# Optimizing for MPI*/OpenMP* on Intel® Xeon Phi™ Coprocessors

Jim Jeffers,  John Pennycook,  Hans Pabst, Heinrich Bockhorst, Intel Corporation

Vince Betro, Paul Peltz, National Institute for Computational Sciences, UTenn

# Agenda

## Aims

- Focus on techniques, not syntax.

- Maximise hands-on programming experience.

- Demonstrate performance gains for real-life applications.

## FAQ Session

- E-mail your questions to: john.pennycook@intel.com

- Most frequent/interesting answered at the end of the day
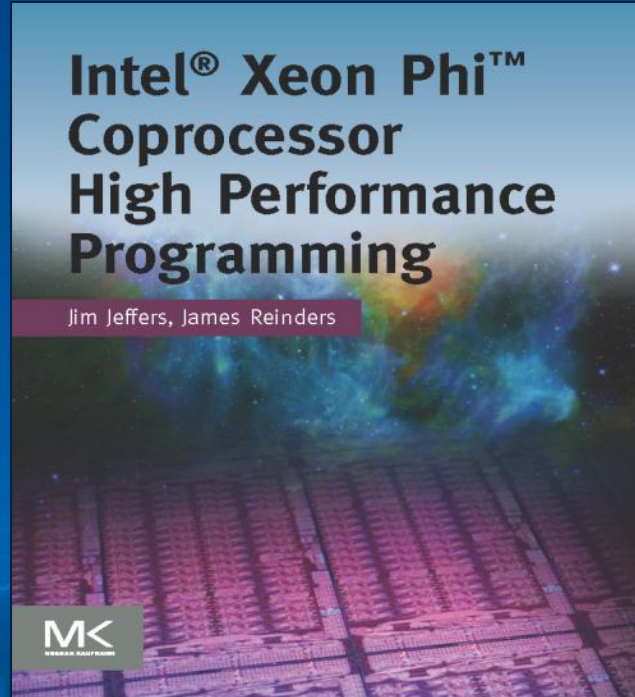
# Agenda

| | | |
|---|---|---|
| 08:30 | Introduction to Intel® Xeon Phi™ Coprocessors | John Pennycook, Jim Jeffers |
| 09:30 | Cluster Administration | Paul Peltz |
| 09:45 | **Hands-on:** Logging in and Running Applications | |
| 10:00 | Break | |
| 10:30 | MPI* and OpenMP* | Vince Betro |
| 11:00 | **Hands-on:** MPI*, OpenMP* and Intel® Trace Analyzer and Collector (ITAC) | |
| 12:00 | Lunch | |
| 13:00 | Compiler-Assisted Offload | Hans Pabst |
| 13:30 | **Hands-on:** Device Extensions in OpenMP* 4.0 | |
| 14:30 | Break | |
| 14:45 | Vectorization | John Pennycook |
| 15:15 | **Hands-on:** Explicit Vectorization with OpenMP* 4.0 | |
| 16:15 | Frequently Asked Questions | All presenters |
| 16:30 | End | |

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

Learn more about this book:
## lotsofcores.com

Intel® Xeon Phi™
Coprocessor
High Performance
Programming

Jim Jeffers, James Reinders

MK
MORGAN KAUFMANN

Available since mid-February 2013.

*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

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

intel®

# Introduction to Intel® Xeon Phi™ Coprocessors

**Parallelism, Parallelism and More Parallelism**

(intel)

# Parallelism in Modern Computer Architecture

## Instruction-Level Parallelism (ILP)

- Pipelining
- Multiple instruction issue
- Out-of-Order execution

## Vectorization (SIMD)

- Single instructions can be applied to more than one piece of data

## Threading (MIMD)

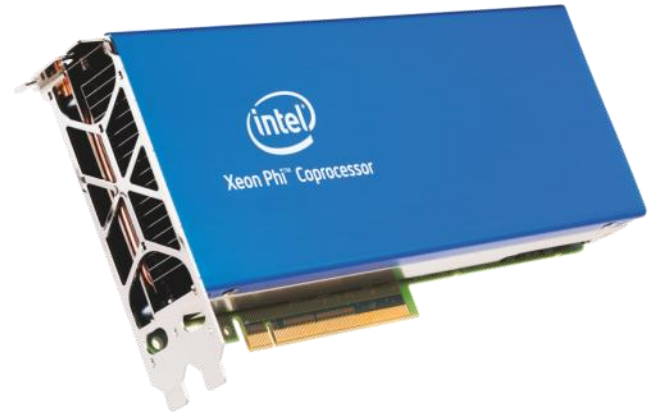- Multiple instances of a program can run simultaneously

# Parallelism in Modern Computer Architecture

Instruction-Level Parallelism (ILP)

- Pipelining
- Multiple instruction issue
- Out-of-Order execution

Vectorization (SIMD)

- Single instructions can be applied to more than one piece of data

Threading (MIMD)

- Multiple instances of a program can run simultaneously

# Heterogeneity in Architecture Designs

## Multi-core

- High single-thread performance.
- Low number of threads.
- Suitable for any workload.

## Many-core

- Low single-thread performance.
- High number of threads.
- Suitable for highly parallel workloads.



Knights Corner – the first Intel® Xeon Phi™ Coprocessor

# Heterogeneity in Architecture Designs

## Multi-core

- High single-thread performance.

- Low number of threads.

- Suitable for any workload.

## Many-core

- Low single-thread performance.

- High number of threads.

- Suitable for highly parallel workloads.

|  | Multi-core | Many-core |
|---|---|---|
| Clock Rate | ≈3 GHz | ≈1 GHz |
| # Cores | 4 – 12 | ≈60 |
| Threads per Core | 1 – 2 | 4 |
| Bandwidth† | ≈60 GB/s | ≈350 GB/s |
| Thermal Design Power | ≈115 W | ≈225 W |

† Peak electrical bandwidth.

# Multi-core vs Many-core – Theoretical Peak Performance

## Peak GFLOP/s in Single Precision

- Clock Rate x Cores x Ops/Cycle x SIMD

## 2 x Intel® Xeon® Processor E5-2670v2

- 2.5 GHz x 20 cores x 2 ops x 8 SIMD
  = 800 GFLOP/s

## Intel® Xeon Phi™ Coprocessor 7120P

- 1.24 GHz x 61 cores x 2 ops x 16 SIMD
  = 2420.48 GFLOP/s



GFLOP/s

|                | Scalar & ST | Vector & ST | Scalar & MT | Vector & MT |
|----------------|-------------|-------------|-------------|-------------|
| 2 x Processor  | 5           | 40          | 100         | 800         |
| Coprocessor    | 1.24        | 19.84       | 151.28      | 2420.48     |

■ 2 x Processor  ■ Coprocessor

Note the logarithmic scale on the y-axis.
ST = Single Thread, MT = Multiple Threads

# Architecture of an Intel® Xeon Phi™ Coprocessor Core

## Two Pipelines

- Scalar unit based on Pentium® processors.
- Dual issue (vector + scalar)

## 512-bit Vector Processing Unit (VPU)

## 4 Hardware Threads

- Cannot issue instructions back to back from same thread
- RR scheduling to hide 4 cycle latency

## 512 KB L2 Cache

# Architecture of an Intel® Xeon Phi™ Coprocessor

## Cache

- 32 KB L1 / 512 KB L2 per core
- Fully coherent

## Core Communication

- Bi-directional ring buffer
- 8 GB GDDR5 shared by all cores

## PCIe*

- Gen2
- 16 channels

# Knights Landing Integrated On-Package Memory



**Cache Model**
Let the hardware automatically manage the integrated on-package memory as an "L3" cache between KNL CPU and external DDR

**Flat Model**
Manually manage how your application uses the integrated on-package memory and external DDR for peak performance

# Consistent Tools & Programming Models



Code

Compiler
Libraries
Parallel Models

Software

Multicore

Manycore

Intel® Xeon®
Processors

Intel® Xeon®
Processor

Intel®
Xeon Phi™
Coprocessor

# Wide Range of Development Options

**Thread Parallelism**

Intel® Math Kernel Library
MPI*

OpenMP*

Intel® Threading Building Blocks
Intel® Cilk™ Plus

Pthreads*

**Vector Parallelism**

Intel® Math Kernel Library

Auto-vectorization

Semi-auto Vectorization
(e.g. #pragma ivdep)

Explicit Vectorization
(e.g. Intel® Cilk™ Plus array notation)

OpenCL*

Intrinsics

Ease of Use

Fine Control

# Utilizing an Intel® Xeon Phi™ Coprocessor

## Native



- Target Code:
  Highly parallel (threaded and vectorized) throughout.

- Potential Bottleneck:
  Serial/scalar code.

## Offload



- Target Code:
  Mostly serial, but with expensive parallel regions.

- Potential Bottleneck:
  PCIe* data transfers.

## Symmetric



- Target Code:
  Highly parallel and performs well on both platforms.

- Potential Bottleneck:
  Load imbalance.

# Hands-on:
# Logging in and Running Applications

**Using Beacon @ NICS**

# Legal Disclaimers

# Legal Disclaimers

## 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 SSE3 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

# Legal Disclaimers

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.  For more information go to http://www.intel.com/performance.

Estimated Results Benchmark Disclaimer:
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.

Software Source Code Disclaimer:
Any software source code reprinted in this document is furnished under a software license and may only be used or copied in accordance with the terms of that license.

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND,  EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF  MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND  NONINFRINGEMENT.  IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

# Legal Disclaimers

(intel)

THE UNIVERSITY *of* TENNESSEE
**NATIONAL LEADERSHIP IN HPC**

SC13
Denver, CO | 2013

# Administering the Phi in a Cluster Environment

## R. Glenn Brook

Director, Application Acceleration Center of Excellence
National Institute for Computational Sciences
glenn-brook@tennessee.edu

## Paul Peltz Jr

HPC Systems Administrator
National Institute for Computational Sciences
ppeltz@tennessee.edu

intel
MIC
Partner

Partnership for Extreme Scale Computing
intel
CRAY
NICS

# The Beacon Project

- Funded by NSF to port and optimize scientific codes to the Intel® Xeon Phi™ coprocessor
- State-funded expansion focuses on energy efficiency, big data applications, and industry
- Example Codes: PSC, H3D, OMEN, ENZO, MADNESS, NWCHEM, Amber, MILC, and MAGMA

| Beacon<br>Cray CS300-AC™ Cluster Supercomputer<br>Peak Performance: 210.1 TFLOP/s | |
|---|---|
| Nodes | 4 service, 6 I/O, 48 compute |
| Interconnect | FDR IB Fat Tree |
| Interconnect Bandwidth | 56 GB/s(total) |
| CPU model | Intel Xeon E5-2670 |
| CPUs per node | 2 8-core, 2.6GHz |
| Memory Bandwidth | 128 GB/s (peak) |
| RAM per node | 256 GB |
| SSD per node | 2 x 480 GB (compute),16 x 300 GB (I/O) |
| Intel® Xeon Phi Coprocessors per node | 4 x 5110P 60-core, 1.053GHz<br>8 GB GDDR5 RAM |

Intel, Xeon, and Intel Xeon Phi are trademarks of Intel Corporation in the U.S. and /or other countries.

