



An Introduction to the Intel® Xeon Phi™ Coprocessor

INFIERI-2013 - July 2013

Leo Borges (leonardo.borges@intel.com)
Intel Software & Services Group

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Intel in High-Performance Computing




Dedicated, Renowned Applications Expertise



Large Scale Clusters for Test & Optimization

Tera-Scale Research




Exa-Scale Labs



Broad Software Tools Portfolio

Defined HPC Application Platform



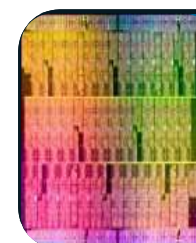
Platform Building Blocks



Manufacturing Process Technologies



Leading Performance, Energy Efficient




Many Integrated Core Architecture

A long term commitment to the HPC market segment

HPC Processor Solutions

Multi-Core



Xeon®
General Purpose Architecture
Leadership Per Core Performance
FP/core CAGR via AVX
Multi-Core CAGR

<u>EN</u> General purpose perf/watt	<u>EP</u> Max perf/watt w/ Higher Memory BW / freq and QPI ideal for HPC	<u>EP 4S</u> Additional compute density	<u>Xeon EX</u> Additional sockets & big memory
--	---	--	---

Many-Core



Intel® Xeon Phi™ Coprocessor
Trades a “big” IA core for multiple lower performance IA cores resulting in higher performance for a subset of highly parallel applications



Highly Parallel Applications Markets, Types, & Hardware



Energy & oil exploration



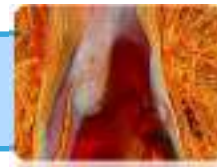
Digital content creation



Climate modeling & weather prediction



Financial analyses, trading



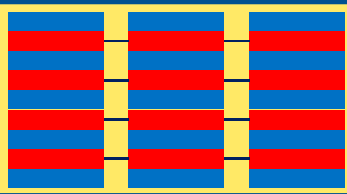
Medical imaging and biophysics



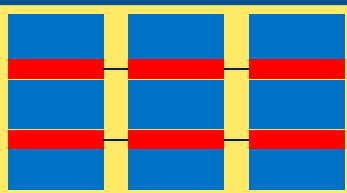
Computer Aided Design & Manufacturing

Parallel Application Types

Fine Grain



Coarse Grain



Embarrassingly Parallel



■ Communication
■ Compute Process

Highly Parallel Compute Kernels

Black-Scholes Sparse/Dense Matrix Mult ...
 FFTs Vector Math LU Factorization

Hardware Options for Parallel Application Types



OR



Intel® MIC Architecture

C/C++
Fortran

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

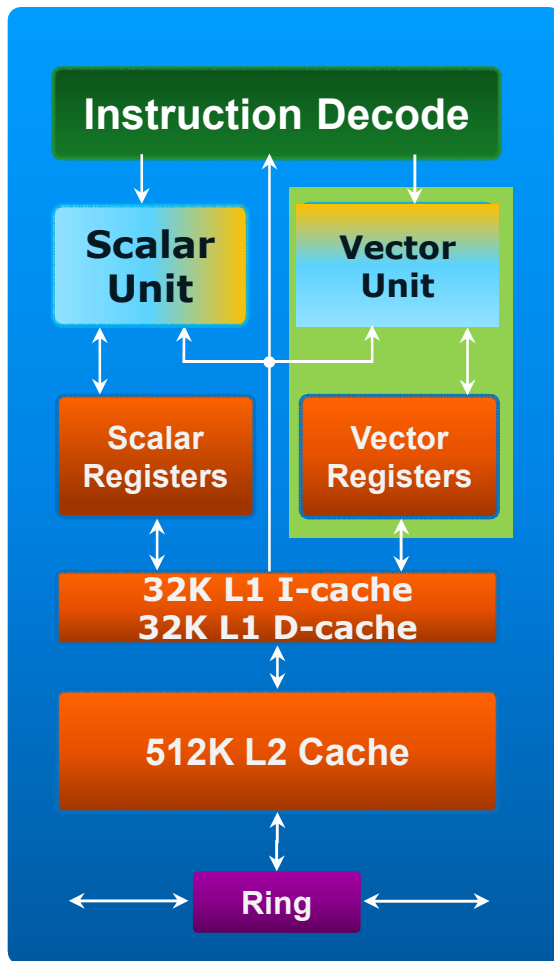
MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread execution unit



>50 in-order cores

- Ring interconnect

64-bit addressing

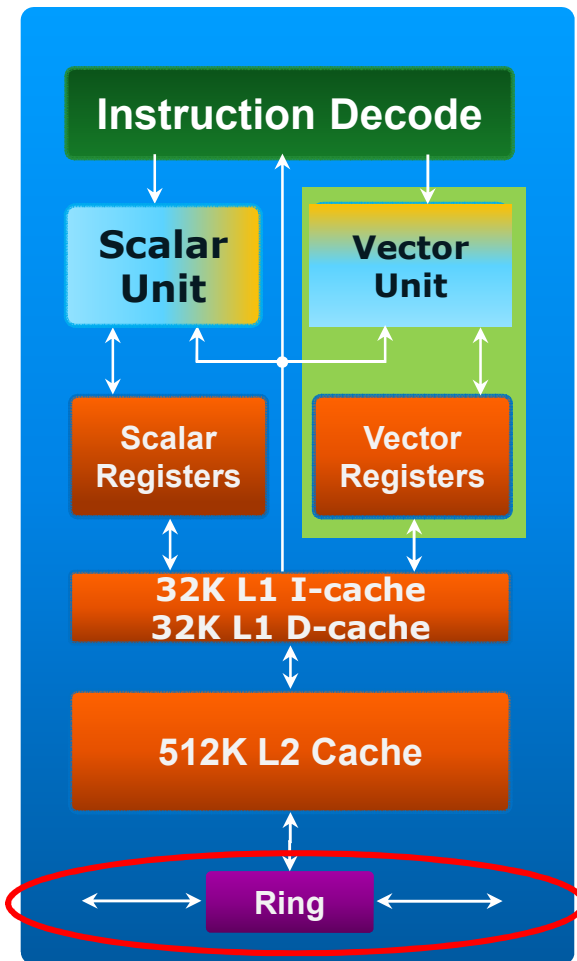
Scalar unit based on Intel® Pentium® processor family

- Two pipelines
 - Dual issue with scalar instructions
- One-per-clock scalar pipeline throughput
 - 4 clock latency from issue to resolution

4 hardware threads per core

- Each thread issues instructions in turn
- Round-robin execution hides scalar unit latency

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread execution unit



>50 in-order cores

- Ring interconnect

64-bit addressing

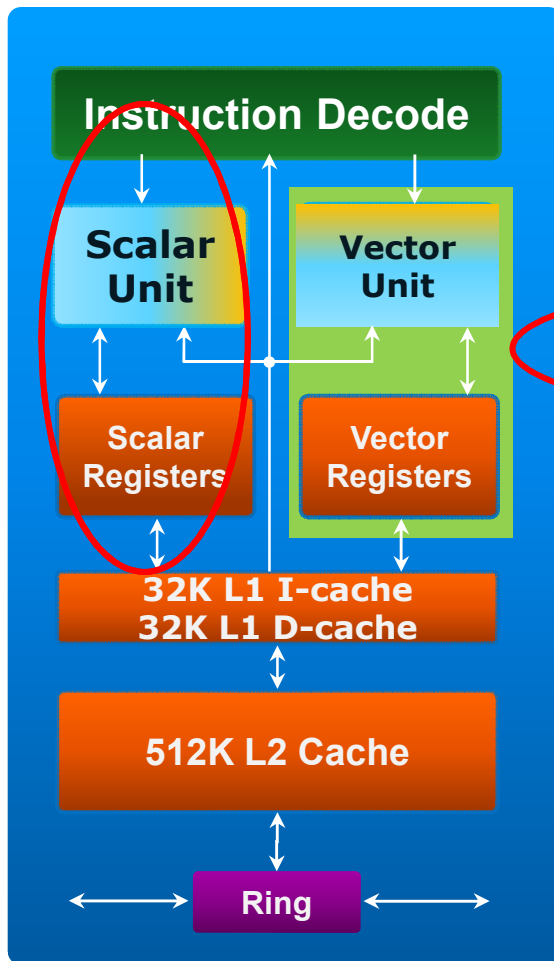
Scalar unit based on Intel® Pentium® processor family

- Two pipelines
 - Dual issue with scalar instructions
- One-per-clock scalar pipeline throughput
 - 4 clock latency from issue to resolution

4 hardware threads per core

- Each thread issues instructions in turn
- Round-robin execution hides scalar unit latency

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread execution unit



>50 in-order cores

- Ring interconnect

64-bit addressing

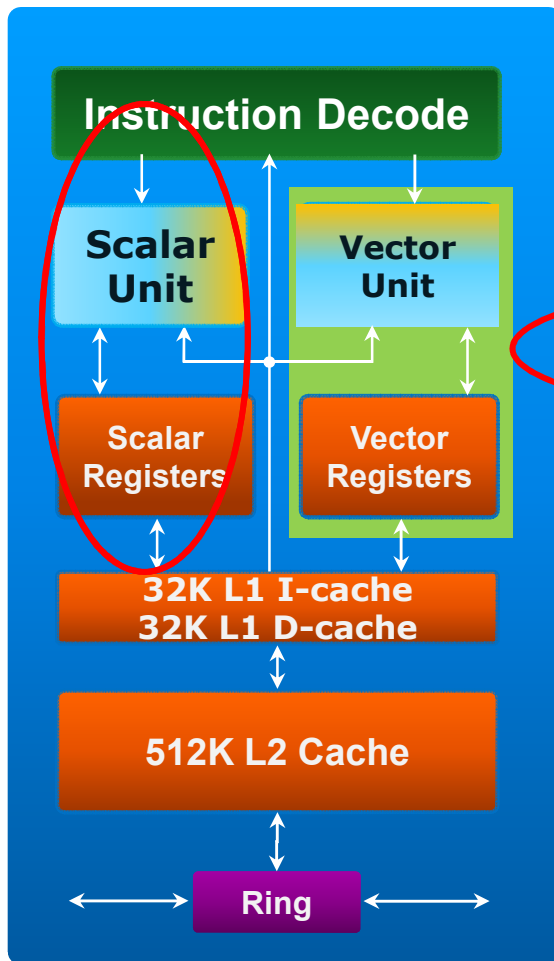
Scalar unit based on Intel® Pentium® processor family

- Two pipelines
 - Dual issue with scalar instructions
- One-per-clock scalar pipeline throughput
 - 4 clock latency from issue to resolution

4 hardware threads per core

- Each thread issues instructions in turn
- Round-robin execution hides scalar unit latency

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread execution unit



>50 in-order cores

- Ring interconnect

64-bit addressing

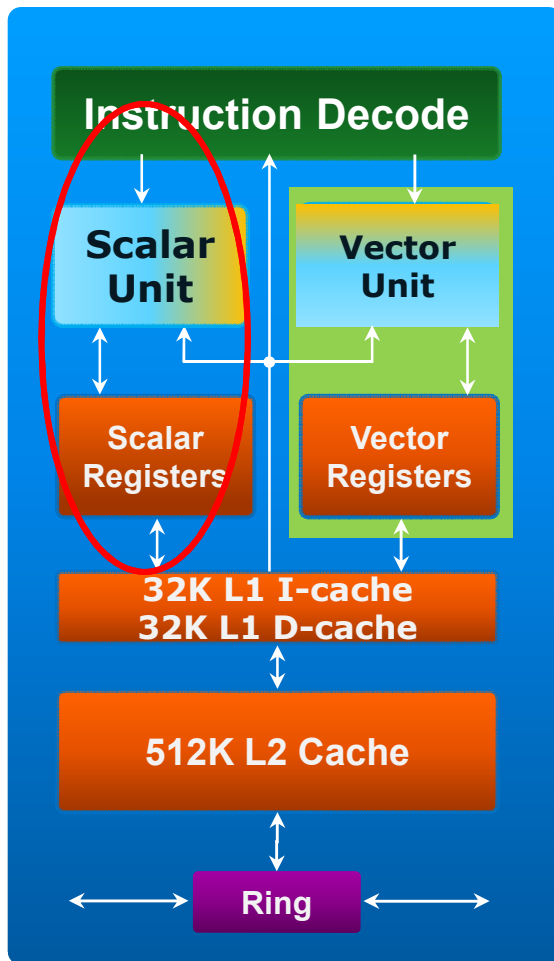
Scalar unit based on Intel® Pentium® processor family

- Two pipelines
 - Dual issue with scalar instructions
- One-per-clock scalar pipeline throughput
 - 4 clock latency from issue to resolution

4 hardware threads per core

- Each thread issues instructions in turn
- Round-robin execution hides scalar unit latency

