism? Why general purpose?

Knights Landing: Intel® Xeon Phi™ Processor Architecture

Performance, Applications, SW Tools and Support

Future Trends and Challenges

# Computing demand continues to grow

**HPC**

Solving bigger and more complex scientific problems to improve day to day lives

**Cloud & online data and services**

Massive growth in online data and services, spurring growth in data centers

**Machine Learning**

Promise of solving problems that are very hard to solve algorithmically

**Genetics & Medical**

Cure for life threatening diseases. Deeper understanding to prevent diseases.

**IoT**

Connected devices slated to grow over 20B by 2020 (Gartner). Drive backend datacenter needs

**Climate/Weather**

**Data Analytics**

Growth across both traditional and emerging usages. Investment both at government and commercial levels

# CPU Compute Growth Trends



Average core count growth in bottom 5 systems in Top-500 list (496-500) vs. Historical growth in per thread performance

Core counts and per thread performance are not from same system

Core counts

Spec Int Speed

Source: Intel, March 2016; Based on historical analysis across the SPEC CPU integer speed and rate benchmarks. Core counts growth based on the data on last five systems from the historical Top-500 lists (www.top500.org)

"Power-wall" slowed frequency increase over last decade
Core counts on exponential growth – much faster than single core performance

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# Exponential Growth in Data Parallelism



**Massive flops per chip with vector and core count growth**

# Exponential Growth in Data Parallelism



**Massive flops per chip with vector and core count growth**

# Exponential Growth in Data Parallelism



**Massive flops per chip with vector and core count growth**

# Parallelism is the way forward

**Trend** → Lots of thread- and data-level parallelism

Systems becoming highly parallel. More vectors, more cores per CPU, more CPUs per system

Single thread performance increasing at slower pace

Significant performance potential for applications that parallelize and vectorize

# Plenty of solutions in play

## Several parallel HW options. Vary with usage

- CPUs
- GPUs
- FPGA solutions
- Application specific accelerators

## Different ways to program them

- MPI/OpenMP/TBB/etc.
- Language extensions with pragmas, etc.
- Different GPU programming models: CUDA, OpenCL, OpenACC, etc.
- Accelerator-specific API
- Research models that try to encompass both CPU and GPU programming

# Software story is important

**Software generally live for decades. Much longer than hardware**

- Important to change software for parallelism in a manner that preserves investment

- They should continue to run and perform well on future hardware

- Choose programming models that lasts long



Chart showing software lifespans from ~1970 to 2016:
- SCRYU/Tetra - CFD
- scSTREAM - CFD
- Dalton Quantum Chemistry
- WRF - Weather
- NWCHEM - Chemistry
- LAPACK - Solvers
- PETSc - Solvers
- UKMO Unified Model - Weather
- Pam-Crash
- Spice
- NASTRAN