# Basic Environment

- Installation

  - Diskfull

  - CentOS 6.2

  - Static IPs

- Workload Manager

  - Torque 4.2.6/Moab 7.2

- MPSS 3.1.2

  - In the process of upgrading cluster to MPSS 3.2.x (June)

- Site wide 1.3PB Lustre File System

  - (/lustre/medusa)

- Local 17TB Lustre on ZFS based SSD file system

  - (/lustre/scratch)

# MPSS Installation Procedures

- Preparing the System

  - Set up /etc/hosts to follow this format.  Many scripts will depend on this later.

    | | | |
    |---|---|---|
    | 10.39.20.12 | beacon001 | beacon001-eth0 |
    | 10.39.20.13 | beacon001-mic0 | |
    | 10.39.20.14 | beacon001-mic1 | |
    | 10.39.20.15 | beacon002 | beacon002-eth0 |
    | 10.39.20.16 | beacon002-mic0 | |
    | 10.39.20.17 | beacon002-mic1 | |

- Generate ssh keys

- Install OFED Software Stack

# MPSS Installation Procedures

- ## MPSS Stack Installation

  - pdsh -w cluster[xxx-xxx] yum -y install --nogpgcheck --noplugins --disablerepo=* $MPSS_LOCATION/*.rpm

  - Install extra RPMS as necessary such as $MPSS_LOCATION/ofed/*.rpm

  - micctrl --initdefaults

    - Copy the /etc/mpss/default.conf and /etc/mpss/mic0.conf files to /etc

    - Use these two files as the template for creating a mic-create-conf.sh script to generate config files for each node

  - Create ifcfg-micbr0 and ifcfg-micX files (External Bridging)

  - micctrl --resetconfig

# MPSS Upgrade Procedures

- MPSS Stack Removal

  - pdsh -w cluster[xxx-xxx] $MPSS_LOCATION/uninstall.sh

  - micctrl --initdefaults on test node

    - update template files if necessary

  - micctrl --resetconfig

  - Upgrade firmware with micflash

# Torque Prologue/Epilogue

- Prologue
  - Set up users home env on MIC
    - Copy ssh keys to $TMP
    - Build .profile
    - Check health of mpssd service and ofed-mic service
    - Set up user account
    - Export/Mount $TMP, $NODE:/opt/intel, /global/opt, and /lustre/$USER
    - Verify MIC's /etc/hosts and fix if necessary
    - copy over mpiexec.hydra and pmi_proxy to $TMP

# Torque Prologue/Epilogue

- Epilogue
  - Unmount all shared file systems
  - Reboot MICs
    - service ofed-mic stop
    - service mpss restart
    - micctrl -Rw
    - micctrl -rwf
    - service ofed-mic start
  - Clean up old processes

# Health Checking and Monitoring

- Ganglia
  - Requires post MIC boot RPM installation currently
  - CPU Metrics
    - Intel disables these by default
    - Our testing determined no measurable effect by having them enabled
    - CPU Metrics can be disabled with "gres=noganglia"

# Scripting Common Commands

- ssh

  - micssh replaces ssh

    - adds -i $ssh_key option to ssh

- mpiexec

  - micmpiexec replaces mpiexec

    - replaces ssh with micssh

    - passes the appropriate environment variables

    - allows for debuggers such as ddt and totalview

# Challenges

- Driver upgrade process (2 to 3.1 and 3.1 to 3.2)
  - The ever changing default.conf and micX.conf files
- The micctrl gamble
  - Intel wants to push everything to using micctrl
  - we need a --dry-run flag
- Intel Compilers and Intel MPI
- Kernel Compatibility
  - rpmbuild --rebuild
    - sufficient for mpss driver
    - difficult for OFED 1.5.4.1 (unless you use OFED-3.5)
- SSH Keys
- OFED 1.5.4.1 or 3.5

# Acknowledgements

We wish to thank the National Science Foundation for their support of the Beacon project, along with the State of Tennessee.

# Contact Information

**Paul Peltz Jr**

HPC Systems Administrator
National Institute for Computational Sciences
ppeltz@utk.edu

# Using Intel MPI and OpenMP 4.0 with the Intel Xeon Phi

**Vincent C. Betro, Ph.D.**

**Hans Pabst, Intel**

**Heinrich Blockhorst, Intel**

**Cray Users Group—May 5, 2014**

# OVERVIEW INFORMATION

# Enabling & Advancing Parallelism
## High Performance Parallel Programming

**Intel tools, libraries and parallel models extend to multicore, many-core and heterogeneous computing**
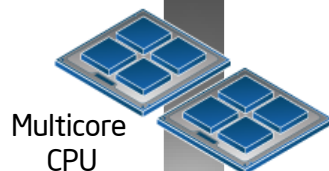
Code

Compiler
Libraries
Parallel Models

(intel)
Software

Multicore

Many-core

Cluster

Multicore CPU

Multicore CPU

Intel® Xeon Phi™ Coprocessor

Multicore Cluster

Multicore & Many-core Cluster

**Use One Software Architecture Today. Scale Forward Tomorrow.**

NICS

3

# Preserve Your Development Investment
## Common Tools and Programming Models for Parallelism

**C/C++**
- OpenCL*
- Intel® Cilk™ Plus
- OpenMP*
- Intel® TBB
- Offload Pragmas
- Intel® C/C++ Compiler

- Intel® MKL
- Intel® MPI

**Fortran**
- Intel® Fortran Compiler
- Coarray
- Offload Directives
- OpenMP*

Multicore

Heterogeneous Computing

Many-core

**Develop Using Parallel Models that Support Heterogeneous Computing**

# Intel® MPI Library Overview

- **Intel is a leading vendor of MPI implementations and tools**

- **Optimized MPI application performance**
  - **Application-specific tuning**
  - **Automatic tuning**

- **Low latency**
  - **Industry leading latency**

- **Interconnect Independence & Runtime Selection**
  - **Multi-vendor interoperability**
  - **Performance optimized support for the latest OFED capabilities through DAPL 2.0**

- **More robust MPI applications**
  - **Seamless interoperability with Intel® Trace Analyzer and Collector**

# Spectrum of Programming Models and Mindsets



Multi-Core Centric

Many-Core Centric

Intel® Xeon ®

Intel® Xeon Phi™

**Multi-Core Hosted**
*General purpose serial and parallel computing*

**Symmetric**
*Codes with balanced needs*

**Many Core Hosted**
*Highly-parallel codes*

**Offload**
*Codes with highly-parallel phases*

Multi-core
(Intel®
Xeon ®)

Many-core
(Intel®
Xeon Phi™)

Main( )
Foo( )
MPI_*( )

Main( )
Foo( )
MPI_*( )

Main( )
Foo( )
MPI_*( )

Foo( )

Main( )
Foo( )
MPI_*( )

Main( )
Foo( )
MPI_*( )

**Range of models to meet application needs**

NICS

# Levels of communication speed

- **Current clusters are not homogenous regarding communication speed:**
  - Inter node (InfiniBand*, Ethernet, etc)
  - Intra node
    - Inter sockets (Quick Path Interconnect)
    - Intra socket

- **Two additional levels to come with Intel® Xeon Phi™ coprocessor:**
  - Host-coprocessor communication
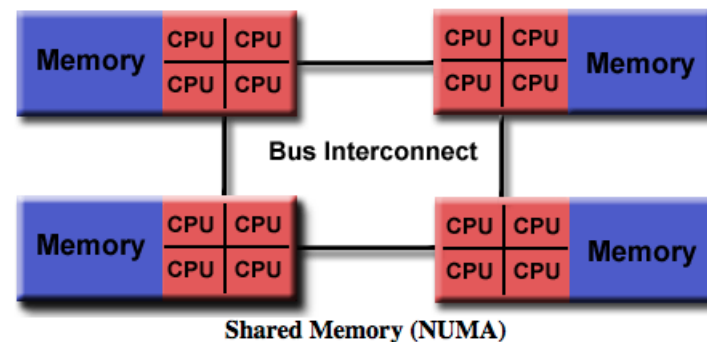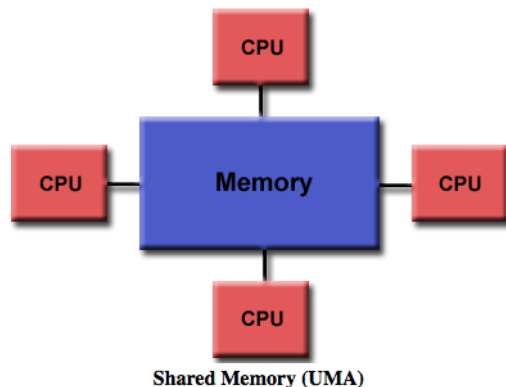  - Inter coprocessor communication

# Selecting network fabrics

- Intel® MPI automatically selects the best available network fabric it can find.
  - Use `I_MPI_FABRICS` to select a different communication device explicitly

- The best fabric is usually based on InfiniBand* (dapl, ofa) for inter node communication and shared memory for intra node

- Available for Intel® Xeon Phi™ coprocessor:
  - `shm, tcp, ofa, dapl`
  - Availability checked in the order `shm:dapl, shm:ofa, shm:tcp` (intra:inter)
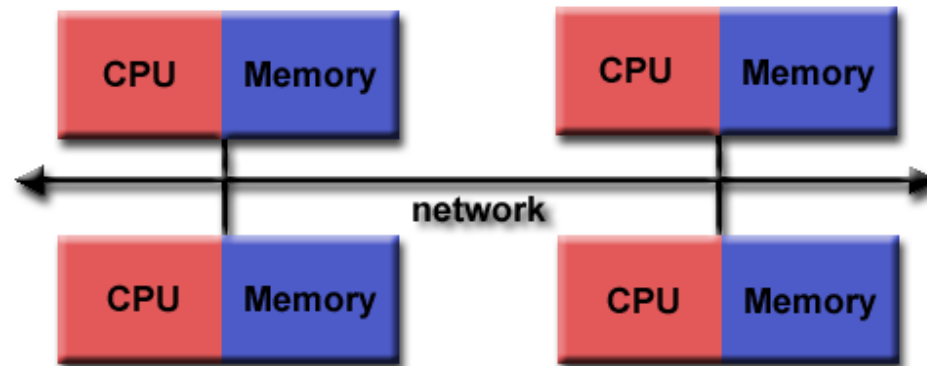
# MPI vs OpenMP: What's the Difference?

- **OpenMP is used in a shared memory space**

- **Often referred to as threading, allows multiple tasks to occur simultaneously which are:**
  - **Embarrassingly parallel**
  - **Are working in disjoint areas of memory**
  - **Can have any "intersections" of threads be caught using atomic or critical sections**



Shared Memory (UMA)



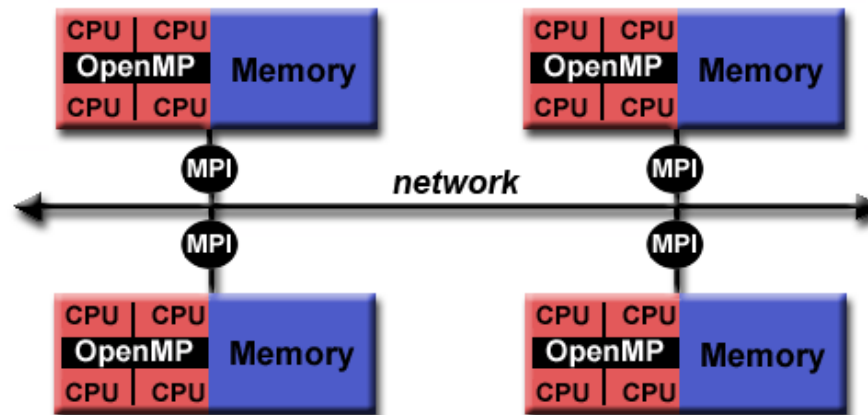Bus Interconnect

Shared Memory (NUMA)

# MPI vs OpenMP: What's the Difference?

- **MPI is used in a distributed memory space**

- **Often referred to as message passing, allows multiple chunks of a problem's domain to be computed on simultaneously**

- **Each computer functions alone with the network acting as the bridges between to update values on process boundaries between iterations**

**NICS**

# MPI vs OpenMP: What's the Difference?

- **The two paradigms can be used together**

- **One example would be having a few MPI ranks on each Xeon Phi and then having many threads spawned by each rank**

- **Another would be to have MPI tasks on CPUs and offload threaded sections to Xeon Phis**

# Xeon Phi Programming Models

- **Native Mode**
  - **Everything runs on the MIC**
  - **All libraries need to be recompiled with -mmic**

- **Offload Mode**
  - **Serial portion runs on host**
  - **Parallel portions are offloaded and run on the MIC**

# Xeon Phi Programming Models

- **Native Mode**
  - **Everything runs on the MIC**
  - **All libraries need to be recompiled with –mmic**

**…is all we will discuss here.  Offload will be covered next by Hans, including OpenMP 4.0!**

**However, we will briefly chat about offload in terms of having many MPI ranks do it.**
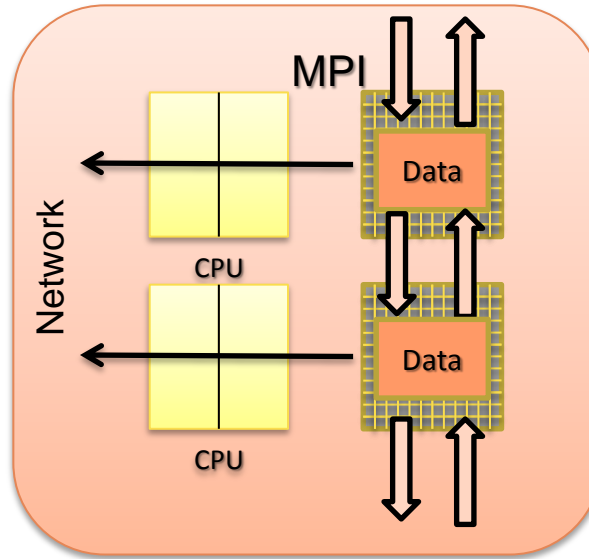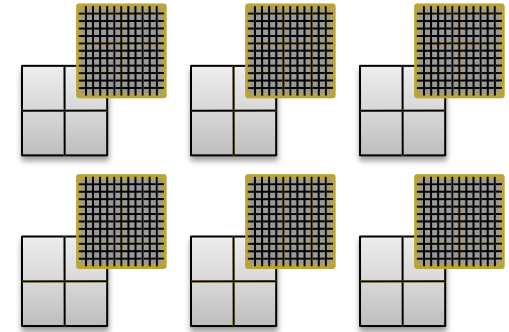
# CO-PROCESSOR OR HOST MPI

# Intel MPI

- **The installed Intel MPI library implements the Message Passing Interface, version 2.2 (MPI-2.2) specifications**

- **3 Programming models are supported**
  - **Co-processor only model**
  - **Symmetric model**
  - **MPI offload model**

- **Intel MPI compilers have an extra 'i' in their name: `mpiicc, mpiicpc, mpiifort`**

- **MPI applications should be launched from the host compute node using `micmpiexec`**

NICS

# Coprocessor-only Programming Model

- **MPI ranks on Intel® Xeon Phi™ coprocessor(only)**

- **All messages into/out of coprocessors**

- **Intel® Cilk™ Plus, OpenMP*, Intel® Threading Building Blocks, Pthreads used directly within MPI processes**



MPI

Network

CPU

Data

CPU

Data

**Homogenous network of many-core CPUs**

Build Intel® Xeon Phi™ binary using the Intel® compiler.