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread execution unit



>50 in-order cores

- Ring interconnect

64-bit addressing

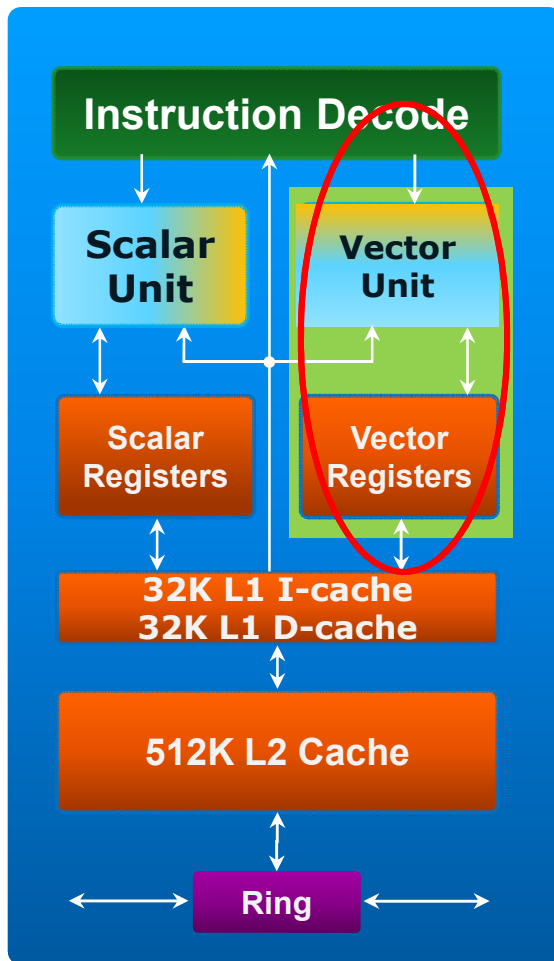
Scalar unit based on Intel® Pentium® processor family

- Two pipelines
 - Dual issue with scalar instructions
- One-per-clock scalar pipeline throughput
 - 4 clock latency from issue to resolution

4 hardware threads per core

- Each thread issues instructions in turn
- Round-robin execution hides scalar unit latency

Each Intel® Xeon Phi™ Coprocessor core is a fully functional multi-thread vector unit



Optimized

- Single and Double precision

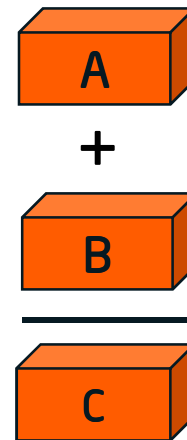
All new vector unit

- 512-bit SIMD Instructions – not Intel® SSE, MMX™, or Intel® AVX
- 32 512-bit wide vector registers
 - Hold 16 singles or 8 doubles per register

Fully-coherent L1 and L2 caches

Reminder: Vectorization, What is it?

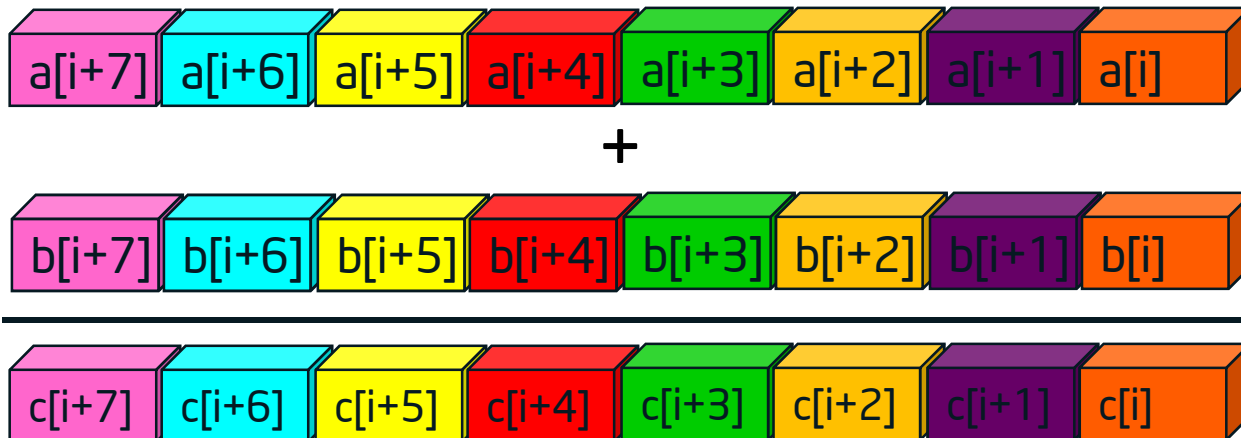
(Graphical View)



Scalar

- One Instruction
- One Mathematical Operation

```
for (i=0; i<=MAX; i++)  
    c[i]=a[i]+b[i];
```



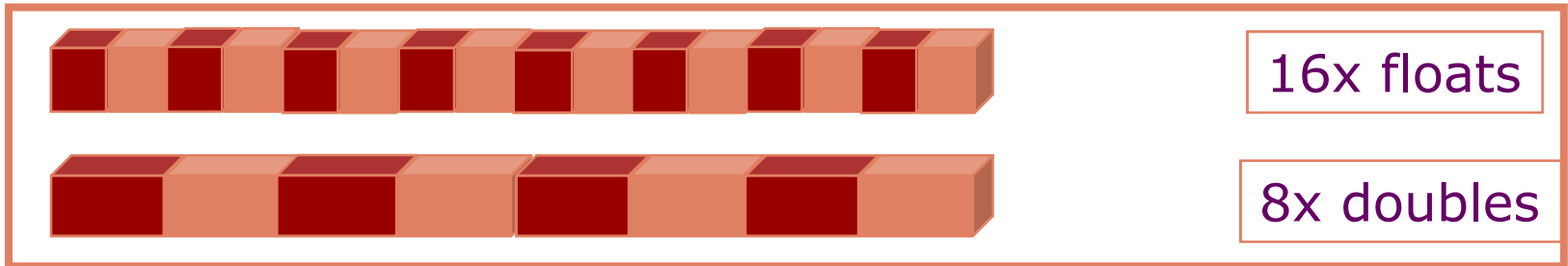
Vector

- One Instruction
- **Eight** Mathematical Operations¹

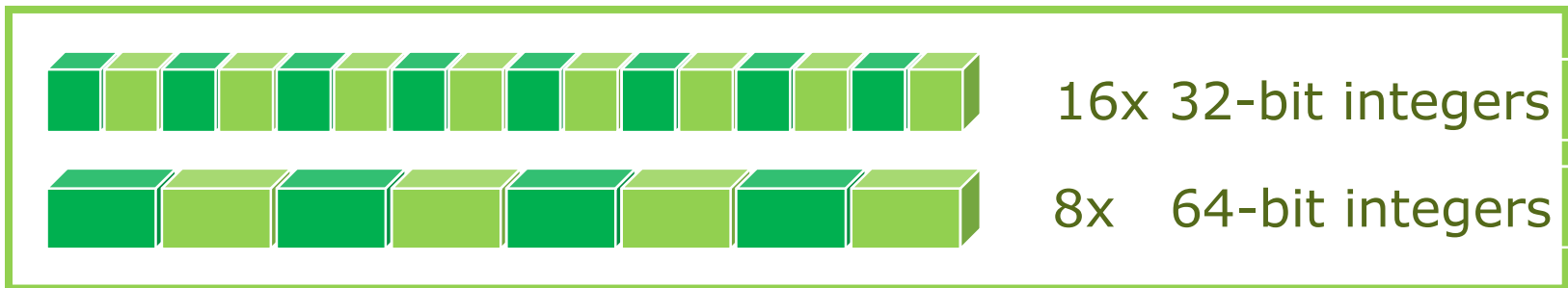
1. Number of operations per instruction varies based on the which SIMD instruction is used and the width of the operands

Data Types for Intel® MIC Architecture

now

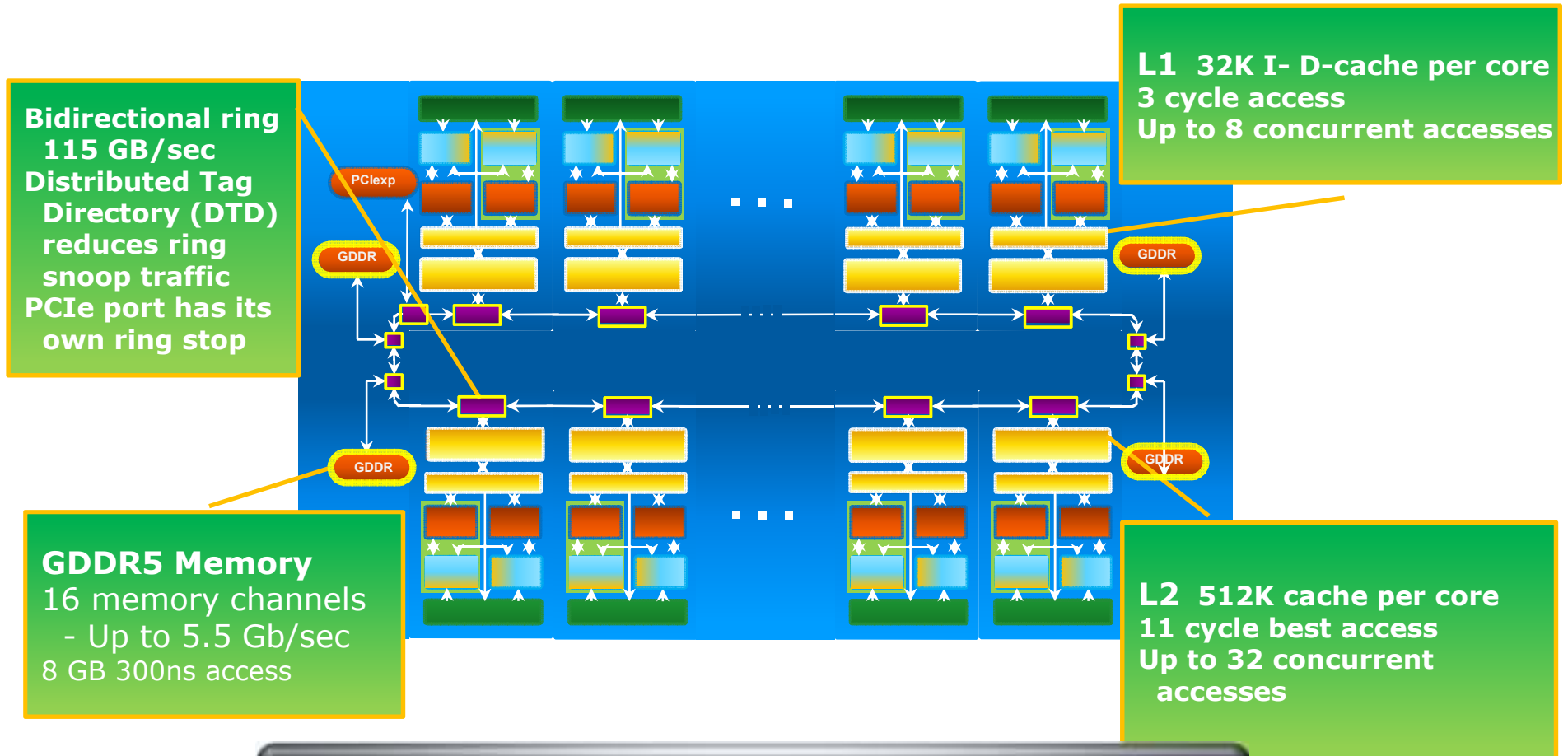


now



Takeaway: Vectorization is very important

Individual cores are tied together via fully coherent caches into a bidirectional ring



Intel® Xeon Phi™ Coprocessor x100 Family Reference Table

Processor Brand Name	Codename	SKU #	Form Factor, Thermal	Board TDP (Watts)	Max # of Cores	Clock Speed (GHz)	Peak Double Precision (GFLOP)	GDDR5 Memory Speeds (GT/s)	Peak Memory BW	Memory Capacity (GB)	Total Cache (MB)	Enabled Turbo	Turbo Clock Speed (GHz)	
Intel® Xeon Phi™ Coprocessor x100	Knights Corner	7120P	PCIe Card, Passively Cooled	300	61	1.238	1208	5.5	352	16	30.5	Y	1.333	
		7120X	PCIe Card, No Thermal Solution	300	61	1.238	1208	5.5	352	16	30.5	Y	1.333	
		5120D	PCIe Dense Form Factor, No Thermal Solution	245	60	1.053	1011	5.5	352	8	30	N	N/A	
		3120P	PCIe Card, Passively Cooled	300	57	1.1	1003	5.0	240	6	28.5	N	N/A	
		3120A	PCIe Card, Actively Cooled	300	57	1.1	1003	5.0	240	6	28.5	N	N/A	
		Previously Launched and Disclosed												
		5110P*	PCIe Card, Passively Cooled	225	60	1.053	1011	5.0	320	8	30	N	N/A	

