### Introduction to Intel scalable architectures

#### Fabio Affinito (SCAI - Cineca)





#### Dennard scaling law (downscaling)



The core frequency Incre and performance do not grow following the arch Moore's law any longer on the

SuperCo

Increase the number of cores to maintain the architectures evolution on the Moore's law





### **Energy trends**



### **Change of paradigm**



# The silicon lattice





50 atoms!

#### Si lattice



#### **Exascale architecture**



System attributes	2001	2010	"2015"		"2018"	
System peak	10 Tera	2 Peta	200 Petaflop/sec		1 Exaflop/sec	
Power	~0.8 MW	6 MW	15 MW		20 MW	
System memory	0.006 PB	0.3 PB	5 PB		32-64 PB	
Node performance	0.024 TF	0.125 TF	0.5 TF	7 TF	1 TF	10 TF
Node memory BW		25 GB/s	0.1 TB/sec	1 TB/sec	0.4 TB/sec	4 TB/sec
Node concurrency	16	12	O(100)	O(1,000)	O(1,000)	O(10,000)
System size (nodes)	416	18,700	50,000	5,000	1,000,000	100,000
Total Node Interconnect BW		1.5 GB/s	150 GB/sec	1 TB/sec	250 GB/sec	2 TB/sec
MTTI		day	O(1 day)		O(1 day)	



Available options...

Right here, right now... two kind of solutions are available on the market:

- IBM+ nVIDIA (Coral-like)
- Intel-based (Xeon/Xeon Phi)



# **IBM+NVIDIA**

News room > News releases >

#### IBM and NVIDIA Launch Supercomputer Centers of Excellence with the U.S. Department of Energy's Oak Ridge and Lawrence Livermore National Labs

Designed to Accelerate Scientific Research and Strengthen National Security Centers Bring Government and Industry Engineers Together to Advance Large-Scale Scientific Applications Leveraging OpenPOWER Innovations for "Summit" and "Sierra" Systems

#### Each node will be based on a Power CPU + 4/6/8 nVIDIA TESLA GPUs connected using an nVIDIA NVlink interconnect



# Intel Xeon and Xeon Phi

Intel will keep on with the production of server processors on the Xeon line, together with the introduction of the Xeon Phi many-core chips

Intel Xeon Phi will not be a co-processor anymore, but a self-standing CPU with a very high number of coresSuch systems are integrated by several vendors in many different configurations (Cray, HP, Lenovo, E4..)



# MARCONI

FERMI, the IBM BlueGene/Q deployed in Cineca ended its lifecycle in 2016

We needed a new HPC machine that could

- increase the computational power
- respect the agreements with PRACE
- satisfy the needs of the italian computing community



## MARCONI

#### Tender proposal

Partition	Installation	CPU	# nodes	# of Racks	Power
A1 – Broadwell (2.1PFlops)	April 2016	E5-2697 v4	1512	25	700KW
A2 - Knight Landing (11 Pflops)	September 2016	KNL	3600	50	1300KW
A3 – Skylake (4.5PFlops)	June 2017	E5-2680 v5	1512	25	700KW

Network: Intel OmniPath



# MARCONI

NeXtScale architecture

nx360M5 nodes:

Supporting Intel HSW & BDW

Able to host both IB network Mellanox EDR & Intel Omni-Path

Twelve nodes are grouped into a Chassis (6 chassis per rack)

The compute node is made of:

2 x Intel Broadwell (Xeon processor E5-2697 v4) 18 cores, 2,3 HGz

8 x 16GB DIMM memory (RAM DDR4 2400 MHz), 128 GB total

1 x 129 GB SATA MLC S3500 Enterprise Value SSD

Further details

1 x link OPA 100GBs

2\*18\*2,3\*16 = 1.325 GFs peak

24 rack in total:

21 rack compute

1 rack service nodes

2 racks core switch



# MARCONI - Network

Network type: new Intel Omnipath Largest Omnipath cluster of the world

Network topology: Fat-tree 2:1 oversubscription tapering at the level of the core switches only

Core Switches: 5 x OPA Core Switch "Sawtooth Forest" 768 ports each

Hdge Switch: 216 OPA Edge Switch "Eldorado Forest" 48 ports each

Maximum system configuration: 5(opa) x 768(ports) x 2(tapering) -> 7680 servers



# **MARCONI - Network**





32 nodes fully interconnected island

# **MARCONI - Storage**

Storage system: 6 x Lenovo GSS-26 Storage

Storage capacity: 17PB(raw capacity)

Storage bandwidth: 100GByte/sec (sustained)

Storage network: Intel Omnipath (directly attached to the OPA switchs)



## **CINECA - Roadmap**



## Intel<sup>®</sup> Xeon Processor Architecture



#### Intel® Xeon® Processor E5-2600 v4 Product Family - TICK



Typically, Increases in Transistor Density Enables New Capabilities, Higher Performance Levels, and Greater Energy Efficiency



CINECA

#### Intel® Xeon® E5-2600 v4 Product Family Overview





#### Intel<sup>®</sup> Xeon<sup>®</sup> Processor E5-2600 v4 Product Family MCC/LCC





### Intel<sup>®</sup> Xeon<sup>®</sup> Processor E5-2600 v4 Product Family HCC

High core count (HCC) die configuration

- Used by SKUs with 16 to 22 cores
  - E5-2699 v4
  - E5-2698 v4
  - E5-2697 v4
  - E5-2697A v4
  - E5-2695 v4
  - E5-2683 v4
- For each core
  - 2.5M last level cache (LLC)
  - Caching agent (CBO)
- For each ring

