# Primer on Computing with Intel Xeon Phi Coprocessors

## Part 1 of 2: Introduction and Programming Models

Vadim Karpusenko and Andrey Vladimirov
Colfax International

Intel HPC Developer Conference at SC14 in New Orleans, LA — November 16, 2014

# Contents

PARALLEL PROGRAMMING AND OPTIMIZATION WITH INTEL XEON PHI COPROCESSORS

HANDBOOK ON THE DEVELOPMENT AND OPTIMIZATION OF PARALLEL APPLICATIONS FOR INTEL XEON PROCESSORS AND INTEL XEON PHI COPROCESSORS

SECOND EDITION

ANDREY VLADIMIROV AND VADIM KARPUSENKO
COLFAX INTERNATIONAL

# §1. MIC Architecture from Developer's Perspective

# Intel Xeon Phi Coprocessors and the MIC Architecture



- PCIe end-point device
- High Power efficiency
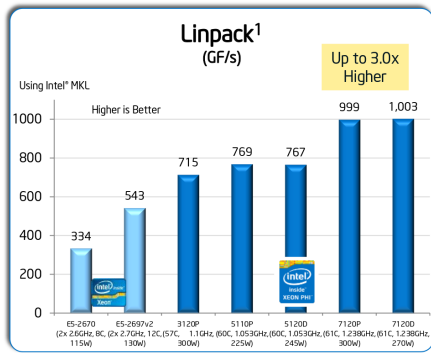- ~ 1 TFLOP/s in DP
- Heterogeneous clustering

For highly parallel applications which reach the scaling limits
on Intel Xeon processors

# Xeon and Xeon Phi Family Product Performance

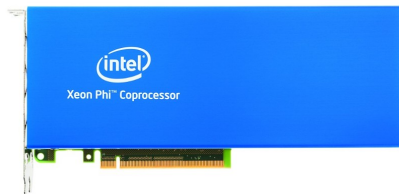Many-core Coprocessors (Xeon Phi) vs Multi-core Processors (Xeon) —

- Better performance per system & performance per watt for parallel applications
- Same programming methods, same development tools.



1. Xeon ran MP Linpack, Xeon Phi ran SMP Linpack. Expected performance difference between the two is estimated in the 3-5% range

Source: "Intel Xeon Product Family: Performance Brief"

# Intel Xeon Phi Coprocessors and the MIC Architecture



- C/C++/Fortran; OpenMP/MPI
- Standard Linux OS
- Up to 768 GB of DDR3 RAM
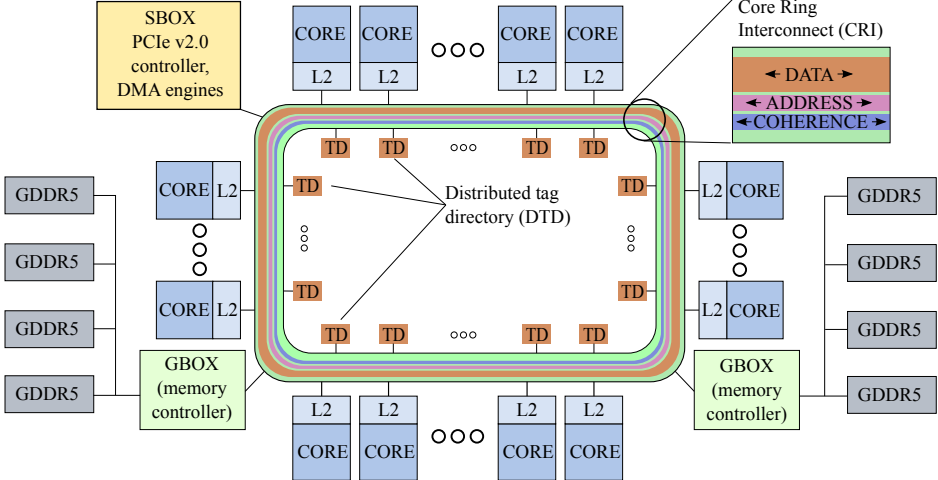- ≤12 cores/socket ≈3 GHz
- 2-way hyper-threading
- 256-bit AVX vectors

- C/C++/Fortran; OpenMP/MPI
- Special Linux $\mu$OS distribution
- 6–16 GB cached GDDR5 RAM
- 57 to 61 cores at ≈1 GHz
- 4 hardware threads per core
- 512-bit IMCI vectors

