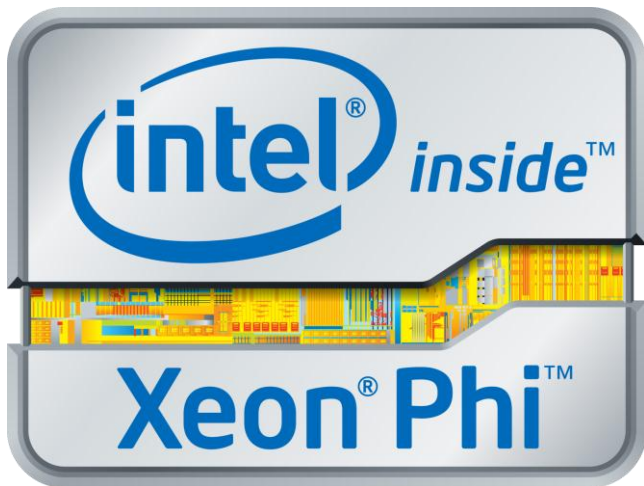# MPI Programming on the Intel® Xeon Phi™ Coprocessor

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

**Intel® Xeon Phi™ Architecture**

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

# 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.

Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance. Buyers should consult other sources of information to evaluate the performance of systems or components they are considering purchasing. For more information on performance tests and on the performance of Intel products, reference www.intel.com/software/products.

| 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.<br><br>Notice revision #20110804 |

# Agenda

- Overview
- Programming Models
- Hybrid Computing
- Load Balancing

**Intel® Xeon Phi™ Architecture**

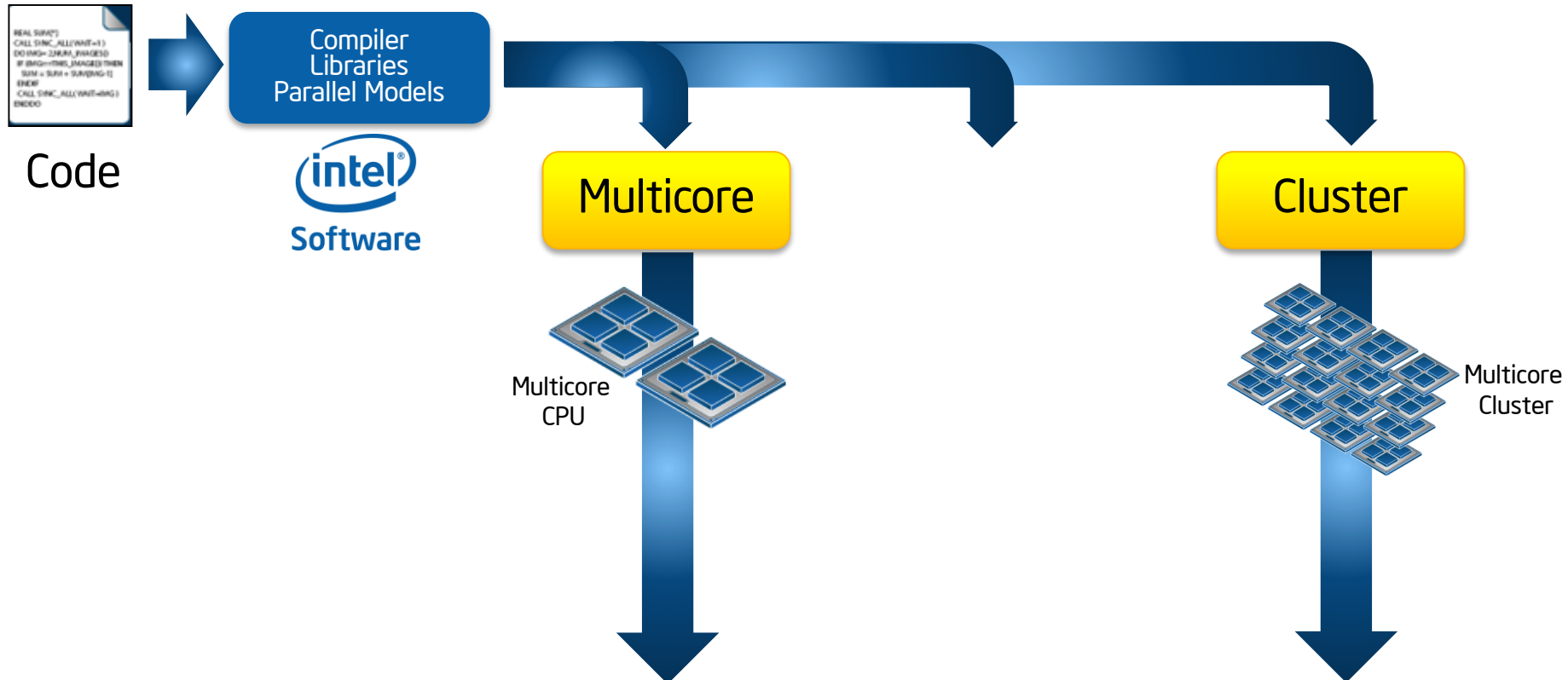**Software & Services Group, Developer Relations Division**