S

SuperComputing Applications and Innovation

CINECA

- Home agent (HA)
- Memory Controller with 2
   DDR4 channels



### What's next ...

- Broadwell (code name) E7 (4-socket server processor models)
- Skylake (code name) server (E5 and E7)
  - Micro-architecture launched in client processors Sep. 2015
  - Intel<sup>®</sup> AVX-512 ( only for server )
  - Expect a lot of additional, key changes
- FPGA and Xeon server integration
- NVM (non-volatile memory) 3D XPoint<sup>™</sup> Technology

# Intel<sup>®</sup> Many Integrated Core Architecture Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor



### Intel<sup>®</sup> Xeon Phi<sup>™</sup> Product Family

based on Intel<sup>®</sup> Many Integrated Core (MIC) Architecture





Continued roadmap

### Knights Landing: Next-Generation Intel® Xeon Phi™



SuperComputing Applications and Innovation

CINECA

# **KNL Mesh Interconnect**



#### **Mesh of Rings**

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

#### **Cache Coherent Interconnect**

- MESIF protocol (F = Forward)
- Distributed directory to filter snoops

#### **Three Cluster Modes**

- (1) All-to-All
- (2) Quadrant
- (3) Sub-NUMA Clustering (SNC)

SuperComputing Applications and Innovation

# Cluster Mode: All-to-All



SuperComputing Applications and Innovation

### Address uniformly hashed across all distributed directories

No affinity between Tile, Directory and Memory

Lower performance mode, compared to other modes. Mainly for fall-back

#### Typical Read L2 miss

- 1. L2 miss encountered
- 2. Send request to the distributed directory
- 3. Miss in the directory. Forward to memory
- 4. Memory sends the data to the requestor

# **Cluster Mode: Quadrant**



#### Chip divided into four virtual Quadrants

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

#### Affinity between the Directory and Memory

Lower latency and higher BW than all-to-all. Software transparent.

L2 miss, 2. Directory access, 3. Memory access,
 Data return



# Cluster Mode: Sub-NUMA Clustering (SNC)



CINECA

SuperComputing Applications and Innovation

Each Quadrant (Cluster) exposed as a separate NUMA domain to OS

#### Looks analogous to 4-Socket Xeon

Affinity between Tile, Directory and Memory

Local communication. Lowest latency of all modes

Software needs to be NUMA-aware to get benefit

L2 Directory access, 3. Memory access, 4. Data return SuperComputing Applications and Innovation

# KNL Core and VPU

Out-of-order core w/ 4 SMT threads VPU tightly integrated with core pipeline 2-wide decode/rename/retire 2x 64B load & 1 64B store port for D\$ L1 prefetcher and L2 prefetcher Fast unaligned and cache-line split support Fast gather/scatter support





## Software Adaption for KNL – Key New Features

Large impact: Intel<sup>®</sup> AVX-512 instruction set

- Slightly different from future Intel<sup>®</sup> Xeon<sup>™</sup> architecture AVX-512 extensions
- Includes SSE, AVX, AVX-2
- Apps built for HSW and earlier can run on KNL (few exceptions like TSX )
- Incompatible with 1st Generation Intel<sup>®</sup> Xeon<sup>™</sup> Phi (KNC)
- Medium impact: New, on-chip high bandwidth memory (MCDRAM) creates heterogeneous (NUMA) memory access
  - can be used transparently too however

Minor impact: Differences in floating point execution / rounding due to FMA and new HW-accelerated transcendental functions - like exp()



### AVX-512 - Greatly increased Register File







	The Intel <sup>®</sup> AVX-512 Subsets [2]
AVX-512DQ	AVX-512 Double and Quad word instructions
	<ul> <li>□ All of (packed) 32bit/64 bit operations AVX-512F doesn't provide</li> <li>□ Close 64bit gaps like VPMULLQ : packed 64x64 → 64</li> <li>□ Extend mask architecture to word and byte (to handle vectors)</li> <li>□ Packed/Scalar converts of signed/unsigned to SP/DP</li> </ul>
AVX-512BW	AVX-512 Byte and Word instructions
	<ul> <li>Extent packed (vector) instructions to byte and word (16 and 8 bit) datatype</li> <li>MMX/SSE2/AVX2 re-promoted to AVX512 semantics</li> <li>Mask operations extended to 32/64 bits to adapt to number of objects in 512bit</li> <li>Permute architecture extended to words (VPERMW, VPERMI2W,)</li> </ul>
AVX-512VL	AVX-512 Vector Length extensions
	<ul> <li>Vector length orthogonality</li> <li>Support for 128 and 256 bits instead of full 512 bit</li> <li>Not a new instruction set but an attribute of existing 512bit instructions</li> </ul>
SuperComputing Applications and Innovation	SuperComputing Applications and Innovation