Timeline: 1970  1980  1990  2000  2010

# Knights Landing: First Intel® Xeon Phi™ <u>Processor</u>

**Enables extreme parallel performance with general purpose programming**



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**

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

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# Knights Landing Overview

**TILE**

| 2 VPU | CHA | 2 VPU |
|-------|-----|-------|
| | 1MB L2 | |
| Core | | Core |



2 x16
1 x4

X4
DMI

MCDRAM  MCDRAM  MCDRAM  MCDRAM

3 DDR4 CHANNELS

EDC  EDC  PCIe Gen 3  DMI  EDC  EDC

Tile

DDR MC        DDR MC

**36 Tiles connected by 2D Mesh Interconnect**

3 DDR4 CHANNELS

EDC  EDC  misc  EDC  EDC

MCDRAM  MCDRAM  MCDRAM  MCDRAM

**Package**

**Chip:** up to **36 Tiles** interconnected by **2D Mesh**

**Tile**: 2 Cores + 2 VPU/core + 1 MB L2

**Memory: MCDRAM:** up to 16 GB on-package; High BW

**DDR4**: 6 channels @ 2400  up to 384GB

**IO:** 36 lanes PCIe Gen3. 4 lanes of DMI for chipset

**Node**: 1-Socket

**Fabric:** Intel® Omni-Path Fabric on-package (not illustrated)

**Vector Peak Perf**: 3+TF DP and 6+TF SP Flops

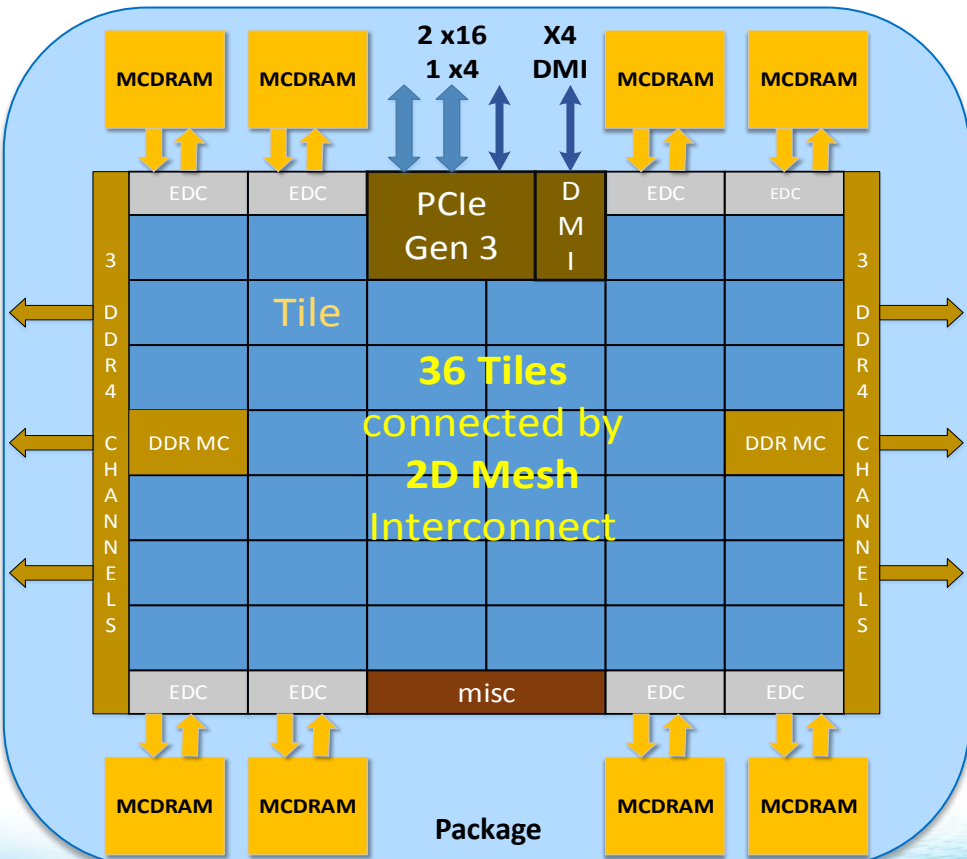**Scalar Perf**: ~3x over Knights Corner

**Streams Triad (GB/s)**: MCDRAM : 450+; DDR: ~90

Note: not all specifications shown apply to all Knights Landing SKUs
Source Intel:  All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. 1Binary Compatible with Intel Xeon processors using Haswell Instruction Set (except TSX). 2Bandwidth numbers are based on STREAM-like memory access pattern when MCDRAM used as flat memory. Results have been estimated based on internal Intel analysis and are provided for informational purposes only.  Any difference in system hardware or software design or configuration may affect actual performance.

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



**Core**: New OoO Core. Balances power efficiency, parallel and single thread performance.

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

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

**CHA**: **C**aching/**H**ome **A**gent. Distributed Tag Directory to keep L2s coherent. MESIF protocol. 2D-Mesh connections for Tile

# Many Trailblazing Improvements in KNL. But why?

| Improvements | What/Why |
|---|---|
| Self Boot Processor | No PCIe bottleneck. Be same as general purpose CPU |
| Binary Compatibility with Xeon | Runs all legacy software. No recompilation. |
| New OoO Core | ~3x higher ST performance over KNC |
| Improved Vector Density | 3+ TFLOPS (DP) peak per chip |
| New AVX 512 ISA | New 512-bit Vector ISA with Masks |
| New memory technology: MCDRAM + DDR | Large High Bandwidth Memory → MCDRAM<br>Huge bulk memory → DDR |
| New on-die interconnect: Mesh | High BW connection between cores and memory |
| Integrated Fabric: Omni-Path | Better scalability to large systems. Lower Cost |

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# Core & VPU

Balanced power efficiency, single thread performance and parallel performance

- 2-wide Out-of-order core
- 4 SMT threads
- 72 in-flight instructions.
- 6-wide execution
- 64 SP and 32 DP Flop/cycle

- Dual ported DL1 → to feed 2 VPU
- Two-level TLB. Large page support
- Gather/Scatter engine
- Unaligned load/store support

- Core resources **shared** or **dynamically repartitioned** between active threads

- General purpose IA core

Thread Selection points

Icache (32KB 8-way)

iTLB

Fetch & Decode

Bpred

Allocate/ Rename

Retire

FP RS (20)

FP RS (20)

Integer Rename Buffer

Integer RF

FP Rename Buffers

FP RF

MEM RS(12)

Recycle Buffer

Int RS (12)

Int RS (12)

Vector ALUS

Vector ALUs

Legacy

TLBs

Dcache (32KB 8-way)

ALU

ALU

# KNL ISA

E5-2600 (SNB[1])    E5-2600v3 (HSW[1])    **KNL (Xeon Phi[2])**

| x87/MMX | x87/MMX | x87/MMX |
| SSE* | SSE* | SSE* |
| AVX | AVX | AVX |
| | AVX2 | AVX2 |
| | BMI | BMI |
| | TSX | |

LEGACY

AVX-512F

AVX-512CD

AVX-512PF

AVX-512ER

No TSX. Under separate CPUID bit

1. Previous Code name Intel® Xeon® processors
2. Xeon Phi = Intel® Xeon Phi™ processor

**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

**C**onflict **D**etection: Improves Vectorization

**P**refetch: Gather and Scatter Prefetch

**E**xponential and **R**eciprocal Instructions

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# AVX-512 CD: Instructions for enhance vectorization

```
for(i=0; i<16; i++) { A[B[i]]++;}
```

```
index = vload &B[i]          // Load 16 B[i]
old_val = vgather A, index   // Grab A[B[i]]
new_val = vadd old_val, +1.0 // Compute new values
vscatter A, index, new_val   // Update A[B[i]]
```