3

# Agenda

- Overview
- Programming Models
- Hybrid Computing
- Load Balancing

# 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

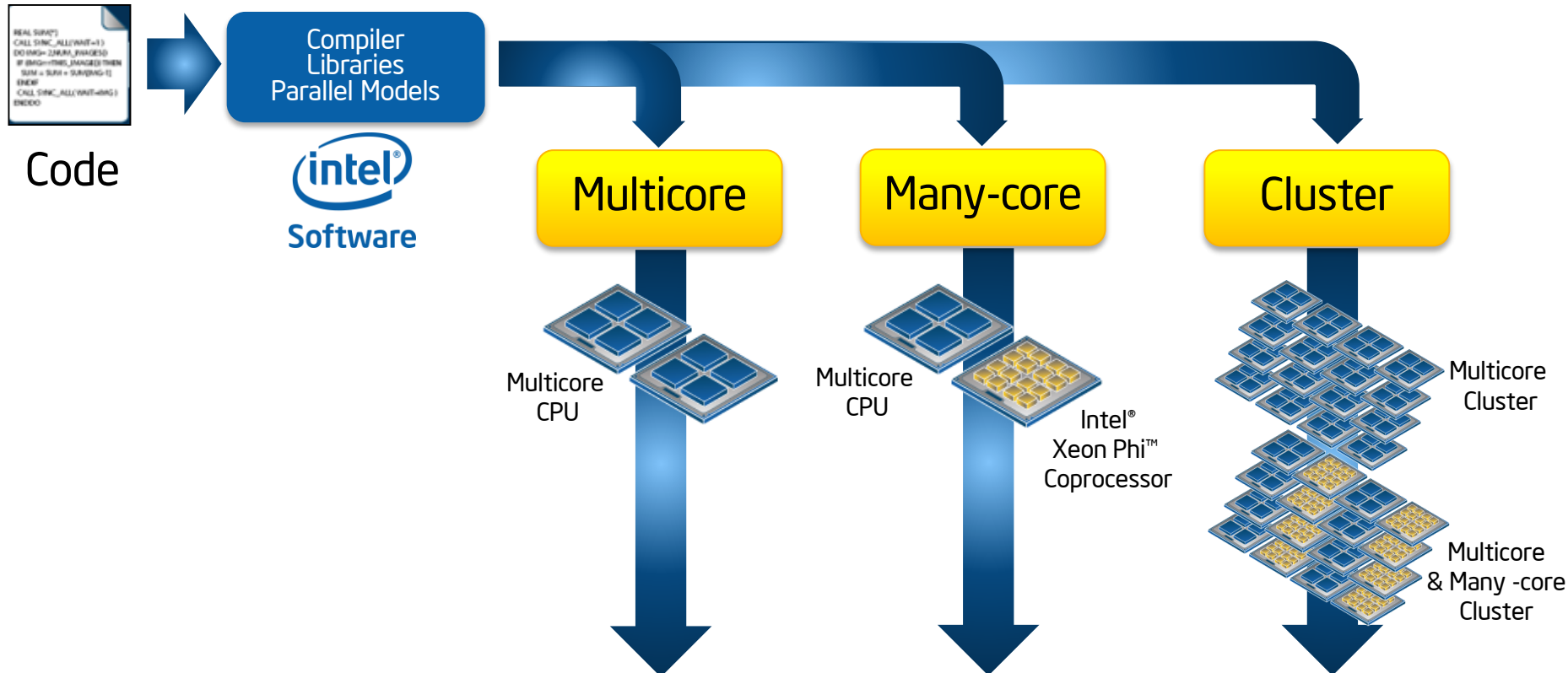Multicore CPU

Cluster

Multicore Cluster

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

# Enabling & Advancing Parallelism
## High Performance Parallel Programming

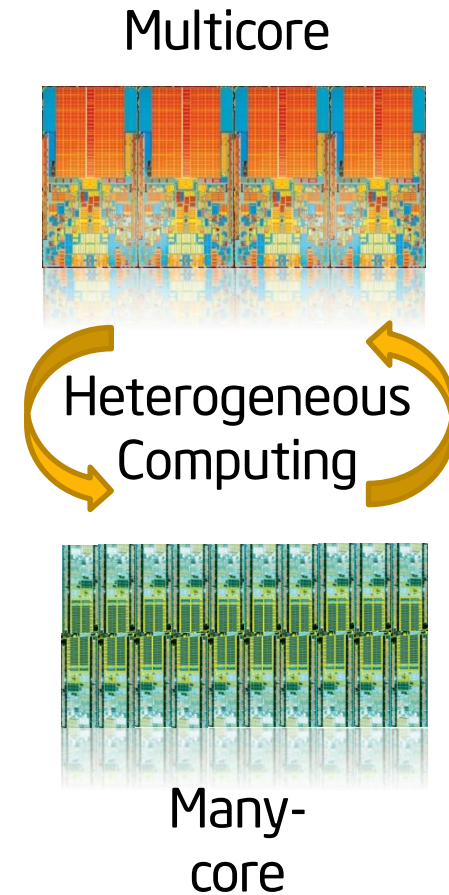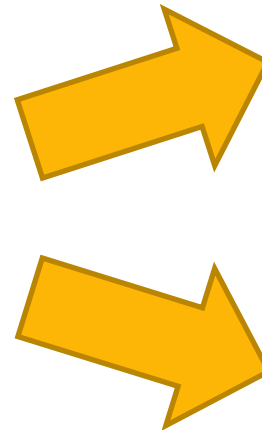**Intel tools, libraries and parallel models extend to multicore, many-core and heterogeneous computing**