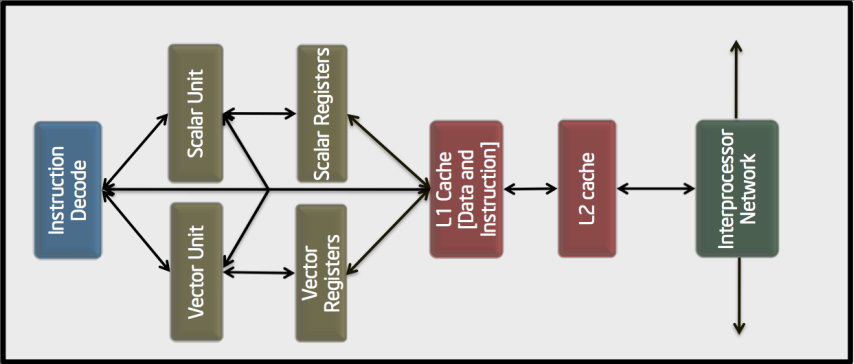# Linux $\mu$OS on Intel Xeon Phi coprocessors (part of MPSS)

```
user@host% lspci | grep -i "co-processor"
06:00.0 Co-processor: Intel Corporation Xeon Phi coprocessor 3120 series (rev 20)
82:00.0 Co-processor: Intel Corporation Xeon Phi coprocessor 3120 series (rev 20)
user@host% sudo service mpss status
mpss is running
user@host% cat /etc/hosts | grep mic
172.31.1.1  host-mic0 mic0
172.31.2.1  host-mic1 mic1
user@host% ssh mic0
user@mic0% cat /proc/cpuinfo | grep proc | tail -n 3
processor : 237
processor : 238
processor : 239
user@mic0% ls /
amplxe   dev     home     lib64    oldroot   proc     sbin     sys     usr
bin      etc     lib      linuxrc  opt       root     sep3.10  tmp     var
```

# Details of the MIC Architecture

# Die Organization Diagram

# Core Topology



| 4 Threads per Core | | Fully Coherent | | |
|---|---|---|---|---|
| 64-bit | 512-wide | L2 Hardware Prefetching | | Ring interconnect |
| In Order | VPU: integer, SP, DP; 3-operand, 16-instruction | 32 KB per core | 512 KB Slice per Core – Fast Access to Local Copy | |
| Specialized Instructions (new encoding) | | | | |
| | HW transcendentals | | | |

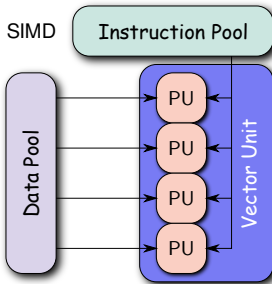# SIMD Operations

## SIMD — Single Instruction Multiple Data

Scalar Loop

```
1  for (i = 0; i < n; i++)
2    A[i] = A[i] + B[i];
```

SIMD Loop

```
1  for (i = 0; i < n; i += 16)
2    A[i:(i+16)] = A[i:(i+16)] + B[i:(i+16)];
```
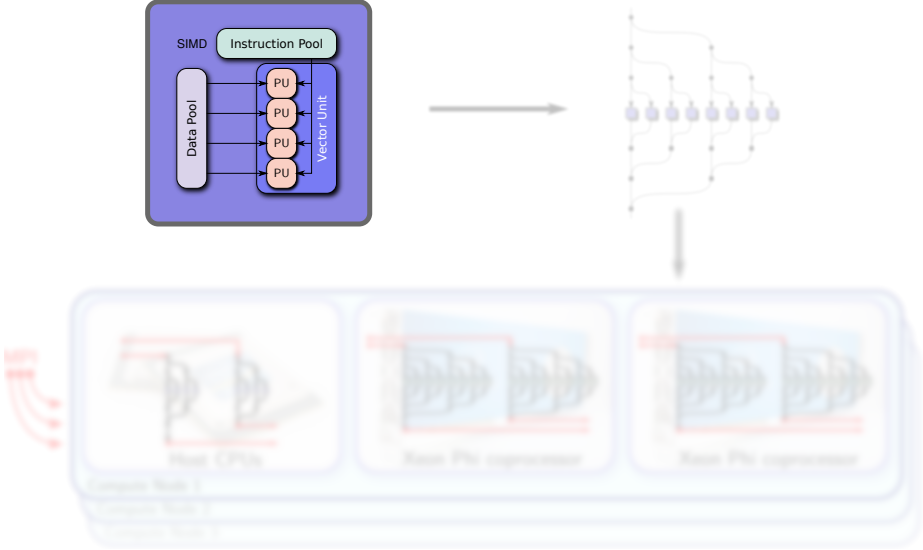
Each SIMD addition operator acts on 16 numbers at a time.
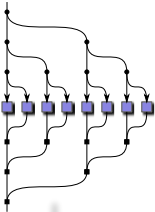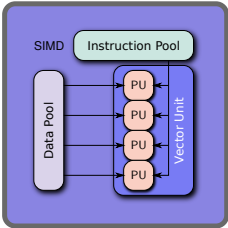
# Instruction Sets in Intel Architectures