✖ Code is wrong if any values within B[i] are duplicated

```
index = vload &B[i]                               // Load 16 B[i]
pending_elem = 0xFFFF;                             // all still remaining
do {
    curr_elem = get_conflict_free_subset(index, pending_elem)
    old_val = vgather {curr_elem} A, index         // Grab A[B[i]]
    new_val = vadd old_val, +1.0                   // Compute new values
    vscatter A {curr_elem}, index, new_val         // Update A[B[i]]
    pending_elem = pending_elem ^ curr_elem        // remove done idx
} while (pending_elem)
```

## AVX-512 Conflict Detection

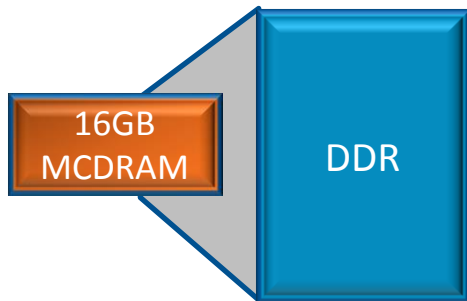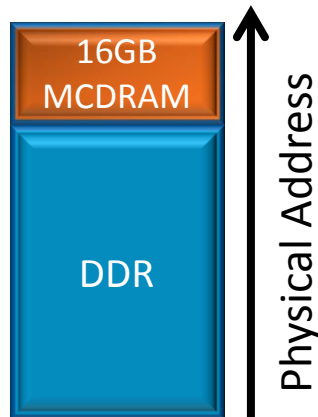| VPCONFLICT{D,Q} zmm1{k1}, zmm2/mem |
| VPBROADCASTM{W2D,B2Q} zmm1, k2 |
| VPTESTNM{D,Q} k2{k1}, zmm2, zmm3/mem |
| VPLZCNT{D,Q} zmm1 {k1}, zmm2/mem |

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# KNL Memory Modes

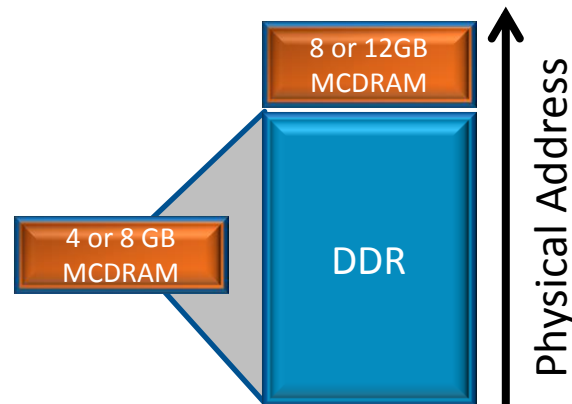**Three** Modes. Selected at boot

## Cache Mode

16GB MCDRAM

DDR

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

## Flat Mode

16GB MCDRAM

DDR

Physical Address

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

## Hybrid Mode
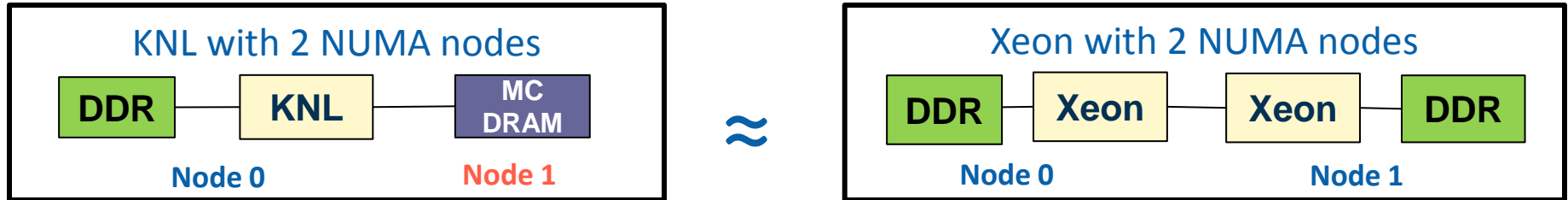
8 or 12GB MCDRAM

4 or 8 GB MCDRAM

DDR

Physical Address

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

# 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 <u>two</u> methods:

- **"Fast Malloc"** functions in High BW library (https://github.com/memkind/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++   (*https://github.com/memkind)