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.**

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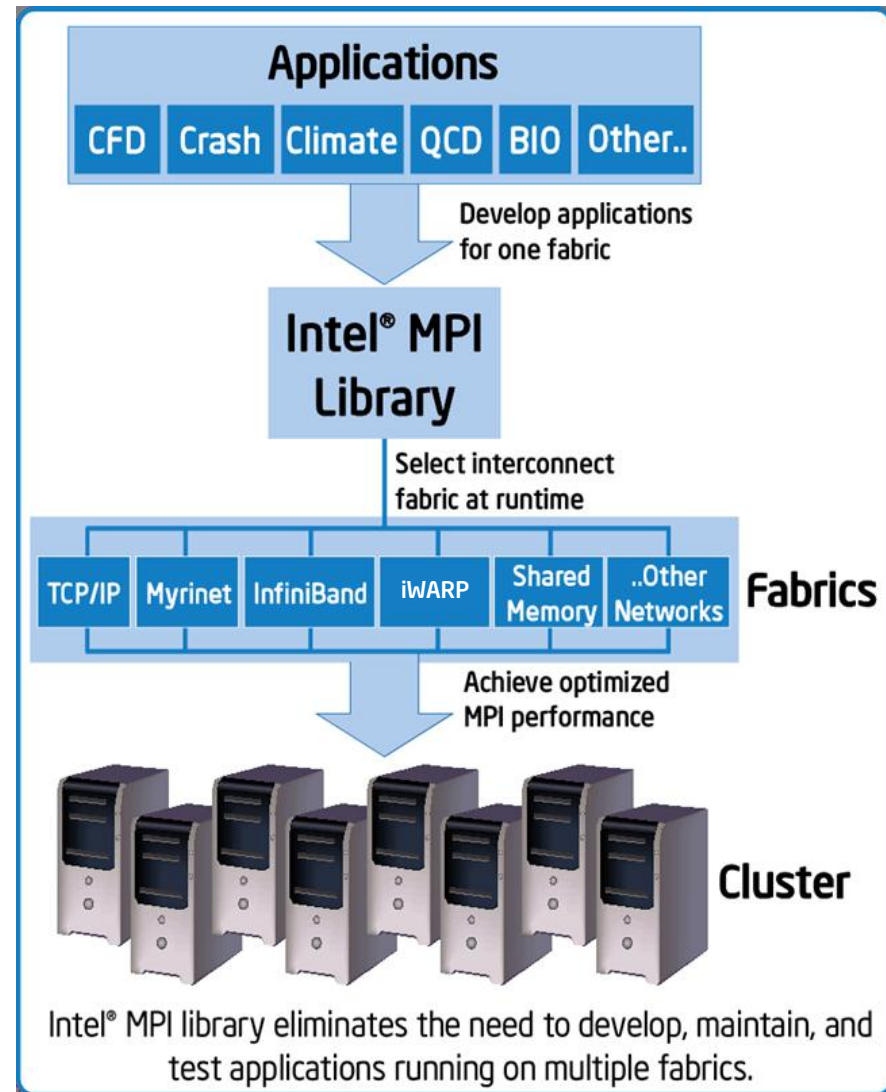## Common Tools and Programming Models for Parallelism

**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
- Lower 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



Applications
CFD | Crash | Climate | QCD | BIO | Other..

Develop applications for one fabric

Intel® MPI Library

Select interconnect fabric at runtime

TCP/IP | Myrinet | InfiniBand | iWARP | Shared Memory | ..Other Networks — Fabrics