### **Other New Instructions**

МРХ	Intel® MPX – Intel Memory Protection Extension
	<ul> <li>Set of instructions to implement checking a pointer against its bounds</li> <li>Pointer Checker support in HW (today a SW only solution of e.g. Intel compilers)</li> <li>Debug and security features</li> </ul>
SGX	Intel® SGX – Intel® Software Guard Extensions
	Intel® Software Guard Extensions enables applications to execute code and protect secrets from within their own protected execution environment, giving developers direct control over their application security
CLFLUSHOPT	Single Instruction – Flush a cache line
	needed for future memory technologies
XSAVE{S,C}	Save and restore extended processor state




SuperComputing Applications and Innovation

CINECA

#### Intel<sup>®</sup> Compiler Processor Switches

Switch	Description
-xmic-avx512	KNL only; already in 14.0
-xcore-avx512	Future XEON only, already in 15.0.1
-xcommon-avx512	AVX-512 subset common to both, already in 15.0.2
-m, -march, /arch	Not yet !
-ax <avx512></avx512>	Same as for "-x <avx512>"</avx512>
-mmic	No – not for KNL







#### MCDRAM: Cache vs Flat Mode



SuperComputing Applications and Innovation

CINECA

#### High Bandwidth On-Chip Memory API

- API is open-sourced (BSD licenses)
  - <u>https://github.com/memkind</u>
  - Uses jemalloc API underneath
    - <u>http://www.canonware.com/jemalloc/</u>
    - <u>https://www.facebook.com/notes/facebook-engineering/scalable-memory-allocation-using-jemall</u> oc/480222803919

Malloc replacement:

SuperComputing Applications and Innovation

```
#include <memkind.h>
    hbw_check_available()
    hbw_malloc, _calloc, _realloc,... (memkind_t kind, ...)
    hbw_free()
    hbw_posix_memalign()
    hbw_get_size(), _psize()
```

ld ... -ljemalloc -lnuma -lmemkind -lpthread

#### HBW API for Fortran, C++

Fortran:

!DIR\$ ATTRIBUTES FASTMEM :: data\_object1, data\_object2

- All Fortran data types supported
- Global, local, stack or heap; scalar, array, ...
- Support in compiler 15.0 update 1 and later versions

#### C++:

standard allocator replacement for e.g. STL like
#include <hbwmalloc.h>
std::vector<int, hbwmalloc::hbw\_allocator>



#### Porting codes on Knights Landing



## Trends that are here to stay

Data parallelism

- Lots of threads, spent on MPI ranks or OpenMP/TBB/pthreads
- Improving support for both peak tput and modest/single thread

Bigger, better, faster memory

- High capacity, high bandwidth, low latency DRAM
- Effective caching and paging
- Increasing support for irregular memory refs, modest tuning

**ISA** innovation

Increasing support for vectorization, new usages



# **Evolution or revolution?**

Incremental changes, significant gains

Parallelization – consistent strategy

- MPI vs. OpenMP already needed to tune and tweak
- Less thread-level parallelism required
- Vectorization more opportunity, more profitable

Enable new features with memory tuning

- Access MCDRAM with special allocation
- Blocking for MCDRAM vs. just cache

### Compatibility





# KNL specific enabling

- Recompilation, with –xMIC-AVX512
- Threading: more MPI ranks, 1 thread/core
- Vectorization: increased efficiency
- MCDRAM and memory tuning: tile, 1GB pages



# What is needed?

• Building

Change compiler switches in make files

• Coding

Parallelization: vectorization, offload

Memory management: MCDRAM enumeration and memory allocation

• Tuning

Potentially fewer threads: more cores but less need for SMT

More memory more MPI ranks



**Take-aways** 

Keep doing what you were doing for KNC and Xeon

#### Some goodness comes for free with a recompile

#### With some extra enabling, use new MCDRAM feature



#### Programming paradigms for the future?

#### Fabio Affinito (SCAI - Cineca)



#### What tools are there?

Differently from hetereogenous architectures (i.e. GPUs), changes can be made without being dramaticaly disruptives.
GPUs require to use CUDA, or in alternative, OpenMP4 directives to address the data copy between host and devices.
Some solutions based on PGAS is avalaible but still far to be optimal





Actually you could use codes already compiled for a x86 arch and they could work (in theory) on a Haswell/Broadwell/KNL processor.

So, at least on principle, the existence of legacy codes is completely preserved.



The fact that a code can run on a KNL architecture doesn't mean that it is really exploiting all the power offered by your machine.

It's like driving a car always with the second gear..



#### Code modernization: a must

If you don't want to waste resources (i.e. money and time) you may need to:

- modernize your code
- re-write your code from scratch (yeah, just like if you were porting it on a GPU!)



#### **Code modernization**



Stage 2: Scalar, Serial Optimization

Stage 3: Vectorization

Stage 4: Parallelization

Stage 5: Scale from Multicore to Many Core



#### Code modernization - 1st stage

At the beginning of your optimization project, select an optimizing development environment. The decision you make at this step will have a profound influence in the later steps. Not only will it affect the results you get, it could substantially reduce the amount of work to do.

The right optimizing development environment can provide you with good compiler tools, optimized, ready-to-use libraries, and debugging and profiling tools to pinpoint exactly what the code is doing at the runtime.



## Code modernization - 2nd stage

Before you begin active parallel programming, you need to make sure your application delivers the right results before you vectorize and parallelize it. Equally important, you need to make sure it does the minimum number of operations to get that correct result.