### Allocate into DDR

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

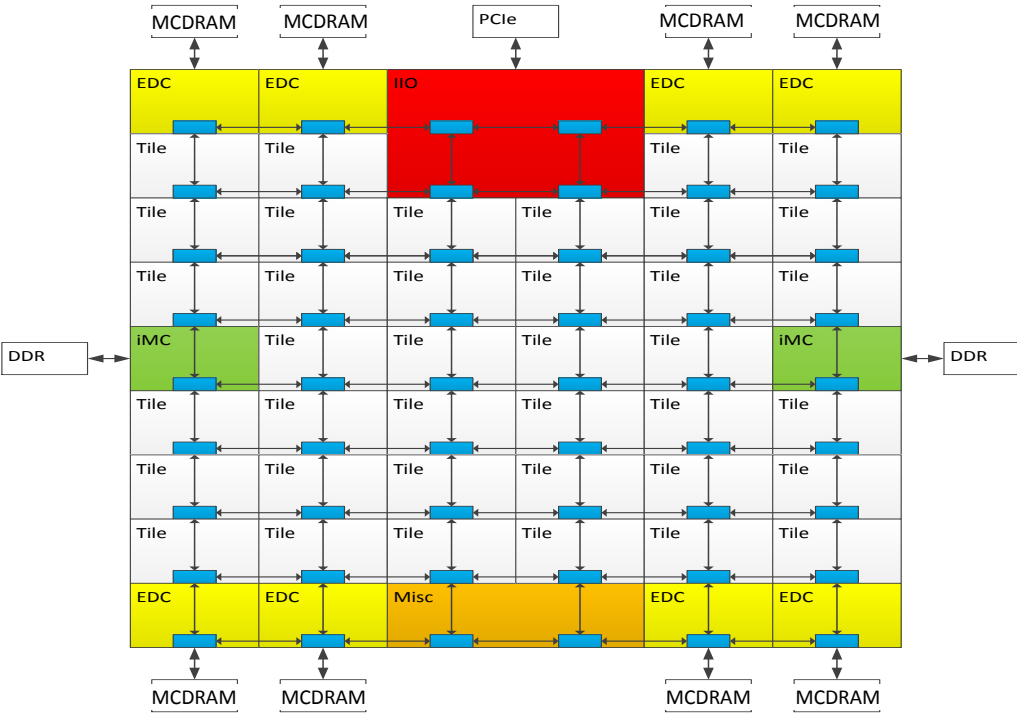### Allocate into MCDRAM

```
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))
```

# KNL Mesh Interconnect



## Mesh of Rings

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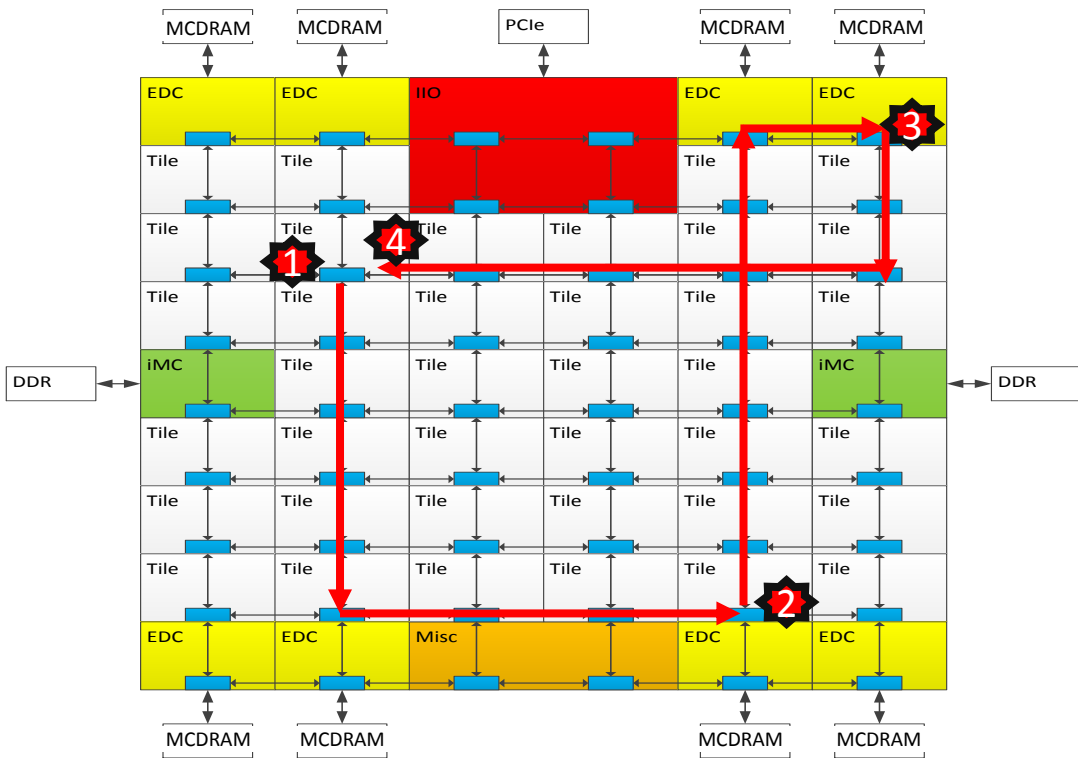- Messages arbitrate at injection and on turn