Achieve optimized MPI performance

Cluster

Intel® MPI library eliminates the need to develop, maintain, and test applications running on multiple fabrics.
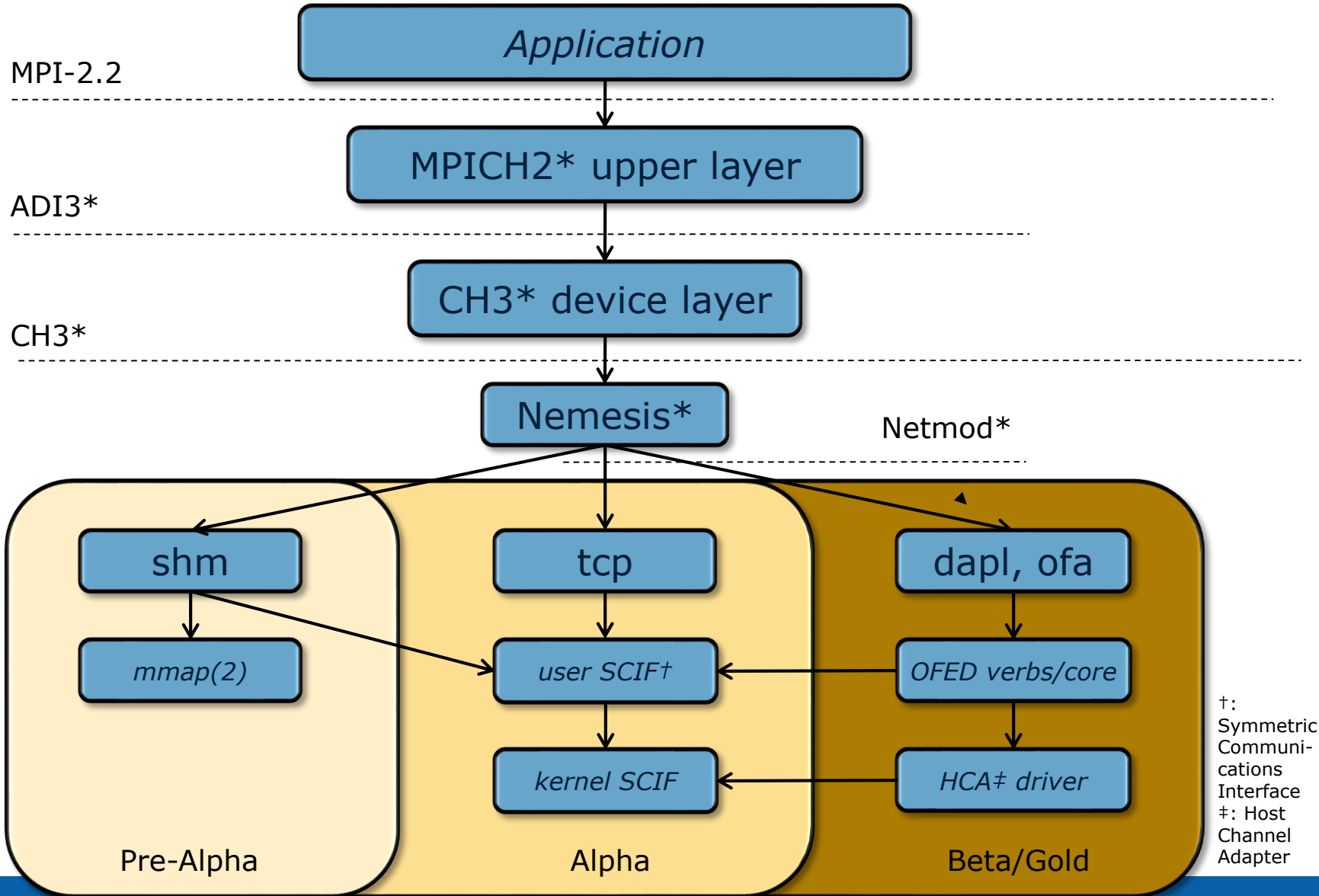
# Levels of communication speed

- 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
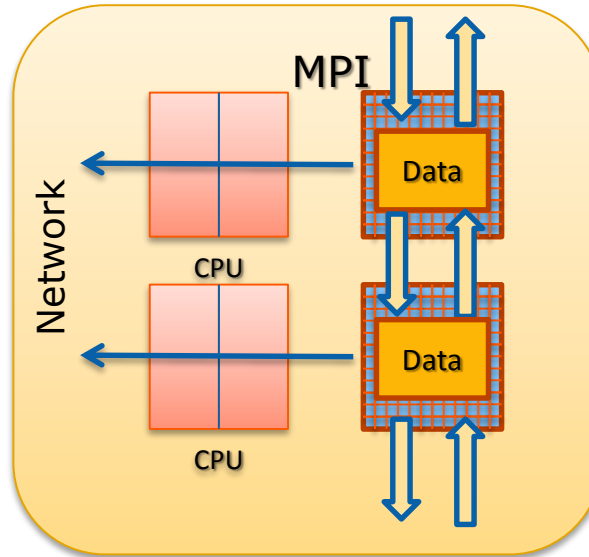
# Intel® MPI Library Architecture & Staging

MPI-2.2

Application

ADI3*

MPICH2* upper layer

CH3*

CH3* device layer

Nemesis*

Netmod*

shm

tcp

dapl, ofa

mmap(2)

user SCIF†

OFED verbs/core

kernel SCIF

HCA‡ driver

Pre-Alpha

Alpha

Beta/Gold

†: Symmetric Communications Interface
‡: Host Channel Adapter

(intel) Software

# Selecting network fabrics

- Intel® MPI selects automatically 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(R) Xeon Phi(TM):
  - shm, tcp, ofa, dapl
  - Availability checked in the order shm:dapl, shm:ofa, shm:tcp (intra:inter)
- Set I_MPI_SSHM_SCIF=1 to enable shm fabric between host and Intel® Xeon Phi™
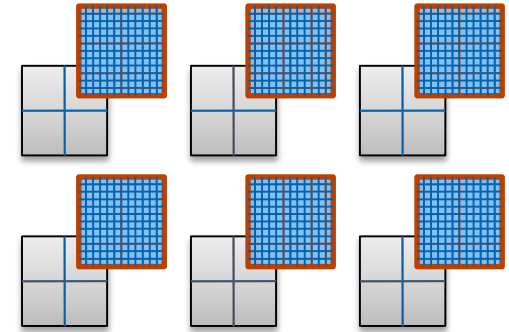
# Agenda