- Choosing the right floating point precision
- Choosing the right approximation method accuracy; polynomial vs. rational
- Avoiding jump algorithms

Computing Applications and Innovatio

- Reducing the loop operation strength by using iteration calculations
- Avoiding or minimizing conditional branches in your algorithms
- Avoiding repetitive calculations, using previously calculated results.

## Code modernization - 3rd stage

Try **vector level parallelism**. First try to vectorize the inner most loop. For efficient vector loops, make sure that there is minimal control flow divergence and that memory accesses are coherent.

Outer loop vectorization is a technique to enhance performance. <u>By default,</u> <u>compilers attempt to vectorize innermost loops in nested loop structures.</u> But, in some cases, the number of iterations in the innermost loop is small. In this case, inner-loop vectorization is not profitable.

However, if an outer loop contains more work, a combination of elemental functions, strip-mining, and pragma/directive SIMD can force vectorization at this outer, profitable level. (we'll see later...)



#### Code modernization - 4th step

Now we get to **thread level parallelization**. Identify the outermost level and try to parallelize it. Obviously, this requires taking care of potential data races and moving data declaration to inside the loop as necessary. It may also require that the data be maintained in a <u>cache efficient manner</u>, to reduce the overhead of maintaining the data across multiple parallel paths.

Since the amount of work needs to compensate for the overhead of parallelization, it helps to have as large a parallel effort in each thread as possible. If the outermost level cannot be parallelized due to unavoidable data dependencies, try to parallelize at the next-outermost level that can be parallelized correctly.



## 4th step: threads

- 1. If the amount of parallel work achieved at the outermost level appears sufficient for the target hardware and likely to scale with a reasonable increase of parallel resources, you are done. Do not add more parallelism, as the overhead will be noticeable (thread control overhead will negate any performance improvement) and the gains are unlikely.
- 2. If the amount of parallel work is insufficient, e.g. as measured by core scaling that only scales up to a small core count and not to the actual core count, attempt to parallelize additional layer, as outmost as possible. Note that you don't necessarily need to scale the loop hierarchy to all the available cores, as there may be additional loop hierarchies executing in parallel.
- 3. If step 2 did not result in scalable code, there may not be enough parallel work in your algorithm. This may mean that partitioning a fixed amount of work among many threads gives each thread too little work, so the overhead of starting and terminating threads swamps the useful work. Perhaps the algorithms can be scaled to do more work, for example by trying on a bigger problem size.
- 4. Make sure your parallel algorithm is cache efficient. If it is not, rework it to be cache efficient, as cache inefficient algorithms do not scale with parallelism.



#### Code modernization - 5th step

Lastly we get to **multi-node (Rank) parallelism**. To many developers message passing interface (MPI) is a black box that "just works" behind the scenes, to transfer data from one MPI task (process) to another.

The beauty of MPI for the developer is that the algorithmic coding is hardware independent. The concern that developers have, is that with the many core architecture with 60+ cores, the communication between tasks may create a communication storm either within a node or across nodes. To mitigate these communication bottlenecks, applications should employ hybrid techniques, employing a few MPI tasks and many OpenMP threads.



#### Leverage your skills!

Some time ago, codes were parallelized using **either** MPI or OpenMP.

Only in some few cases both of them were used.

#### Nowadays, using MPI+OpenMP is a must!

Two parallelization levels are necessary to exploit both the intra-node (OpenMP) and inter-node (MPI) parallelism.



#### Leverage your skills!

Plus, a lot of new features of the MPI and OpenMP standards permit to successfully deal with the new HPC architectures, in particular to exploit data-level parallelism, reducing the latencies, etc.

**MPI**: Non-blocking collectives, RMA, Shred memory parallelism, etc.

#### OpenMP: SIMD, tasks, etc.



#### Data-level parallelism

In particular for the new Intel architectures, a lot of performances depend on how much you're able to use the data-level parallelism, in terms of SIMD vectorization.

We recall, for example, that in MARCONI A1, a factor 16 of the peak-performance is given by the vectorization...



#### Overview of the vectorization techniques. Getting ready for AVX-512

#### Fabio Affinito (SCAI - Cineca)



# The need for SIMD vectorization

Is the Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessor right for me?



"Is the Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessor right for me?", by Eric Gardner -



https://software.intel.com/en-us/articles/is-the-intel-xeon-phi-coprocessor-right-for-me

SuperComputing Applications and Innovation

Single thread (ST) performance is limited in today's CPUs

- Clock frequency constraints
- Difficult to discover "near" Instruction level parallelism (ILP) by hardware

More transistors dedicated to exploit "distant" parallelism

- Task level parallelism (TLP)
  - Improves Multi Thread performance (MT)
- Data level parallelism (DLP)
  - Improves Single Thread performance (ST)
  - Enabled by using SIMD vectors

#### How to enable SIMD vectorization?

Enabling parallelism with Intel<sup>®</sup> Parallel Studio XE 2015 tool suite



Single programming model for all your code

SuperComputing Applications and Innovation

- Based on standards: OpenMP/MPI, C/C++/Fortran
- Programmers/tools responsibility to expose DLP/TLP parallelism

#### Exposing TLP/DLP in your application will benefit today and future Intel<sup>®</sup> Xeon<sup>®</sup> processors and Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessors

Including SIMD vectorization on future Intel<sup>®</sup> AVX-512 products

# Single Instruction Multiple Data (SIMD)

#### Technique for exploiting DLP on a single thread

- Operate on more than one element at a time
- Might decrease instruction counts significantly

#### Elements are stored on SIMD registers or vectors

#### Code needs to be vectorized

Vectorization usually on *inner* loops

SuperComputing Applications and Innovation

Main and *remainder* loops are generated



#### Past, present, and future of Intel SIMD types





SuperComputing Applications and inpovation VX-512 instructions, check out James Reinders' initial and updated post for this topic.