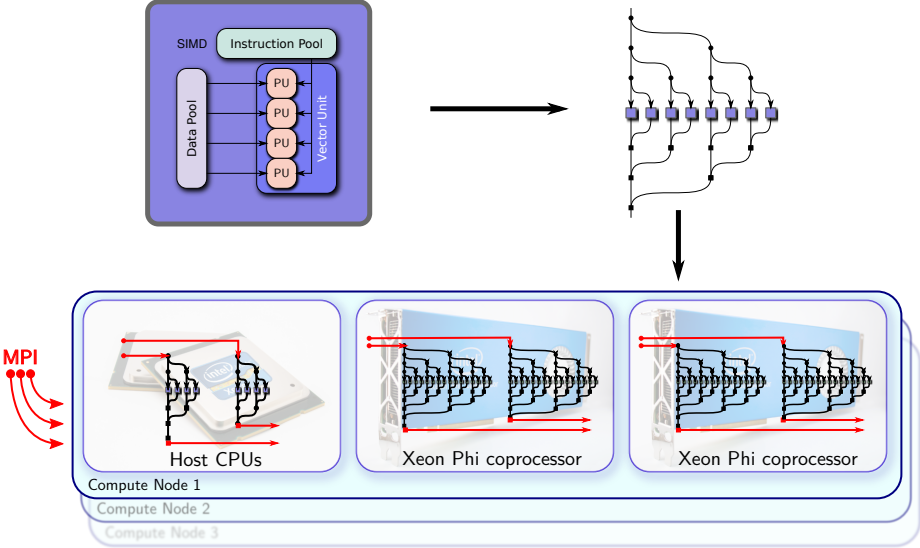| Instruction Set | Year and Intel Processor | Vector registers | Packed Data Types |
|---|---|---|---|
| MMX | 1997, Pentium | 64-bit | 8-, 16- and 32-bit integers |
| SSE | 1999, Pentium III | 128-bit | 32-bit single precision FP |
| SSE2 | 2001, Pentium 4 | 128-bit | 8 to 64-bit integers; SP & DP FP |
| SSE3–SSE4.2 | 2004 – 2009 | 128-bit | (additional instructions) |
| AVX | 2011, Sandy Bridge | 256-bit | single and double precision FP |
| AVX2 | 2013, Haswell | 256-bit | integers, additional instructions |
| IMCI | 2012, Knights Corner Intel Xeon Phi coproc. | 512-bit | 32- and 64-bit integers; single & double precision FP |
| AVX-512 | (future) Knights Landing | 512-bit | 32- and 64-bit integers; single & double precision FP |