- Overview
- **Programming Models**
- Hybrid Computing
- Load Balancing

# Coprocessor only Programming Model

- MPI ranks on Intel® Xeon Phi™ (only)
- All messages into/out of coprocessors
- Intel® Cilk™ Plus, OpenMP*, Intel® Threading Building Blocks, Pthreads used directly within MPI processes

**Homogenous network of many-core CPUs**

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

Upload the binary to the Intel® Xeon Phi™.

Run instances of the MPI application on Intel® Xeon Phi™ 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™ Architecture

```
# mpiicc -mmic -o test_hello.MIC test.c
```
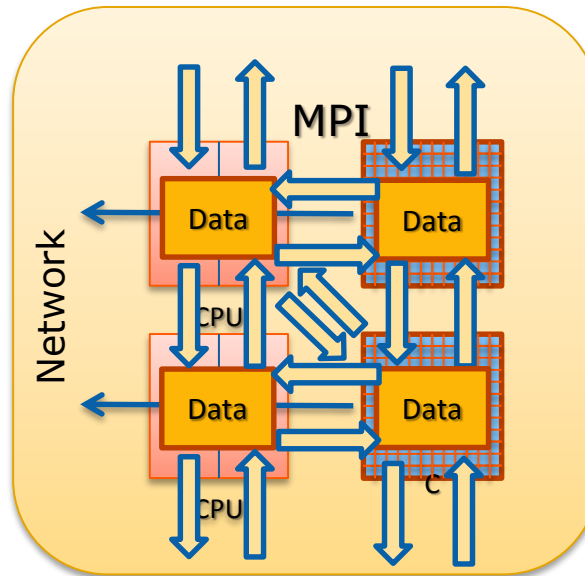
- Upload the coprocessor executable

```
# sudo micput 172.31.1.1 ./test_hello.MIC
    /tmp/test_hello.MIC
```

- Remark: If NFS available no explicit uploads required (just copies)!

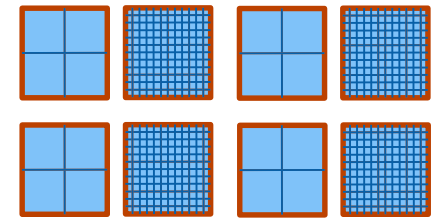- Launch the application on the coprocessor from host

```
# mpiexec -n 2 -wdir /tmp -host 172.31.1.1
    /tmp/test_hello.MIC
```

- Alternatively: login to the coprocessor and execute the already uploaded mpiexec.hydra there!

# 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**



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

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

Run instances of the MPI application on different mixed nodes.

# Symmetric Programming 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
```
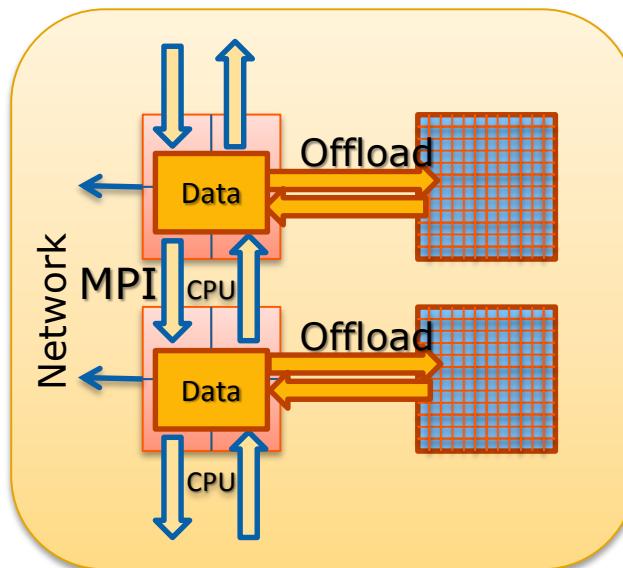
- Upload the Intel® Xeon Phi™ executable

```
# sudo micput 172.31.1.1 ./test_hello.MIC
    /tmp/test_hello.MIC
```

- Launch the application on the host and the coprocessor from the host