# Intel® AVX2/IMCI/AVX-512 differences

	Intel <sup>®</sup> Initial Many Core Instructions	Intel <sup>®</sup> Advanced Vector Extensions 2 <b>AVX2</b>	Intel <sup>®</sup> Advanced Vector Extensions 512 <b>AVX-512</b>
Introduction	2012	2013	2015
Products	Knights Corner	Haswell, Broadwell	Knights Landing, future Intel® Xeon® and Xeon® Phi <sup>™</sup> products
Register file	SP/DP/int32/int64 data types 32 x 512-bit SIMD registers 8 x 16-bit mask registers	SP/DP/int32/int64 data types 16 x 256-bit SIMD registers No mask registers (instr. blending)	SP/DP/int32/int64 data types 32 x 512-bit SIMD registers 8 x (up to) 64-bit mask
ISA features	Not compatible with AVX*/SSE* No unaligned data support Embedded broadcast/cvt/swizzle MVEX encoding	Fully compatible with AVX/SSE* Unaligned data support (penalty) VEX encoding	Fully compatible with AVX*/SSE* Unaligned data support (penalty) Embedded broadcast/rounding EVEX encoding
Instruction features	Fused multiply-and-add (FMA) Partial gather/scatter Transcendental support	Fused multiply-and-add (FMA) Full gather	Fused multiply-and-add (FMA) Full gather/scatter Transcendental support (ERI only) Conflict detection instructions PFI/BWI/DQI/VLE (if applies)

AVSoperComputingApplicationsandulingoration of Intel® MIC and Intel® Xeon® architecture

## Vectorization on Intel<sup>®</sup> compilers





### Auto vectorization

#### Relies on the compiler for vectorization

- No source code changes
- Enabled with -vec compiler knob (default in -02 and -03 modes)

Option	Description
-00	Disables all optimizations.
-01	Enables optimizations for speed which are know to not cause code size increase.
-02/-0 (default)	<ul> <li>Enables intra-file interprocedural optimizations for speed, including:</li> <li>Vectorization</li> <li>Loop unrolling</li> </ul>
-03	<ul> <li>Performs O2 optimizations and enables more aggressive loop transformations such as:</li> <li>Loop fusion</li> <li>Block unroll-and-jam</li> <li>Collapsing IF statements</li> <li>This option is recommended for applications that have loops that heavily use floating-point calculations and process large data sets. However, it might incur in slower code, numerical stability issues, and compilation time increase.</li> </ul>

Compiler smart enough to apply loop transformations

It will allow to vectorize more loops

5

SuperComputing Applications and Innovation

CINECA


## Vectorization: target architecture options

#### On which architecture do we want to run our program?

Option	Description			
-mmic	Builds an application that runs natively on Intel <sup>®</sup> MIC Architecture.			
<u>-xfeature</u> <u>-xHost</u>	<ul> <li>Tells the compiler which processor features it may target, referring to which instruction sets and optimizations it may generate (not available for Intel® Xeon Phi<sup>TM</sup> architecture). Values for <i>feature</i> are:</li> <li>COMMON-AVX512 (includes AVX512 FI and CDI instructions)</li> <li>MIC-AVX512 (includes AVX512 FI, CDI, PFI, and ERI instructions)</li> <li>CORE-AVX512 (includes AVX512 FI, CDI, BWI, DQI, and VLE instructions)</li> <li>CORE-AVX2</li> <li>CORE-AVX-I (including RDRND instruction)</li> <li>AVX</li> <li>SSE4.2, SSE4.1</li> <li>ATOM_SSE4.2, ATOM_SSSE3 (including MOVBE instruction)</li> <li>SSE3, SSE3, SSE2</li> </ul> When using -*Host, the compiler will generate instructions for the highest instruction set available on the compilation host processor.			
<u>-axfeature</u>	Tells the compiler to generate multiple, feature-specific auto-dispatch code paths for Intel <sup>®</sup> processors if there is a performance benefit. Values for <i>feature</i> are the same described for <code>-xfeature</code> option. Multiple features/paths possible, e.g.: <code>-axSSE2</code> , AVX. It also generates a baseline code path for the default case.			



### Auto vectorization: not all loops will vectorize

Data dependencies between iterations

- Proven Read-after-Write data (i.e., loop carried) dependencies
- Assumed data dependencies
  - Aggressive optimizations (e.g., IPO) might help

Vectorization won't be efficient

- Compiler estimates how better the vectorized version will be
- Affected by data alignment, data layout, etc.