# Three Layers of Parallelism

# Three Layers of Parallelism

# Three Layers of Parallelism

# Examples of Solutions with the Intel MIC Architecture
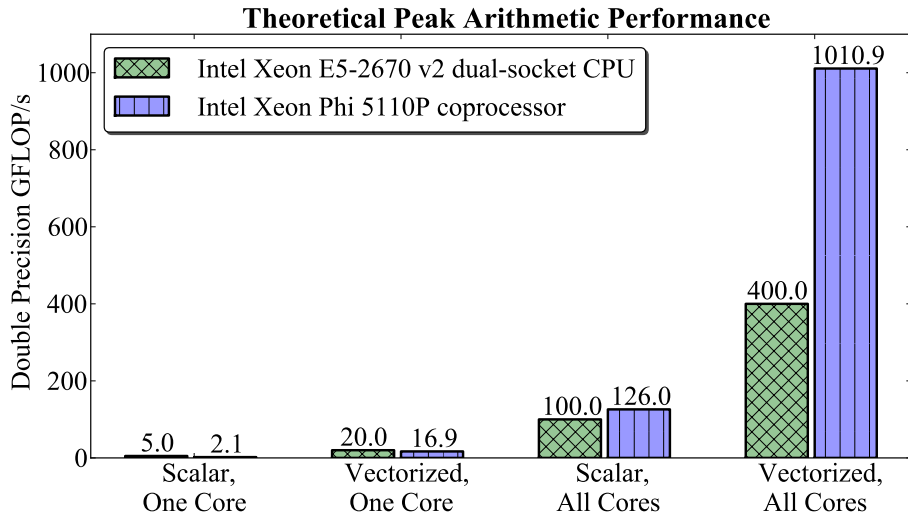


Colfax's CXP7450 workstation with two Intel Xeon Phi coprocessors
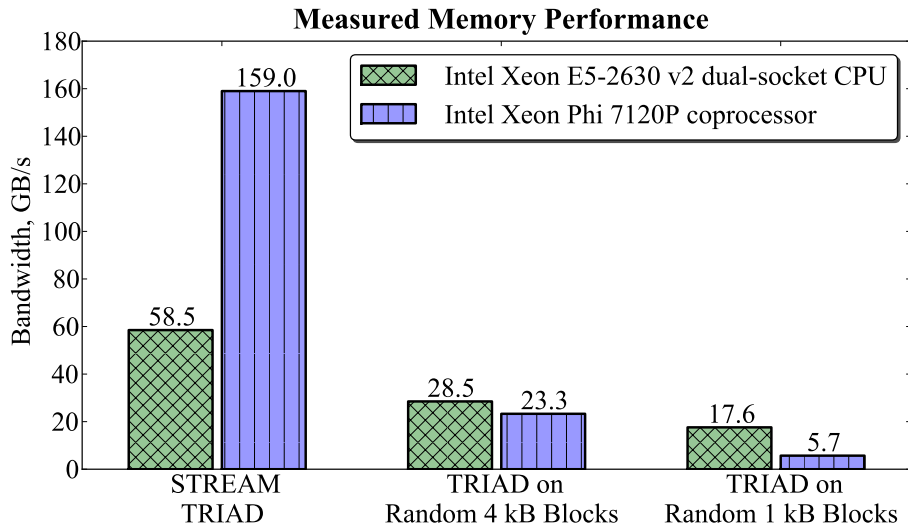
Colfax's CXP9000 server with eight Intel Xeon Phi coprocessors

`xeonphi.com`

# Data and Task Parallelism (Vectors and Cores)



Theoretical Peak Arithmetic Performance

# Memory Access Pattern



Measured Memory Performance

# §2. Programming Models and Application Porting

# Offload and Native modes

- Explicit offload mode:

| Host | Coprocessor |
|------|-------------|
| ```
main() {

#pragma offload target(mic)


}
``` | ```


    myFunction();

``` |

- Native mode:

| Host | Coprocessor |
|------|-------------|
| | ```
main() {

    myFunction();

}
``` |

# Native Programming

# Native Execution

### "Hello World" application:

```c
#include <stdio.h>
#include <unistd.h>
int main(){
    printf("Hello world! I have %ld logical cores.\n",
    sysconf(_SC_NPROCESSORS_ONLN ));
}
```

### Compile and run on host:

```
user@host% icc hello.c
user@host% ./a.out
Hello world! I have 32 logical cores.
user@host%
```

# Native Execution

Compile and run the same code on the coprocessor in the native mode:

```
user@host% icc hello.c -mmic
user@host% scp a.out mic0:~/
a.out 100% 10KB 10.4KB/s 00:00
user@host% ssh mic0
user@mic0% pwd
/home/user
user@mic0% ls
a.out
user@mic0% ./a.out
Hello world! I have 240 logical cores.
user@mic0%
```

- Use -mmic to produce executable for MIC architecture
- Must transfer executable to coprocessor (or NFS-share) and run from shell
- Native MPI applications work the same way (need Intel MPI library)

# Porting User Applications for Native Execution

Simple CPU applications can be compiled for native execution on Xeon Phi coprocessors by supplying the flag "-mmic" to the Intel compiler:

```
user@host% icpc -c myobject1.cc -mmic
user@host% icpc -c myobject2.cc -mmic
user@host% icpc -o myapplication myobject1.o myobject2.o -mmic
```

Same for coprocessor-only MPI applications:

```
user@host% mpiicpc -c myobject1.cc -mmic
user@host% mpiicpc -c myobject2.cc -mmic
user@host% mpiicpc -o myapplication myobject1.o myobject2.o -mmic
```

# Native Applications with Autotools

- Use the Intel compiler with flag `-mmic`
- Eliminate assembly and unnecessary dependencies
- Use `--host=x86_64` to avoid "program does not run" errors

Example, the GNU Multiple Precision Arithmetic Library (GMP):

```
user@host% wget https://ftp.gnu.org/gnu/gmp/gmp-5.1.3.tar.bz2
user@host% tar -xf gmp-5.1.3.tar.bz2
user@host% cd gmp-5.1.3
user@host% ./configure CC=icc CFLAGS="-mmic"  --disable-assembly --host=x86_64
...
user@host% make
...
```

# Running Native MPI Applications on Coprocessors

```
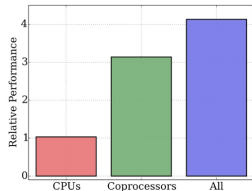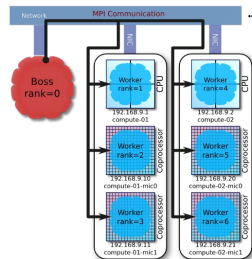user@host% export I_MPI_MIC=1
user@host% mpiicpc -mmic -o HelloMPI.MIC HelloMPI.c
user@host% scp HelloMPI.MIC mic0:~/
user@host% mpirun -host mic0 -np 2 ~/HelloMPI.MIC
Hello World from rank 1 running on host-mic0!
Hello World from rank 0 running on host-mic0!
MPI World size = 2 processes
```

- Enable the MIC architecture in Intel MPI: `I_MPI_MIC=1`
- Copy or NFS-share MPI library & executables with coprocessor
- Use `mpiicpc` with `-mmic` to compile
- Launch as if `mic0` is a remote host

# Native Heterogeneous Clustering in Action

**Heterogeneous Clustering
with Homogeneous Code:
Asian Option Pricing**

- Monte Carlo method

- MPI + OpenMP + automatic vectorization

- The same C code for clusters of
  a) CPUs
  b) Coprocessors
  c) CPUs+Coprocessors (heterogeneous)

- **No code modification** to run on
  the Intel MIC architecture

- No platform-specific tunning

- Bridged network configuration





More information in paper at `xeonphi.com/papers/heterogeneous`

# Offload Programming

# Explicit Offload Programming Model

"Hello World" in the explicit offload model:

```c
#include <stdio.h>

int main(int argc, char * argv[] ) {

    printf("Hello World from host!\n");

#pragma offload target(mic)
    {
        printf("Hello World from coprocessor!\n"); fflush(0);
    }

    printf("Bye\n");
}
```

Application launches and runs on the host, but some parts of code and data are moved ("offloaded") the coprocessor.

# Compiling and Running an Offload Application

```
user@host% icpc hello_offload.cpp -o hello_offload
user@host% ./hello_offload
Hello World from host!
Bye
Hello World from coprocessor!
```

- Code inside of #pragma offload is offloaded automatically
- Console output on Intel Xeon Phi coprocessor is buffered and mirrored to the host console
- If coprocessor is not installed, code inside #pragma offload may fall back to the host system

Computing with Xeon Phi

# Offloading Functions and Data

```
1  int* __attribute__((target(mic))) data;
2
3  __attribute__((target(mic))) void MyFunction(int* foo) {
4      // ... implement function as usual
5  }
6
7  int main(int argc, char * argv[] ) {
8      // ...
9  #pragma offload target(mic) inout(data : length(N))
10     {
11         MyFunction(data);
12     }
13 }
```

- Functions and data used on coprocessor must be marked with the specifier __attribute__((target(mic)))

# Data Marshalling for Dynamically Allocated Data

```
1  double *p1=(double*)malloc(sizeof(double)*N);
2  // ...
3
4  #pragma offload target(mic)                        \
5          in   (p1 : length(N) alloc_if(1) free_if(0)) \
6          out  (p2 : length(N) align(4096))          \
7          inout(p3[0:N/2]  : into(p4[N/2:N/2]))
8    {
9      // ...
10   }
```

- #pragma offload recognizes clauses in, out, inout and nocopy
- Data size (length), alignment, redirection, retention, and other properties may be specified
- Marshalling is required for pointer-based data

# Multiple Coprocessors and Fallback to Host

```
1  const int nDevices = _Offload_number_of_devices();
2  #pragma omp parallel num_threads(nDevices) if(nDevices>0)
3    {
4      const int iDevice = omp_get_thread_num();
5  #pragma offload target(mic: iDevice) if(nDevices>0)
6        {
7          MyFunction(/*...*/);
8        }
9    }
```

- Up to 8 coprocessors, up to 32 host threads
- All offloads start simultaneously and block the respective thread

# Asynchronous Offload

- By default, #pragma offload blocks until offload completes
- Use clause "signal" with any pointer to avoid blocking
- Use #pragma offload_wait to block where needed

```
char* offload0;
#pragma offload target(mic:0) signal(offload0) in(data : length(N))
 { /* ... will not block code execution because of clause "signal" */ }

DoSomethingElse();

/* Now block until offload signalled by pointer "offload0" completes */
#pragma offload_wait target(mic:0) wait(offload0)
```

- Used to overlap communication with computation

# Alternative: Virtual-shared Memory Model