```
# mpiexec -n 2 -host <hostname> ./test_hello : -wdir
/tmp -n 2 -host 172.31.1.1 /tmp/test_hello.MIC
```
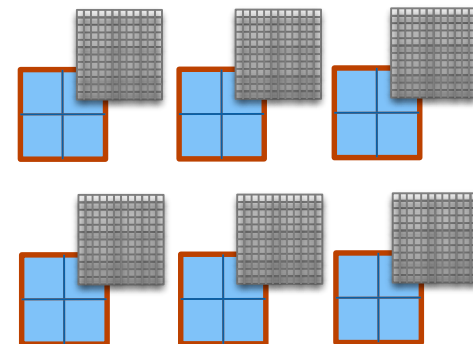
# MPI+Offload Programming Model

- 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™



**Homogenous network of heterogeneous nodes**

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

Run instances of the MPI application on the host, offloading code onto 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™ Architecture 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
  - Previous compilers needed `–offload-build` flag
- Execute on host(s) as usual
  ```
  # mpiexec -n 2 ./test
  ```
- MPI processes will offload code for acceleration

# Agenda

- Overview
- Programming Models
- **Hybrid Computing**
- Load Balancing

# 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 (cont'd)

- 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 (including coprocessors)

# 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™ Architecture!

# 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 (takes precedence):

    **I_MPI_PIN_DOMAIN =** *<size>[:<layout>]*

    <size> =

    | | |
    |---|---|
    | **omp** | Adjust to OMP_NUM_THREADS |
    | **auto** | #CPUs/#MPIprocs |
    | **<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 mpiexec of Intel® MPI reads argument sets from the command line:
  - Sections between ":" define an argument set (also lines in a configfile, but not yet available in Beta)
  - Host, number of nodes, but also environment can be set independently in each argument set

```
# mpiexec –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

# Co-Processor only and Symmetric Support

- Full hybrid support on Intel® Xeon from Intel ® MPI extends to Intel® Xeon Phi(TM)
- `KMP_AFFINITY=balanced` (only on coprocessor) in addition to `scatter` and `compact`
- Recommendations:
  - Explicitly control where MPI processes and threads run in a hybrid application
  (according to threading model and application)
  - Avoid splitting cores among MPI processes, i.e. `I_MPI_PIN_DOMAIN` should be a multiple of 4
  - Try different `KMP_AFFINITY` settings for your application

# MPI+Offload Support

- Define thread affinity manually per single MPI process (pseudo syntax!):

```
# export OMP_NUM_THREADS=4

# mpiexec -env KMP_AFFINITY=[1-4] -n 1 -host myMIC ... :
           -env KMP_AFFINITY=[5-8] -n 1 -host myMIC ... :
           ...
```

# Agenda

- Overview
- Programming Models
- Hybrid Computing
- Load Balancing

# Intel® Xeon Phi™ Coprocessor Becomes a Network Node



Intel® Xeon® Processor

Intel® Xeon Phi™ Coprocessor

Virtual Network Connection

Intel® Xeon® Processor

Intel® Xeon Phi™ Coprocessor

Virtual Network Connection

## Intel® Xeon Phi™ Architecture + Linux enables IP addressability

# Load Balancing

- Situation
  - Host and coprocessor computation performance are different
  - Host and coprocessor internal communication speed is different
- MPI in symmetric mode is like running on a heterogenous cluster
- Load balanced codes (on homogeneous cluser) may get imbalanced!
- Solution? No general solution!
  - Approach 1: Adapt MPI mapping of (hybrid) code to performance characteristics: #m processes per host, #n process per coprocessor(s)
  - Approach 2: Change code internal mapping of workload to MPI processes
    - Example: uneven split of calculation grid for MPI processes on host vs. coprocessor(s)
  - Approach 3: ...
- Analyze load balance of application with ITAC
  - Ideal Interconnect Simulator