*Please refer to our technical documentation for Silicon stepping information



Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

MPI Programming Models

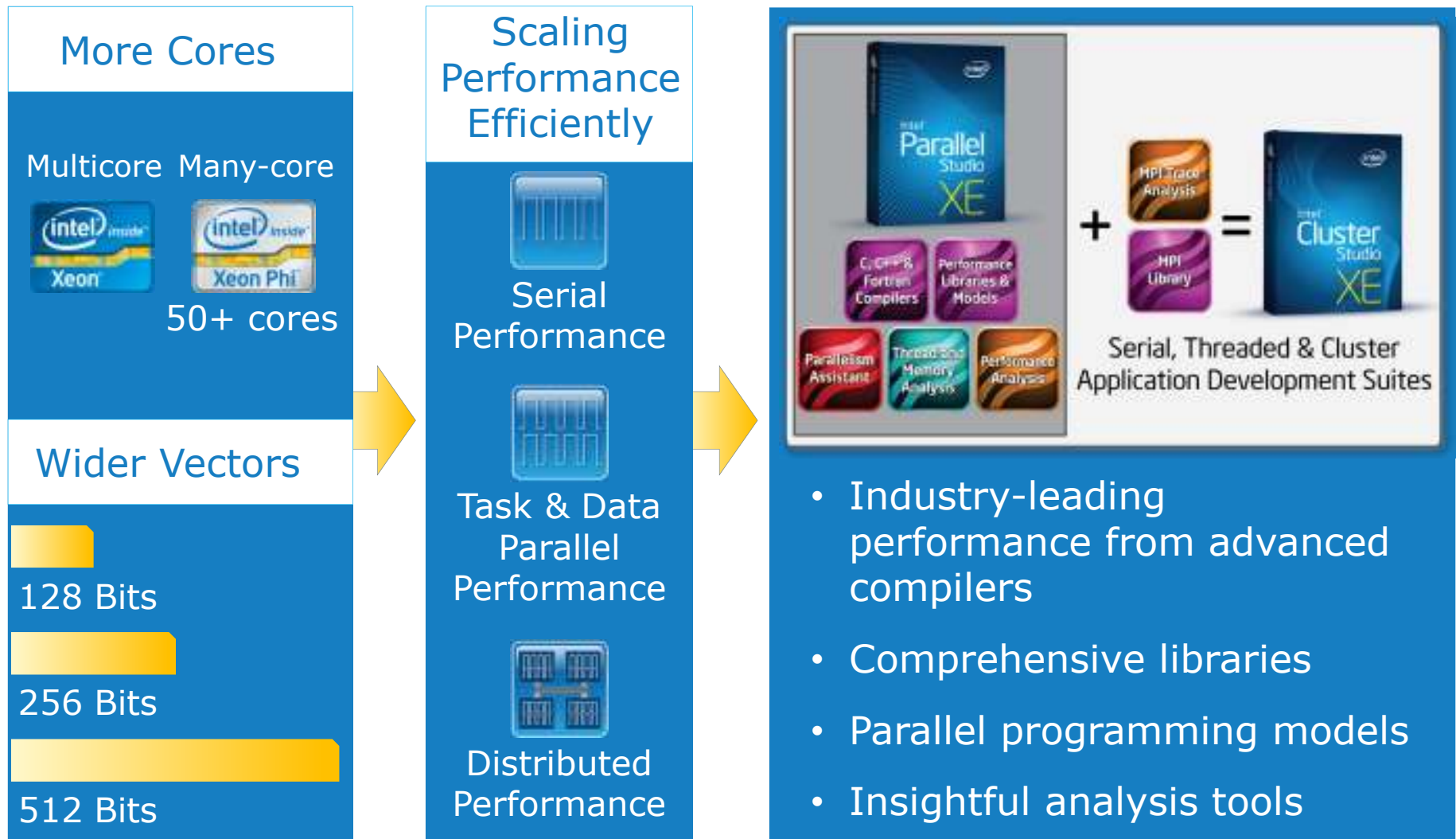
Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

More Cores. Wider Vectors. Performance Delivered.

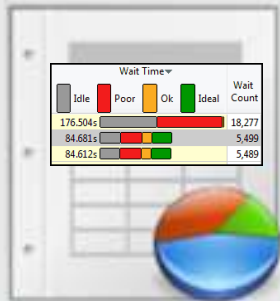
Intel® Parallel Studio XE 2013 and Intel® Cluster Studio XE 2013



Comprehensive set of SW tools

Code Analysis

Advisor XE
VTune Amplifier XE
Inspector XE
Trace Analyzer



Libraries & Compilers

Math Kernel Library
Integrated Performance Primitives
Intel Compilers



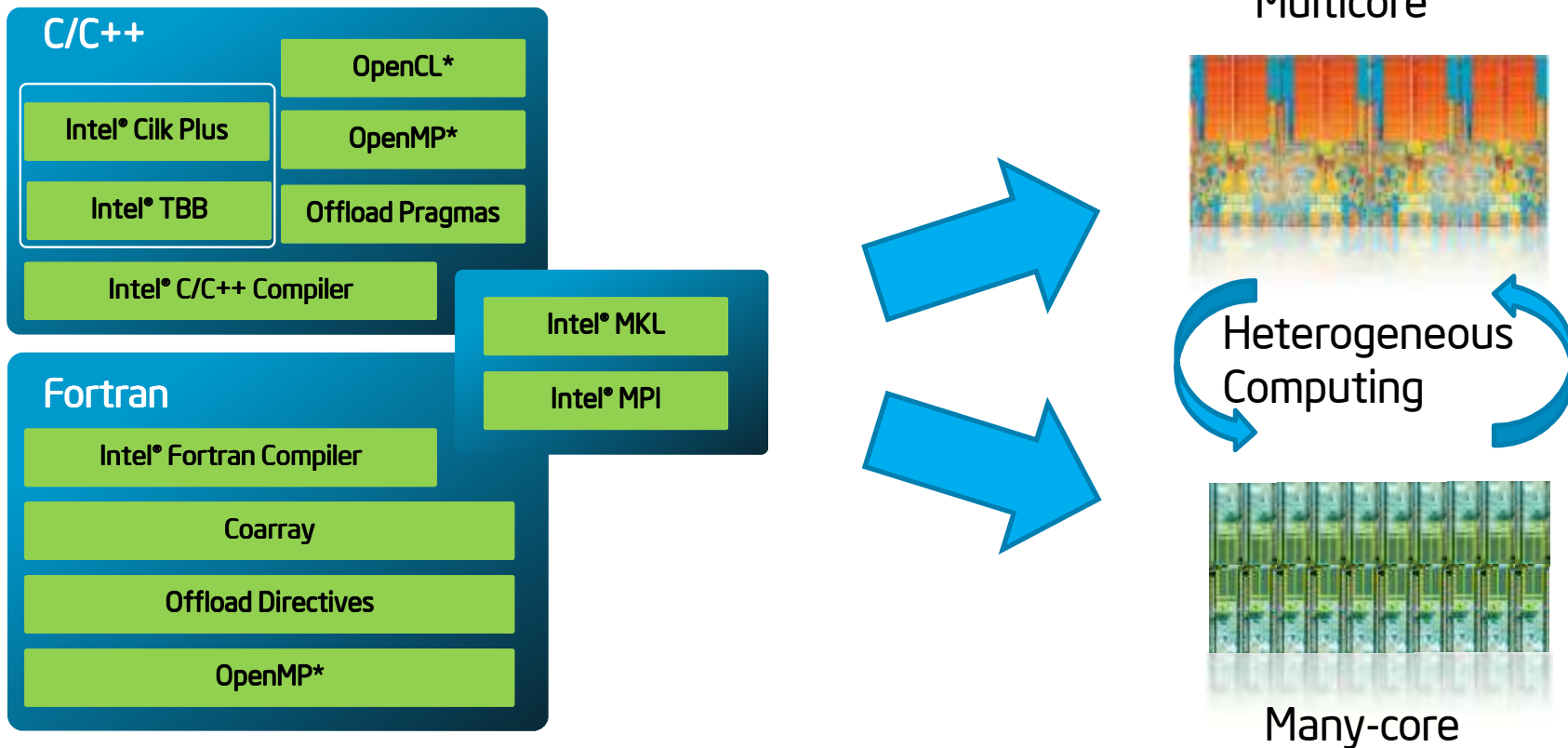
Programming Models

Intel Cilk Plus
Threading Building Blocks
OpenMP
OpenCL
MPI
Offload/Native/MYO



Preserve Your Development Investment

Common Tools and Programming Models for Parallelism



Develop Using Parallel Models that Support Heterogeneous Computing

Introduction

High-level overview of the Intel® Xeon Phi™ platform: Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:

- Native
- Offload
 - Explicit block data transfer
 - Offloading with Virtual Shared Memory