```
1  _Cilk_shared int arr[N]; // This is a virtual-shared array
2
3  _Cilk_shared void Compute() { // This function may be offloaded
4      // ... function uses array arr[]
5  }
6
7  int main() {
8      // arr[] can be initialized on the host
9      _Cilk_offload Compute(); // and used on coprocessor
10     // and the values are returned to the host
11 }
```

- Alternative to Explicit Offload
- Data synced from host to coprocessor before the start of offload
- Data synced from coprocessor to host at the end of offload

# §3. Heterogeneous Computing with the MIC Architecture

# Inter-Operation of Offload, MPI and OpenMP

# Teaming Xeon and Xeon Phi Coprocessors

Programming models allow a range of CPU+MIC coupling modes

Xeon - Multi-Core Centric        *Breadth*        MIC - Many-Core Centric

| Multi-Core Hosted | Offload | Symmetric | Many Core Hosted |
|---|---|---|---|
| General serial and parallel computing | Code with highly-parallel phases | Codes with balanced needs | Highly-parallel codes |

# Heterogeneous Distributed Computing with Xeon Phi

**Option 1: MPI+OpenMP with Offload.**

- MPI processes are multi-threaded with OpenMP.

- MPI runs only on CPUs.

- MPI processes offload to coprocessor(s).

- OpenMP in offload regions.

# Heterogeneous Distributed Computing with Xeon Phi

**Option 2: Symmetric pure MPI (native mode).**

- MPI processes on hosts
- Native MPI processes on the coprocessor.
- No OpenMP.
- xeonphi.com/papers/p2p

# Heterogeneous Distributed Computing with Xeon Phi

**Option 3: Symmetric hybrid MPI+OpenMP.**

- MPI processes on hosts
- Native MPI processes on the coprocessor.
- Multi-threading with OpenMP.
- xeonphi.com/papers/p2p

# File I/O in MPI Applications on Coprocessors

# RAM Filesystem

- Files are stored in the coprocessor RAM
- Does not survive MPSS restart or host reboot
- Fastest method
- Good for local pre-staged input or runtime scratch data

# Virtio Transfer to Local Host Drives

- Files are stored on a physical or virtual drive on the host
- Written data persistent across reboots
- Fast method
- Cannot share a drive between coprocessors
- Good for distributed checkpointing

# Network Storage

- Files are stored on a remote file server
- Can share a mount point across the cluster
- Lustre has scalable performance
- NFS is slow but easy to set up
- More information: xeonphi.com/papers/io

# §4. Future-Proofing and Knights Landing

# Future-Proofing: Reliance on Compiler and Libraries

Ease of use

*Threading Options*                    *Vector Options*

| Threading Options | Vector Options |
|---|---|
| Intel® Math Kernel Library API* | Intel® Math Kernel Library |
| Intel® Threading Building Blocks / Intel® Cilk™ Plus | Array Notation: Intel® Cilk™ Plus |
| OpenMP* | Auto vectorization |