Unsupported loop structure

- While-loop, for-loop with unknown number of iterations
- Complex loops, unsupported data types, etc.
- (Some) function calls within loop bodies
  - Not the case for SVML functions

**RaW dependency** 

for (int i = 0; i < N; i++)
 a[i] = a[i-1] + b[i];</pre>

Inefficient vectorization

Function call within loop body



### Auto vectorization on Intel<sup>®</sup> compilers

Vectorization breakdown for loop candidates in Polyhedron benchmark suite



#### SuperComputing Applications and Innovation

100%

CINECA

### Validating vectorization success

#### Generate compiler report about optimizations

- -qopt-report[=n]
- -qopt-report-file=<fname>
- -qopt-report-phase=<phase>

Generate report (level [1..5], default 2) Optimization report file (stderr, stdout also valid) Info about opt. phase:



```
loop
        Loop nest optimizations
        Auto-parallelization
par
        Vectorization
vec
openmpOpenMP
offload Offload
        Interprocedural optimizations
ipo
        Profile Guided optimizations
pqo
        Code generation optimizations
cq
tcollect
                Trace analyzer (MPI) collection
        All optimizations (default)
all
```

Vectorized loop



## Guided vectorization: disambiguation hints

Get rid of assumed vector dependencies

Assume function arguments won't be aliased

- C/C++: Compile with -fargument-noalias
- C99 "restrict" keyword for pointers
  - Compile with -restrict otherwise

Ignore assumed vector dependencies (compiler directive)

- C/C++: #pragma ivdep
- Fortran: !dir\$ ivdep



### Some Intel<sup>®</sup> compiler directives

Directive	Description
distribute, distribute_point	Instructs the compiler to prefer loop distribution at the location indicated.
inline	Instructs the compiler to inline the calls in question.
ivdep	Instructs the compiler to ignore assumed vector dependencies.
loop_count	Indicates the loop count is likely to be an integer.
optimization_level	Enables control of optimization for a specific function.
parallel/noparallel	Facilitates auto-parallelization of an immediately following loop; using keyword always forces the compiler to auto-parallelize; noparallel pragma prevents auto-parallelization.
[no]unroll	Instructs the compiler the number of times to unroll/not to unroll a loop
[no]unroll_and_jam	Prevents or instructs the compiler to partially unroll higher loops and jam the resulting loops back together.
unused	Describes variables that are unused (warnings not generated).
[no]vector	Specifies whether the loop should be vectorised. In case of forcing vectorization that should be according to the given <u>clauses</u> .

#### Guided vectorization: #pragma simd

#### Force loop vectorization ignoring **all** dependencies – Additional <u>clauses</u> for specify reductions, etc.

SIMD loop

SIMD function

```
void v_add(float *c, float *a, float *b)
{
    #pragma simd
    for (int i = 0; i < N; i++)
        c[i] = a[i] + b[i];
    }
    for</pre>
```

```
__declspec(vector)
void v_add(float c, float a, float b)
{
    c = a + b;
}
...
for (int i = 0; i < N; i++)
    v_add(C[i], A[i], B[i]);</pre>
```



#### Guided vectorization: #pragma simd

#### Also supported in OpenMP

- Almost same functionality/syntax
  - Use #pragma omp simd [clauses] for SIMD loops
  - Use #pragma omp declare simd [clauses] for SIMD functions
- See <u>OpenMP 4.0 specification</u> for more information



### Explicit vectorization with array notation

Express high-level vector parallel array operations

- Valid notation in Fortran since Fortran 90
- − Supported in C/C++ by Intel<sup>®</sup> compiler (<u>Cilk<sup>™</sup> Plus</u>) and GCC 4.9
  - Enabled by default on Intel<sup>®</sup> compiler, use -fcilkplus option on GCC
- No additional modifications to source code
- Most arithmetic and logic operations already overloaded
- Also built-in reducers for array sections

#### Vectorization becomes explicit

- C/C++ syntax: array-expression[lower-bound:length[:stride]]

#### Samples

a[:]	// All elements
a[2:6]	// Elements 2 to 7
a[:][5]	// Column 5
a[0:3:2]	// Elements 0,2,4

SIMD function invoked with array notation

```
__declspec(vector)
void v_add(float c, float a, float b)
{
    c = a + b;
}
...
v_add(C[:], A[:], B[:]);
```



# Improving vectorization: data layout

#### Vectorization more efficient with unit strides

- Non-unit strides will generate gather/scatter
- Unit strides also better for data locality
- Compiler might refuse to vectorize

AoS vs SoA

- Layout your data as Structure of Arrays (SoA)

Traverse matrices in the right direction

- C/C++: a[i][:], Fortran: a(:,i)
- Loop interchange might help
  - Usually the compiler is smart enough to apply it
  - Check compiler optimization report

#### Array of Structures vs Structure of Arrays

```
// Array of Structures (AoS)
struct coordinate {
   float x, y, z;
} crd[N];
   ...
for (int i = 0; i < N; i++)
   ... = ... f(crd[i].x, crd[i],y, crd[i].z);</pre>
```

Consecutive elements in memory \_\_\_\_\_

#### x0 y0 z0 x1 y1 z1 ...x(n-1) y(n-1) z(n-1)

```
// Structure of Arrays (SoA)
struct coordinate {
   float x[N], y[N], z[N];
} crd;
...
for (int i = 0; i < N; i++)
... = ... f(crd.x[i], crd.y[i], crd.z[i]);</pre>
```

Consecutive elements in memory

-x0 x1 ...x(n-1) <mark>y0 y1 ...y(n-1)</mark> z0 z1 ...z(n-1) —



# Improving vectorization: data alignment

Unaligned accesses might cause significant performance degradation

- Two instructions on current Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessor
- Might cause "false sharing" problems
  - Consumer/producer thread on the same cache line