## 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

# Cluster Mode: All-to-All



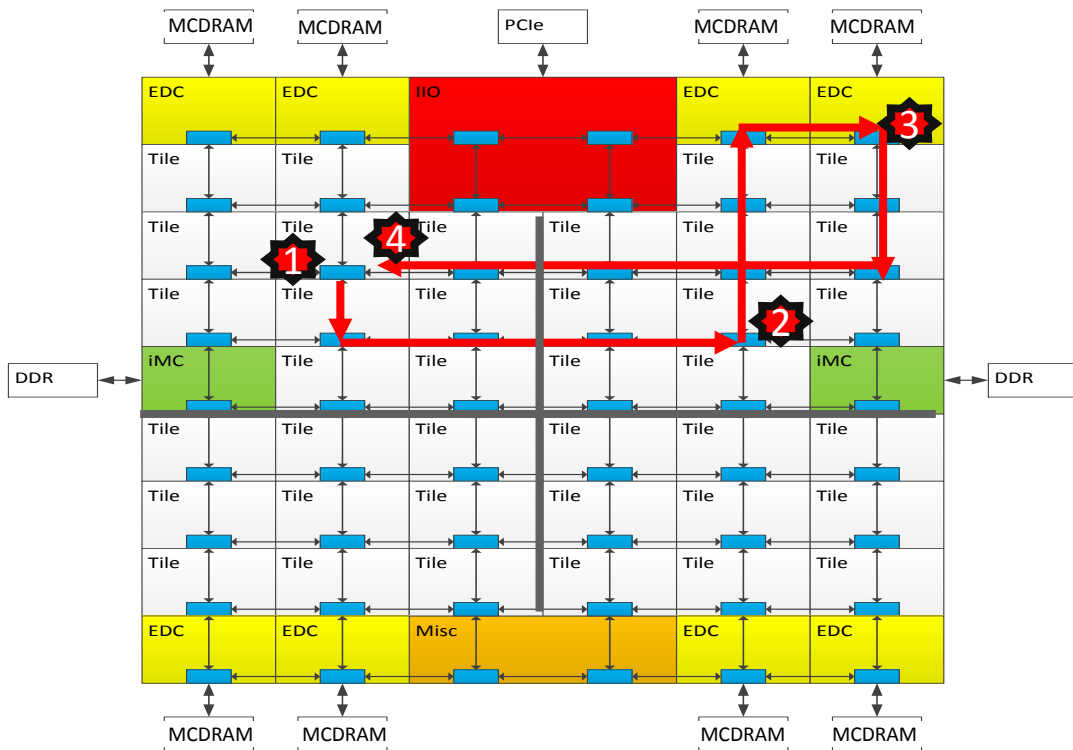**Address uniformly hashed across all distributed directories**

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: Quadrant



Chip divided into four virtual Quadrants

Address hashed to a Directory in the same quadrant as the Memory

Affinity between the Directory and Memory

Lower latency and higher BW than all-to-all.  SW Transparent.

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

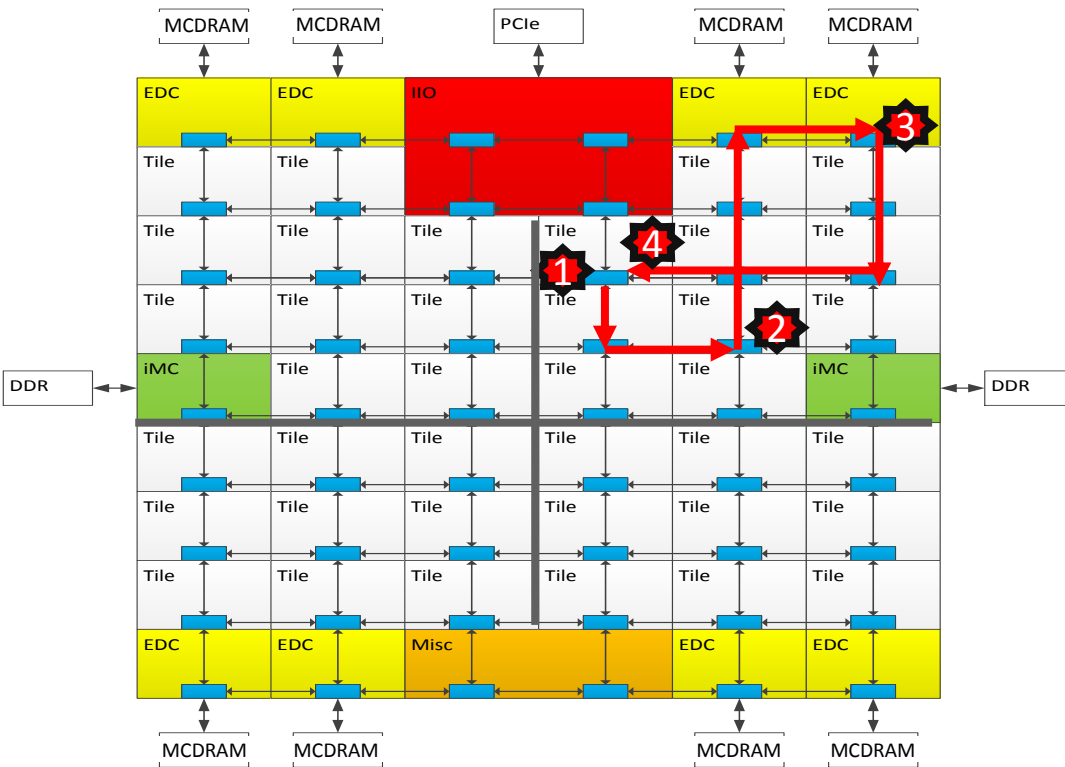# Cluster Mode: Sub-NUMA Clustering (SNC)



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

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# KNL w/ Intel® Omni-Path Fabric