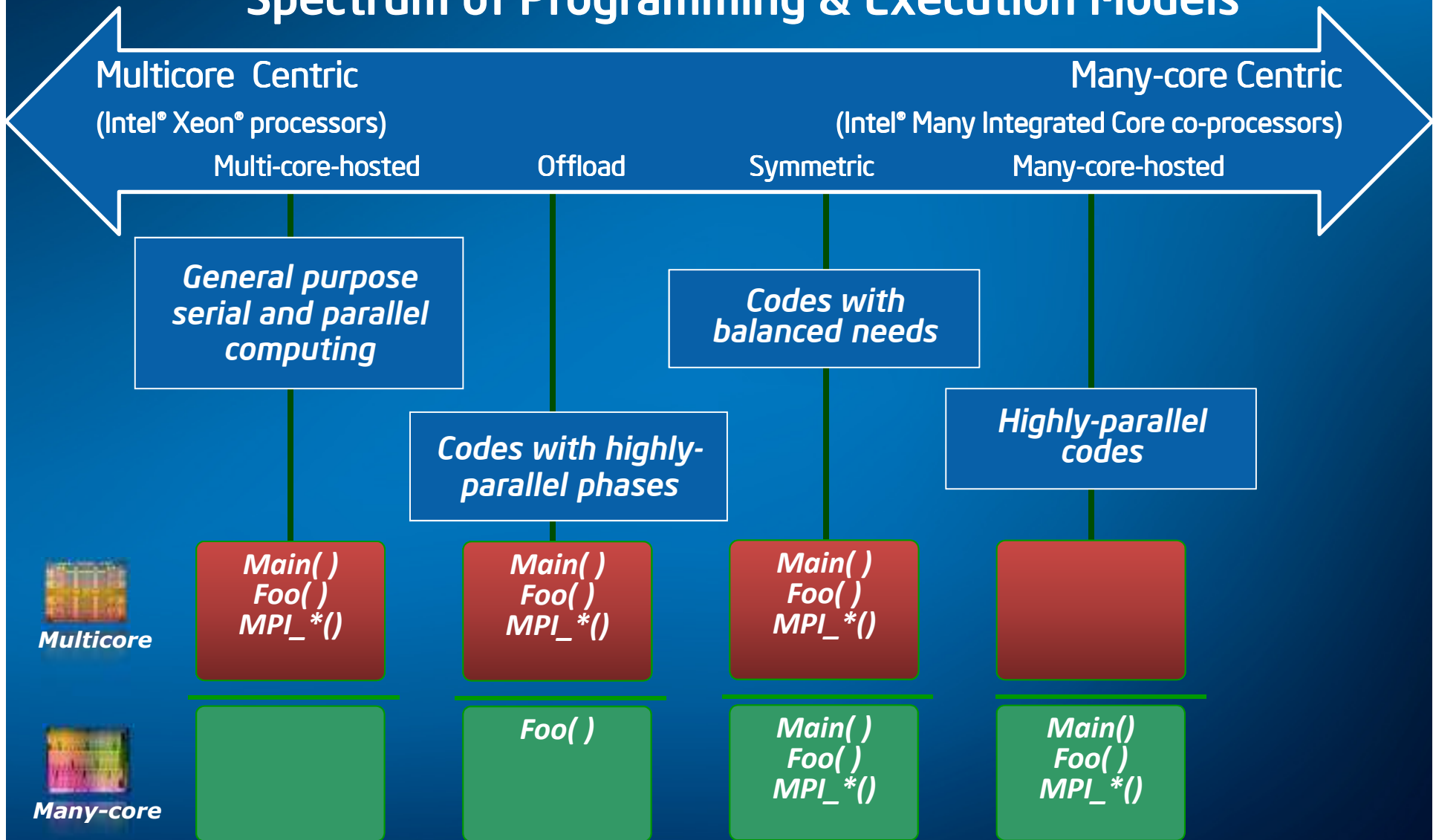
Performance and Thread Parallelism

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

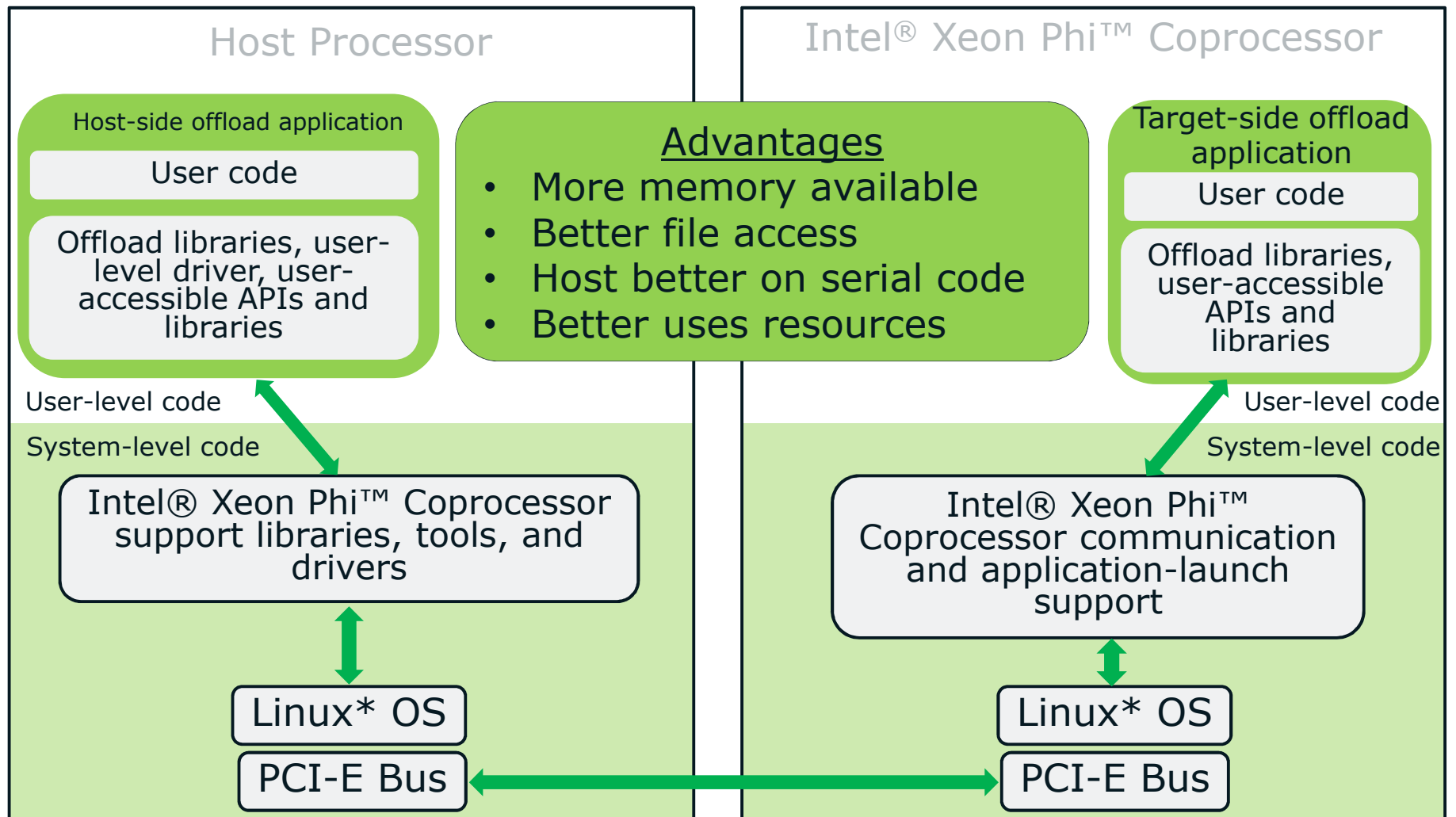
22 Profiling: Intel® Trace Analyzer and Collector

Spectrum of Programming & Execution Models

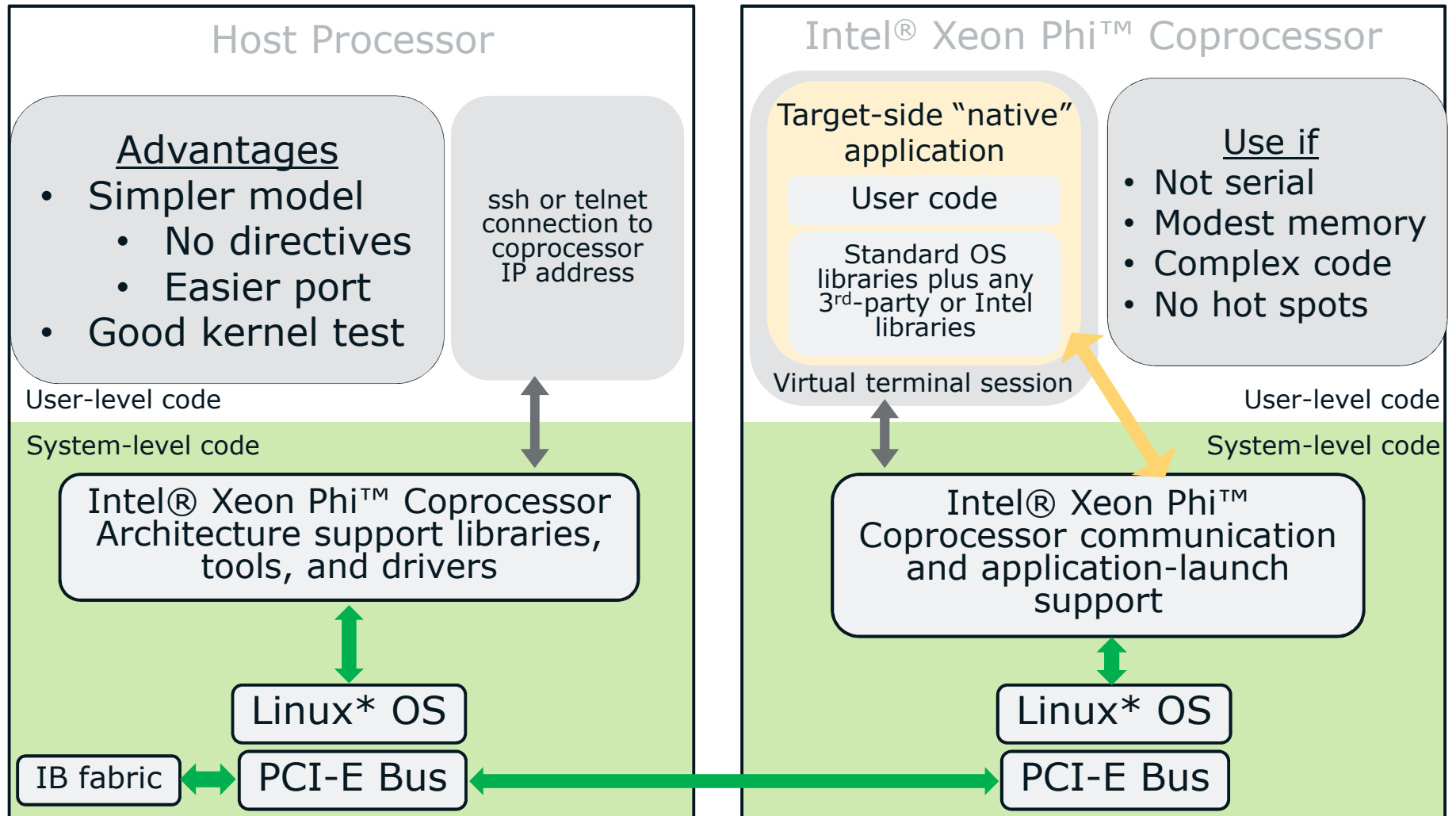


Range of Models to Meet Application Needs

Intel® Xeon Phi™ Coprocessor runs *either* as an accelerator for offloaded host computation



Or Intel® Xeon Phi™ Coprocessor runs as a native or MPI* compute node via IP or OFED



The Intel® Manycore Platform Software Stack (Intel® MPSS) provides Linux* on the coprocessor

Authenticated users can treat it like another node

```
ssh mic0 top
```

```
Mem: 298016K used, 7578640K free, 0K shrd, 0K buff, 100688K cached
CPU:  0.0% usr  0.3% sys  0.0% nic 99.6% idle  0.0% io  0.0% irq  0.0% sirq
Load average: 1.00 1.04 1.01 1/2234 7265
  PID  PPID  USER      STAT   VSZ  %MEM  CPU  %CPU  COMMAND
  7265  7264  fdkew     R      7060  0.0   14   0.3   top
    43     2  root     SW           0  0.0   13   0.0   [ksoftirqd/13]
  5748     1  root     S      119m  1.5  226   0.0   ./sep_mic_server3.8
  5670     1  micuser  S     97872  1.2    0   0.0   /bin/coi_daemon --coiuser=micuser
  7261  5667  root     S     25744  0.3    6   0.0   sshd: fdkew [priv]
  7263  7261  fdkew     S     25744  0.3  241   0.0   sshd: fdkew@notty
  5667     1  root     S     21084  0.2    5   0.0   /sbin/sshd
  5757     1  root     S      6940  0.0   18   0.0   /sbin/getty -L -l /bin/noauth 1152
     1     0  root     S      6936  0.0   10   0.0   init
  7264  7263  fdkew     S      6936  0.0    6   0.0   sh -c top
```

Intel MPSS supplies a virtual FS and native execution

```
sudo scp /opt/intel/composerxe/lib/mic/libiomp5.so root@mic0:/lib64
scp a.out mic0:/tmp
ssh mic0 /tmp/a.out my-args
```

Add `-mmic` to compiles to create native programs

```
icc -O3 -g -mmic -o nativeMIC myNativeProgram.c
```

Alternately, use the offload capabilities of Intel® Composer XE to access coprocessor

Offload directives in source code trigger Intel Composer to compile objects for both host and coprocessor

```
#pragma offload target(mic) inout(A:length(2000))           C/C++  
!DIR$ OFFLOAD TARGET(MIC) INOUT(A: LENGTH(2000))          Fortran
```

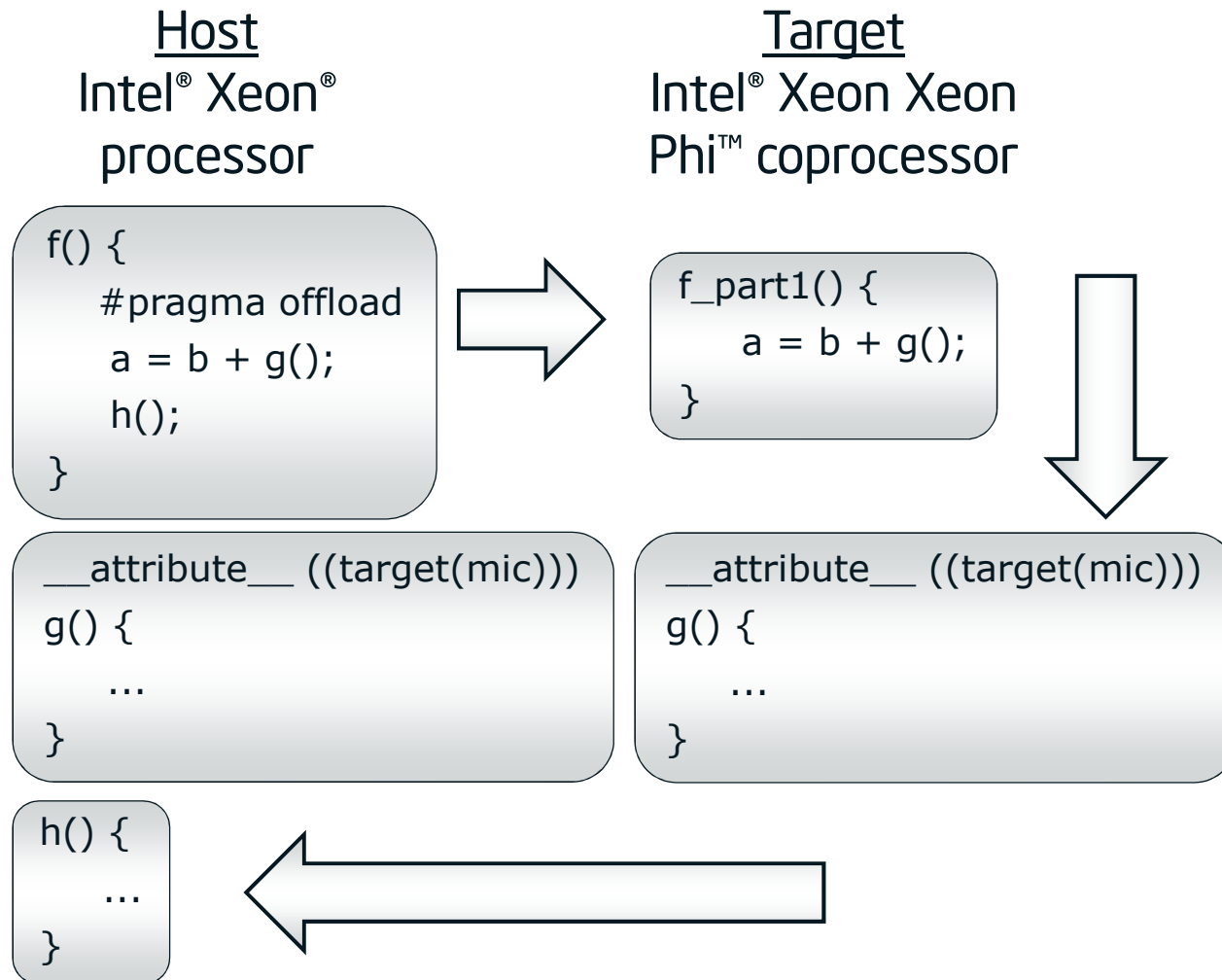
When the program is executed and a coprocessor is available, the offload code will run on that target

- Required data can be transferred explicitly for each offload
- Or use Virtual Shared Memory (`_Cilk_shared`) to match virtual addresses between host and target coprocessor

Offload blocks initiate coprocessor computation and can be synchronous or asynchronous

```
#pragma offload_transfer target(mic) in(a: length(2000)) signal(a)  
!DIR$ OFFLOAD_TRANSFER TARGET(MIC) IN(A: LENGTH(2000)) SIGNAL(A)  
_Cilk_spawn _Cilk_offload asynch-func()
```

Offload directives are independent of function boundaries



Execution

- If at first offload the target is available, the target program is loaded
- At each offload if the target is available, statement is run on target, else it is run on the host
- At program termination the target program is unloaded

Example: Compiler Assisted Offload

- Offload section of code to the coprocessor.

```
float pi = 0.0f;
#pragma offload target(mic)
#pragma omp parallel for reduction(+:pi)
for (i=0; i<count; i++) {
    float t = (float)((i+0.5f)/count);
    pi += 4.0f/(1.0f+t*t);
}
pi /= count;
```