Alignment is generally unknown at compile time

- Every vector access is potentially an unaligned access
  - Vector access size = cache line size (64-byte)
- Compiler might "peel" a few loop iterations
- In general, only one array can be aligned, though When possible, we have to
  - Align our data
    - Tell the compiler data is aligned
      - Might not be always the case





## Improving vectorization: data alignment

How to	Language	Syntax	Semantics
	C/C++	<pre>void* _mm_malloc(int size, int n)</pre>	Allocate memory on hean aligned to <i>n</i>
	C/C++	<pre>int posix_memalign   (void **p, size_t n, size_t size)</pre>	byte boundary.
	C/C++	declspec(align(n)) array	
align data	Fortran (not in common section)	<pre>!dir\$ attributes align:n::array</pre>	Alignment for variable declarations.
	Fortran (compiler option)	-align <i>n</i> byte	
	C/C++	<pre>#pragma vector aligned</pre>	Vectorize assuming all array data
tell the compiler	Fortran	!dir\$ vector aligned	accessed are aligned (may cause fault otherwise).
about it	C/C++	assume_aligned( <i>array</i> , <i>n</i> )	Compiler may assume array is aligned
	Fortran	!dir\$ assume_aligned array:n	to <i>n</i> byte boundary.



### Vectorization with multi-version loops



SuperComputing Applications and Innovation

## **Other considerations**

Loop tiling/blocking to improve data locality

- Square tiles so elements can be reused

Use streaming loads/stores to save bandwidth

- #pragma vector [non]temporal(list)
- -qopt-streaming-stores=[always|never|auto]
- -qopt-streaming-cache-evict[=n] (Intel<sup>®</sup> MIC only)

#### Tune software prefetcher

- -qopt-prefetch[=n]
- -qprefetch-distance=n1[,n2]
- #pragma [no]prefetch [clauses]

(Intel<sup>®</sup> MIC only) (Intel<sup>®</sup> MIC only)



## Low level (explicit) vectorization

#### A.k.a "ninja programming"

Vectorization relies on the programmer with some help from the compiler

Might be convenient for low level performance tuning of critical hotspots

Not portable among different SIMD architectures

SIMD C++ class		Intrinsics	Assembly		
finclude <fvec.h> 732vec4 a,b,c; a = b + c;</fvec.h>		<pre>#include <xmmintrin.h>m128 a,b,c; a = _mm_add_ps(b,c);</xmmintrin.h></pre>	<pre>m128 a,b,c; asm { movaps xmm0,b movaps xmm1,c addps xmm0,xmm1 movaps a, xmm0 }</pre>		
ntel) Intrinsics Guide	The Intel Intrinsic many Intel instruc	s Guide is an interactive reference tool for Intel intrinsic instru tions - including Intel® SSE, AVX, AVX-512, and more - withou	ctions, which are C style functions that provide access to $\times$ t the need to write assembly code.		
cnnologies	<u></u>				
	(sqrt x)?				
SSE SSE					
© 55E2	m512d _mm512_mask_rsqrt14_pd (m512d src,mmask8 k,m512d a) vrsqrt14pd				
= SSE3	m512d _mm512_maskz_rsqrt14_pd (mmask8 k,m512d a) vrsqrt14pd				
SSE4 1	m512d _mm512_rsqrt14_pd (m512d a) vrsqrt14pd				
SSE4.1	Synopsis				
	stropers				
	moldmmold_rsqrt14_pd_(mold_a) #include "zmmintrin.h"				
EMA	Instruction: vrsqrt14pd zmm {k}, zmm				
= AVX-512	CPUID Flags:	AVASIZE			
KNC	Description				
SVMI	Compute the approximate reciprocal square root of packed double-precision (64-bit) floating-point elements in a, and store the results in dst. The				
Other	maximum relative error for this approximation is less than 2^-14.				
- ould	Operation				
<b>tegories</b> Application-Targeted Arithmetic Bit Manipulation Cast	FOR j := 0 to i := dst[i ENDFOR dst[MAX:512]	7  *64 +63:i] := APPROXIMATE(1.0 / SQRT(a[i+63:i])) := 0			
Compare	m512 mm513	mask reart14 ps ( m512 src mmask16 k m	512 a) vreart14ac		
Convert	molzmmolz_mask_rsqrt14_ps (molz_src,mmask16 K,m512 a) vrsqrt14ps				
Cryptography	m512 _mm512_maskz_rsqrt14_ps (mmask16 k,m512 a)		vrsqrt14ps		
Elementary Math	m512 _mm512	<pre>2_rsqrt14_ps (m512 a)</pre>	vrsqrt14ps		



# How to get ready for Intel<sup>®</sup> AVX-512?

#### BKM: Start optimizing your application today for current generation of Intel<sup>®</sup> Xeon<sup>®</sup> processors and Intel<sup>®</sup> Xeon<sup>™</sup> Phi coprocessors

Tune your AVX-512 kernels on non-existing silicon

- Compile with latest compiler toolchains
  - Intel<sup>®</sup> compiler (v15.0): -xCOMMON-AVX512, -xMIC-AVX512, -xCORE-AVX512
  - GNU compiler (v4.9): -mavx512f, -mavx512cd, -mavx512er, -mavx512pf
- Run Intel<sup>®</sup> Software Development emulator (<u>SDE</u>)
  - Emulate (future) Intel<sup>®</sup> Architecture Instruction Set Extensions (e.g. Intel<sup>®</sup> MPX, ...)
  - Tools available for detailed analysis
    - Instruction type histogram
    - Pointer/misalignment checker
  - Also possible to debug the application while emulated