| | Semi-auto vectorization: #pragma (vector, simd) |
| Pthreads* | OpenCL* |
| | C/C++ Vector Classes (F32vec16, F64vec8) |

Depth

Fine control

# Next Generation MIC: Knights Landing (KNL)

- 2nd generation MIC product: code name Knights Landing (KNL)
- Intel's 14 nm manufacturing process
- A processor (running the OS) or a coprocessor (PCIe device)
- On-package high-bandwidth memory w/flexible memory models: flat, cache, & hybrid
- Intel Advanced Vector Extensions AVX-512 (public)



*Source: Intel Newsroom*

# Getting Ready for the Future

- Porting HPC applications to today's MIC architecture makes them ready for future architectures, such as KNL

- Xeon, KNC and KNL are not binary compatible, therefore assembly-level tuning will not scale forward.

- Reliance on compiler optimization and using optimized libraries (such as Intel MKL) ensures future-readiness.



*Today*

*22nm process PCIe coprocessor*

*Tomorrow*

*14nm Standalone CPU*

*Source: Intel Newsroom*

Source: https://www.brighttalk.com/webcast/10773/116329

Source: https://www.brighttalk.com/webcast/10773/116329

# Colfax Developer Training

Colfax runs one-day and four-day trainings for organizations and individuals on parallel programming with Intel Xeon Phi coprocessors.



Information: xeonphi.com/training

# Additional Reading

It all comes down to
PARALLEL
PROGRAMMING !
(applicable to processors
and Intel® Xeon Phi™
coprocessors both)

Forward, Preface
Chapters:
1. Introduction
2. High Performance Closed
   Track
   Test Drive!
3. A Friendly Country Road Race
4. Driving Around Town:
   Optimizing A Real-World
   Code Example
5. Lots of Data (Vectors)
6. Lots of Tasks (not Threads)
7. Offload
8. Coprocessor Architecture
9. Coprocessor System Software
10. Linux on the Coprocessor
11. Math Library
12. MPI
13. Profiling and Timing
14. Summary
Glossary, Index

Learn more about this book:
## lotsofcores.com



**Intel® Xeon Phi™
Coprocessor
High Performance
Programming**

Jim Jeffers, James Reinders

MK

*This book belongs on the
bookshelf of every HPC
professional. Not only does it
successfully and accessibly
teach us how to use and
obtain high performance on
the Intel MIC architecture, it is
about much more than that. It
takes us back to the universal
fundamentals of high-
performance computing
including how to think and
reason about the performance
of algorithms mapped to
modern architectures, and it
puts into your hands powerful
tools that will be useful for
years to come.*
—Robert J. Harrison
Institute for Advanced
Computational Science,
Stony Brook University

Available since February 2013.

Intel® Xeon Phi™ Coprocessor High Performance Programming,
Jim Jeffers, James Reinders, (c) 2013, publisher: Morgan Kaufmann

*©2013, James Reinders & Jim Jeffers, book image used with permission