Upload the binary to the Intel® Xeon Phi™ coprocessor.

Run instances of the MPI application on Intel® Xeon Phi™ coprocessor nodes.

# Coprocessor-only Programming Model

- **MPI ranks on the Intel® Xeon Phi™ coprocessor(s) only**

- **MPI messages into/out of the coprocessor(s)**

- **Threading possible**

- Build the application for the Intel® Xeon Phi™ coprocessor

  ```
  # mpiicc -mmic -o test_hello.MIC test.c
  ```

- Launch the application on the coprocessor from host

  ```
  # export I_MPI_MIC=enable
  # mpirun -n 2 -host host-mic0 ./test_hello.MIC
  ```

- Alternatively: login to the Intel® Xeon Phi™ coprocessor and start the MPI run there!

- No NFS: Upload the Intel® Xeon Phi™ executable and add working directory flag

  ```
  # scp ./test_hello.MIC host-mic0:/my_mic_dir/
  # mpirun … -wdir /my_mic_dir/ …
  ```

# Using Intel MPI: MIC 2 MIC

- If one MIC card is not sufficient for your domain decomposition, you may use all MIC cards on the node or even multiple MIC cards on multiple nodes.

- If you need to use MICs on multiple nodes, you must request multiple nodes with qsub and check which ones they are with "cat $PBS_NODEFILE"
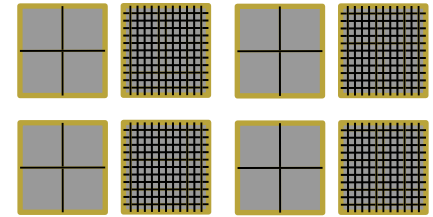
- Note that the MIC requires the –wdir argument

```
micmpiexec -n 2 -wdir $TMPDIR -host beacon#-mic0
$TMPDIR/mpi_hello.MIC : -n 2 -wdir $TMPDIR  -host
beacon#-mic1 $TMPDIR/mpi_hello.MIC
```

# Symmetric Programming Model

- **MPI ranks on Intel® Xeon Phi™ Architecture and host CPUs**

- **Messages to/from any core**

- **Intel® Cilk™ Plus, OpenMP*, Intel® Threading Building Blocks, Pthreads* used directly within MPI processes**

**Heterogeneous network of homogeneous CPUs**

MPI

Network

Data · Data

CPU

Data · Data

CPU

**Build binaries by using the resp. compilers targeting Intel® 64 and Intel® Xeon Phi™ Architecture.**

**Upload the binary to the Intel® Xeon Phi™ coprocessor.**

**Run instances of the MPI application on different mixed nodes.**

# Symmetric model

- **MPI ranks on the coprocessor(s) and host CPU(s)**

- **MPI messages into/out of the coprocessor(s)  and host CPU(s)**

- **Threading possible**

- Build the application for Intel®64 and the Intel® Xeon Phi™ Architecture separately

```
# mpiicc -o test_hello test.c
# mpiicc –mmic -o test_hello.MIC test.c
```

- Launch the application on the host and the coprocessor

```
# export I_MPI_MIC=enable
# mpirun -n 2 -host <hostname> ./test_hello :
        -n 2 -host host-mic0 ./test_hello.MIC
```

- No NFS: Upload the Intel® Xeon Phi™ executable and add flag

```
# scp ./test_hello.MIC host-mic0:/my_mic_dir/
# mpirun … : –wdir /my_mic_dir/ …
```

# Utilize the POSTFIX env variable
## *Support for NFS-shared cards*

- **Assumption: The current working directory is available with identical path on the coprocessor (e.g. mounted)**
- **Specify the suffix of the coprocessor binary**

```
# export I_MPI_MIC_POSTFIX=.MIC
```

- **Specify the node names in a file**

```
# cat mpi_hosts
host
host-mic0
```

- **Execute in symmetric mode on the host and the coprocessor**

```
# export I_MPI_MIC=enable

# mpirun -f mpi_hosts -n 4 ./test_hello
```

- **The binary** `./test_hello${I_MPI_MIC_POSTFIX}` **will be used by mpirun on the coprocessor**

# Utilize the PREFIX env variable
## *Support for NFS-shared cards*

- **Assumption: The current working directory is available with identical path on the coprocessor (e.g. mounted)**
- **Place the coprocessor binary in a separate directory, but with identical basename of the host**
  ```
  # mpiicc –mmic -o ./MIC/test_hello test.c
  ```
- **Specify the prefix of the coprocessor binary**

  ```
  # export I_MPI_MIC_PREFIX=./MIC/
  ```

- **Execute in symmetric mode on the host and the coprocessor**

  ```
  # export I_MPI_MIC=enable

  # mpirun –f mpi_hosts –n 4 ./test_hello
  ```

- **The binary `${I_MPI_MIC_PREFIX}./test_hello` will be used by mpirun on the coprocessor**

# Launching an MPI Application with Manually Specified Hosts

- **Launch an MPI application on mic0 of the current compute node with**
  - `micmpiexec -n 1 -wdir $TMPDIR -host beacon#-mic0 $TMPDIR/application.MIC`

- **Launch an MPI application on both mic0 and mic1 of the current compute node with**
  - `micmpiexec -n 1 -wdir $TMPDIR -host beacon#-mic0 $TMPDIR/application.MIC : -n 1 -wdir $TMPDIR -host beacon#-mic1 $TMPIDR/application.MIC`

- **Launch an MPI application on both MICs and the compute node with**
  - `micmpiexec -n 1 -wdir $TMPDIR -host beacon#-mic0 $TMPDIR/application.MIC : -n 1 -wdir $TMPDIR -host beacon#-mic1 $TMPIDR/application.MIC : -n 1 -host beacon# ./application`

# Launching an MPI Application with a Machine File

- **The machine file needs to be of the form `<host>:<number of ranks>`**

- **Sample machine file named hosts_file:**
  - `beacon11:8`
  - `beacon12:8`
  - `beacon11-mic0:2`
  - `beacon11-mic1:2`
  - `beacon12-mic0:2`
  - `beacon12-mic1:2`

- **Launch the MPI application using**
  - `micmpiexec -machinefile hosts_file -n 16 ./application : -n 8 -wdir $TMPDIR -genv LD_LIBRARY_PATH $TMPDIR/lib:/lib64:/lib $TMPDIR/application.MIC`

# Launching an MPI Application with Process Pinning

- **export I_MPI_PIN=1**

- **export I_MPI_PIN_PROCESSOR_LIST="0-6,8-14"**

- **export I_MPI_DEBUG=4 or 5**

- **Allows you to pin MPI processes to sockets on the host.**

**NICS**

# Example of Native Mode MPI application experiences





Both GROMACS and BLAST are bioinformatics packages used to look at proteins, lipids, and nucleic acids, often for the purpose of Genomics research or drug docking research.

Both use the MPI-OpenMP hybrid model of execution, where the problem is decomposed in an embarrassingly parallel fashion over several MPI ranks and then the consequent divisions are threaded in execution.

This approach allows the developers to take advantage of the large number of threads available on the Intel Xeon Phi along with allowing them to decompose the problem well so it fits on a given coprocessor.

# MPI on HOST + OFFLOAD on PHI

# MPI+Offload Programming Model

- **MPI ranks on Intel® Xeon® processors (only)**

- **All messages into/out of host CPUs**

- **Offload models used to accelerate MPI ranks**

- **Intel® Cilk™ Plus, OpenMP\*, Intel® Threading Building Blocks, Pthreads\* within Intel® Xeon Phi™ coprocessors**



**Homogenous network of heterogeneous nodes**

Build Intel® 64 executable with included offload by using the Intel® compiler.

Run instances of the MPI application on the host, offloading code onto the coprocessor.

Advantages of more cores and wider SIMD for certain applications

# MPI+Offload Programming Model

- **MPI ranks on the host CPUs only**

- **MPI messages into/out of the host CPUs**

- **Intel® Xeon Phi™ coprocessor as an accelerator**

- Compile for MPI and internal offload

  `# mpiicc –o test test.c`

- Latest compiler compiles by default for offloading if offload construct is detected!
  - Switch off by `-no-offload` flag

- Execute on host(s) as usual

  `# mpirun -n 2 ./test`

- MPI processes will offload code for acceleration

# Launching an MPI Application with Each Rank Offloading to a Different Card

- **When specifying the target for offload by a given rank, one can use mic:N, where N=0,1,2,3… for as many Xeon Phis as are on the system.**

- **Despite the fact that the pragmas are pre-processed, one can use N as a variable to allow it to depend on which MPI rank is calling the offload at runtime.**

# MPI+Offload Support

- **How to control mapping of threads on the coprocessor?**
  - **How do I avoid that offload of first MPI process interferes with offload of second MPI process, i.e. by using identical cores/threads on the coprocessor?**
  - **Default: No special support (now). Offloads from MPI processes handled by system like offloads from independent processes (or users).**

- **Define thread affinity manually per single MPI process:**

```
# export OMP_NUM_THREADS=8
# mpirun -env KMP_AFFINITY=proclist=[1-8],explicit -n 1
        -host myHost ./test_mpioffload :
        -env KMP_AFFINITY=proclist=[9-16],explicit -n 1
        -host myHost ./test_mpioffload : ...
        ...
```

# MPI+Offload Support (ctd.)

- **Alternative: Use `KMP_PLACE_THREADS` and Intel® MPI rank number `PMI_RANK` in a wrapper script:**

```
# cat ./wrapoffload.sh
  cores=$(( (OMP_NUM_THREADS+3)/4 ))
  offset=$(( cores*PMI_RANK ))
  export KMP_PLACE_THREADS=${cores}Cx4T,${offset}O
  ./test_mpioffload


# export OMP_NUM_THREADS=8

# mpirun -n 4 -host myHost ./wrapoffload.sh
```

- **The mapping will be:**

```
MPI rank 0: KMP_PLACE_THREADS=2Cx4T,0O == [1-8]

MPI rank 1: KMP_PLACE_THREADS=2Cx4T,2O == [9-16]
```

# OFFLOADING TO MULTIPLE CARDS FROM ONE OR MORE HOST PROCESSES

# Simulataneous Computing using OpenMP

- **Any OMP call blocks until the statement completes, unless the "nowait" modifier is used**

- **To use both the host and MIC simultaneously, multiple threads need to be executed on host**
  - One or more threads that contain an offload call
  - Other threads have the host do some work

- **With OpenMP, this achieved using OpenMP task calls**

NICS

# Offloading to multiple cards from one host process

- **Unoptimized fashion: http://software.intel.com/en-us/forums/topic/393649**

```
//Allocate memory on each of N mics
for(N=0; N<=1; N++){
#pragma offload_transfer target(mic:N) nocopy(pi: alloc_if(1) free_if(0)) signal(&test1[N])
}
for(N=0; N<=1; N++){
#pragma offload target(mic:N) in(N,num_steps,step,sum) inout(pi[N]) signal(&test[N]) wait(&test1[N])
 {
 #pragma omp parallel for reduction(+:sum)
 for (i=0;i<iter; i++)
 <parallel loop body>
 }
}
//Free memory on mic
for(N=0; N<=1; N++){
#pragma offload_transfer target(mic:N) nocopy(pi: alloc_if(0) free_if(1)) wait(&test[N])
}
```

NICS

# Offloading to multiple cards from one host process

- **According to OpenMP 4.0 (and in my opinion to be pedantic and make sure it behaves):**

```
#pragma omp parallel
#pragma omp single
{
#pragma omp task
  #pragma omp target(mic) OR #pragma offload target(mic)
  {
      <various serial code>
      #pragma omp parallel for
      for (int i=0; i<limit; i++)
         <parallel loop body>
  }
#pragma omp task
  {<host code or another offload>}
}
```

# Example of Offload Mode OpenMP application experiences

- **NAMD** (**NA**noscale **M**olecular **D**ynamics program)

  - freeware molecular dynamics simulation package
  - uses Charm++ parallel programming model
  - simulates large systems (millions of atoms)

  Charm++ has been ported to the Intel Xeon Phi, and NAMD is being run across several Intel Xeon Phis in offload mode on Beacon and other resources.