#### Summary

Programmers are mostly responsible of exposing DLP (SIMD) parallelism Intel<sup>®</sup> compilers provide sophisticated/flexible support for vectorization

- Auto, guided (assisted), and low-level (explicit) vectorization
- Based on OpenMP standards and specific directives
- Easily portable across different Intel<sup>®</sup> SIMD architectures
- Fine-tuning of generated code is key to achieve the best performance
  - Check whether code is actually vectorized
  - Data layout, alignment, remainder loops, etc.

Get ready for Intel<sup>®</sup> AVX-512 by optimizing your application today on current generation of Intel<sup>®</sup> Xeon<sup>®</sup> processors and Intel<sup>®</sup> Xeon<sup>™</sup> Phi coprocessors



### **Online resources**

Intel<sup>®</sup> Xeon Phi<sup>™</sup>

uperComputing Applications and Innovation

- <u>Developer portal</u> Programming guides, tools, trainings, case studies, etc.
- <u>Solutions catalog</u> Existing Intel<sup>®</sup> Xeon Phi<sup>™</sup> solutions for known codes

Intel<sup>®</sup> software development tools, performance tuning, etc.

- <u>Documentation library</u>
   All available documentation about Intel software
- <u>Learning lab</u>
   Learning material with Intel<sup>®</sup> Parallel Studio XE
- <u>Performance</u>
   Resources about performance tuning on Intel hardware
- <u>Forums</u> Public discussions about Intel SIMD, threading, ISAs, etc.

Other resources (white papers, benchmarks, case studies, etc.)

- <u>Go parallel</u> BKMs for Intel multi- and many-core architectures
- <u>Colfax research</u> Public
  - Publications and material on parallel programming
- <u>Bayncore</u> <u>labs</u> Research and development activities (WIP)

#### **Recommended books**





Parallel Programming with Intel Parallel Studio XE

Stephen Blair-Chappell, Andrew Stokes

SuperComputing Applications and Innovation

High performance parallelism pearls: multi-core and many-core approaches, by James Reinders and Jim Jeffers, Morgan Kaufmann, 2014 Intel® Xeon Ph



exander Supalov, Andrey Semin, Michael Klemm

open

Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessor high-performance programming, by Jim Jeffers and James Reinders, Morgan Kaufmann, 2013

Optimizing HPC applications with Intel<sup>®</sup> cluster tools, by Alexander Supalov et al, Apress, 2014

The software optimization handbook, by Aart Bik, Intel<sup>®</sup> press, 2004

*Parallel programming with Intel® Parallel Studio XE*, by Stephen Blair-Chappell and Andrew Stokes, Wrox press,

2012 Pell, Andrew Stokes 2012 SuperComputing Applications and Innovation Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor High Performance Programming



#### The Software Vectorization Handbook

Applying Multimedia Extensions for Maximum Performance Aart I.C. Bik



91

#### Introduction to the hands-on



### General layout of a HPC cluster





In order to permit to all users to work, a job scheduler is in charge to define the execution of the jobs.

On the Cineca HPC clusters, the job scheduler is PBS. PBS can work

- with the submission of a script
- or in an interactive session



#### qsub

```
[faffinit@r000u18l02 ~]$ qsub -I -A cin priorit -l select=1:ncpus=36:mpiprocs=36 -l walltime=15:00
gsub: waiting for job 62238.r000u17l01 to start
gsub: job 62238.r000u17l01 ready
                                             [faffinit@r000u18l02 ausurf]$ cat job.marconi
[faffinit@r041c03s02 ~]$
                                             #!/bin/bash
                                             #PBS -l walltime=6:00:00
                                             #PBS -l select=2:ncpus=16:mpiprocs=8
                                             #PBS -o job.out
                                             #PBS -e job.err
                                             #PBS - A cin staff
    interactive
                                             cd SPBS O WORKDIR
                                             export OMP NUM THREADS=1
                                             module load profile/phys autoload ge
                                             mpirun -np 16 $QE HOME/bin/pw.x -input input > log
               batch
                                             #module load autoload openmpi
                                             #module load profile/advanced
                                             #module load fftw/3.3.4--openmpi--1-10.3--gnu--6.1.0
                SuperComputing Applications and Inn([faffinit@r000u18l02 ausurf]$ qsub job.marconi
```

#### Workbench: Poisson

We will use as a workbench to test the Intel software tools a small program that simulate the solution of the Poisson equation on a grid.

Exercise: familiarize with the code, understanding (just a little bit) what is it doing, what are the computational kernels, etc.

Exercise: compile the code and run on one single node. Check the execution times changing the number of MPI processes, OpenMP threads, etc



**Graphic sessions** 

In order to work with Intel tools, it can be necessary to work with a graphic session rather than text-only ssh.A tool is avalaible from Cineca to work with a VNC session: RCM You can download RCM from:

http://www.hpc.cineca.it/software/rcm



Credentials etc.

# GALILEO login: login.galileo.cineca.it usernames: a08traXX where XX=21..24 account\_no = train\_cintel16

