#### Multi- and many-core computing for Physics applications

S. F. Schifano

University of Ferrara and INFN-Ferrara

X Seminar on Nuclear, Subnuclear and Applied Physics

June 2-8, 2013

Porto Conte, Alghero, Italy

June 2-8, 2013 1 / 101

#### Outline

Efficient programming of high-performance processors for physics applications.

- aspects of parallel computing (in short)
- e multi- and many-core architectures: "classic" CPUs, GP-GPUs, Xeon-Phi
- programming issues and performance results: Lattice-Boltzmann as case study
- Multicore and FPGA-based systems for L0-trigger processors



#### Background: Let me introduce myself

Development of computing systems optimized for computational physics:

- APEmille and apeNEXT: LQCD-machines
- AMchip: pattern matching processor, installed at CDF
- Janus: FPGA-based system for spin-glass simulations
- QPACE: Cell-based machine, mainly LQCD
- AuroraScience: multi-core based machine

#### APEmille e apeNEXT (2000 and 2004)





イロト イヨト イヨト イヨト

S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

## Janus (2007)

- 256 FPGAs
- 16 boards
- 8 host PC



#### From custom to commodity computing systems

- making the own processor is not a easy job
- commodity system have evolved in the right directions

< 6 ×

#### Designing an ideal LQCD-machine

Assuming computation and I/O can be performed in parallel, an **ideal** LQCD-machine should satisfy the equation:

w(n)/F = I(n,m)/b

where:

- *w*(*n*) is the number of floating point operations performed by the processor when execute a computation of size *n* (i.e. when execute the Dirac operator on a sublattice size of *n*)
- *F* is the number of floating point operations performed by the processor per clock cycle
- *I*(*n*, *m*) is the information exchange function
- *b* is the number of bits exchanged by the processor with the rest of the system per clock cycle

G. Bilardi, A. Pietracaprina, G. Pucci, F. Schifano, R. Tripiccione, *The Potential of On-Chip Multiprocessing for QCD Machines*, HiPC 2005, LNCS vol. 3769.

#### Information Exchange Function

The I(n, m) function is the **information exchange function** which defines:

the amount of of bit echanged by the processor with the rest of the system, assuming it has on board a storage element of m bits (register file, memory on-chip, cache, ...) and performs a computation of size n.

On a parallel machine the I(n, m) function is composed by two parts:

$$I(n,m) = I_{\rm lc}(n,m) + I_{\rm nb}(n,m)$$

- $I_{lc}(n, m)$ : the quantity of bits exchanged with the **local** memory
- Inb(n, m): the quantity of bits exchanged with the remote memories accessed via network

#### **Processors Comparison**

Using the balance equation we have evaluated different processor architecture

	apeNEXT	BG/L	Cell	ITANIUM2
frequency	200Mhz	700Mhz	3.2Ghz	1.6Ghz
$\lambda$	180nm	130nm	90nm	90nm
L <sub>w</sub>	64	32/64	32/64	64
F	8	4	64/8	4
m	32kb	32Mb	20Mb	72Mb
b <sub>lc</sub>	128	62.85	64	x
b <sub>nb</sub>	48	24	192	32 – <i>x</i>
ξlm	n.a.	6.07/2.28	2.27/6.06	4.04
ξмм	0.76	9.53/4.77	<b>0.61</b> /2.42	2.20

 $\xi$  is a measure of how well the processor is balanced between computing and I/O for QCD application

## QPACE Machine (2008)

- 8 backplanes per rack
- 256 nodes (2048 cores)
- 16 root-cards
- 8 cold-plates
- 26 Tflops peak double-precision
- 35 KWatt maximum power consumption
- 750 MFLOPS / Watt
- TOP-GREEN 500 in Nov.'09 and July'10



#### Use of recent processors

 QPACE has been the first attempt (in our community) to use a commodity processor interconnected by a custom network

 what I would like to discuss now is how and how well we can use recent developed processors for our applications

which issues we have to face out ?

how to program them ?

#### A "modern" CPU architecture: my point of view !



 $\ldots$  YES  $\ldots$  (the core of) a modern CPU is still based on the 1950 Von Neumann model !!