**NICS**

# Useful Environment Variables

- **The following applies only to offload mode execution**

- **All environment variables defined on the host are replicated on the MIC in offload mode**

- **To modify specific MIC vaules, `MIC_ENV_PREFIX` must be defined**

```
OMP_NUM_THREADS=8
OMP_STACKSIZE=16M
MIC_ENV_PREFIX=MIC_
```

```
MIC_OMP_NUM_THREADS=96
MIC_OMP_STACKSIZE=4M
```

- **For `csh`: `setenv ENV_VARIABLE VALUE`**

- **For `sh`: `export ENV_VARIABLE=VALUE`**

**NICS**

# Useful Environment Variables part 2

- `OFFLOAD_REPORT` **can be useful when trying to debug code that offloads**
  - `OFFLOAD_REPORT=1`
    - **Gives basic information (e.g. CPU time) about whether code blocks marked for offload are running on the host or coprocessor**
  - `OFFLOAD_REPORT=2`
    - **Gives detailed information (e.g. CPU time and data transfer) about the offload process**

- **Use** `MIC_HOST_LOG` **to output traces to a file**
  - `MIC_HOST_LOG=~/app/mic.log`

# MPI+OPENMP ON CO-PROCESSOR (OR HOST)

# Running Hybrid Code in Native Mode

- **One major advantage of running on the Xeon Phi is that despite its lower clock speed per core, there are an egregious number of threads that can be spawned from LOCAL MPI ranks on a card**

- **Here, we will discuss a paradigm for doing this as well as how to determine the optimal combination of MPI ranks and threads such that your domain fits on the cards and the threads remain saturated throughout computation**

**NICS**

# Traditional Cluster Computing

- **MPI is »the« portable cluster solution**

- **Parallel programs use MPI over cores inside the nodes**

  - Homogeneous programming model

  - "Easily" portable to different sizes of clusters

  - No threading issues like »False Sharing« (common cache line)

  - Maintenance costs only for one parallelization model

# Traditional Cluster Computing (contd.)

- **Hardware trends**
  - **Increasing number of cores per node - plus cores on co-processors**
  - **Increasing number of nodes per cluster**
- **Consequence: Increasing number of MPI processes per application**
- **Potential MPI limitations**
  - **Memory consumption per MPI process, sum exceeds the node memory**
  - **Limited scalability due to exhausted interconnects (e.g. MPI collectives)**
  - **Load balancing is often challenging in MPI**

# Hybrid Computing

- **Combine MPI programming model with threading model**

- **Overcome MPI limitations by adding threading:**

  - **Potential memory gains in threaded code**

  - **Better scalability (e.g. less MPI communication)**

  - **Threading offers smart load balancing strategies**

- **Result: Maximize performance by exploitation of hardware (incl. co-processors)**

# Example: MPI Load Imbalance



i ⟶

Difficult to implement load balancing in nodes with MPI

Nodes

Proc 4  Proc 5  ...

j

Proc 0  Proc 1  Proc 2  Proc 3

4 Cores per Node ⟶

Dark red = high load

# Example: Hybrid Load Balance



4 Threads per Node on 4 Cores

# Options for Thread Parallelism

Intel® Math Kernel Library

OpenMP*

Intel® Threading Building Blocks
Intel® Cilk™ Plus

Pthreads* and other threading libraries

Ease of use / code maintainability

Programmer control

Choice of unified programming to target Intel® Xeon and Intel® Xeon Phi™!

# Intel® MPI Support of Hybrid Codes

- **Intel® MPI is strong in mapping control**

- **Sophisticated default or user controlled**
  - `I_MPI_PIN_PROCESSOR_LIST` **for pure MPI**
  - **For hybrid codes (default, takes precedence):**
    `I_MPI_PIN_DOMAIN` **=*<size>[:<layout>]***
      - **<size> =**

        | | |
        |---|---|
        | `omp` | **Adjust to OMP_NUM_THREADS** |
        | `auto` | **#CPUs/#MPIprocs (default)** |
        | `<n>` | **Number** |

      - **<layout> =**

        | | |
        |---|---|
        | `platform` | **According to BIOS numbering** |
        | `compact` | **Close to each other** |
        | `scatter` | **Far away from each other** |

- **Naturally extends to hybrid codes on Intel® Xeon Phi™**

# Intel® MPI Support of Hybrid Codes

- **Define `I_MPI_PIN_DOMAIN` to split logical processors into non-overlapping subsets**

- **Mapping rule:** 1 MPI process per 1 domain



Pin OpenMP threads inside the domain with `KMP_AFFINITY` (or in the code)

# Intel® MPI Environment Support

- **The execution command mpirun of Intel® MPI reads argument sets from the command line:**
  - Sections between „ : " define an argument set (alternatively a line in a configfile specifies a set)
  - Host, number of nodes, but also environment can be set independently in each argument set

```
# mpirun –env I_MPI_PIN_DOMAIN 4 –host myXEON ...
        : -env I_MPI_PIN_DOMAIN 16 –host myMIC
```

- **Adapt the important environment variables to the architecture**
  - `OMP_NUM_THREADS`, `KMP_AFFINITY` for OpenMP
  - `CILK_NWORKERS` for Intel® Cilk™ Plus

# Coprocessor-only and Symmetric Support

- **Full hybrid support on Intel® Xeon from Intel® MPI extends to the Intel® Xeon Phi™ coprocessor**

- `KMP_AFFINITY=balanced` **(only on the coprocessor) in addition to** `scatter` **and** `compact`

- `KMP_PLACE_THREADS=<n>Cx<m>T,<o>O` **(**`<n>`**-Cores times** `<m>`**-Threads with** `<o>`**-cores Offset, only on coprocessor) in addition to** `KMP_AFFINITY` **for exact but still generic thread placement**

- **Recommendations:**
  - **Explicitly control where MPI processes and threads run in a hybrid application (according to threading model)**
  - **Avoid splitting cores among MPI processes, i.e.** `I_MPI_PIN_DOMAIN` **should be a multiple of 4**
  - **Try different** `KMP_AFFINITY` **and/or** `KMP_PLACE_THREADS` **settings for your application**

# OS Thread Affinity Mapping

- **The Intel® Xeon Phi™ coprocessor has N cores, each with 4 hardware thread contexts, for a total of M=4*N threads**

- **The OS maps "procs" to the M hardware threads:**

| MIC core | 0 | | | | 1 | | … | (N-2) | (N-1) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIC HW thread | 0 | 1 | 2 | 3 | 0 | 1 | … | 3 | 0 | 1 | 2 | 3 |
| OS "proc" | 1 | 2 | 3 | 4 | 5 | 6 | … | (M-4) | 0 | (M-3) | (M-2) | (M-1) |

- **The OS runs on proc 0, which lives on core (N-1)!**
  - **Rule of thumb: Avoid using OS procs 0, (M-3), (M-2), and (M-1) to avoid contention with the OS**
    - **Only less than 2% resources unused (1/#cores)**
  - **Especially important when using the offload model due to data transfer activity!**
  - **But: Non-offload applications may slightly benefit from running on core (N-1)**

# OS Thread Affinity Mapping (ctd.)

- **OpenMP library maps to the OS "procs"**
- **Examples (for non-offload apps which benefit from core N-1):**
  - `KMP_AFFINITY=compact,granularity=thread,compact`

| MIC core | 0 | | | | 1 | | … | (N-2) | (N-1) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIC HW thread | 0 | 1 | 2 | 3 | 0 | 1 | … | 3 | 0 | 1 | 2 | 3 |
| OS "proc" | 1 | 2 | 3 | 4 | 5 | 6 | … | (M-4) | 0 | (M-3) | (M-2) | (M-1) |
| OpenMP thread | 0 | 1 | 2 | 3 | 4 | 5 | … | (M-5) | (M-4) | (M-3) | (M-2) | (M-1) |

  - `KMP_AFFINITY=balanced,granularity=thread`
    `OMP_NUM_THREADS=n=M/2`

| MIC core | 0 | | | | 1 | | … | (N-2) | (N-1) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIC HW thread | 0 | 1 | 2 | 3 | 0 | 1 | … | 3 | 0 | 1 | 2 | 3 |
| OS "proc" | 1 | 2 | 3 | 4 | 5 | 6 | … | (M-4) | 0 | (M-3) | (M-2) | (M-1) |
| OpenMP thread | 0 | 1 | | | 3 | 4 | … | | (n-2) | (n-1) | | |

# OS Thread Affinity Mapping (ctd.)

- **Use `balanced` affinity to minimize False Sharing!**
  - `KMP_PLACE_THREADS=2Cx2T,0O`
  - **but still with implicit default mapping `scatter,granularity=thread`**

| MIC core | 0 | | | | 1 | | … | (N-2) | (N-1) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIC HW thread | 0 | 1 | 2 | 3 | 0 | 1 | … | 3 | 0 | 1 | 2 | 3 |
| OS "proc" | 1 | 2 | 3 | 4 | 5 | 6 | … | (M-4) | 0 | (M-3) | (M-2) | (M-1) |
| OpenMP thread | 0 | 2 | | | 1 | 3 | | | | | | |

  - `KMP_PLACE_THREADS=2Cx2T,0O` **and `KMP_AFFINITY=balanced`**

| MIC core | 0 | | | | 1 | | … | (N-2) | (N-1) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIC HW thread | 0 | 1 | 2 | 3 | 0 | 1 | … | 3 | 0 | 1 | 2 | 3 |
| OS "proc" | 1 | 2 | 3 | 4 | 5 | 6 | … | (M-4) | 0 | (M-3) | (M-2) | (M-1) |
| OpenMP thread | 0 | 1 | | | 2 | 3 | | | | | | |

# Additional Resources

- **Other documentation, presentations, and even a community forum can be found at**
  - **http://software.intel.com/mic-developer**
  - **http://www.openmp.org/mp-documents/OpenMP4.0.0.pdf**

  **Let's do a Trace Analyzer Lab!**

  **http://software.intel.com/en-us/articles/intelr-xeon-phitm-advanced-workshop-labs**

# Contact

## Vincent Betro
vbetro@tennessee.edu

## Help
help@nics.utk.edu

**NICS**

# OpenMP* 4.0 for HPC in a Nutshell

**C U G   2 0 1 4   ·   L u g a n o**

Hans Pabst, May 5th 2014
[Slides by Dr. M. Klemm, Intel Corp.]