Fabric integrated *on package*

First product with integrated fabric
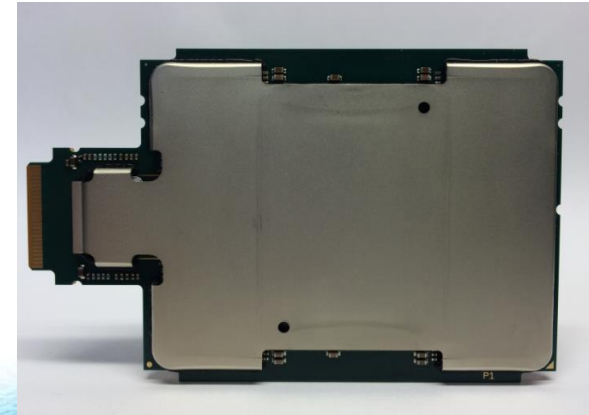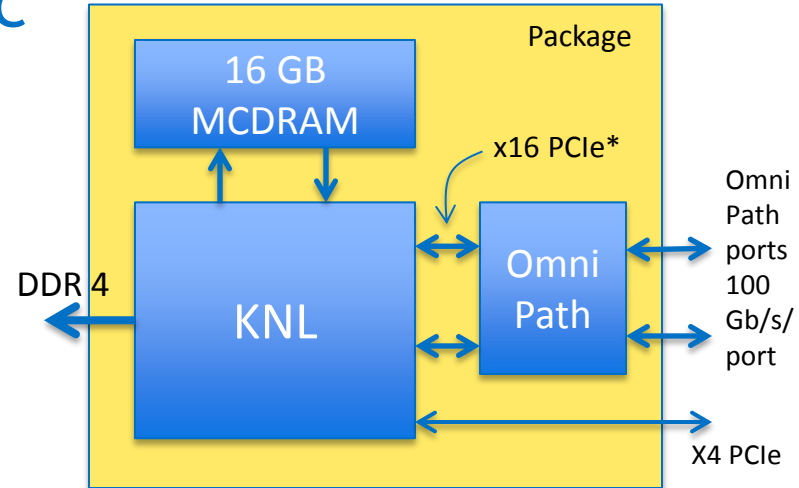
Connected to KNL die via 2 x16 PCIe* ports
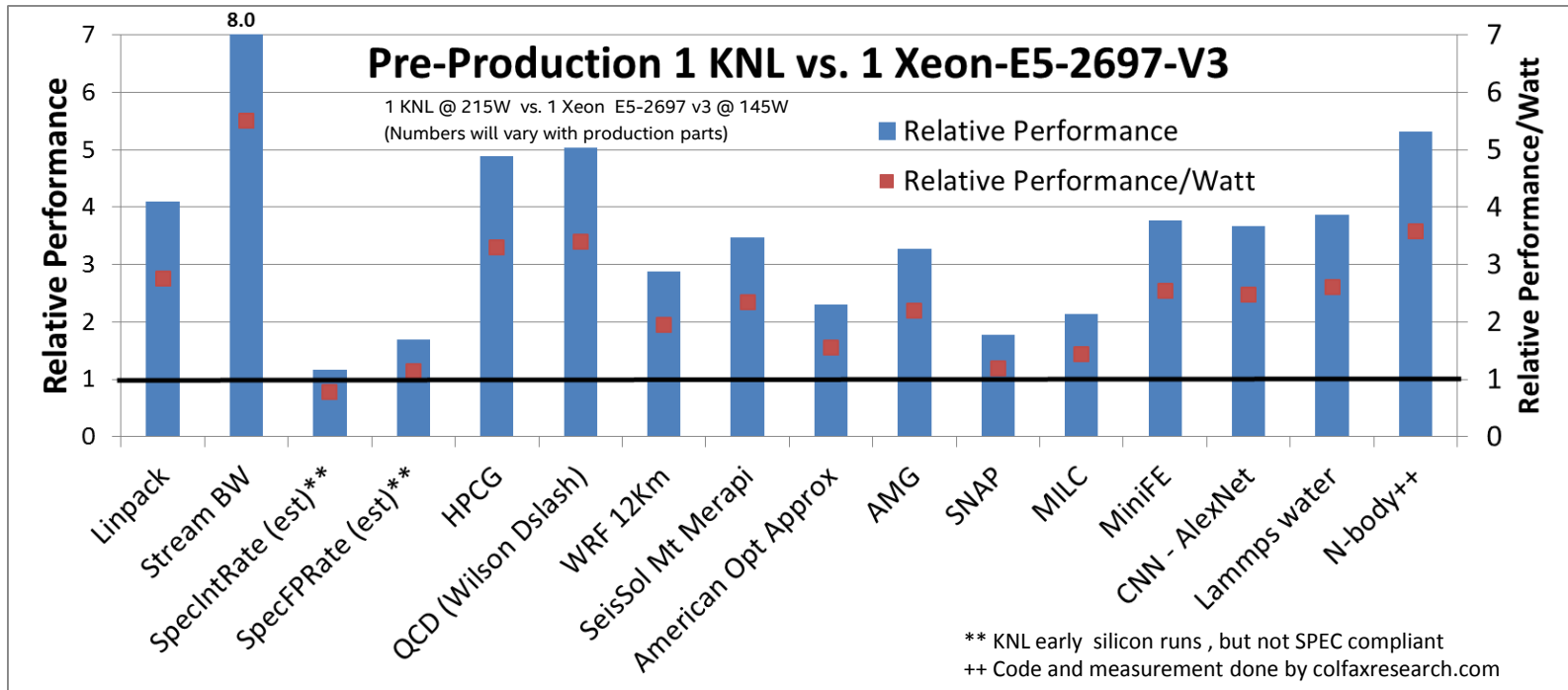Output: 2 Omni-Path ports
- 25 GB/s/port (bi-dir)