S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

June 2-8, 2013 12 / 101

#### The Von-Neumann Architecture

- Instruction fetch: obtain instruction from program-storage
- Instruction decode: identify the operation to execute
- operand fetch: locate and get operand-data
- execute: compute result value (or status)
- result store: write result into storage (for later use)
- next instruction: identify next instruction to fetch
- Ø go to step 1

#### J. Backus

... thus programming is basically planning and detailing the enormous traffic of words through the von Neumann bottleneck, and much of that traffic concerns not significant data itself, but where to find it.

< 回 > < 回 > < 回 >

#### **CPU** performances

At beginning, CPU performances have heavily relied on hardware:

- clock frequency
- supports to optimized memory time access: e.g. levels of caches, ROB,
  ...
- supports to increase ILP: brach-predictors, out-of-order execution, ...

#### Hardware Evolution



S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

#### The Multi-core processors era begins !

Multi-core architecture allows CPU performances to scale according to Moore's law.

- Increase frequency beyond
  ≈ 3 GHz is not possible
- assembly more CPUs in a single silicon device ✓
- great impact on application performance and design X
- move challenge to exploit high-performance computing from HW to SW X



June 2-8, 2013

16/101

#### Microprocessor Transistor Counts 1971-2011 & Moore's Law

# Many Different Cores



- all large core: multi-core Intel x86 CPUs
- many small core: NVIDIA GPUs accelerators
- all small cores: MIC architectures, Intel Xeon Phi accellerator
- mixed large and small cores: Cell, AMD-Fusion, NVIDIA-Denver

3 > < 3 >

#### "Classic" CPU Architectures

The most recent Intel *classic* CPU micro-architectures is the *Sandybridge*:



8 cores, 1 shared L3-cache

< 🗇 🕨

#### "Classic" CPU Architectures

Main features:

- 6-10 (and soon more) cores
- frequency  $\approx$  3 GHz
- 3 levels of caches, 2 withing a core and 1 shared
- support for SIMD execution: AVX 256-bits
- e.g.: Xeon E5-2680 Sandybridge: 691.2/345.6 GFlops SP/DP

Programming issues:

- core parallelism
- data parallelism
- cache optimizations
- Non Uniform Memory Architecture (NUMA)

3

14 TH 14

< 6 ×

#### "Classic" CPU Architectures: Performances

- c = 8 cores
- SIMD instructions on 256-bit operands: each vector register can pack n = 4(8) double (single) precision numbers
- each core can execute two operations per clock-cycle: one add and one mul

$$P = f \times 2 \times n \times c$$



S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

#### Numa SMP Multi-socket Multi-core Systems



- Symmetric Multi-processor Architecture (SMP)
- Non Uniform Memory Architecture (NUMA)

< 6 ×

#### Accelerator: Is this a really new concept ?







June 2-8, 2013 22 / 101

#### Accelerator: today they appear much better !



June 2-8, 2013 23 / 101

イロト イポト イヨト イヨト

#### Accelerator-based systems: hardware view

At top-level a PCIe accelarator-card (GP-GPU, ...) sits inside a standard commodity server with one or two multi-core CPUs.



・ 何 ト ・ ヨ ト ・ ヨ ト

#### The Sausage Machine Model

A computer is like a sausage machine:



- ... no input-meat ... no output-sausage !!
- ... it produces results if you provide enough input-data !!

< 6 N

#### Are accellerators good sausage machines ? FPS-164 and VAX (1976):

- Floating Point: F = 11 Mflop/s, IO Rate: B = 44 MB/s
- Ratio of flops to bytes of data movement: R = 0.25 Flops / Byte
- Host-device latency: O(1) clock-cycle

Nvidia Kepler K20 and PciE (2012):

- Floating Point: F = 1170 Gflop/s (DP), IO Rate: B = 8 GB/s
- Ratio of flops to bytes of data movement: R = 146.25 Flops / Byte
- Host-device latency: 𝒪(10 − 100) clock-cycles
- Flop/s are cheap, so are provisioned in excess,
- data needs to be re-used and processed several times by the FPUs,
- smart programming techniques to hide data movement latency, e.g. recompute data instead of access memory.