# Legal Disclaimer & Optimization Notice

SOFTWARE AND SERVICES

# OpenMP API

- De-facto standard, OpenMP 4.0 out since July 2013

- API for C/C++ and Fortran for shared-memory parallel programming

- Based on directives (pragmas in C/C++)

- Portable across vendors and platforms

- Supports various types of parallelism

# Evolution of Hardware (at Intel)

*Images not intended to reflect actual die sizes*

|  | 64-bit Intel® Xeon® processor | Intel® Xeon® processor 5100 series | Intel® Xeon® processor 5500 series | Intel® Xeon® processor 5600 series | Intel® Xeon® processor E5-2600v2 series | Intel® Xeon Phi™ Co-processor 7120P |
|---|---|---|---|---|---|---|
| Frequency | 3.6GHz | 3.0GHz | 3.2GHz | 3.3GHz | 2.7GHz | 1.238MHz |
| Core(s) | 1 | 2 | 4 | 6 | 12 | 61 |
| Thread(s) | 2 | 2 | 8 | 12 | 24 | 244 |
| SIMD width | 128 (2 clock) | 128 (1 clock) | 128 (1 clock) | 128 (1 clock) | 256 (1 clock) | 512 (1 clock) |

# Levels of Parallelism in OpenMP 4.0

| Level | | Description | Annotation |
|---|---|---|---|
| Cluster | 🟥 | Group of computers communicating through fast intercon... | OpenMP 4.0 for Devices |
| Coprocessors/Accelerators | 🟩 | Special compute devices attached to the local node through special interconnect | |
| Node | 🟩 | Group of processors communicating through shared memory | OpenMP 4.0 Affinity |
| Socket | 🟩 | Group of cores communicating through shared cache | |
| Core | 🟩 | Group of functional units communicating through registers | |
| Hyper-Threads | 🟩 | Group of thread contexts sharing functional units | |
| Superscalar | 🟧 | Group of instructions sharing functional units | |
| Pipeline | 🟧 | Sequence of instructions sharing functional unit... | OpenMP 4.0 SIMD |
| Vector | 🟩 | Single instruction using multiple functional units | |

**SOFTWARE AND SERVICES**

# OpenMP Intro in Three Slides (1)

```
#pragma omp parallel
{
    #pragma omp for
    for (i = 0; i<N; i++)
    {…}

    #pragma omp for
    for (i = 0; i< N; i++)
    {…}
}
```

fork

distribute work

barrier

distribute work

barrier

join

# OpenMP Intro in Three Slides (2)

```
double a[N];
double l,s = 0;
#pragma omp parallel for reduction(+:s) private(l) \
                         schedule(static,4)

for (i = 0; i<N; i++)
{
    l = log(a[i]);
    s += l;
}
```

s=0

s'=0  s''=0  s'''=0  s''''=0

distribute work

barrier

s'+= s''        s'''+= s''''

s'+= s'''

s = s'

# OpenMP Intro in Three Slides (3)

```
#pragma omp parallel
#pragma omp single
for(e = l->first; e ; e = e->next)
    #pragma omp task
        process(e);
```

fork

join

# OpenMP 4.0 for Devices

# Device Model

- OpenMP 4.0 supports accelerators/coprocessors
- Device model:
  - One host
  - Multiple accelerators/coprocessors of the same kind



Coprocessors

Host

# OpenMP 4.0 for Devices - Constructs

- Transfer control [and data] from the host to the device

- Syntax (C/C++)
  ```
  #pragma omp target [data] [clause[[,] clause],…]
  structured-block
  ```

- Syntax (Fortran)
  ```
  !$omp target [data] [clause[[,] clause],…]
  structured-block
  !$omp end target [data]
  ```

- Clauses
  ```
  device(scalar-integer-expression)
  map(alloc | to | from | tofrom: list)
  if(scalar-expr)
  ```

SOFTWARE AND SERVICES

# Execution Model

- The `target construct` transfers the control flow to the target device
    - Transfer of control is sequential and synchronous
    - The transfer clauses control direction of data flow
    - Array notation is used to describe array length
- The `target data` construct creates a scoped device data environment
    - Does not include a transfer of control
    - The transfer clauses control direction of data flow
    - The device data environment is valid through the lifetime of the target data region
- Use `target update` to request data transfers from within a target data region

# Execution Model

- ## Data environment is lexically scoped
  - Data environment is destroyed at closing curly brace
  - Allocated buffers/data are automatically released

# Example

```
#pragma omp target data device(0) map(alloc:tmp[:N]) map(to:input[:N]) map(from:res)
   {
#pragma omp target device(0)
#pragma omp parallel for
    for (i=0; i<N; i++)
      tmp[i] = some_computation(input[i], i);


    update_input_array_on_the_host(input);


#pragma omp target update device(0) to(input[:N])


#pragma omp target device(0)
#pragma omp parallel for reduction(+:res)
    for (i=0; i<N; i++)
      res += final_computation(input[i], tmp[i], i)
   }
```

host

target

host

target

host

# `teams` Construct

- Support multi-level parallel devices

- Syntax (C/C++):
  ```
  #pragma omp teams [clause[[,] clause],…]
  structured-block
  ```

- Syntax (Fortran):
  ```
  !$omp teams [clause[[,] clause],…]
  structured-block
  ```

- Clauses
  ```
  num_teams(integer-expression)
  num_threads(integer-expression)
  default(shared | none)
  private(list), firstprivate(list)
  shared(list), reduction(operator : list)
  ```

# Offloading SAXPY to a Coprocessor

```
int main(int argc, const char* argv[]) {
  float *x = (float*) malloc(n * sizeof(float));
  float *y = (float*) malloc(n * sizeof(float));
  // Define scalars n, a, b & initialize x, y

#pragma omp target data map(to:x[0:n])
  {
#pragma omp target map(tofrom:y)
#pragma omp teams num_teams(num_blocks) num_threads(nthreads)
```



**all do the same**

```
  for (int i = 0; i < n; i += num_blocks){
    for (int j = i; j < i + num_blocks; j++) {
        y[j] = a*x[j] + y[j];
  } }
  }
  free(x); free(y); return 0;
}
```

# Offloading SAXPY to a Coprocessor

```
int main(int argc, const char* argv[]) {
  float *x = (float*) malloc(n * sizeof(float));
  float *y = (float*) malloc(n * sizeof(float));
  // Define scalars n, a, b & initialize x, y


#pragma omp target data map(to:x[0:n])
{
#pragma omp target map(tofrom:y)
#pragma omp teams num_teams(num_blocks) num_threads(bsize)
```

**all do the same**

```
#pragma omp distribute
  for (int i = 0; i < n; i += num blocks){
```

**workshare (w/o barrier)**

```
#pragma omp parallel for
    for (int j = i; j < i + num blocks; j++) {
```

**workshare (w/ barrier)**

```
        y[j] = a*x[j] + y[j];
  } }
} free(x); free(y); return 0; }
```

# Offloading SAXPY to a Coprocessor

```c
int main(int argc, const char* argv[]) {
  float *x = (float*) malloc(n * sizeof(float));
  float *y = (float*) malloc(n * sizeof(float));
  // Define scalars n, a, b & initialize x, y

#pragma omp target map(to:x[0:n]) map(tofrom:y)
  {
#pragma omp teams distribute parallel for \
        num_teams(num_blocks) num_threads(bsize)
    for (int i = 0; i < n; ++i){
      y[i] = a*x[i] + y[i];
    }
  }

  free(x); free(y); return 0;
}
```

# OpenMP 4.0 Affinity

# NUMA is here to Stay…

- (Almost) all multi-socket compute servers are NUMA systems
  - Different access latencies for different memory locations
  - Different bandwidth observed for different memory locations
- Example: Intel® Xeon E5-2600v2 Series processor

# Thread Affinity – Why It Matters?



STREAM Triad, Intel® Xeon E5-2697v2

Y-axis: GB/sec [higher is better] (0.00 to 100.00)

X-axis: # of threads/cores (1 to 24)

Legend: compact, par — scatter, par — compact, seq — scatter, seq

# Thread Affinity – Processor Binding

Binding strategies depends on machine and the app

- Putting threads far, i.e. on different packages
  - (May) improve the aggregated memory bandwidth
  - (May) improve the combined cache size
  - (May) decrease performance of synchronization constructs
- Putting threads close together, i.e. on two adjacent cores which possible share the cache
  - (May) improve performance of synchronization constructs
  - (May) decrease the available memory bandwidth and cache size (per thread)

# Thread Affinity in OpenMP* 4.0

- OpenMP 4.0 introduces the concept of places…
  - set of threads running on one or more processors
  - can be defined by the user
  - pre-defined places available:
    - threads      one place per hyper-thread
    - cores        one place exists per physical core
    - sockets      one place per processor  package

… and affinity policies…
  - spread      spread OpenMP threads evenly among the places
  - close       pack OpenMP threads near master thread
  - master      collocate OpenMP thread with master thread
- … and means to control these settings
  - Environment variables OMP_PLACES and OMP_PROC_BIND
  - clause proc_bind for parallel regions

# Thread Affinity Example

- Example (Intel® Xeon Phi™ Coprocessor): Distribute outer region, keep inner regions close

```
OMP_PLACES=cores(8); OMP_NUM_THREADS=4,4

#pragma omp parallel proc_bind(spread)
    #pragma omp parallel proc_bind(close)
```

# We're Almost Through

- OpenMP 4.0 is a major leap for OpenMP
  - New kind of parallelism has been introduced
  - Support for heterogeneous systems with coprocessor devices

- OpenMP 4.0 has more to offer!
  - Improved Fortran 2003 support
  - User-defined reductions
  - Task dependencies
  - Cancellation

- Video series including this content
  - http://software.intel.com/en-us/videos/part-1-of-5-openmp-40-for-simd-and-affinity-features-with-intel-xeon-processors-and-intel

# The last Slide...

- OpenMP 4.0 support in Intel® Compiler
  - Introduced Intel® Composer XE 2013 SP1
    - SIMD Constructs (except combined constructs)
    - OpenMP for devices (except combined constructs)
    - OpenMP Affinity
  - Feature-complete in Intel® Composer XE 2015

- Intel® Software Development Tools 2015 Beta
  - Try out: http://bit.ly/sw-dev-tools-2015-beta

# Optimization Techniques for Implicit and Explicit Vectorization

**CUG 2014, Lugano, Switzerland**

# What is SIMD?

## Scalar Code

- Executes code one element at a time.

## Vector Code

- Executes code multiple elements at a time.
- Single Instruction Multiple Data.

[Scalar] 1 elem at a time
```
addss xmm1, xmm2
```

[SSE] 4 elems at a time
```
addps xmm1, xmm2
```

[AVX] 8 elems at a time
```
vaddps ymm1, ymm2, ymm3
```

[MIC / AVX-512] 16 elems at a time
```
vaddps zmm1, zmm2, zmm3
```

# Preparing Code for SIMD



Identify Hotspots

Integer or FP?

FP

Integer

Can convert to SP?

Yes → Change to SP

No

Precision is important: impacts the SIMD width.

Re-layout data for SIMD efficiency

Align data structures

Convert code to SIMD form

Follow SIMD coding guidelines

Optimize memory access patterns and prefetch (if appropriate)

Further optimization

# Session Plan

## Hands-off

- Data Layout and Alignment

- Implicit Vectorization

- Explicit Vectorization

- SIMD Intrinsics

## Hands-on

- Explicit Vectorization with OpenMP* 4.0

# Session Plan

## Hands-off

- **Data Layout and Alignment**
- Implicit Vectorization
- Explicit Vectorization
- SIMD Intrinsics

## Hands-on

- Explicit Vectorization with OpenMP* 4.0

# Data Layout – Why It's Important

## Instruction-Level

- Hardware is optimized for contiguous loads/stores.

- Support for non-contiguous accesses differs with hardware.
  (e.g. AVX2/KNC gather)

## Memory-Level

- Contiguous memory accesses are cache-friendly.

- Number of memory streams can place pressure on prefetchers.

# Data Layout – Common Layouts

## Array-of-Structs (AoS)

| x | y | z | x | y | z |
|---|---|---|---|---|---|
| x | y | z | x | y | z |
| x | y | z | x | y | z |

- Pros:
  Good locality of {x, y, z}.
  1 memory stream.

- Cons:
  Potential for gather/scatter.

## Struct-of-Arrays (SoA)

| x | x | x | x | x | x |
|---|---|---|---|---|---|
| y | y | y | y | y | y |
| z | z | z | z | z | z |

- Pros:
  Contiguous load/store.

- Cons:
  Poor locality of {x, y, z}.
  3 memory streams.

## Hybrid (AoSoA)

| x | x | y | y | z | z |
|---|---|---|---|---|---|
| x | x | y | y | z | z |
| x | x | y | y | z | z |

- Pros:
  Contiguous load/store.
  1 memory stream.

- Cons:
  Not a "normal" layout.

# Data Alignment – Why It's Important

**Cache Line 0**

| 0 | 1 | 2 | 3 | … | … | 6 | 7 |
|---|---|---|---|---|---|---|---|

**Cache Line 1**

| 8 | 9 | … | … | … | … | … | … |
|---|---|---|---|---|---|---|---|

| 0 | 1 | 2 | 3 |
|---|---|---|---|

| 6 | 7 | 8 | 9 |
|---|---|---|---|

**Aligned Load**
- Address is aligned.
- One cache line.
- One instruction.

**Unaligned Load**
- Address is not aligned.
- Potentially multiple cache lines.
- Potentially multiple instructions.

# Data Alignment – Why It's Important

## Cache Associativity

- L1 and L2 on Knights Corner are 8-way associative.
    - L1: 8 cache lines that are 32KB/8 = 4KB apart.
    - L2: 8 cache lines that are 512KB/8 = 64KB apart.