- Offload any function call to the coprocessor.

```
#pragma offload target(mic) \
    in(transa, transb, N, alpha, beta) \
    in(A:length(matrix_elements)) \
    in(B:length(matrix_elements)) \
    in(C:length(matrix_elements)) \
    out(C:length(matrix_elements) alloc_if(0))
{
    sgemm(&transa, &transb, &N, &N, &N, &alpha, A, &N, B, &N,
        &beta, C, &N);
}
```

Example: Compiler Assisted Offload

- An example in Fortran:

```
!DEC$ ATTRIBUTES OFFLOAD : TARGET( MIC ) :: SGEMM
!DEC$ OMP OFFLOAD TARGET( MIC ) &
!DEC$ IN( TRANSA, TRANSB, M, N, K, ALPHA, BETA, LDA, LDB, LDC ), &
!DEC$ IN( A: LENGTH( NCOLA * LDA ) ), &
!DEC$ IN( B: LENGTH( NCOLB * LDB ) ), &
!DEC$ INOUT( C: LENGTH( N * LDC ) )
CALL SGEMM( TRANSA, TRANSB, M, N, K, ALPHA, &
           A, LDA, B, LDB BETA, C, LDC )
```

Example – share work between coprocessor and host using OpenMP*

```
omp_set_nested(1);
#pragma omp parallel private(ip)
{
#pragma omp sections
{
#pragma omp section
/*      use pointer to copy back only part of potential array,
to avoid overwriting host */
#pragma offload target(mic) in(xp) in(yp) in(zp) out(ppot:length(np1))
#pragma omp parallel for private(ip)
  for (i=0;i<np1;i++) {
    ppot[i] = threed_int(x0,xn,y0,yn,z0,zn,nx,ny,nz,xp[i],yp[i],zp[i]);
  }
#pragma omp section
#pragma omp parallel for private(ip)
  for (i=0;i<np2;i++) {
    pot[i+np1] =
    threed_int(x0,xn,y0,yn,z0,zn,nx,ny,nz,xp[i+np1],yp[i+np1],zp[i+np1]);
  }
}
}
```

Top level, runs on host
Runs on coprocessor
Runs on host

Pragmas and directives mark data and code to be offloaded and executed

	C/C++ Syntax
Offload pragma	#pragma offload <clauses> <statement> Allow next statement to execute on coprocessor or host CPU
Variable/function offload properties	__attribute__((target(mic))) Compile function for, or allocate variable on, both host CPU and coprocessor
Entire blocks of data/code defs	#pragma offload_attribute(push, target(mic)) #pragma offload_attribute(pop) Mark entire files or large blocks of code to compile for both
	Fortran Syntax
Offload directive	!dir\$ omp offload <clauses> <statement> Execute OpenMP* parallel block on coprocessor
	!dir\$ offload <clauses> <statement> Execute next statement or function on coproc.
Variable/function offload properties	!dir\$ attributes offload:<mic> :: <ret-name> OR <var1, var2, ...> Compile function or variable for CPU and coprocessor
Entire code blocks	!dir\$ offload begin <clauses> !dir\$ end offload

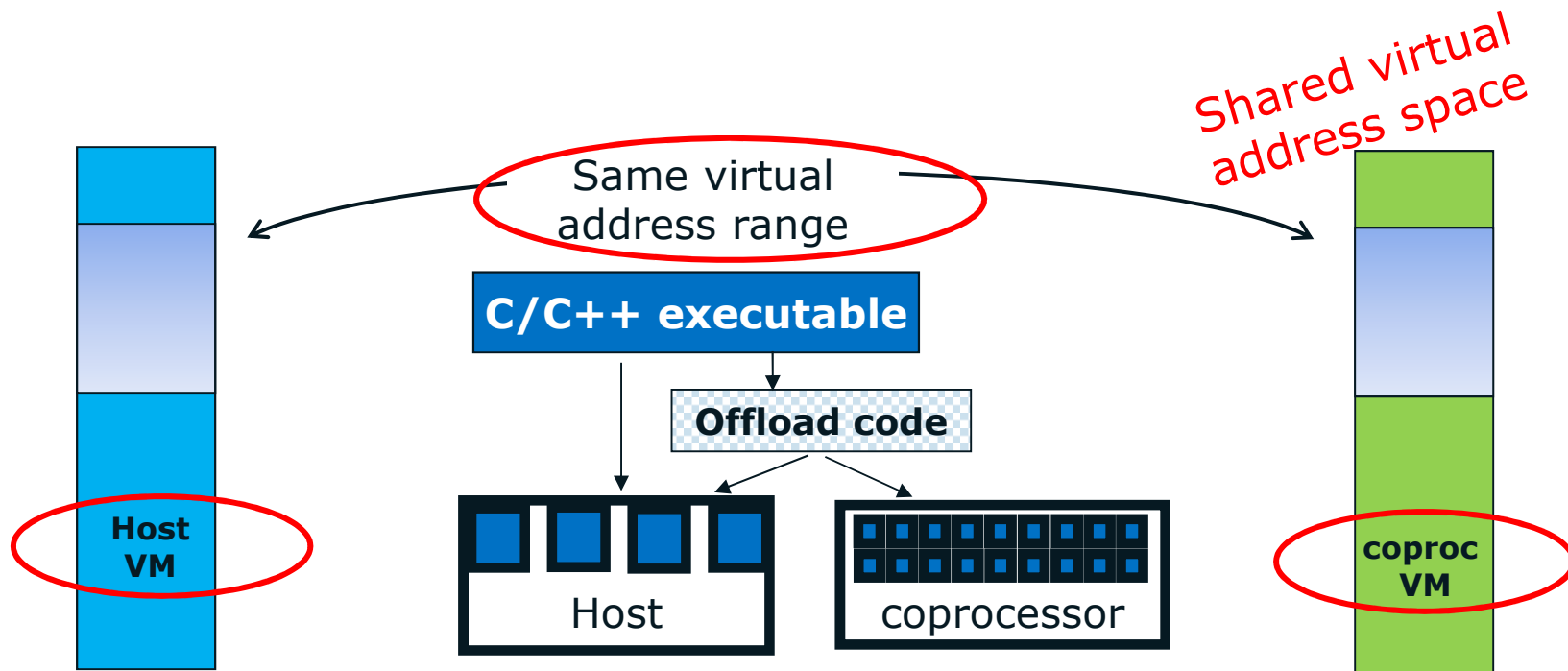
Options on offloads can control data copying and manage coprocessor dynamic allocation

Clauses	Syntax	Semantics
Multiple coprocessors	<code>target (mic[:unit])</code>	Select specific coprocessors
Conditional offload	<code>if (condition) / mandatory</code>	Select coprocessor or host compute
Inputs	<code>in (var-list modifiers_{opt})</code>	Copy from host to coprocessor
Outputs	<code>out (var-list modifiers_{opt})</code>	Copy from coprocessor to host
Inputs & outputs	<code>inout (var-list modifiers_{opt})</code>	Copy host to coprocessor and back when offload completes
Non-copied data	<code>nocopy (var-list modifiers_{opt})</code>	Data is local to target
Modifiers		
Specify copy length	<code>length (N)</code>	Copy N elements of pointer's type
Coprocessor memory allocation	<code>alloc_if (bool)</code>	Allocate coprocessor space on this offload (default: TRUE)
Coprocessor memory release	<code>free_if (bool)</code>	Free coprocessor space at the end of this offload (default: TRUE)
Control target data alignment	<code>align (N bytes)</code>	Specify minimum memory alignment on coprocessor
Array partial allocation & variable relocation	<code>alloc (array-slice) into (var-expr)</code>	Enables partial array allocation and data copy into other vars & ranges

To handle more complex data structures on the coprocessor, use Virtual Shared Memory

An identical range of virtual addresses is reserved on both host and coprocessor: changes are shared at offload points, allowing:

- Seamless sharing of complex data structures, including linked lists
- Elimination of manual data marshaling and shared array management
- Freer use of new C++ features and standard classes



Example: Virtual Shared Memory

- Shared between host and Xeon Phi

```
// Shared variable declaration
_Cilk_shared T in1[SIZE];
_Cilk_shared T in2[SIZE];
_Cilk_shared T res[SIZE];

_Cilk_shared void compute_sum()
{
    int i;
    for (i=0; i<SIZE; i++) {
        res[i] = in1[i] + in2[i];
    }
}

(...)

// Call compute sum on Target
_Cilk_offload compute_sum();
```

Virtual Shared Memory uses special allocation to manage data sharing at offload boundaries

Declare virtual shared data using `_Cilk_shared` allocation specifier

Allocate virtual dynamic shared data using these special functions:

```
_Offload_shared_malloc(), _Offload_shared_aligned_malloc(),  
_Offload_shared_free(), _Offload_shared_aligned_free()
```

Shared data copying occurs automatically around offload sections

- Memory is only synchronized on entry to or exit from an offload call
- Only modified data blocks are transferred between host and coprocessor

Allows transfer of C++ objects

- Pointers are transportable when they point to “shared” data addresses

Well-known methods can be used to synchronize access to shared data and prevent data races within offloaded code

- E.g., locks, critical sections, etc.

This model is integrated with the Intel® Cilk™ Plus parallel extensions

Note: Not supported on Fortran - available for C/C++ only

Data sharing between host and coprocessor can be enabled using this Intel® Cilk™ Plus syntax

What	Syntax
Function	<pre>int _Cilk_shared f(int x){ return x+1; }</pre> <p>Code emitted for host and target; may be called from either side</p>
Global	<pre>_Cilk_shared int x = 0;</pre> <p>Datum is visible on both sides</p>
File/Function static	<pre>static _Cilk_shared int x;</pre> <p>Datum visible on both sides, only to code within the file/function</p>
Class	<pre>class _Cilk_shared x {...};</pre> <p>Class methods, members and operators available on both sides</p>
Pointer to shared data	<pre>int _Cilk_shared *p;</pre> <p><i>p</i> is local (not shared), can point to shared data</p>
A shared pointer	<pre>int *_Cilk_shared p;</pre> <p><i>p</i> is shared; should only point at shared data</p>
Entire blocks of code	<pre>#pragma offload_attribute(push, _Cilk_shared)</pre> <pre>#pragma offload_attribute(pop)</pre> <p>Mark entire files or blocks of code <code>_Cilk_shared</code> using this pragma</p>

Intel® Cilk™ Plus syntax can also specify the offloading of computation to the coprocessor

Feature	Example
Offloading a function call	<pre>x = _Cilk_offload func(y);</pre> <p>func executes on coprocessor if possible</p>
	<pre>x = _Cilk_offload_to (card_num) func(y);</pre> <p>func must execute on specified coprocessor or an error occurs</p>
Offloading asynchronously	<pre>x = _Cilk_spawn _Cilk_offload func(y);</pre> <p>func executes on coprocessor; continuation available for stealing</p>
Offloading a parallel for-loop	<pre>_Cilk_offload _Cilk_for(i=0; i<N; i++){ a[i] = b[i] + c[i]; }</pre> <p>Loop executes in parallel on coprocessor. The loop is implicitly “un-inlined” as a function call.</p>

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

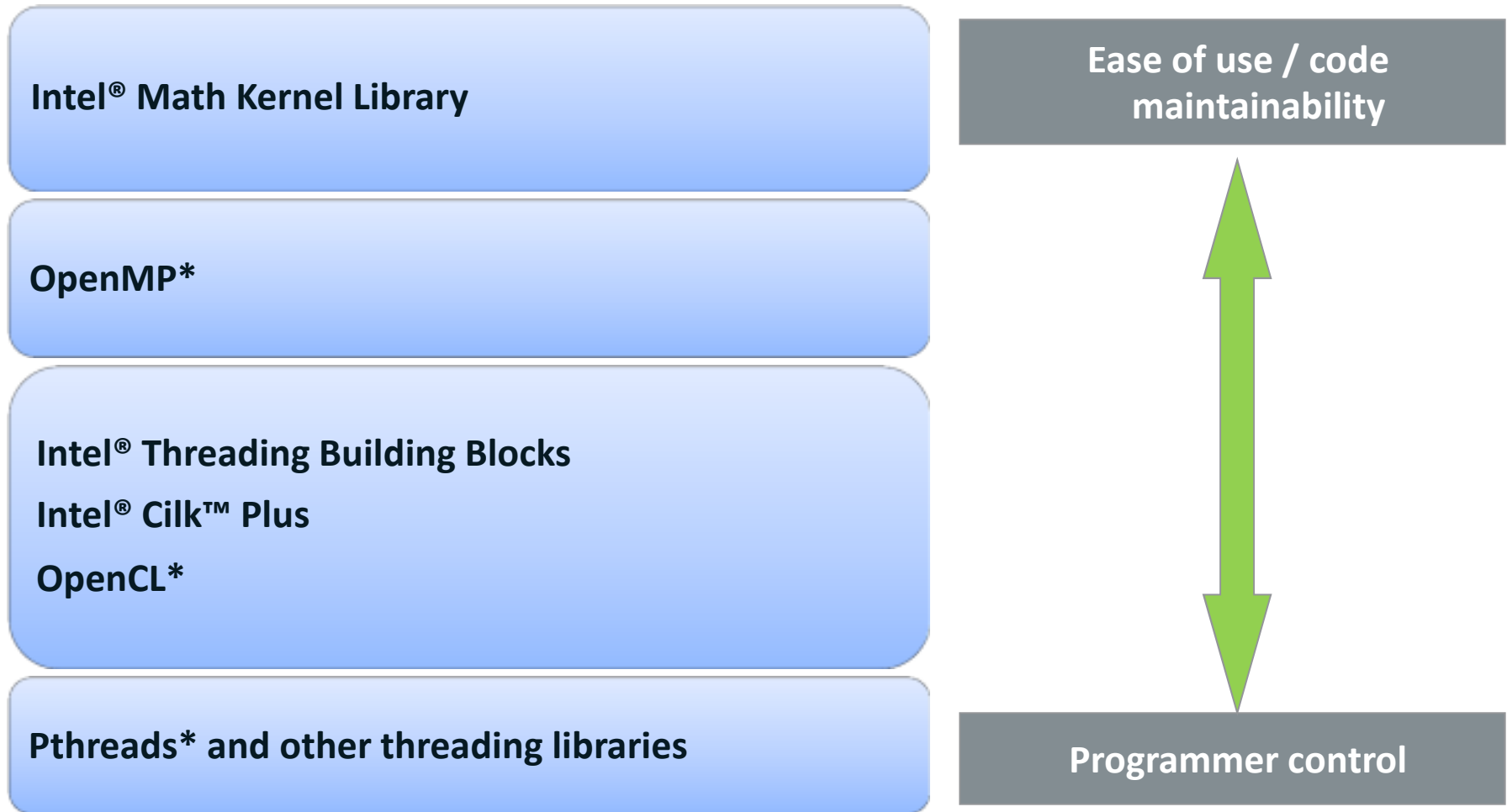
MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Options for Thread Parallelism