э.

イロト イポト イヨト イヨト

#### Performance Evaluation: Amdhal's Law

How much can I accelerate my application ?

Amdahl's Law approximately states:

Suppose a car is traveling between two cities 60 miles apart, and has already spent one hour traveling half the distance at 30 mph. No matter how fast you drive the last half, it is impossible to achieve 90 mph average before reaching the second city. Since it has already taken you 1 hour and you only have a distance of 60 miles total; going infinitely fast you would only achieve 60 mph.

#### Accelerator and the Amdahl's Law



The speedup of an accelerated program is limited by the time needed for the host fraction of the program.



June 2-8, 2013 28 / 101

#### Accelerator Issues: the Amdahl's law

Let assume that:

• a computation has a execution time:

$$T_s = t_A + t_B, \quad t_A = P \cdot T_s, \quad t_B = (1 - P) \cdot T_s$$

 execution time of portion A can be improved by a factor N using an accelerated version of the code

• execution time of portion *B* is run on the host and remain un-parallelized.

Under this assumptions the execution time of the new code is

$$T_{\rho} = t_{A}/N + t_{B} = (P \cdot T_{s})/N + (1 - P) \cdot T_{s}$$

Then the **speedup** (a measure of how fast is the new code) is:

$$S(n) = T_s/T_p = \frac{T_s}{((P \cdot T_s)/N + (1-P) \cdot T_s)} = \frac{1}{\binom{P}{N} + (1-P)}$$

where P is the fraction of code accelerated, and N is the improving factor.

# Accelerator Issues: the Amdahl's law

Plotting the speed-up as function of *N*:



even if I improve the 3/4 of my code by large values of *N* the maximum speedup I can achieve is limited to 4!!!

S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

Accelerator Issues: Host-Device Latency

... bandwidth problems can be cured with money. Latency problems are harder because the speed of light is fixed - you can't bribe God.

Anonymous

Moving data between Host and GPU is limited by **bandwidth** and **latency**:

T(n) = l + n/B

- accelerator processor clock period is O(1)ns
- PciE latency is O(1)µs

3

不得る 不足る 不足る

## So ... what's better ? Multi-core CPUs or Accelerators

what's better to plow a ground ?





It depends on what do we need. As rule of thumb:

- Iow-latency and reasonable throughput: left
- high-througput and reasonable-latency: right

Better if you can use both !!! May be hard to program and get good efficiency !

#### Let's look on GPUs in more detail: GPU evolution



- GPUs evolve much faster in terms of raw-computing power
- Fast-growing video-game market forces innovation

June 2-8, 2013 33 / 101

( ) < ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) <

Image: A matrix

#### GPUs vs CPUs architecture



- GPUs specialized for highly data-parallel and intensive computation (exactly what rendering is about)
- more transistors devoted to data-processing rather than data caching and flow-control

June 2-8, 2013 34 / 101

#### Architecture of GP-GPU: NVIDIA GT200



S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

June 2-8, 2013 35 / 101

#### NVIDIA GT200: SM Architecture

- 64KB register file (2KB for each SP)
- 16KB shared memory
- 8KB constant cache
- 8 32-bit ALU/FPU
- 1 64-bit FMAD
- 8 branch units



#### 1 SM executes in parallel up to 8 thread and manages up to 1024 threads

June 2-8, 2013 36 / 101

14 E 5
#### NVIDIA GT200: SM Architecture



S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

#### June 2-8, 2013 37 / 101

#### **NVIDIA GT200 Specs**

Processor:

- 1.3 GHz, 936 Gflops SP, 78 Gflops DP
- 10 Texture Processor Cluster (TPC)
- 1 TPC includes 3 Streaming Multiprocessor (SM) + 1 Texture Memory
- 1 SM include 8 Streaming Processor (SP) loosely corresponding to a modern CPU-core with 8-way SIMD computing-unit
- a total of  $10 \times 3 \times 8 = 240$  threads

Memory:

- 4GB Global Memory, 512-bit, 102.4 GB/s, 400 ... 600 cycles of latency
- Constant memory 64 KB (RO)
- Texture memory 256 KB (RO, two dimensional locality)

160 Watt (typical, w/Memory), < 6 Gflops/Watt (SP, w/MUL), 1.5 K €

3

イロト イポト イヨト イヨト

#### NVIDIA GPU Architecture Evolution

	GT200	GF100	GK110
#SM	15	14	14
#cuda-core	240	448	2688
GHz	1.2	1.15	800
GF/s (SP/DP)	936 / 78	1030 / 515	3950 / 1320
mem. bits	512	384	384
mem. MHz	1107	1500	1500
mem. GB/s (ECC-off)	141.7	144	250
Watt	160	215	235

イロト イポト イヨト イヨト

### **GPU Programming Model**

- execution has an hierarchical structure:
  - a grid of blocks
  - each block is a 1-2-3 D array of threads
- host launches a grid of thread-blocks
- a CUDA kernel (program executed on the device) is executed by an array of threads





< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### Vector sum example

```
// device code
global void vadd ( double * A. double * B. double * C ) {
 int i = threadIdx.x + blockIdx.x * blockDim.x;
 C[i] = A[i] + B[i];
}
int main () {
 double A h[N], B h[N], C h[N];
 double * A d. * B d. * C d:
 srand48();
 vinit(double *A. double *B. double *C);
  // allocate and copy data on the device
 cudaMalloc((void**) &A d); cudaMalloc((void**) &B d); cudaMalloc((void**) &C d);
 cudaMemcpy(A_d, A_h, N, H2D); cudaMemcpy(B_d, B_h, N, H2D); cudaMemcpy(C_d, C_h, N, H2D);
 dim3 dimBlock(64, 1); // size of thread-block
 dim3 dimGrid(N/64. 1) : // size of block-arid
 // run kernel
 vadd <<< dimGrid, dimBlock >>> (A,B,C);
 cudaThreadSynchronize(); // wait until kernel terminates !!!!
  // copy results back to host
 cudaMemcpy(C h, C d, N, D2H);
  // feee memory device
 cudaFree(A d): cudaFree(B d): cudaFree(C d):
```

### **GPU Programming Issues**

- host-to-device latency: Amdhal's law
- memory access latency:
   \$\mathcal{O}(10^3)\$ processor cycles, run many threads to hide memory-latency
- high-data parallelism: many threads-per-block and many blocks-per-grid

### Where we are going ?



... towards a convergence between CPU and GPU architectures

S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

June 2-8, 2013 43 / 101

イロト イポト イヨト イヨト

# First attempt to merge GPU and CPU concepts: MIC architectures

MIC: Many Integrated Core Architecture



- Knights Ferry: development board
- Knights Corners: production board
- Intel Xeon-Phi: commercial board

#### Intel MIC Systems: Knights Corners

- Yet another accelerator board
- PCIe interface
- Knights Corners: 61 x86 core @ 1.2 GHz
- each core has 32KB L1 instruction cache, 32KB L1 data cache, and 256KB L2 cache
- 512-bit SIMD unit: 16 SP, 8 DP
- multithreading: 4 threads / core
- 8 MB L3 shared coherent cache
- 4-6 GB GDDR5

#### **MIC Architectures**



- cores based on Pentium architecures
- $\approx$  60 cores
- in-order architecture
- wide SIMD instructions

イロト イポト イヨト イヨト

#### **Core Architectures**

- Scalar pipeline derived from the dual-issue Pentium processor
- Fully coherent cache structure
- 4 execution threads per core
- Separate register sets per thread
- Fast access to its 256KB local subset of a coherent L2 cache.
- 32KB instruction-cache and 32KB data-cache per core
- 3-operand, 16-wide vector processing unit (VPU)
- VPU executes integer, single-precision float, and double precision
- 1024 bits wide, bi-directional (512 bits in each direction)



э.

### MIC Programming Model

#### native:

```
icc -mmic pippo.c -o pippo
```

offload:

using approriate  ${\tt pragmas}$  to mark code that will be transparently executed onto the MIC board

Programming is well integrated with many languages:

- openMP
- TBB
- Oilk
- ...