## Set Conflicts

- Occur when references are a multiple of 4K (L1) or 64K (L2) apart.
- Look for high cache miss rate, even though working set < cache capacity.
- Solution is to pad arrays appropriately.

# Data Alignment – Sample Applications

## 1) Align Memory

- _mm_malloc(bytes, 64)        /        !dir$ attributes align:64

## 2) Access Memory in an Aligned Way

- for (i = 0; i < N; i++) { array[i] … }

## 3) Tell the Compiler

- #pragma vector aligned        /        !dir$ vector aligned
- __assume_aligned(p, 16)        /        !dir$ assume_aligned (p, 16)
- __assume(i % 16 == 0)        /        !dir$ assume (mod(i, 16) .eq. 0)

# Data Alignment – Real-life Applications



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 8 | 9 | … |   |   |   |   |   |
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |   |

☐ Data

# Data Alignment – Real-life Applications



Data

# Data Alignment – Real-life Applications



| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Data

Halo

# Data Alignment – Real-life Applications



Data

Halo

# Data Alignment – Real-life Applications



Data

Halo

Padding

Not strictly necessary…

0 1 2 3 4 5 6 7 8 9 …

# Data Alignment – Real-life Applications



Data

Halo

Padding

Not strictly necessary…

# Session Plan

## Hands-off

- Data Layout and Alignment
- **Implicit Vectorization**
- Explicit Vectorization
- SIMD Intrinsics

## Hands-on

- Explicit Vectorization with OpenMP* 4.0

# Implicit Vectorization

- Very powerful, but a compiler cannot make unsafe assumptions.

```
int* g_size;

void not_vectorizable
(float* a, float* b, float* c, int* ind) {
    for (int i = 0; i < *g_size; i++) {
        int j = ind[i];
        c[j] += a[i] + b[i];
    }
}
```

- Unsafe Assumptions:
  - a, b and c point to different arrays.
  - Value of global g_size is loop-invariant.
  - ind[i] is a one-to-one mapping.

# Implicit Vectorization

- Very powerful, but a compiler cannot make unsafe assumptions.

```
int* g_size;

void vectorizable
(float* restrict a, float* restrict b, float* restrict c, int* restrict ind) {
    int size = *g_size;
    #pragma ivdep
    for (int i = 0; i < size; i++) {
        int j = ind[i];
        c[j] += a[i] + b[i];
    }
}
```

- Safe Assumptions:
  - a, b and c point to different arrays. (restrict)
  - Value of global g_size is loop-invariant. (pointer dereference outside loop)
  - ind[i] is a one-to-one mapping. (#pragma ivdep)

# Implicit Vectorization – Improving Performance

Getting code to vectorize is only half the battle

- "LOOP WAS VECTORIZED" != "the code is optimal"
- Vectorized code can be slower than the scalar equivalent.

Compiler will <u>always</u> choose correctness over performance

- "Hints" and pragmas can't possibly cover all the situations…
- … but we can usually rewrite loop bodies to assist the compiler.

# Implicit Vectorization – Common Code Transformations

| | |
|---|---|
| Pattern: | Error checking within a vectorizable loop. |
| Issue: | Number of loop iterations unknown at compile time. |
| Original: | if (error condition) { exit } |
| Transform: | if (error condition) { error = true; } … if (error) exit |

| | |
|---|---|
| Pattern: | Memory access guarded by conditional. |
| Issue: | Compiler cannot assume memory is safe to access. |
| Original: | if (condition) a[i] += result; |
| Transform: | a[i] += (condition) ? result : 0; |

# Implicit Vectorization – Common Code Transformations

| | |
|---|---|
| **Pattern**: | Loading neighbour values, with a branch for boundary conditions. |
| **Issue(s)**: | Compiler may generate a gather;  array[position-1] may be unsafe. |
| **Original**: | double left = (column > 0) ? array[position -1] : boundary |
| **Transform**: | Pad array with appropriate halo data: avoids branch and ensures load is safe. |

| | |
|---|---|
| **Pattern**: | Loop contains OpenMP atomics, intrinsics, inline assembly |
| **Issue**: | Compiler cannot vectorize these things. |
| **Original**: | for (...) { <br> // vectorizable <br> // non-vectorizable <br> } |
| **Transform**: | Separate vectorizable and non-vectorizable code into two loops. |

# Session Plan

## Hands-off

- Data Layout and Alignment
- Implicit Vectorization
- **Explicit Vectorization**
- SIMD Intrinsics

## Hands-on

- Explicit Vectorization with OpenMP* 4.0

# Explicit Vectorization

## Compiler Responsibilities

- Allow programmer to declare that code <u>can</u> and <u>should</u> be run in SIMD.
- Generate the code the programmer asked for.

## Programmer Responsibilities

- Correctness (e.g. no dependencies, no invalid memory accesses).
- Efficiency (e.g. alignment, loop order, masking).

# Explicit Vectorization – Motivating Example 1

```
float sum = 0.0f;
float *p = a;
int step = 4;

#pragma omp simd reduction(+:sum) linear(p:step)
for (int i = 0; i < N; ++i) {
        sum += *p;
        p += step;
}
```

- The two += operators have different meaning from each other.
- The programmer should be able to express those differently.
- The compiler has to generate different code.
- The variables *i, p* and *step* have different "meaning" from each other.

# Explicit Vectorization – Motivating Example 2

```
#pragma omp declare simd simdlen(16)
uint32_t mandel(fcomplex c)
{
    uint32_t count = 1; fcomplex z = c;
    for (int32_t i = 0; i < max_iter; i += 1) {
        z = z * z + c;
        int t = cabsf(z) < 2.0f;
        count += t;
        if (!t) { break;}
    }
    return count;
}
```

- mandel() function is called from a loop over X/Y points.

- We would like to vectorize that outer loop.

- Compiler creates a vectorized function that acts on a vector of 16 c values.

# Explicit Vectorization – Performance Impact



M. Klemm, A. Duran, X. Tian, H. Saito, D. Caballero, and X. Martorell, "Extending OpenMP with Vector Constructs for Modern Multicore SIMD Architectures. In Proc. of the Intl. Workshop on OpenMP", pages 59-72, Rome, Italy, June 2012. LNCS 7312.

# Explicit Vectorization – Array Notation

"Long Form"

```
C[0:N] = A[0:N] + B[0:N];
D[0:N] = C[0:N] * C[0:N];
```

"Short Form"

```
for (i = 0; i < N; i += V) {
    C[i:V] = A[i:V] + B[i:V];
    D[i:V] = C[i:V] * C[i:V];
}
```

- Long form is more elegant, but short form is better for performance.
- Imagine each array assignment as a for loop:
  - For large N, long form will not keep C[0:N] in cache.
  - For appropriate V, short form keeps C[i:V] in registers.

- Some differences between Intel® Cilk™ Plus notation (C/C++) and Fortran 90.

# Explicit Vectorization – OpenMP* SIMD Loops

#pragma omp simd / !$omp simd => for/do loop is a SIMD loop.

| | |
|---|---|
| safelen (*length*) | Maximum distance between two iterations executed concurrently by a SIMD instruction. |
| linear (*list[:linear-step]*) | List items are private and have a linear relationship with respect to the iteration space. |
| aligned (*list[:alignment]*) | List items are aligned to a platform-dependent value (or the value of the optional parameter). |

private (*list*), lastprivate (*list*), reduction (*reduction-identifier:list*) and collapse (*n*) are also supported, with functionality matching that of omp for.

# Explicit Vectorization – OpenMP* SIMD Functions

#pragma omp declare simd / !$omp declare simd => function will be called from a SIMD loop.

| | |
|---|---|
| simdlen (*length*) | Maximum number of concurrent arguments to the function (i.e. maximum SIMD width). |
| uniform (*argument-list*) | List items have the same value for all SIMD lanes, and can therefore be broadcast. |
| inbranch<br>notinbranch | Function always called inside a conditional.<br>Function never  called inside a conditional. |

linear (*argument-list[:linear-step]*) and aligned (*argument-list[:alignment]*) are also supported, with functionality matching that of omp simd.

# Implicit vs Explicit Vectorization – Important Differences

## Implicit

- Automatic dependency analysis.
  (e.g. recognises SIMD reductions)

- Recognizes idioms with data dependencies.
  (e.g. array[i++] = x; -> vcompress)

- Non-inline functions will be scalarized.

- Limited support for outer-loop vectorization
  (only with –O3).

## Explicit

- No dependency analysis.
  (e.g. SIMD reductions <u>must</u> be declared)

- Recognizes idioms without data dependencies.

- Non-inline functions can be vectorized.

- Outer loops can be vectorized.

# Session Plan

## Hands-off

- Data Layout and Alignment
- Implicit Vectorization
- Explicit Vectorization
- **SIMD Intrinsics**

## Hands-on

- Explicit Vectorization with OpenMP* 4.0

# SIMD Intrinsics – An Introduction

## Intrinsic Functions

- Substituted by compiler for specific instruction(s).
- Higher level of abstraction than assembly.

## SIMD Intrinsics

- Explicit vectorization at an instruction level.
- Direct manipulation of SIMD and mask registers.

```
// Scalar code for vector add.
for (int i = 0; i < N; i++) {
    c[i] = a[i] + b[i];
}

// SIMD intrinsics for vector add.
for (int i = 0; i < N; i += 8) {
    __m512 ai = _mm512_load_pd(&a[i]);
    __m512 bi = _mm512_load_pd(&b[i]);
    __m512 ci = _mm512_add_pd(ai, bi);
    _mm512_store_pd(&c[i], ci);
}
```

# SIMD Intrinsics – When and How?

## When to Use Intrinsics

- You think you can beat the compiler (after asm inspection).
- Your code absolutely needs to run as fast as possible.
- Nothing else works -- intrinsics should be your last resort!

## Common Use-cases

- Utilizing particular instructions (e.g. bit-manipulation, cryptography).
- "Horizontal" vectorization (e.g. computing dot products of SIMD registers).
- Exploiting memory layout knowledge (e.g. AoS to SoA transpose).

# SIMD Intrinsics – Intel® Intrinsics Guide



Filter by ISA.

Filter by functionality.

Expand any intrinsic for a detailed description.

Available at: http://software.intel.com/sites/landingpage/IntrinsicsGuide/

# SIMD Intrinsics – Performance Impact

## "Ninja" Performance Gap

- The gap between naïve code and the <u>best</u>-optimized code.
- Average speed-up of 24x on Sandy Bridge and Knights Ferry.

## Closing the Gap

- Apply well-known algorithmic techniques (e.g. cache blocking).
- Use a recent compiler, and explicit vectorization features.
- Average speed-up of only ~1.3x from intrinsics.

N. Satish, C. Kim, J. Chhugani, H. Saito, R. Krishnaiyer, M. Smelyanskiy, M. Girkar and P. Dubey, "Can Traditional Programming Bridge the Ninja Performance Gap for Parallel Computing Applications?", in Proceedings of the International Symposium on Computer Architecture (ISCA), Portland, OR, 2012

# Hands-on:
# Explicit Vectorization with OpenMP* 4.0

**Accelerating N-Body Codes with SIMD Instructions**

# Optimization Techniques for Implicit and Explicit Vectorization

## Summary

- Intel® Composer XE provides many alternative methods for vectorization.

- Generating vector code is <u>not</u> complicated:
    - Compiler reports / tools help with auto-vec.
    - Explicit vectorization an industry standard (OpenMP* 4.0).


- Extensions to existing languages allow you to target vector architectures while keeping your source code readable, maintainable, and familiar.

# Legal Disclaimers

# Legal Disclaimers

## 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 SSE3 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

# Legal Disclaimers

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. For more information go to http://www.intel.com/performance.

Estimated Results Benchmark Disclaimer:
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.

Software Source Code Disclaimer:
Any software source code reprinted in this document is furnished under a software license and may only be used or copied in accordance with the terms of that license.