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

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism: OpenMP

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

OpenMP* on the Coprocessor

- The basics work just like on the host CPU
 - For both native and offload models
 - Need to specify `-openmp`
- There are 4 hardware thread contexts per core
 - Need at least $2 \times n_{\text{core}}$ threads for good performance
 - For all except the most memory-bound workloads
 - Often, 3x or 4x (number of available cores) is best
 - Very different from hyperthreading on the host!
 - `-opt-threads-per-core=n` advises compiler how many threads to optimize for
 - If you don't saturate all available threads, be sure to set `KMP_AFFINITY` to control thread distribution

OpenMP defaults

- OMP_NUM_THREADS defaults to
 - 1 x ncore for host (or 2x if hyperthreading enabled)
 - 4 x ncore for native coprocessor applications
 - 4 x (ncore-1) for offload applications
 - one core is reserved for offload daemons and OS
- Defaults may be changed via environment variables or via API calls on either the host or the coprocessor

Target OpenMP environment (offload)

Use target-specific APIs to set for coprocessor target only, e.g.

`omp_set_num_threads_target()` (called from host)

`omp_set_nested_target()` etc

- Protect with `#ifdef __INTEL_OFFLOAD`, undefined with `-no-offload`
- Fortran: `USE MIC_LIB` and `OMP_LIB` C: `#include <offload.h>`

Or define MIC – specific versions of env vars using

`MIC_ENV_PREFIX=MIC` (no underscore)

- Values on MIC **no longer default to values on host**
- Set values specific to MIC using

`export MIC_OMP_NUM_THREADS=120` (all cards)

`export MIC_2_OMP_NUM_THREADS=180` for card #2, etc

`export MIC_3_ENV="OMP_NUM_THREADS=240|KMP_AFFINITY=balanced"`

Stack Sizes for Coprocessor

For the main thread, (thread 0), default stack limit is 12 MB

- In offloaded functions, stack is used for local or automatic arrays and compiler temporaries
- To increase limit, export MIC_STACKSIZE (e.g. =100M)
 - default unit is K (Kbytes)
- For native apps, use ulimit -s (default units are Kbytes)

For worker threads: default stack size is 4 MB

- Space only needed for those local variables or automatic arrays or compiler temporaries for which each thread has a private copy
- To increase limit, export OMP_STACKSIZE=10M (or as needed)
- Or use dynamic allocation (may be less efficient)

Typical error message if stack limits exceeded:

offload error: process on the device 0 was terminated by SEGFAULT

Thread Affinity Interface

Allows OpenMP threads to be bound to physical or logical cores

- export environment variable `KMP_AFFINITY=`
 - `physical` use all physical cores before assigning threads to other logical cores (other hardware thread contexts)
 - `compact` assign threads to consecutive h/w contexts on same physical core (eg to benefit from shared cache)
 - `scatter` assign consecutive threads to different physical cores (eg to maximize access to memory)
 - `balanced` blend of compact & scatter (currently only available for Intel® MIC Architecture)
- Helps optimize access to memory or cache
- Particularly important if all available h/w threads not used
 - else some physical cores may be idle while others run multiple threads
- See compiler documentation for (much) more detail

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism: TBB

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Intel® Threading Building Blocks

Widely used C++ template library for parallelism

C++ Library for parallel programming

- Takes care of managing multitasking

Runtime library

- Scalability to available number of threads

Cross-platform

- Windows, Linux, Mac OS* and others

<http://threadingbuildingblocks.org>

Intel® Threading Building Blocks

Generic Parallel Algorithms

Efficient scalable way to exploit the power of multi-core without having to start from scratch

Concurrent Containers

Common idioms for concurrent access

- a scalable alternative serial container with a lock around it

TBB Flow Graph

Task scheduler

The engine that empowers parallel algorithms that employs task-stealing to maximize concurrency

Thread Local Storage

Scalable implementation of thread-local data that supports infinite number of TLS

Synchronization Primitives

User-level and OS wrappers for mutual exclusion, ranging from atomic operations to several flavors of mutexes and condition variables

Miscellaneous

Thread-safe timers

Threads

OS API wrappers

Memory Allocation

Per-thread scalable memory manager and false-sharing free allocators

parallel_for usage example

```
#include <tbb/blocked_range.h>
#include <tbb/parallel_for.h>
using namespace tbb;
```

```
class ChangeArray{
    int* array;
public:
    ChangeArray(int* a): array(a) {}
    void operator()(const blocked_range<int>& r) const {
        for (int i = r.begin(); i != r.end(); i++) {
            Foo (array[i]);
        }
    }
};
```

```
int main (){
    int a[n];
    // initialize array here...
    parallel_for (blocked_range<int>(0, n), ChangeArray(a));
    return 0;
}
```

ChangeArray class defines a for-loop body for parallel_for

blocked_range – TBB template representing 1D iteration space

As usual with C++ function objects the main work is done inside operator()

A call to a template function parallel_for<Range, Body>: with arguments
Range → blocked_range
Body → ChangeArray

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism: MKL

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

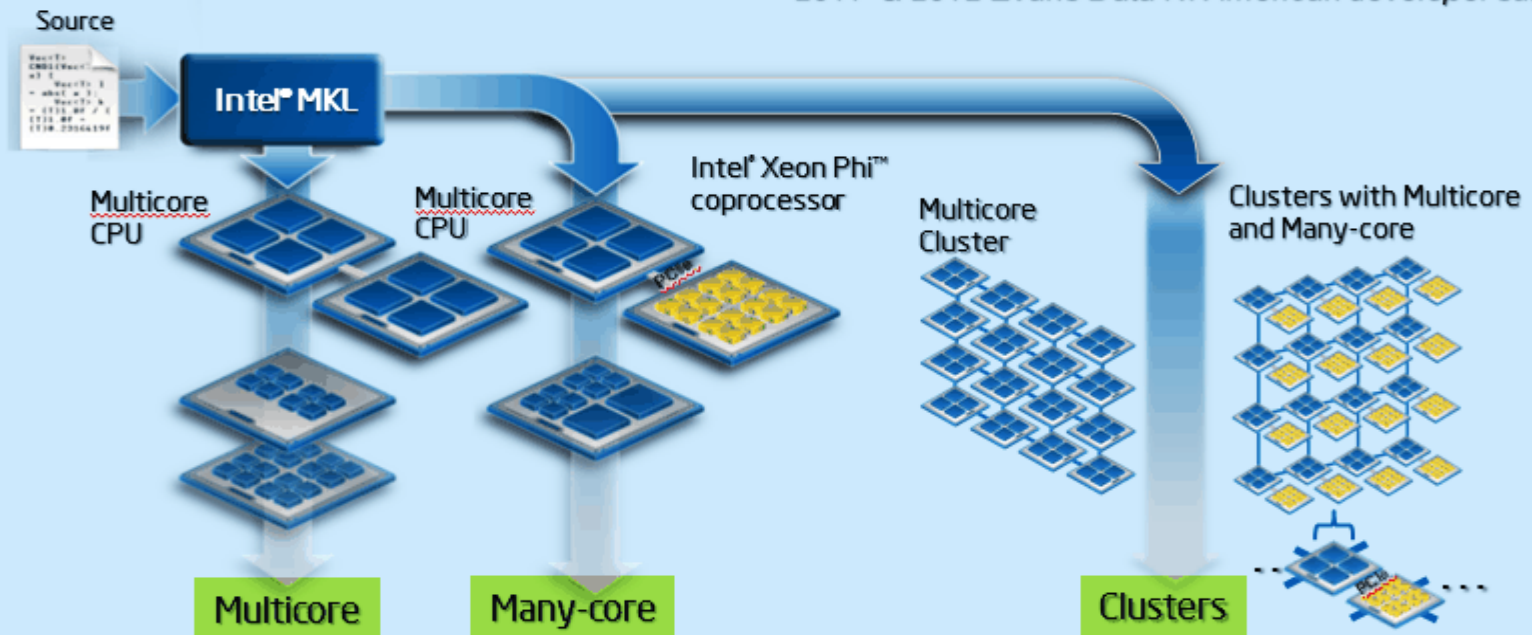
Profiling: Intel® Trace Analyzer and Collector

Conclusions

Intel® MKL is industry's leading math library *

Linear Algebra	Fast Fourier Transforms	Vector Math	Vector Random Number Generators	Summary Statistics	Data Fitting
<ul style="list-style-type: none"> • BLAS • LAPACK • Sparse solvers • ScaLAPACK 	<ul style="list-style-type: none"> • Multidimensional (up to 7D) • FFTW interfaces • Cluster FFT 	<ul style="list-style-type: none"> • Trigonometric • Hyperbolic • Exponential, Logarithmic • Power / Root • Rounding 	<ul style="list-style-type: none"> • Congruential • Recursive • Wichmann-Hill • Mersenne Twister • Sobol • Neiderreiter • Non-deterministic 	<ul style="list-style-type: none"> • Kurtosis • Variation coefficient • Quantiles, order statistics • Min/max • Variance-covariance • ... 	<ul style="list-style-type: none"> • Splines • Interpolation • Cell search

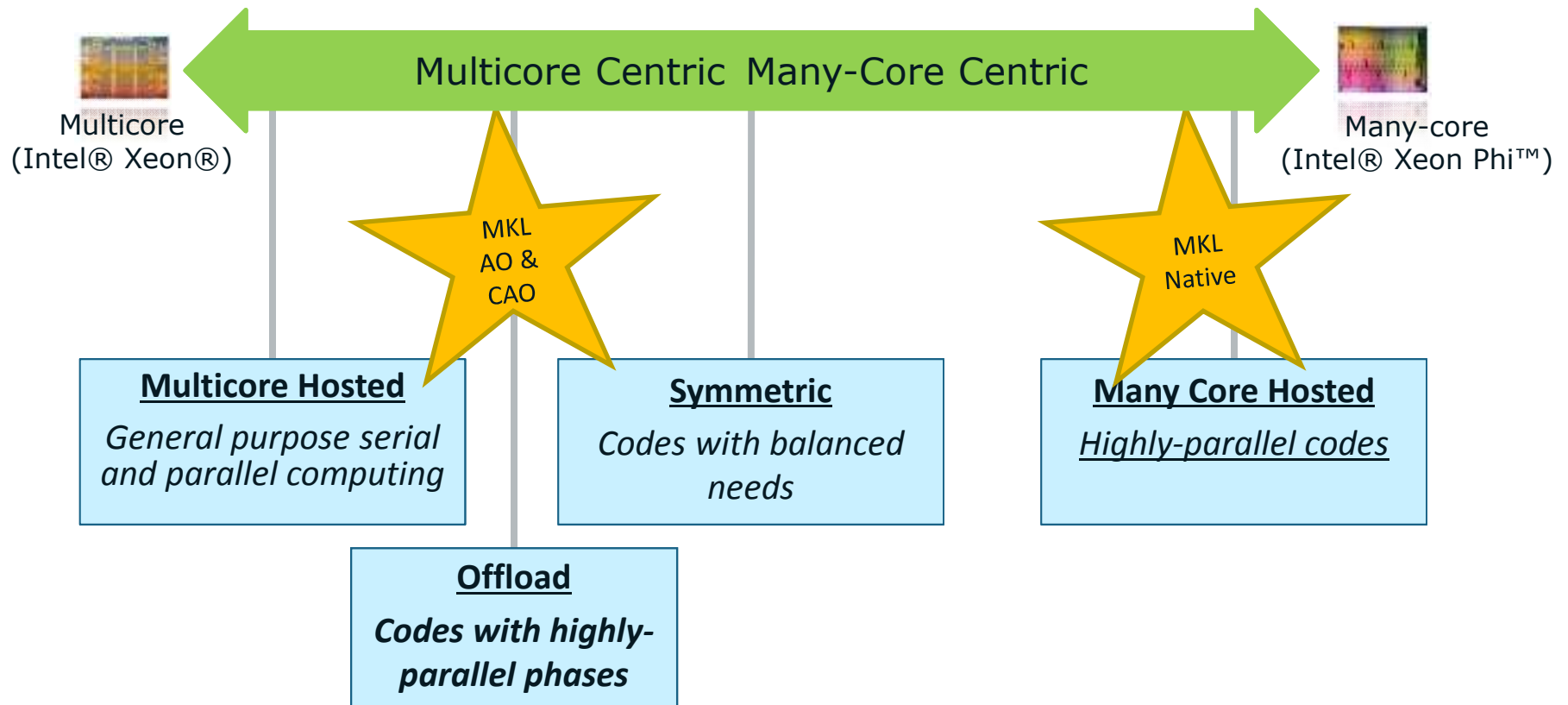
* 2011 & 2012 Evans Data N. American developer surveys



MKL Usage Models on Intel® Xeon Phi™ Coprocessor

- **Automatic Offload**
 - No code changes required
 - Automatically uses both host and target
 - Transparent data transfer and execution management
- **Compiler Assisted Offload**
 - Explicit controls of data transfer and remote execution using compiler offload pragmas/directives
 - Can be used together with Automatic Offload
- **Native Execution**
 - Uses the coprocessors as independent nodes
 - Input data is copied to targets in advance

MKL Execution Models



Work Division Control in MKL Automatic Offload

Examples	Notes
<code>mkl_mic_set_Workdivision(MKL_TARGET_MIC, 0, 0.5)</code>	Offload 50% of computation only to the 1 st card.