3

4 **A** N A **A** N A **A** N

#### Parallelism management

- offload a code that spanws threads
- use openMP

```
for(t = 0; t < NTHREAD; t++) {
    pthread_create(&threads[t], NULL, threadFunc, (void *) &tData[t]);
}
for(t = 0; t < NTHREAD; t++) {
    pthread_join(threads[t], NULL);
}</pre>
```

```
#pragma omp parallel private(tid)
{
   tid = omp_get_thread_num();
   theadFunc((void *) &targv[tid]);
}
```

э.

イロト イポト イヨト イヨト

#### Example: saxpy

```
C = s \times A + B
```

```
#define N 1717
void vinit (double *A, double *C) {
 int i:
 for (i=0; i<N; i++){</pre>
   A[i] = drand48(); B[i] = drand48(); C[i] = 0.0;
}
int main () {
 double A[N], B[N], C[N];
 double s;
 srand48();
 vinit(double *A, double *B, double *C);
 s = rand48();
 #pragma offload target(mic:-1) in(A,B:lenght(N)) in(s) inout(C:lenght(N))
    #pragma omp parallel for private(i)
   for ( i=0; i < N; i++ )
      C[i] = s * A[i] + B[i]
```

June 2-8, 2013 50 / 101

э.

#### Example: vector sum

```
#define N 1717
void attribute ((target(mic))) vadd ( double *A , double *B , double *C )
void vinit (double *A. double *B. double *C) {
 int i;
 for (i=0; i<N; i++){</pre>
   A[i] = drand48(); B[i] = drand48(); C[i] = 0.0;
int main () {
 double A[N], B[N], C[N];
 srand48():
 vinit(double *A, double *B, double *C);
 #pragma offload target(mic:0) in(A,B:lenght(N)) inout(C:lenght(N))
    vadd (A,B,C);
void vadd (double *A. double *B. double *C) {
#ifdef MIC
 int i;
 for (i=0: i<N: i++)</pre>
    C[i] = A[i] + B[i];
#else
 fprint(stderr, "This code is running on the host\n");
#endif
```

### **MIC Programming Issues**

#### • core parallelism:

- keep all 60 cores (1 reserver for OS) busy
- runs 2-3 (up-to) 4 threads/core is necessary to hide memory latency

#### vector parallelism:

- enable data-parallelism
- enable use of 512-bit vector instructions

#### Amdhal's law:

- transfer time between host and MIC-board not negligible
- hide transfer time overlapping computation and processing

### Case Study

Lattice Boltzmann application

- "classic" multi-core: Sandybridge
- GP-GPU
- Xeon-Phi

< ロ > < 同 > < 回 > < 回 >

#### The Lattice Boltzmann Method

- Lattice Boltzmann method (LBM) is a class of computational fluid dynamics (CFD) methods.
- Simulation of synthetic dynamics described by the discrete Boltzmann equation, instead of the Navier-Stokes equations.
- The key idea:
  - a set of virtual particles called populations arranged at edges of a discrete and regular grid
  - interacting by propagation and collision reproduce after appropriate averaging – the dynamics of fluids.

Relevant features:

- "Easy" to implement complex physics.
- Good computational efficiency on MPAs.

-

く 同 ト く ヨ ト く ヨ ト

#### The D2Q37 Lattice Boltzmann Model

- Correct treatment of:
  - Navier-Stokes equations of motion
  - heat transport equations
  - perfect gas state equation ( $P = \rho T$ )
- D2 model with 37 components of velocity
- Suitable to study behaviour of compressible gas and fluids
- optionally in presence of **combustion**<sup>1</sup> effects.

<sup>1</sup>chemical reactions turning cold-mixture of reactants into hot-mixture of burnt product.

June 2-8, 2013

55 / 101

S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

### Simulation of the Rayleigh-Taylor (RT) Instability

Instability at the interface of two fluids of different densities triggered by gravity.



A cold-dense fluid over a less dense and warmer fluid triggers an instability that mixes the two fluid-regions (till equilibrium is reached).