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

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

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# Pre-production KNL Performance and Performance/Watt



**Pre-Production 1 KNL vs. 1 Xeon-E5-2697-V3**

1 KNL @ 215W vs. 1 Xeon E5-2697 v3 @ 145W
(Numbers will vary with production parts)

■ Relative Performance
■ Relative Performance/Watt

** KNL early silicon runs , but not SPEC compliant
++ Code and measurement done by colfaxresearch.com

Categories: Linpack, Stream BW, SpecIntRate (est)**, SpecFPRate (est)**, HPCG, QCD (Wilson Dslash), WRF 12Km, SeisSol Mt Merapi, American Opt Approx, AMG, SNAP, MILC, MiniFE, CNN - AlexNet, Lammps water, N-body++

**Significant performance improvement for compute and bandwidth sensitive workloads, while still providing good general purpose out-of-box throughput performance.**

Avinash Sodani CGO PPoPP HPCA Keynote 2016

# MCDRAM Cache Hit Rate



MCDRAM performs well as cache for many workloads
Enables good out-of-box performance without memory tuning

# Deep Learning Training on KNL



Significant boost in deep learning training performance with KNL
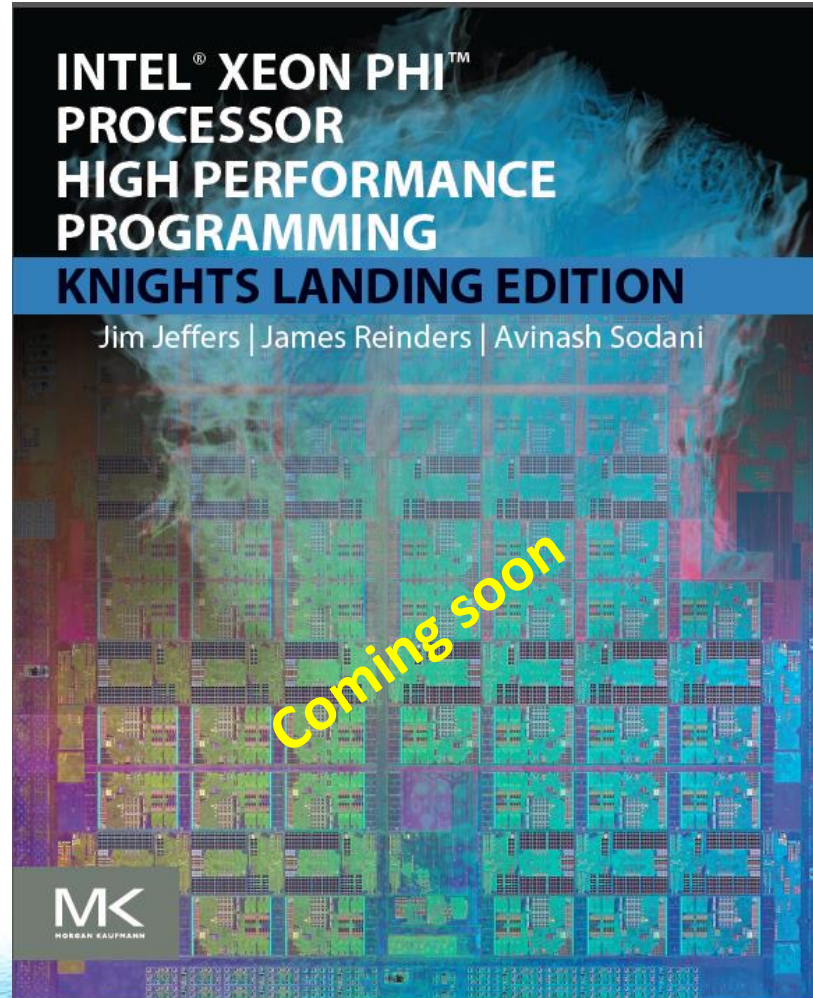Setting a trend for future increase with same programming model

# Programming for KNL

No different than programming a CPU

Same basics apply
- Exploit thread parallelism – Use all cores
  - Using parallel runtimes like MPI, OpenMP, TBB, etc.
  - Not always necessary to use all threads/core to get best performance
- Exploit the data parallelism – Vectorize!

- Utilize high bandwidth memory

Similar optimizations help both Intel®
Xeon® and Xeon Phi™ processors



INTEL® XEON PHI™
PROCESSOR
HIGH PERFORMANCE
PROGRAMMING
KNIGHTS LANDING EDITION

Jim Jeffers | James Reinders | Avinash Sodani

Coming soon

MK
MORGAN KAUFMANN

# Tools support evolving rapidly

**Auto-Vectorize**
- New instructions that help vectorize loops, e.g., Vconflict
- Aggressive vectorization and multi-versioning
- Masking and predication

**Language constructs to express parallelism**
- OpenMP pragmas
- Task level parallelism
- Higher level language constructs/libraries

**Compiler hints to guide Optimizations**
- Compiler pragmas as hints for vectorization
- Aliasing/alignment directives

**Feedback on code changes for parallelization**
- Meaningful and actionable compiler feedback about optimizations
- Profiling tools to better understand the program behavior
- Drive compiler optimization through runtime metrics