(intel)

# Additional Reading

### Available November 17, 2014

- 28 chapters
- 69 expert contributors
- Numerous "Real World" Code "Recipes" and examples using OpenMP, MPI, TBB, C, C++, OpenCL, Fortran.
- Successful techniques, tips for vectorization, scalable parallel coding, load balancing, data structure and memory tuning, applicable to both processors and coprocessors!
- All figures, diagrams and code freely downloadable. (Nov'14)

# Next Session Starting in 15 Minutes

- Example: N-body simulation
- Optimization essentials
- Node-level tuning
- Scalability in a cluster





**N-Body Simulation Performance**

Processor: Intel Xeon E5-2697 v2
Coprocessor: Intel Xeon Phi 7120P

# Thank you for tuning in,
# and
# have a wonderful journey
# to the Parallel World!

http://research.colfaxinternational.com/

P.S.: We are hiring! xeonphi.com/jobs

# §5. Resources/Backup Slides

# Reference Guides

- Intel C++ Compiler 14.0 User and Reference Guide
- Intel VTune Amplifier XE User's Guide
- Intel Trace Analyzer and Collector Reference Gude
- Intel MPI Library for Linux* OS Reference Manual
- Intel Math Kernel Library Reference Manual
- Intel Software Documentation Library
- MPI Routines on the ANL Web Site
- OpenMP Specifications

# Intel's Top 10 List

1. Download programming books: "Intel Xeon Phi Coprocessor High Performance Programming" by Jeffers & Reinders, and "Parallel Programming and Optimization with Intel Xeon Phi Coprocessors" by Colfax.

2. Watch the parallel programming webinar

3. Bookmark and browse the mic-developer website

4. Bookmark and browse the two developer support forums: "Intel MIC Architecture" and "Threading on Intel Parallel Architectures".

5. Consult the "Quick Start" guide to prepare your system for first use, learn about tools, and get C/C++ and Fortran-based programs up and running

Intel's Top 10 resources list: `http://xeonphi.com/top10`

# Intel's Top 10 List (continued)

6. Try your hand at the beginning lab exercises
7. Try your hand at the beginner/intermediate real world app exercises
8. Browse the case studies webpage to view examples from many segments
9. Begin optimizing your application(s); consult your programming books, the ISA reference manual, and the support forums for assistance.
10. Hone your skills by watching more advanced video workshops

Intel's Top 10 resources list: `http://xeonphi.com/top10`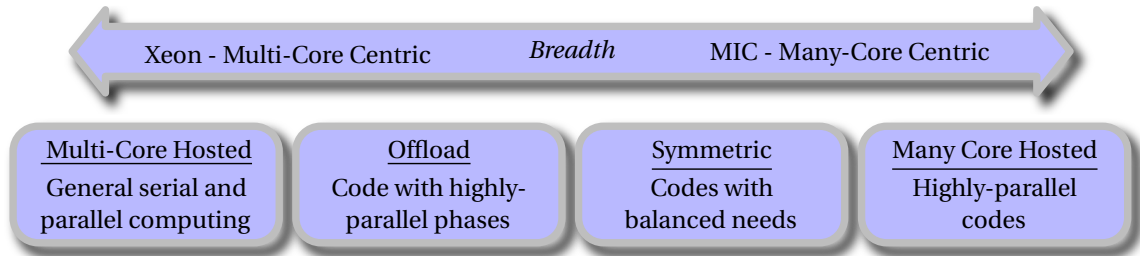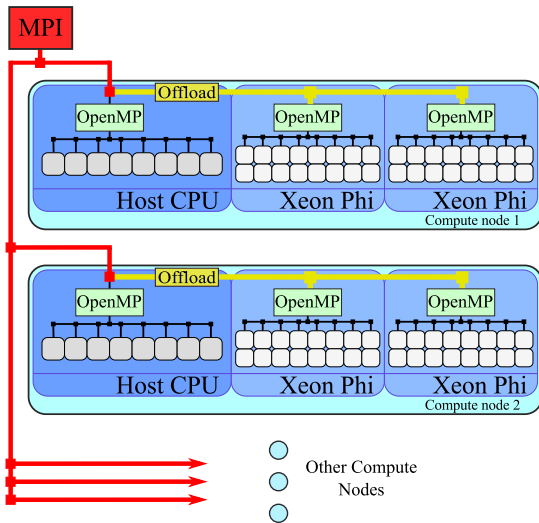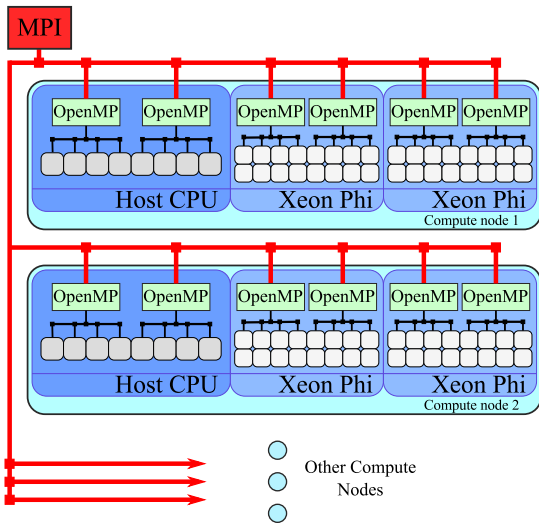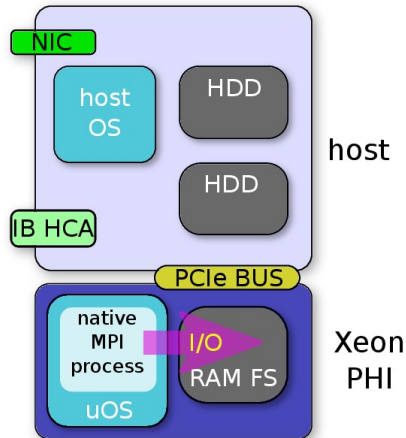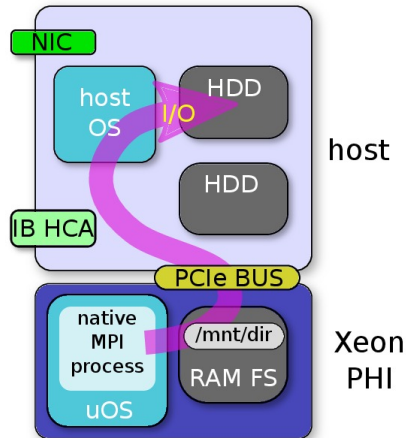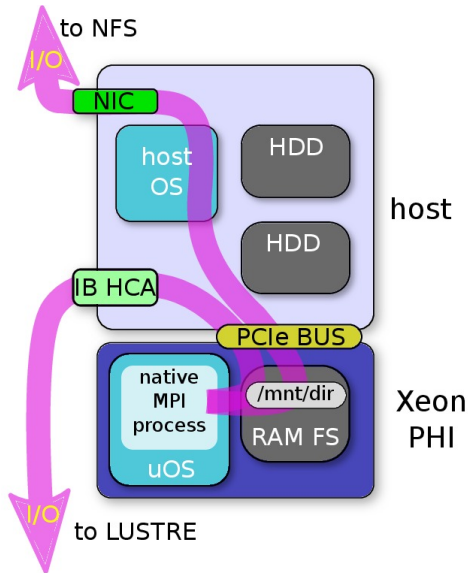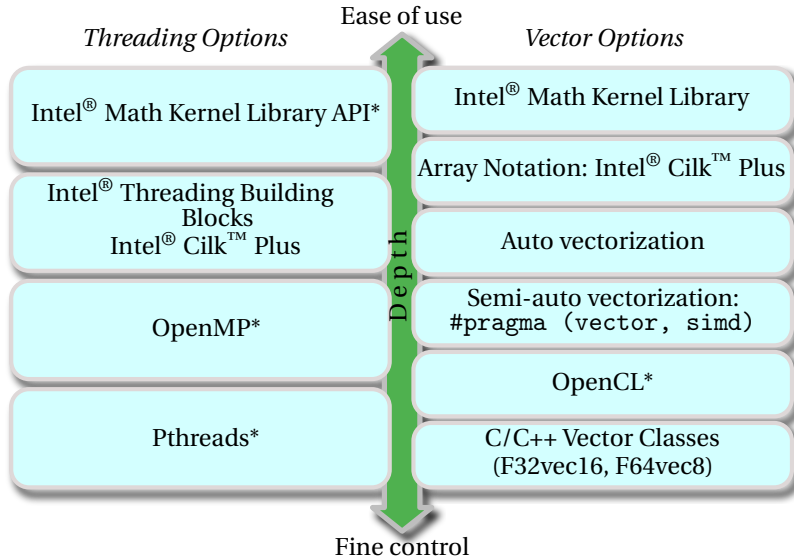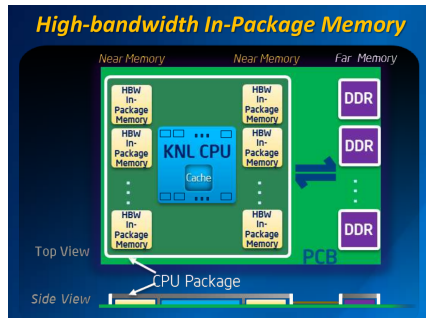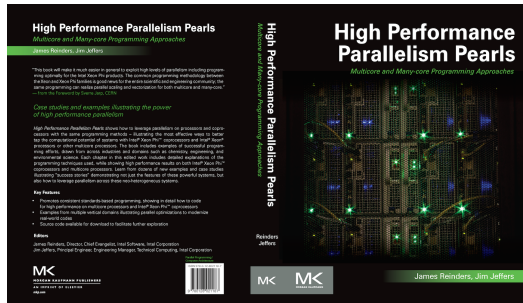