June 2-8, 2013 56 / 101

イロト イポト イヨト イヨト

### LBM Computational Scheme

```
foreach time_step
  foreach lattice_point
    propagate();
    collide();
    endfor
endfor
```

#### Embarrassing parallelism

All sites can be processed in parallel applying in sequence propagate and collide.

#### Challenge

Efficient implementation on computing systems to exploit a large fraction of peak performance.

14 TH 14

#### D2Q37 propagation scheme



#### D2Q37 propagation scheme



Gather 37 populations from 37 different lattice-sites,

### D2Q37 propagation



- applies to each lattice-cell,
- requires to access cells at distance 1,2, and 3,
- gathers populations at the edges of the arrows at the center point,
- performs memory accesses with sparse addressing patterns.

4 A 1

### D2Q37: boundary-conditions



- we simulate a 2D lattice with periodic-boundaries along x-direction
- top and the bottom boundary conditions are enforced:
  - ▶ to adjust some values at sites y = 0...2 and  $y = N_y 3...N_y 1$
  - e.g. set vertical velocity to zero

This step (bc) is computed before the collision step.

- collision is computed to each lattice-cell
- computational intensive: for the D2Q37 model, and requires > 7600 DP operations
- completely local: arithmetic operations require only the populations associate to the site

イモトイモト

< 6 N

#### D2Q37: version 1 and 2

We have developed two versions of the code:

#### • Version 1:

- computes propagation and collision in two separate steps;
- is used if reactive dynamics is enable
- requires computing of the divergence of the velocity field between the two steps; to do so, we need a further step in which data is gathered from memory.

#### Version 2:

- merges computation of propagation and collision in just one single step;
- saves to access memory twice and improves performances.

3

### Implementation on Sandybridge CPUs

N. sockets	2	
CPU family	Xeon E5-2680	
frequency	2.7 GHz	
cores/socket	8	
L3-cache/socket	20 MB	
Peak Perf. DP	345.6 GFlops	
Peak Memory Bw	85.3 GBytes	

- Advanced Vector Extensions (256-bit)
- Symmetric Multi-Processor (SMP) system:
  - programming view: single processor with 16-24 cores
  - memory address space shared among cores

 Non Uniform Memory Access (NUMA) system: memory access time depends on relative position of thread and data allocation.

$$T_{exe} \ge \max\left(\frac{W}{F}, \frac{I}{B}\right) = \max\left(\frac{7666}{345.2}, \frac{592}{85.312}\right) \text{ ns} = \max(22.2, 6.94) \text{ ns}$$

June 2-8, 2013

64 / 101

#### **Relevant Optimization**

Applications approach peak performance if hardware features are exploited by the code:

- core parallelism: all cores has to work in parallel, e.g. running different functions or working on different data-sets (MIMD/multi-task or SPMD parallelism);
- vector programming: each core has to process data-set using vector (streaming) instructions (SIMD parallelism);
- cache data reuse: data loaded into cache has to be reused as long as possible to save memory access;
- **NUMA control**: time to access memory depends on the relative allocation of data and threads.

3

#### Memory layout for LB : AoS vs SoA

Iattice stored as AoS:

```
typedef struct {
    double p1; // population 1
    double p2; // population 2
    ...
    double p37; // population 37
} pop_t;
pop_t lattice2D[SIZEX*SIZEY];
```

Iattice stored as SoA:

```
typedef struct {
    double p1[SIZEX*SIZEY]; // population 1 array
    double p2[SIZEX*SIZEY]; // population 2 array
    ...
    double p37[SIZEX*SIZEY]; // population 37 array
} pop_t;
pop_t lattice2D;
```

• AoS exploits cache-locality of populations of as site: relevant for computing collision

SoA exploits data locality of corresponding populations of sites: suitable for GPUs.

Two copies of the lattice are kept in memory: each step read from prv and write onto nxt.