# Future Trends

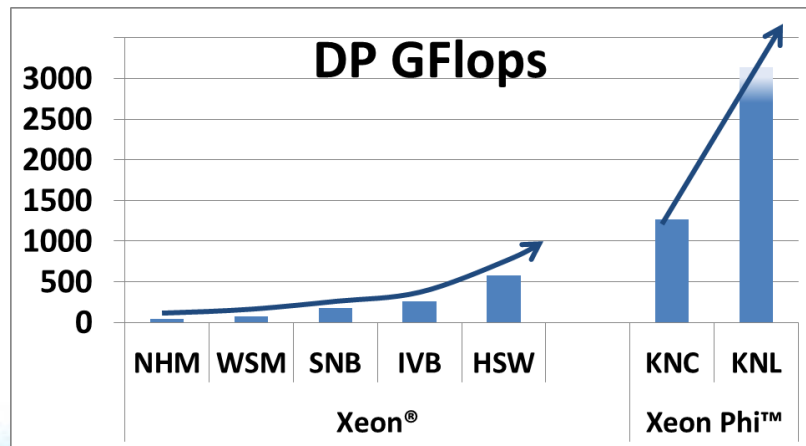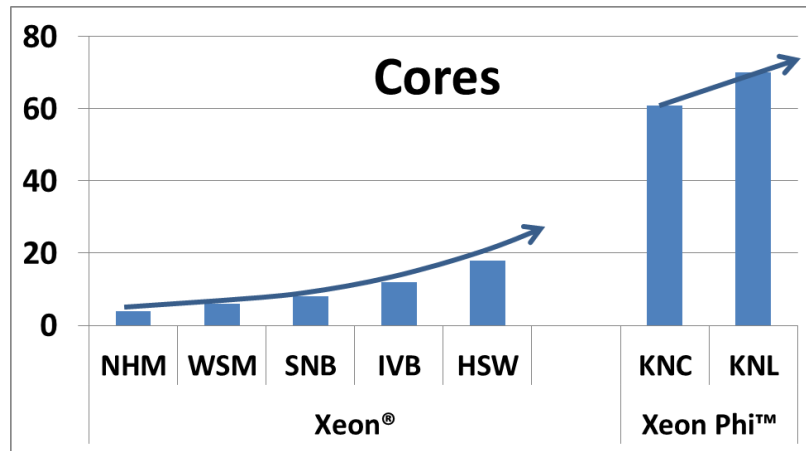Transistor density will increase

→ more cores and flops

→ more integration of system components

Power will continue to be a big challenge

- Intense focus on power efficient designs
- System power efficiency via integration
- More intelligent power management to better share power among components
- Usage-specific instructions and functionality for power efficiency

## More parallel solutions in future
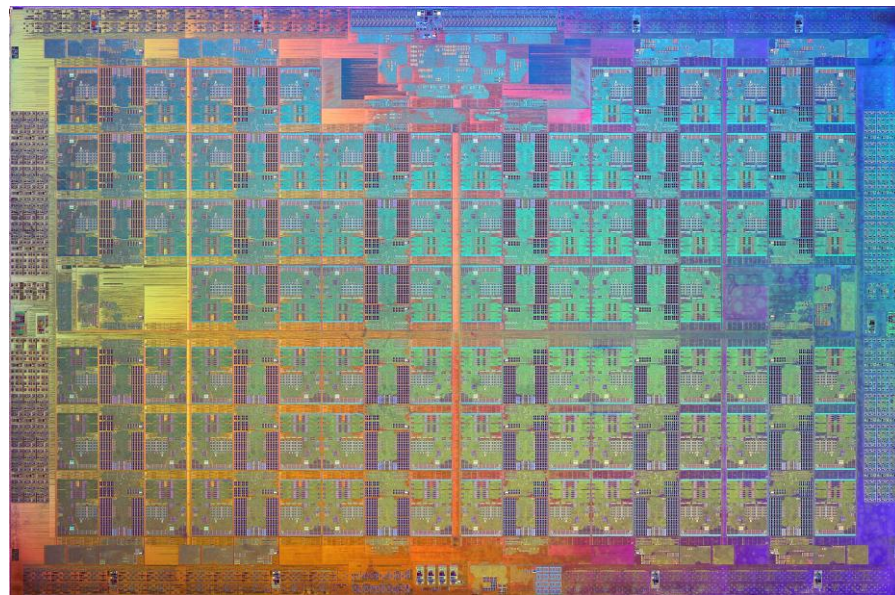
# Some Future SW Challenges

- Better **load balancing** between different threads
  - More task based parallelization, instead of bulk synchronous model

- **Data locality** conscious coding
  - Utilize caches well. Good for both performance and power

- **Reducing memory capacity** per thread
  - This can limit utilizing all cores in a CPU due to capacity constraints

- **Algorithms** that minimize global communications

- Continue to **improve tools** that provide relevant and actionable feedback to programmer on parallelization opportunities

# Summary

Knights Landing Xeon Phi™ processor
Massively parallel **CPU** with **general purpose programming**

- More parallel machines in future

- Parallelizing applications critical for performance

- Choice of "how" to parallelize is important → Software has a long life time



CPU + general purpose programming
provides a stable base for parallel software