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

(intel)

# Legal Disclaimers

# Optimizing for MPI*/OpenMP* on Intel® Xeon Phi™ Coprocessors

**CUG 2014, Lugano, Switzerland**

# Frequently Asked Questions

# Levels of Parallelism

| | |
|---|---|
| Cluster | Group of computers<br>communicating through fast interconnect |
| Coprocessors/Accelerators | Special compute devices<br>attached to the local node through special interconnect |
| Node | Group of processors<br>communicating through shared memory |
| Socket | Group of cores<br>communicating through shared cache |
| Core | Group of functional units<br>communicating through registers |
| Hyper-Threads | Group of thread contexts sharing functional units |
| Superscalar | Group of instructions sharing functional units |
| Pipeline | Sequence of instructions sharing functional units |
| Vector | Single instruction using multiple functional units |

(intel)

# Levels of Parallelism in OpenMP* 4.0

| | |
|---|---|
| Cluster | Group of computers communicating through fast interconnect |
| Coprocessors/Accelerators | Special compute devices attached to the local node through special interconnect |
| Node | Group of processors communicating through shared memory |
| Socket | Group of cores communicating through shared cache |
| Core | Group of functional units communicating through registers |
| Hyper-Threads | Group of thread contexts sharing functional units |
| Superscalar | Group of instructions sharing functional units |
| Pipeline | Sequence of instructions sharing functional units |
| Vector | Single instruction using multiple functional units |

(intel)

# Levels of Parallelism in MPI* + OpenMP* 4.0

| | |
|---|---|
| **Cluster** | Group of computers communicating through fast interconnect |
| **Coprocessors/Accelerators** | Special compute devices attached to the local node through special interconnect |
| **Node** | Group of processors communicating through shared memory |
| **Socket** | Group of cores communicating through shared cache |
| **Core** | Group of functional units communicating through registers |
| **Hyper-Threads** | Group of thread contexts sharing functional units |
| Superscalar | Group of instructions sharing functional units |
| Pipeline | Sequence of instructions sharing functional units |
| **Vector** | Single instruction using multiple functional units |

(intel)

# Summary

## Compilation

- Single toolchain for CPU and MIC; can often just add "-mmic".
- Support for industry standards (including MPI* and OpenMP*).

## Optimization

- Biggest benefits from tuning memory and vector behaviour.
- Straightforward parallel tuning techniques still apply.
- Intel tools can help to identify optimization opportunities (and are always improving).
- "Dual-Transforming-Tuning Advantage"; optimize for processor <u>and</u> coprocessor.

# Legal Disclaimers

# Legal Disclaimers

(intel)

# Legal Disclaimers

(intel)

# Legal Disclaimers

The above statements and any others in this document that refer to plans and expectations for the third quarter, the year and the future are forward-looking statements that involve a number of risks and uncertainties. Words such as "anticipates," "expects," "intends," "plans," "believes," "seeks," "estimates," "may," "will," "should" and their variations identify forward-looking statements. Statements that refer to or are based on projections, uncertain events or assumptions also identify forward-looking statements. Many factors could affect Intel's actual results, and variances from Intel's current expectations regarding such factors could cause actual results to differ materially from those expressed in these forward-looking statements. Intel presently considers the following to be the important factors that could cause actual results to differ materially from the company's expectations. Demand could be different from Intel's expectations due to factors including changes in business and economic conditions; customer acceptance of Intel's and competitors' products; supply constraints and other disruptions affecting customers; changes in customer order patterns including order cancellations; and changes in the level of inventory at customers. Uncertainty in global economic and financial conditions poses a risk that consumers and businesses may defer purchases in response to negative financial events, which could negatively affect product demand and other related matters.  Intel operates in intensely competitive industries that are characterized by a high percentage of costs that are fixed or difficult to reduce in the short term and product demand that is highly variable and difficult to forecast. Revenue and the gross margin percentage are affected by the timing of Intel product introductions and the demand for and market acceptance of Intel's products; actions taken by Intel's competitors, including product offerings and introductions, marketing programs and pricing pressures and Intel's response to such actions; and Intel's ability to respond quickly to technological developments and to incorporate new features into its products. The gross margin percentage could vary significantly from expectations based on capacity utilization; variations in inventory valuation, including variations related to the timing of qualifying products for sale; changes in revenue levels; segment product mix; the timing and execution of the manufacturing ramp and associated costs; start-up costs; excess or obsolete inventory; changes in unit costs; defects or disruptions in the supply of materials or resources; product manufacturing quality/yields; and impairments of long-lived assets, including manufacturing, assembly/test and intangible assets.  Intel's results could be affected by adverse economic, social, political and physical/infrastructure conditions in countries where Intel, its customers or its suppliers operate, including military conflict and other security risks, natural disasters, infrastructure disruptions, health concerns and fluctuations in currency exchange rates. Expenses, particularly certain marketing and compensation expenses, as well as restructuring and asset impairment charges, vary depending on the level of demand for Intel's products and the level of revenue and profits. Intel's results could be affected by the timing of closing of acquisitions and divestitures. Intel's results could be affected by adverse effects associated with product defects and errata (deviations from published specifications), and by litigation or regulatory matters involving intellectual property, stockholder, consumer, antitrust, disclosure and other issues, such as the litigation and regulatory matters described in Intel's SEC reports. An unfavorable ruling could include monetary damages or an injunction prohibiting Intel from manufacturing or selling one or more products, precluding particular business practices, impacting Intel's ability to design its products, or requiring other remedies such as compulsory licensing of intellectual property. A detailed discussion of these and other factors that could affect Intel's results is included in Intel's SEC filings, including the company's most recent reports on Form 10-Q, Form 10-K and earnings release.

Rev. 7/17/13

# Intel® Xeon Phi™ Product Family
Intel ® ITAC

Heinrich Bockhorst, May 5th 2014

Software and Services Group
Intel Corporation

# Introduction – What is Tracing?

- Record program execution
  - Program events such as function enter/exit, communication

- 1:1 protocol of the actual program execution
  - Sampling gathers statistical information

- Accurate data

- Easily get loads of data

# Event based approach

- Event = time stamp + thread ID + description
  - Function entry/exit
  - Messages
  - Collective operations
  - Counter samples
- Strengths:
  - Predict detailed program behavior
  - Record exact sequence of program states – keep timing consistent
  - Collect information about exchange of messages: at what times and in which order
  - Detect temporal dependencies

# Intel® Trace Analyzer and Collector



Compare the event timelines of two communication profiles

Blue = computation
Red = communication

Chart showing how the MPI processes interact

# Intel® Trace Analyzer and Collector Overview

- **Intel® Trace Analyzer and Collector helps the developer:**
  - Visualize and understand parallel application behavior
  - Evaluate profiling statistics and load balancing
  - Identify communication hotspots

- **Features**
  - Event-based approach
  - Low overhead
  - Excellent scalability
  - Comparison of multiple profiles
  - Powerful aggregation and filtering functions
  - Fail-safe MPI tracing
  - Provides API to instrument user code
  - MPI correctness checking
  - Idealizer

# Profiles: Flat Function Profile

- Statistics about functions

Optimization Notice

# Timelines: Event Timeline

- Get impression of program structure
- Display functions, messages and collective operations for each process/thread along time-axis
- Retrieval of detailed event information

Optimization Notice

# Communication Profiles

- Statistics about point-to-point or collective communication
- Generic matrix supports grouping by several attributes in each dimension
  Sender, Receiver, Data volume per msg, Tag, Communicator, Type
- Available attributes: Count, Bytes transferred, Time, Transfer rate

# View - zooming

# Aggregation Example

Optimization Notice

# Full ITAC Functionality on Intel® Xeon Phi™

# ITAC Prerequisites

- Set ITAC environment (per user)

  `# source /opt/intel/itac/8.1.2.033/intel64/bin/itacvars.sh impi4`

  – Identical for host and the coprocessor

- No NFS: Upload ITAC library manually

  `# sudo scp /opt/intel/itac/8.1.2.033/mic/slib/libVT.so host-mic0:/lib64/`

# ITAC Usage with Intel® Xeon Phi™ Coprocessor

- Run with –trace flag (without linkage) to create a trace file
  - MPI+Offload
    ```
    # mpirun –trace -n 2 ./test
    ```
  - Coprocessor only
    ```
    # mpirun –trace -n 2 -wdir /tmp
        -host host-mic0 ./test_hello.MIC
    ```
  - Symmetric
    ```
    # mpirun –trace -n 2 -host michost./test_hello :
        -wdir /tmp -n 2 -host host-mic0
        ./test_hello.MIC
    ```
- Flag „-trace" will implicitly pre-load libVT.so (which finally calls libmpi.so to execute the MPI call)
- Set `VT_LOGFILE_FORMAT=stfsingle` to create a single trace

# ITAC Usage with Intel® Xeon Phi™: Compilation Support

- Compile and link with „–trace" flag

  ```
  # mpiicc -trace -o test_hello test.c

  # mpiicc –trace –mmic -o test_hello.MIC test.c
  ```

  - Linkage of libVT library

- Compile with –tcollect flag

  ```
  # mpiicc –tcollect -o test_hello test.c

  # mpiicc –tcollect –mmic -o test_hello.MIC test.c
  ```

  - Linkage of libVT library
  - Will do a full instrumentation of your code, i.e. All user functions will be visible in the trace file
  - Maximal insight, but also maximal overhead

- Use the VT API of ITAC to manually instrument your code.

- Run Intel® MPI program as usual without „-trace" flag

  ```
  # mpirun ...
  ```

# Improving Load Balance: Real World Case



Collapsed data per node and coprocessor

Host
16 MPI procs x
1 OpenMP thread

Coprocessor
8 MPI procs x
28 OpenMP threads

Too high load on Host
= too low load on Coprosseor

# Improving Load Balance: Real World Case



Collapsed data per node and coprocessor

Host
16 MPI procs x
1 OpenMP thread

Coprocessor
24 MPI procs x
8 OpenMP threads

Too low load on Host
= too high load on Coprocessor

# Improving Load Balance: Real World Case



Collapsed data per node and coprocessor

<u>Host</u>
16 MPI procs x
1 OpenMP thread

<u>Coprocessor</u>
16 MPI procs x
12 OpenMP thrds

Perfect balance
Host load = Coprocessor load

# Ideal Interconnect Simulator (Idealizer)

- What is the Ideal Interconnect Simulator?
  - Using a ITAC trace of an MPI application, simulate it under ideal conditions
    - Zero network latency
    - Infinite network bandwidth
    - Zero MPI buffer copy time
    - Infinite MPI buffer size
  - Only limiting factors are concurrency rules, e.g.,
    - A message can not be received before it is sent
    - An All-to-All collective may end only when the last thread starts

# Ideal Interconnect Simulator (Idealizer)



**Actual trace**

**Idealized Trace**

# Building Blocks: Elementary Messages



Early Send / Late Receive

zero duration

P1  MPI_Isend

P2  MPI_Recv PI_Recv

zero duration

Late Send / Early Receive

zero duration

P1  MPI_Isend

P2  MPI_Recv

Load imbalance

# Building Blocks: Collective Operations



Actual trace (Gigabit Ethernet)

Simulated trace (Ideal interconnect)

Same MPI_Alltoallv

Legend:
257 = MPI_Alltoallv
506 = User_Code

Same timescale in both figures

# Application Imbalance Diagram: Total

# Application Imbalance Diagram: Breakdown

**Intel® Xeon Phi™ Coprocessor**

**Software & Services Group, Developer Products Division**

# Online Resources

- Intel® MPI Library product page

  www.intel.com/go/mpi

- Intel® Trace Analyzer and Collector product page

  www.intel.com/go/traceanalyzer

- Intel® Clusters and HPC Technology forums

  http://software.intel.com/en-us/forums/intel-clusters-and-hpc-technology/

- Intel® Xeon Phi™ Coprocessor Developer Community

  http://software.intel.com/en-us/mic-developer

# Summary

- The ease of use of Intel® MPI and related tools like the Intel Trace Analyzer and Collector extends from the Intel Xeon architecture to the Intel® Xeon Phi™ coprocessor.

Optimization Notice

25

# Legal Disclaimer & Optimization Notice

**Intel® Xeon Phi™ Coprocessor**

**Software & Services Group, Developer Products Division**

Optimization Notice

# Agenda

Start tuning on host

Overview of Intel® VTune™ Amplifier XE

Efficiency metrics

Problem areas

Optimization
Notice

(intel)

# Performance Analysis Methodology
## Optimization: A Top-down Approach

- Use top down approach
- Understand application and system characteristics
  - Use appropriate tools at each level

| System Config, BIOS, OS, Network I/O, Disk I/O, Database Tuning, etc. | — **System** |

| Application Design Algorithmic Tuning Driver Tuning Parallelization | — **Application** |

| Cache/Memory Instructions SIMD others | — **Processor** |

Optimization Notice
(intel)

# Performance Analysis Methodology
## Optimization: A Top-down Approach

- Use top down approach
- Understand application and system characteristics
  - Use appropriate tools at each level

System Config, BIOS, OS, Network I/O, Disk I/O, Database Tuning, etc.