#### **Code Optimizations**

#### • core parallelism:

- lattice split over the cores
- pthreads library to handle parallelism
- NUMA library to control allocations of data and threads

#### • instruction parallelism:

- exploiting vector instructions (AVX)
- process 4 lattice-sites in parallel

< 6 ×

#### **Core Parallelism**

Standard POSIX Linux pthread library is used to manage parallelism:

```
for ( step = 0; step < MAXSTEP; step++ ) {</pre>
 if ( tid == 0 || tid == 1 ) {
    comm():
           // exchange borders
    propagate(); // apply propagate to left - and right-border
  } else {
    propagate(); // apply propagate to the inner part
 pthread barrier wait(...):
 if ( tid == 0 )
   bc(); // apply bc() to the three upper row-cells
 if ( tid == 1 )
   bc(); // apply bc() to the three lower row-cells
 pthread barrier wait(...);
 collide(): // compute collide()
 pthread barrier wait(..);
```

## Vector Programming

Components of 4 cells are combined/packed in a AVX vector of 4-doubles

GCC and ICC vectorization by

- enabling auto-vectorization flags,
   e.g. -mAVX, -mavx
- using the \_mm256 vector type and intrinsics functions (\_mm256\_add\_pd(),...)
- using the vector\_size attribute (only GCC)



```
typedef double fourD __attribute__ ((vector_size(4*sizeof(double))));
typedef struct {
  fourD p1; // population 1
  fourD p2; // population 2
    ...
   fourD p37; // population 37
} v_pop_type;
```

June 2-8, 2013 69 / 101

#### **Collide Performance**



GCC no-autovec: 18% of peak

- GCC autovec: 31% of peak
- GCC intrinsics: 62% of peak

( ) < ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) < )
 ( ) <

< 🗇 🕨

#### Propagate Performance



memcopy: 80% of peak

- GCC autovec: 22% of peak
- GCC intrinsics: 40% of peak

-

- T - N

#### **Optimization of Propagate**



- cache data-reuse: reordering of populations allows a better CACHE-reuse and improves performances of propagate;
- NUMA control: using the NUMA library to control data and thread allocation avoids overheads in accessing memory;
- cache blocking: load the cache with a small data-subset and work on it as long as possible;
- non-temporal instructions: store data directly to memory without request of read-for-ownership, and save time.

< ロ > < 同 > < 回 > < 回 >
## **Optimization of Propagate**

Sandybridge Execution Time (ms)							
$L_x \times L_y$	Size (GB)	Base	+NUMA Ctrl	+New-Labelling	+Cache-Blocking	+NT	
256 × 8000	0.56	116.08	47.84	36.22	27.54	20.86	
256  imes 16000	1.13	234.44	95.90	72.14	55.16	41.62	
256  imes 32000	2.26	414.32	190.95	143.13	110.34	82.97	
480 imes 8000	1.06	215.92	89.83	67.76	51.42	39.04	
480  imes 16000	2.12	338.96	178.34	134.77	103.18	77.99	
480  imes 32000	4.23	711.62	356.87	269.64	205.28	156.23	
$1680\times16000$	7.41	1376.55	625.31	472.54	372.34	279.16	

Sandybridge Bandwidth (GB/s))							
L <sub>x</sub> × L <sub>y</sub> Size (GB) Base +NUMA Ctrl +New-Labelling +Cache-Blocking +N							
256  imes 8000	0.56	10.44	25.34	33.48	44.16	58.31	
256  imes 16000	1.13	10.34	25.29	33.61	44.02	58.35	
256  imes 32000	2.26	11.71	25.40	33.88	43.99	58.49	
480  imes 8000	1.06	10.53	25.31	33.55	44.34	58.40	
480  imes 16000	2.12	13.41	25.49	33.73	44.13	58.38	
480 × 32000	4.23	12.78	25.48	33.72	44.33	58.25	
1680  imes 16000	7.41	11.56	25.45	33.68	43.83	58.46	

June 2-8, 2013 73 / 101

2

イロト イポト イヨト イヨト

## **Optimization of Propagate**



version including all optimizations performs at  $\approx$  58 MB/s,  $\approx$  67% of peak and very close to memory-copy (68.5 MB/s).