# Improving Load Balance: Real World Case

Collapsed data per node and coprocessor card

Host
16 MPI procs x
1 OpenMP thread

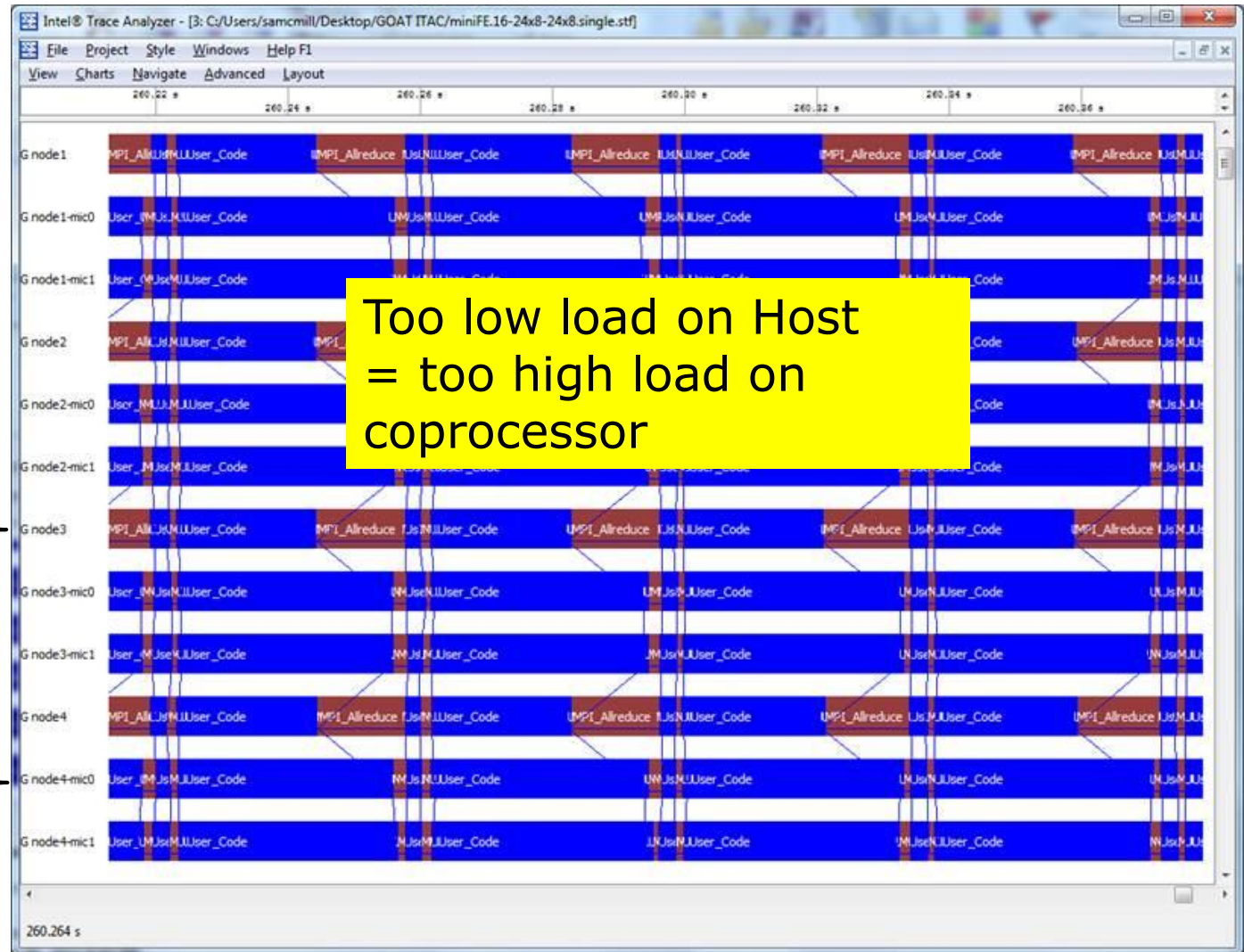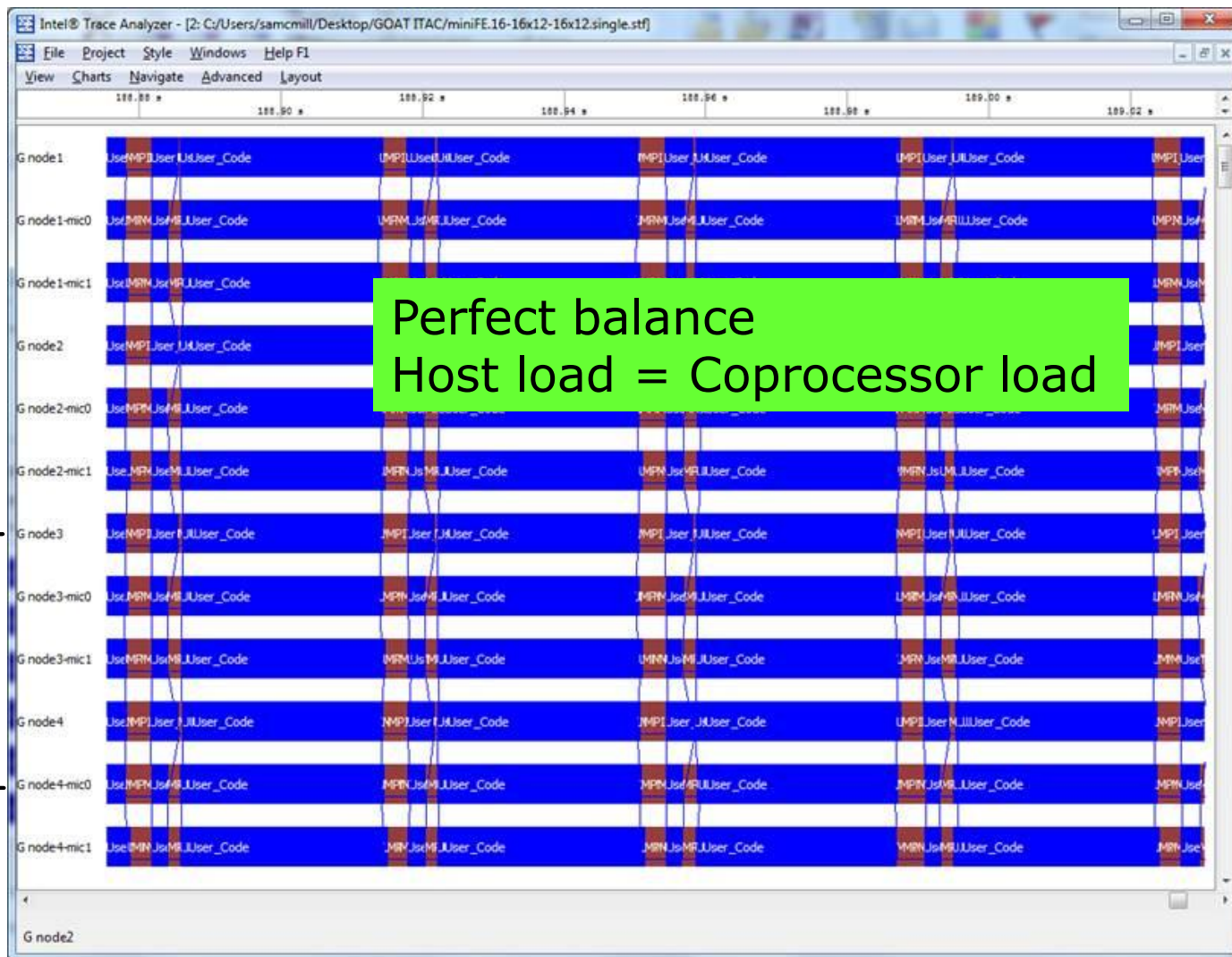Coprocessor
8 MPI procs x
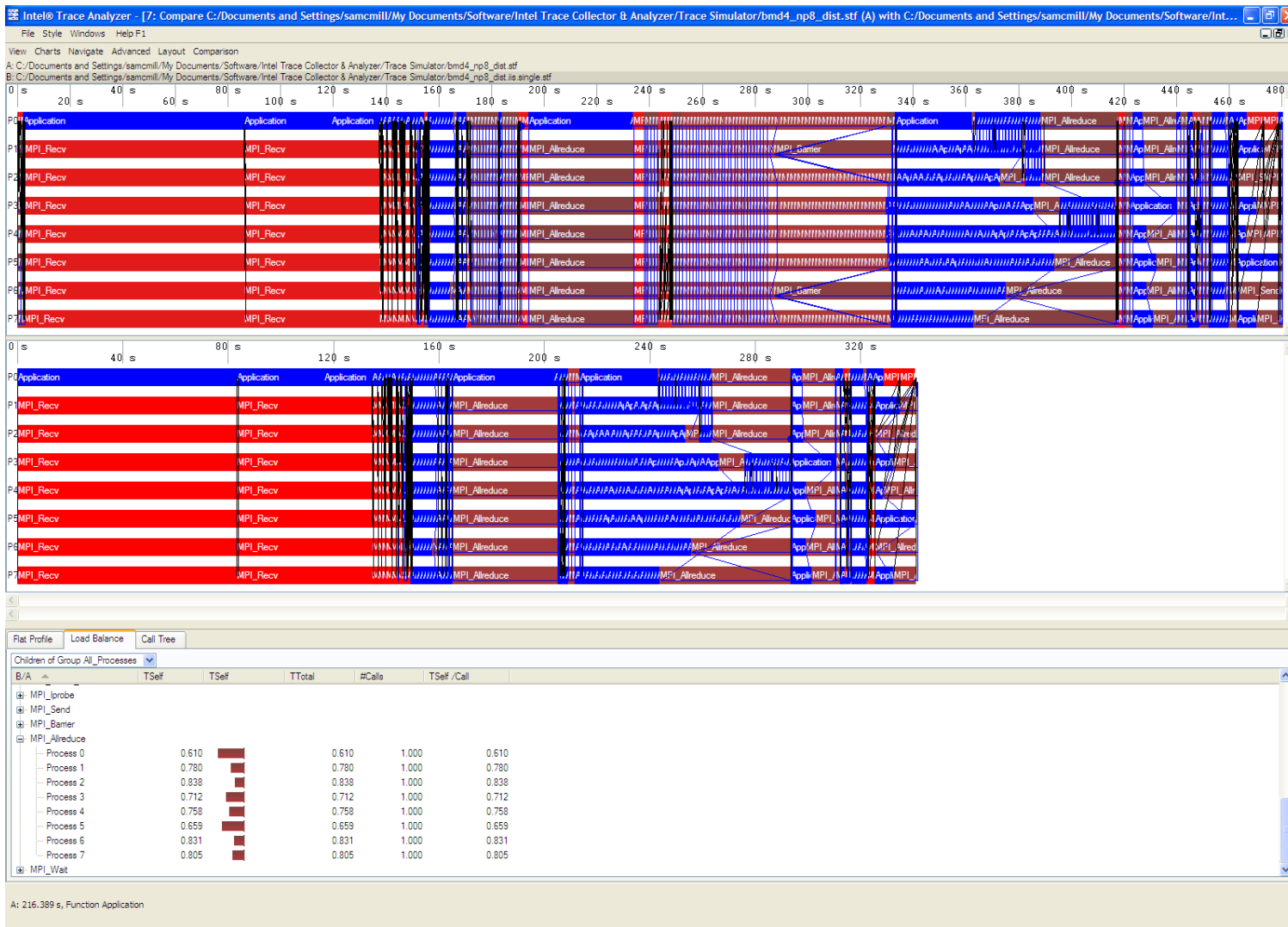28 OpenMP threads



Too high load on Host = too low load on coprocessor

# Improving Load Balance: Real World Case

Collapsed data
per node and
coprocessor card

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

<u>Coprocessor</u>
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 card

<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 (IIS)

- What is the Ideal Interconnect Simulator (IIS)?
  - 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™ Architecture**

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

# Questions?