Examples	Notes
<code>MKL_MIC_0_WORKDIVISION=0.5</code>	Offload 50% of computation only to the 1 st card.

How to Use MKL with Compiler Assisted Offload

- The same way you would offload any function call to the coprocessor.
- An example in C:

```
#pragma offload target(mic) \  
  in(transa, transb, N, alpha, beta) \  
  in(A:length(matrix_elements)) \  
  in(B:length(matrix_elements)) \  
  in(C:length(matrix_elements)) \  
  out(C:length(matrix_elements) alloc_if(0))  
{  
    sgemm(&transa, &transb, &N, &N, &N, &alpha, A, &N, B, &N,  
          &beta, C, &N);  
}
```


Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

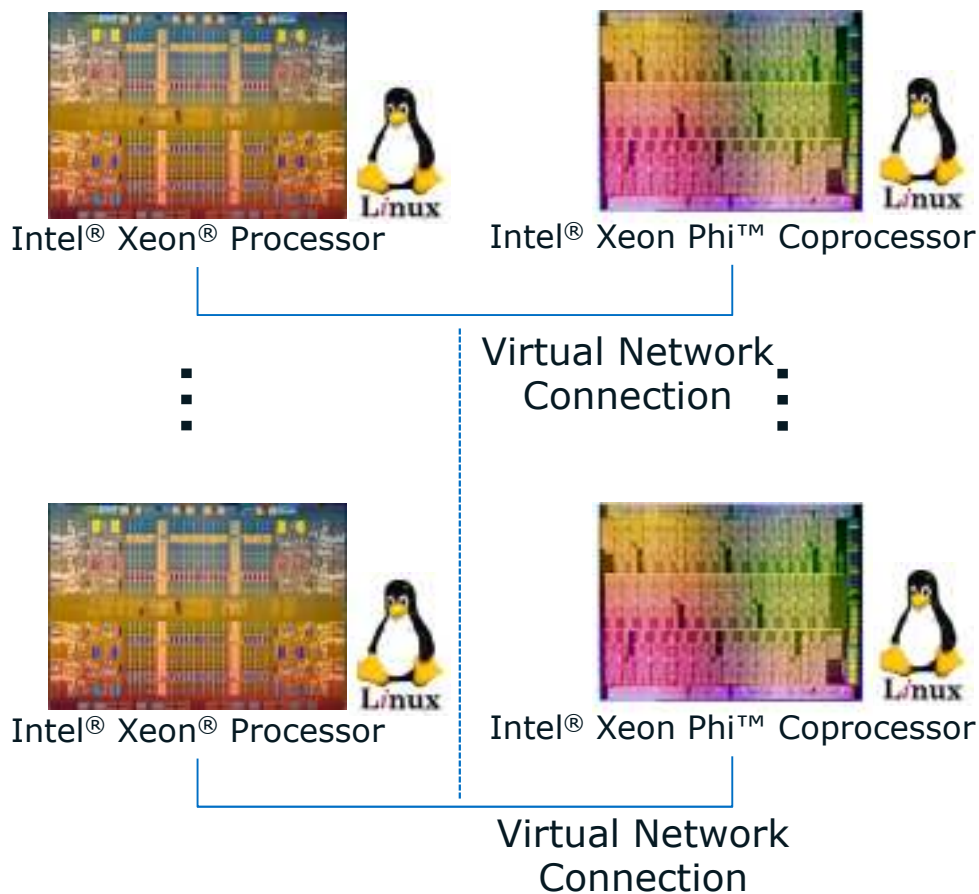
MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Intel® Xeon Phi™ Coprocessor Becomes a Network Node



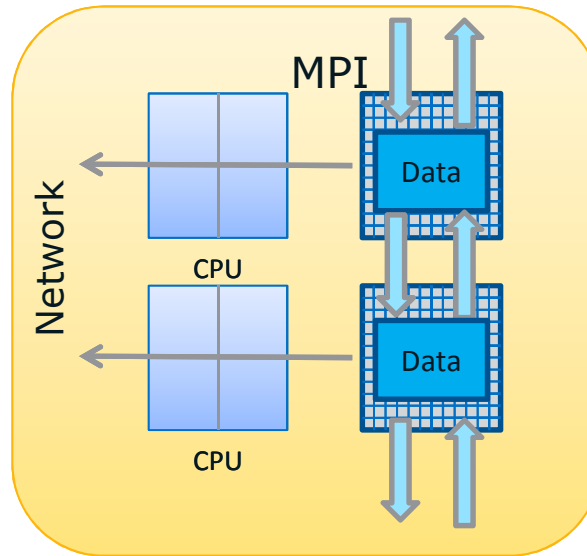
Intel® Xeon Phi™ Architecture + Linux enables IP addressability

Coprocessor only Programming Model

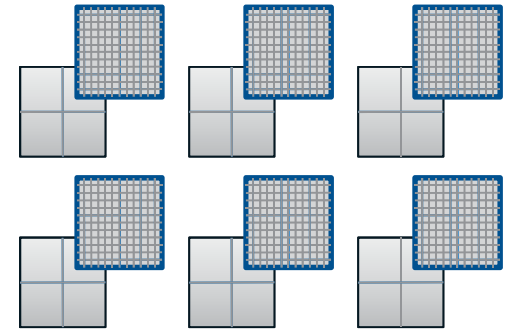
MPI ranks on Intel® Xeon Phi™ coprocessor (only)

All messages into/out of the 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 coprocessor binary using the Intel® compiler

Upload the binary to the Intel Xeon Phi coprocessor

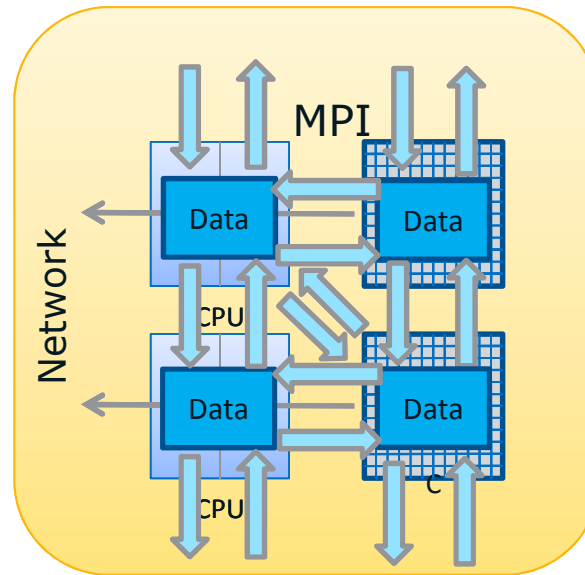
Run instances of the MPI application on Intel Xeon Phi coprocessor nodes

Symmetric Programming Model

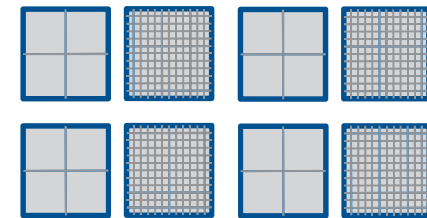
MPI ranks on Intel® Xeon Phi™ Architecture and host CPUs

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 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

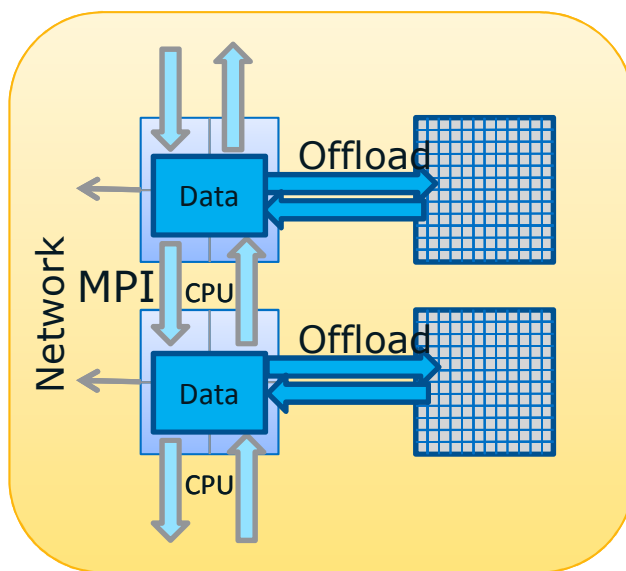
MPI+Offload Programming Model

MPI ranks on Intel® Xeon® processors (only)

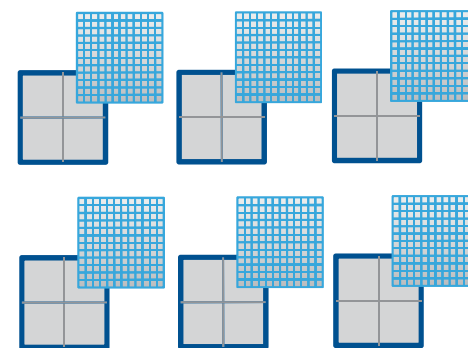
All messages into/out of host CPUs

Offload models used to accelerate MPI ranks

Intel® Cilk™ Plus, OpenMP*, Intel® Threading Building Blocks, Pthreads* within Intel® Xeon Phi™ coprocessor



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 coprocessor

Advantages of more cores and wider SIMD for certain applications

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

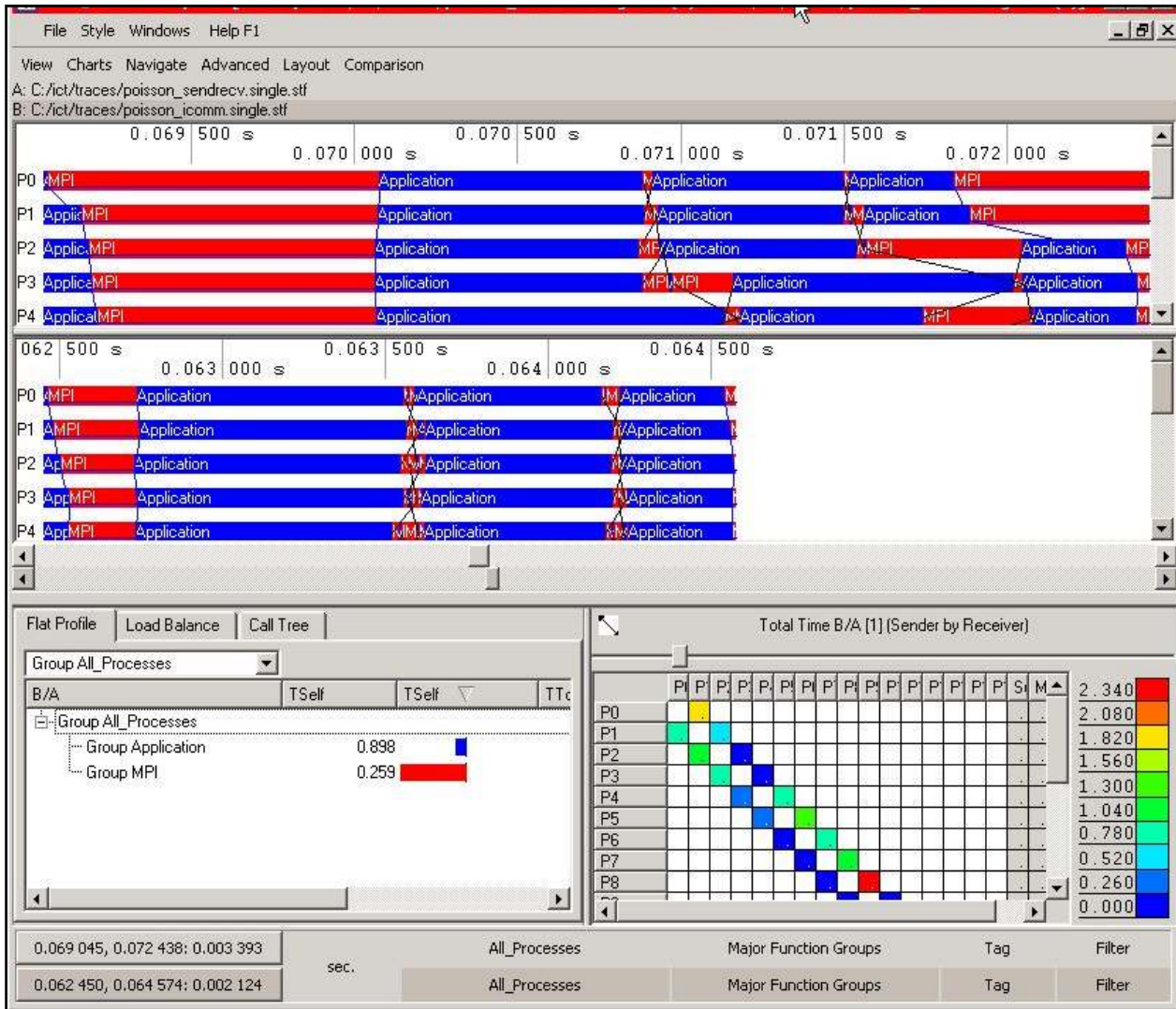
MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Intel® Trace Analyzer and Collector