System

Application Design
Algorithmic Tuning
Driver Tuning
Parallelization

Application

Cache/Memory
Instructions
SIMD
others

Processor

VTune™ Amplifier XE can help here

Optimization
Notice

(intel)

# Start with host-based profiling to identify vectorization/ parallelism/ offload candidates

Start with representative/reasonable workloads!

Use Intel® VTune™ Amplifier XE to gather hot spot data

- Tells what functions account for most of the run time

- Often, this is enough

  - But it does not tell you much about program structure

Optimization Notice

# Start with **host-based** profiling to identify vectorization/ parallelism/ offload candidates

Start with representative/reasonable workloads!

Use Intel® VTune™ Amplifier XE to gather hot spot data

- Tells what functions account for most of the run time

- Often, this is enough

  - But it does not tell you much about program structure

Alternately, profile functions & loops using Intel® Composer XE
- Build with options

  `-profile-functions -profile-loops=all -profile-loops-report=2`
- Run the code (which may run slower) to collect profile data
- Look at the resulting `dump` files, or open the `xml` file with the data viewer `loopprofileviewer.sh` located in the compiler `./bin` directory
- Tells you
  - which loops and functions account for the most run time
  - how many times each loop executes (min, max and average)

Optimization Notice

# Correctness/Performance Analysis of Parallel code

Intel® Inspector XE and thread-reports in VTune™ Amplifier XE are not available on the Intel® Xeon Phi™ coprocessor

So…

Optimization Notice

# Correctness/Performance Analysis of Parallel code

Intel® Inspector XE and thread-reports in VTune™ Amplifier XE are not available on the Intel® Xeon Phi™ coprocessor

So…

- Use Intel Inspector XE on your code with **offload disabled** (on host) to identify correctness errors (e.g., deadlocks, races)
  - Once fixed, then enable offload and continue debugging on the coprocessor

Optimization Notice

(intel)

# Correctness/Performance Analysis of Parallel code

Intel® Inspector XE and thread-reports in VTune™ Amplifier XE are not available on the Intel® Xeon Phi™ coprocessor

So…

- Use Intel Inspector XE on your code with **offload disabled** (on host) to identify correctness errors (e.g., deadlocks, races)
  - Once fixed, then enable offload and continue debugging on the coprocessor
- Use VTune Amplifier XE's parallel performance analysis tools to find issues on the host by running your program with **offload disabled**
  - Fix everything you can
  - Then study scaling on the coprocessor using lessons from host tuning to further optimize parallel performance
    - Be wary of synchronization across more than a handful of threads
    - Pay attention to load balance.

Optimization Notice

(intel)

# Agenda

Start tuning on host

Overview of Intel® VTune™ Amplifier XE

Efficiency metrics

Problem areas

Optimization Notice

(intel)

# Intel® VTune™ Amplifier XE
## Tune Applications for Scalable Multicore Performance

- **Fast, Accurate Performance Profiles**
  - Hotspot (Statistical call tree)
  - Hardware-Event Based Sampling
- **Thread Profiling**
  - Visualize thread interactions on timeline
  - Balance workloads
- **Easy set-up**
  - Pre-defined performance profiles
  - Use a normal production build
- **Compatible**
  - Microsoft*, GCC*, Intel compilers
  - C/C++, Fortran, Assembly, .NET*
  - Latest Intel processors and compatible processors[1]
- **Find Answers Fast**
  - Filter out extraneous data
  - View results tied to source/assembly lines
  - Event multiplexing
- **Windows* or Linux***
  - Visual Studio* Integration (Windows)
  - Standalone user interface and command line
  - 32 and 64-bit

[1] IA-32 and Intel® 64 architectures.
Many features work with compatible processors.
Event based sampling requires a genuine Intel Processor.

Optimization Notice

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance



Result Display Tabs

Optimization Notice

# VTune™ Amplifier XE visualizes performance



Result Analysis Type

# VTune™ Amplifier XE visualizes performance



Result Viewpoint

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance



Result Components

# VTune™ Amplifier XE visualizes performance



Grid Pane

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance



Source View /
Per line localization

Optimization
Notice

(intel)

# VTune™ Amplifier XE visualizes performance



Source View / View / Hot spot Navigation controls

# VTune™ Amplifier XE visualizes performance

# VTune™ Amplifier XE visualizes performance

# For event collection the coprocessor is treated as a special HW architecture

# Project properties provides the means to invoke data collection by target type

# Launch Application serves many uses, from host/offload to native execution

# Search directories have been reorganized to speed symbol resolution during finalization

- Enumerate source directories under this tab



- Put library paths here

Notable coprocessor library paths:
/lib/firmware/mic
/usr/linux-k1om-4.7/linux-k1om/lib64
/opt/intel/composerxe/lib/mic
/opt/intel/composerxe/tbb/lib/mic
/opt/intel/composerxe/mkl/lib/mic
/opt/intel/mpi-rt/4.1.0/mic

Optimization Notice

(intel)

# General Exploration runs a set of events to drive top-down analysis

# Agenda

Start tuning on host

Overview of Intel® VTune™ Amplifier XE

Efficiency metrics

Problem areas

Optimization
Notice

(intel)

# Cycles Per Instruction (CPI), a standard measure, has some special kinks

- Threads on each Intel® Xeon™ Phi core share a clock
  - If all 4 HW threads are active, each gets ¼ total cycles
- Multi-stage instruction decode requires two threads to utilize the whole core – one thread only gets half
- With two ops/per cycle (U-V-pipe dual issue):

| Threads per Core | | Best CPI per Core | Best CPI per Thread |
|---|---|---|---|
| 1 | x | 1.0 | = 1.0 |
| 2 | x | 0.5 | = 1.0 |
| 3 | x | 0.5 | = 1.5 |
| 4 | x | 0.5 | = 2.0 |

- To get thread CPI, multiply by the active threads

Optimization Notice

(intel)

# As an efficiency metric, CPI must be considered carefully: it IS a ratio

- Changes in CPI absent major code changes can indicate general latency gains/losses

| Metric | Formula | Investigate if |
|---|---|---|
| CPI per Thread | CPU_CLK_UNHALTED/ INSTRUCTIONS_EXECUTED | > 4.0, or increasing |
| CPI per Core | (CPI per Thread) / Number of hardware threads used | > 1.0, or increasing |

- Note the effect on CPI from applied optimizations
- Reduce high CPI through optimizations that target latency
  - Better prefetch
  - Increase data reuse through better blocking

Optimization Notice

(intel)

# Two more examples why absolute CPI value is less important than changes

- Scaling data from a typical lab workload:

| Metric | 1 hardware thread / core | 2 hardware threads / core | 3 hardware threads / core | 4 hardware threads / core |
|---|---|---|---|---|
| CPI per Thread | 5.24 | 8.80 | 11.18 | 13.74 |
| CPI per Core | 5.24 | 4.40 | 3.73 | 3.43 |

- Observed CPIs from several tuned workloads:

Optimization Notice

# Efficiency Metric: Compute to Data Access Ratio

- Measures an application's computational density, and suitability for Intel® Xeon Phi™ coprocessors

| Metric | Formula | Investigate if |
|---|---|---|
| Vectorization Intensity | VPU_ELEMENTS_ACTIVE / VPU_INSTRUCTIONS_EXECUTED | |
| L1 Compute to Data Access Ratio | VPU_ELEMENTS_ACTIVE / DATA_READ_OR_WRITE | < Vectorization Intensity |
| L2 Compute to Data Access Ratio | VPU_ELEMENTS_ACTIVE / DATA_READ_MISS_OR_ WRITE_MISS | < 100x L1 Compute to Data Access Ratio |

- Increase computational density through vectorization and reducing data access (see cache issues, also, DATA ALIGNMENT!)

Optimization Notice

(intel)

# Agenda

Start tuning on host

Overview of Intel® VTune™ Amplifier XE

Efficiency metrics

Problem areas*

*tuning suggestions requiring deeper understanding of architectural tradeoffs and application data handling details are highlighted with this "ninja" notation

Optimization Notice

(intel)

# Problem Area: L1 Cache Usage

- Significantly affects data access latency and therefore application performance

| Metric | Formula | Investigate if |
|--------|---------|----------------|
| L1 Misses | DATA_READ_MISS_OR_WRITE_MISS + L1_DATA_HIT_INFLIGHT_PF1 | |
| L1 Hit Rate | (DATA_READ_OR_WRITE – L1 Misses) / DATA_READ_OR_WRITE | < 95% |

- Tuning Suggestions:
  - Software prefetching
  - Tile/block data access for cache size
  - Use streaming stores

  If using 4K access stride, may be experiencing conflict misses

  Examine Compiler prefetching (Compiler-generated L1 prefetches should not miss)

40

Optimization Notice

(intel)

# Problem Area: Data Access Latency

- Significantly affects application performance

| Metric | Formula | Investigate if |
|--------|---------|----------------|
| Estimated Latency Impact | (CPU_CLK_UNHALTED<br> – EXEC_STAGE_CYCLES<br> – DATA_READ_OR_WRITE)<br>   / DATA_READ_OR_WRITE_MISS | >145 |

- Tuning Suggestions:
  - Software prefetching
  - Tile/block data access for cache size
  - Use streaming stores

Check cache locality – turn off prefetching and use CACHE_FILL events - reduce sharing if needed/possible

If using 64K access stride, may be experiencing conflict misses

Optimization Notice

# Problem Area: TLB Usage

- Also affects data access latency and therefore application performance

| Metric | Formula | Investi-gate if: |
|--------|---------|------------------|
| L1 TLB miss ratio | DATA_PAGE_WALK/DATA_READ_OR_WRITE | > 1% |
| L2 TLB miss ratio | LONG_DATA_PAGE_WALK / DATA_READ_OR_WRITE | > .1% |
| L1 TLB misses per L2 TLB miss | DATA_PAGE_WALK / LONG_DATA_PAGE_WALK | > 100x |

- Tuning Suggestions:
  - Improve cache usage & data access latency
  - If L1 TLB miss/L2 TLB miss is high, try using large pages
  - For loops with multiple streams, try splitting into multiple loops
  - If data access stride is a large power of 2, consider padding between arrays by one 4 KB page

Optimization Notice

(intel)

# Problem Area: VPU Usage

- Indicates whether an application is vectorized successfully and efficiently

| Metric | Formula | Investigate if |
|--------|---------|----------------|
| Vectorization Intensity | VPU_ELEMENTS_ACTIVE / VPU_INSTRUCTIONS_EXECUTED | <8 (DP), <16(SP) |

- Tuning Suggestions:
  - Use the Compiler vectorization report!
  - For data dependencies preventing vectorization, try using Intel® Cilk™ Plus #pragma SIMD (if safe!)
  - Align data and tell the Compiler!
  - Restructure code if possible: Array notations, AOS->SOA

Optimization Notice 📖

# Problem Area: Memory Bandwidth

- Can increase data latency in the system or become a performance bottleneck

| Metric | Formula | Investigate if |
|---|---|---|
| Memory Bandwidth | (UNC_F_CH0_NORMAL_READ + UNC_F_CH0_NORMAL_WRITE+ UNC_F_CH1_NORMAL_READ + UNC_F_CH1_NORMAL_WRITE) X 64/time | < 80GB/sec (practical peak 140GB/sec)<br><br>(with 8 memory controllers) |

- Tuning Suggestions:
  - Improve locality in caches
  - Use streaming stores
  - Improve software prefetching

44

Optimization Notice

(intel)

# Final caution: coprocessor collections can generate dense volumes of data
Example: DGEMM on 60+ cores



Tip: Use a CPU Mask to reduce data volume while maintaining equivalent accuracy.

Optimization Notice

# Summary

- Vectorization, Parallelism, and Data locality are critical to good performance for the Intel® Xeon Phi™ Coprocessor

- Event names can be misleading – we recommend using the metrics given in this presentation or our tuning guide at http://software.intel.com/en-us/articles/optimization-and-performance-tuning-for-intel-xeon-phi-coprocessors-part-2-understanding

- Intel® VTune™ Amplifier XE supports collecting all of the above metrics, as well as providing special analysis types like General Exploration and Memory Bandwidth

Optimization
Notice

(intel)

# Please return your evaluation forms!

Optimization
Notice

(intel)

# 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