Compare the event timelines of two communication profiles

Blue = computation
Red = communication

Chart showing how the MPI processes interact

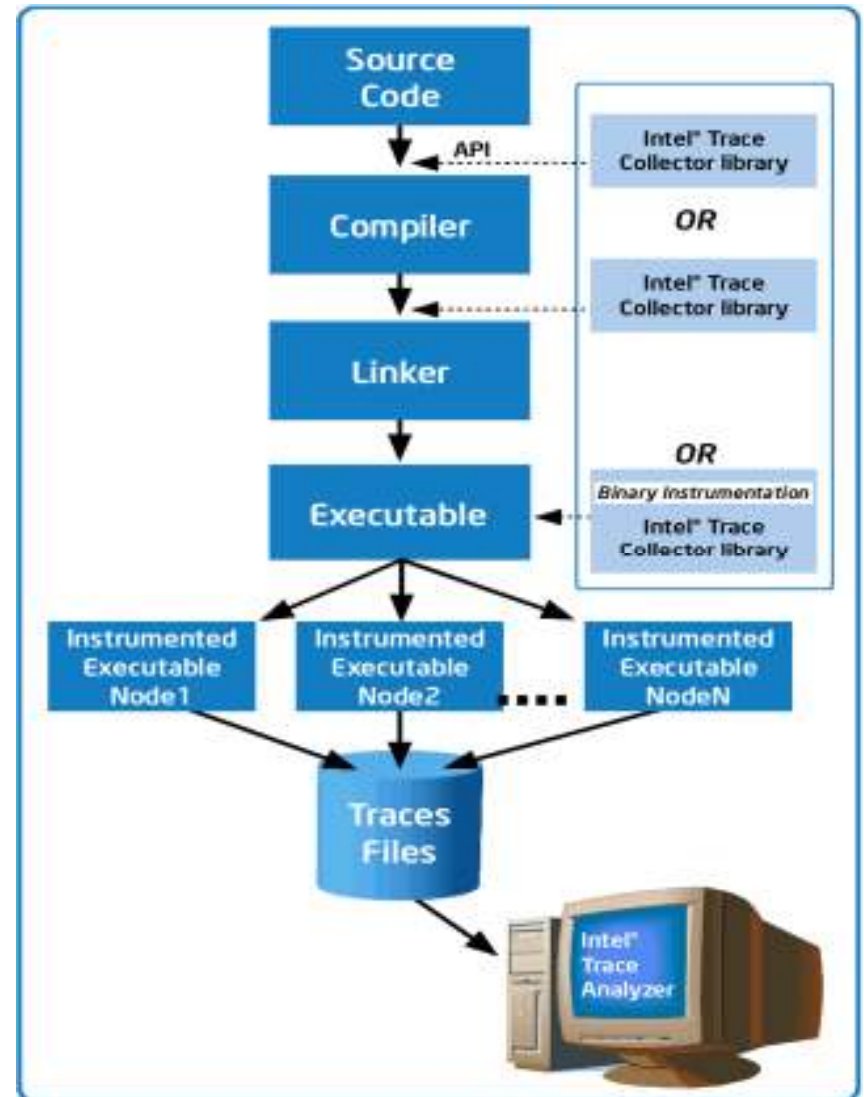
Intel® Trace Analyzer and Collector Overview

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



Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Collecting Hardware Performance Data

Hardware counters and events

- 2 counters in core, most are thread specific
- 4 outside the core (uncore) that get no thread or core details
- See PMU documentation for a full list of events

Collection

- Invoke from Intel® VTune™ Amplifier XE
- If collecting more than 2 core events, select multi-run for more precise results or the default multiplexed collection, all in one run
- Uncore events are limited to 4 at a time in a single run
- Uncore event sampling needs a source of PMU interrupts, e.g. programming cores to CPU_CLK_UNHALTED

Output files

- Intel VTune Amplifier XE performance database

Intel® VTune™ Amplifier XE offers a rich GUI

The screenshot displays the Intel VTune Amplifier XE 2011 interface. The main window shows a 'Hotspots' view with a table of CPU time usage for various functions. A call stack pane is visible on the right, showing the current stack selected. The bottom section features a timeline area with multiple tracks for threads, CPU usage, and frames over time. A filter area at the bottom allows for filtering by module, thread, and call stack mode.

Callouts from the right side of the image point to the following features:

- Menu and Tool bars
- Analysis Type
- Viewpoint currently being used
- Tabs within each result
- Grid area
- Current grouping
- Stack Pane
- Timeline area
- Filter area

Function	CPU Time	Overhead Time	Module
FireObject::checkCollision	6.542s	0ms	SystemProcedur...
dllStopPlugin	6.346s	0ms	RenderSystem_D...
TaskManagerTBB::WaitForSyst...	6.155s	0ms	Smoke.exe
FireObject::ProcessFireCollisio...	5.118s	0ms	SystemProcedur...
TaskManagerTBB::ParallelFor	2.905s	0ms	Smoke.exe
BaseThreadInitThunk	2.832s	0ms	kernel32.dll
OpenFileStreamDataStream...	2.602s	0ms	OpenMain.dll
Selected 1 row(s):		6.542s	

Intel® VTune™ Amplifier XE on Intel® Xeon Phi™ coprocessors

Adjust Data Grouping

- Function - Call Stack
- Module - Function - Call Stack
- Source File - Function - Call Stack
- Thread - Function - Call Stack
- ... (Partial list shown)

No Call Stacks Yet

Double Click Function to View Source

Filter by Timeline Selection (or by Grid Selection)



Filter by Module & Other Controls

Intel VTune Amplifier XE 2013

Lightweight Hotspots - Hotspots

Grouping: Function / Thread / H/W Context / Call Stack

Function / Thread / H/W Context / Call Stack	CPU Time	Instructions Retired	CPI Rate
p [libiomp5.so]	726.807s	145,160,000,000	5.4
p diffusion_tiled\$omp\$parallel@239	546.193s	81,720,000,000	7.2
p _do_softirq	47.174s	2,200,000,000	23.3
p _raw_spin_unlock_irq	23.651s	5,510,000,000	4.6
p diffusion_baseline	8.495s	4,250,000,000	2.1
p _raw_spin_unlock_irqrestore	3.349s	160,000,000	22.8
p _ticket_spin_lock	3.220s	1,410,000,000	2.4
p system_call_after_swaps	2.009s	110,000,000	19.9
p run_timer_softirq	1.835s	0	0.0
p hrtimer_run_pending	1.670s	0	0.0
p weighted_cpupload	1.046s	20,000,000	57.0
p libc-2.14.90.so	0.872s	130,000,000	7.3
Selected 1 row(s):	0.404s	60,000,000	7.3

Hardware Event Sample

Viewing: 1 of 1 | selected stack

- 100.0% (0.404s of 0.404s)
- vmlinux[Unknown] - sched.c
- vmlinux[Unknown] - sched.c
- vmlinux[Unknown] - sched.c

Thread

- Thread (0x0)
- Thread (0x17)
- Thread (0x16)
- Thread (0x1a)
- Thread (0x2d)
- Thread (0x1a)
- Thread (0x1a)
- Thread (0x1a)
- CPU Time

Legend:

- Running
- CPU Time
- CPU Time
- CPU Time

No filters are applied. Process: Any Process Thread: Any Thread Module: Any Module

Call Stack Mode: Only user functions. Inline Mode: on

Intel® VTune™ Amplifier XE displays event data at function, source & assembly levels

The screenshot displays the Intel VTune Amplifier XE 2013 interface. The top menu bar includes 'Analysis Target', 'Analysis Type', 'Collection Log', 'Summary', 'Bottom-up', 'Top-down Tree', 'Tasks', and 'diffusion.c'. Below the menu, there are tabs for 'Source' and 'Assembly'. The main window is divided into two panes: 'Source' on the left and 'Assembly' on the right. The 'Source' pane shows C code with columns for 'CPU Time' and 'Instructions'. The 'Assembly' pane shows assembly instructions with columns for 'Code Loc...', 'Sou...', 'Assembly', 'CP...', and 'Ins...'. A vertical scroll bar between the panes features a 'Heat Map' showing hotspots. Several callout boxes provide instructions: 'Time on Source / Asm' points to the CPU Time and Instructions columns; 'Quick Asm navigation: Select source to highlight Asm' points to the source code; 'Quickly scroll to hot spots. Scroll Bar "Heat Map" is an overview of hot spots' points to the scroll bar; 'Right click for instruction reference manual' points to a right-click action on an assembly instruction; and 'Click jump to scroll Asm' points to a blue link in the assembly pane.

Time on Source / Asm

Quick Asm navigation: Select source to highlight Asm

Quickly scroll to hot spots. Scroll Bar "Heat Map" is an overview of hot spots

Right click for instruction reference manual

Click jump to scroll Asm

Intel® VTune™ Amplifier XE 2013

Introduction

High-level overview of the Intel® Xeon Phi™ platform:
Hardware and Software

Intel Xeon Phi Coprocessor programming considerations:
Native or Offload

Performance and Thread Parallelism

MPI Programming Models

Tracing: Intel® Trace Analyzer and Collector

Profiling: Intel® Trace Analyzer and Collector

Conclusions

Conclusions: Intel® Xeon Phi™ Coprocessor supports a variety of programming models

The familiar Intel development environment is available:

- Intel® Composer: C, C++ and Fortran Compilers
- OpenMP*
- Intel® MPI Library support for the Intel® Xeon Phi™ Coprocessor
 - Use as an MPI node via TCP/IP or OFED
- Parallel Programming Models
 - Intel® Threading Building Blocks (Intel® TBB)
 - Intel® Cilk™ Plus
- Intel support for gdb on Intel Xeon Phi Coprocessor
- Intel Performance Libraries (e.g. Intel Math Kernel Library)
 - Three versions: host-only, coprocessor-only, heterogeneous
- Intel® VTune™ Amplifier XE for performance analysis
- Standard runtime libraries, including pthreads*

Intel® Xeon Phi™ Coprocessor Developer site:

<http://software.intel.com/mic-developer>



One Stop Shop for:

Tools & Software Downloads

Getting Started Development Guides

Video Workshops, Tutorials, & Events

Code Samples & Case Studies

Articles, Forums, & Blogs

Associated Product Links

<http://software.intel.com/mic-developer>

Resources

<http://software.intel.com/mic-developer>

- Developer's Quick Start Guide
- Programming Overview
- New User Forum at

<http://software.intel.com/en-us/forums/intel-many-integrated-core>

<http://software.intel.com/en-us/articles/programming-and-compiling-for-intel-many-integrated-core-architecture>

<http://software.intel.com/en-us/articles/advanced-optimizations-for-intel-mic-architecture>

Intel® Composer XE 2013 for Linux* User and Reference Guides

Intel Premier Support <https://premier.intel.com>



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 © , 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

Offloaded data have some restrictions and directives to channel their transfer

Offload data are limited to scalars, arrays, and “bitwise-copyable” structs (C++) or derived types (Fortran)

- No structures with embedded pointers (or allocatable arrays)
- No C++ classes beyond the very simplest
- Fortran 2003 object constructs also off limits, mostly
- Data exclusive to the coprocessor has no restrictions

Offload data includes all scalars & named arrays in lexical scope, which are copied both directions automatically

- IN, OUT, INOUT, NOCOPY are used to limit/channel copying
- Data not automatically transferred:
 - Local buffers referenced by local pointers
 - Global variables in functions called from the offloaded code
- Use IN/OUT/INOUT to specify these copies – use LENGTH

alloc_if() and free_if() provide a means to manage coprocessor memory allocs

Both default to true: normally coprocessor variables are created/destroyed with each offload

A common convention is to use these macros:

```
#define ALLOC alloc_if(1)
#define FREE free_if(1)
#define RETAIN free_if(0)
#define REUSE alloc_if(0)
```

To allocate a variable and keep it for the next offload

```
#pragma offload target(mic) in(p:length(n) ALLOC RETAIN)
```

To reuse that variable and keep it again:

```
#pragma offload target(mic) in(p:length(n) REUSE RETAIN)
```

To reuse one more time, then discard:

```
#pragma offload target(mic) in(p:length(n) REUSE FREE)
```