## Full Code Performance Version 1

Lattice size:  $\approx$  250  $\times$  16000 cells

	NVIDIA C2050	2-WS	2-SB
propagate	29.11 ms	140.00 ms	42.12 ms
collide	154.10 ms	360.00 ms	146.00 ms
propagate	84 GB/s	17.5 GB/s	60 GB/s
collide	205.4 GF/s	88 GF/s	220 GF/s
T/site	44 ns	130 ns	46 ns
MLUps	22	7.7	21.7
Р	172 GF/s	60 GF/s	166 GF/s
R <sub>max</sub>	33%	<b>38%</b>	48%
ξ (collide)	_	1.19	1.27

- NVIDIA Tesla C2050,  $\approx$  500 GF DP,  $\approx$  144 GB/s peak (PARCFD'11)
- 2-WS: Intel dual 6-core (Westmere),  $\approx$  160 GF DP,  $\approx$  60 GB/s peak (ICCS'11)
- 2-SB: Intel dual 8-core (Sandybridge),  $\approx$  345 GF DP,  $\approx$  85.3 GB/s peak

$$\xi = \frac{P}{N_c \times \mathbf{v} \times f}$$

## Full Code Performance Version 2

Execution of propagate and collide performed in a single step.

Lattice size:  $\approx 250 \times 16000$  cells.

	NVIDIA C2050	2-WS	2-SB
propagateCollide	167.2 ms	410.0 ms	144.0 ms
propagateCollide	190 GF/s	77 GF/s	224 GF/s
T/site	40 ns	110 ns	35 ns
MLUps	25	9.3	28.2
Р	188 GF/s	72 GF/s	216 GF/s
R <sub>max</sub>	36%	45%	<b>62%</b>
$\xi$ (propColl)	-	1.05	1.29

• NVIDIA Tesla C2050,  $\approx$  500 GF DP,  $\approx$  144 GB/s peak (PARCFD'11)

- 2-WS: Intel dual 6-core (Westmere),  $\approx$  160 GF DP,  $\approx$  60 GB/s peak (ICCS'11)
- 2-SB: Intel dual 8-core (Sandybridge),  $\approx$  345 GF DP,  $\approx$  85.3 GB/s peak

Difference of  $\xi$  might be accounted to different speed of memory-controllers.

э.

不同 トイラトイラト

### Results

LBM code on CPUs supporting the new AVX instructions carefully exploiting:

• core parallelism

- vector/streaming parallelism
- cache blocking, cache data-reuse and not-temporal instruction

Results:

- AVX version improves performances of collide and <code>propagate</code> by a factor  $\approx$  2X w.r.t. the SSE
- efficiency is high: 45% 62% for the dual-socket

э.

不同 トイラトイラト

# **GPU** Implementation

- JUDGE JÜlich Dedicated Gpu Environment
- PLX system at CINECA

### Compute Nodes:

- ► 54 Compute nodes IBM System x iDataPlex dx360 M3
- node: 2 Intel Xeon X5650(Westmere) 6-core processor 2,66 GHz
- Main memory: 96 GB
- Network: IB QDR HBA
- GPU: 2 NVIDIA Tesla M2050 (Fermi) 1,15 GHz (448 cores), 3 GB memory

### Complete System:

- 648 CPU cores
- 108 GPU cards
- 5,1 TB main memory
- 62,5 Teraflops peak performance

### Node Parallelism



- lattice  $L_x \times L_y$  is split into sub-lattices  $\frac{L_x}{N_0} \times L_y$  along X-direction
- on each node *borders* of neighbor sub-lattices are replicated (ghost borders)
- on each node the sub-lattice is further split on GPUs
- nodes have been logically arranged in a ring
- X-splitting requires to exchange Y-borders with neighbours
- make easy parallelization w/o bad impacts on performance
- ... other splitting can be used.

## Host Program

The lattice is stored as a Structure of Arrays (SOA) to exploit data-coalescing.

## **CUDA Grids Layouts**

Physical lattice of  $8 \times 16$  sites: each CUDA-thread process a lattice-point.



propagate() and collide()



### **Execution Times**

Lattice size:  $252 \times 16384$ 

#threads/block	32	64	128	256	512		
ECC Enabled							
propagate() ms	33.6	28.5	29.1	29.5	29.2		
propagate() GB/s	72.7	85.8	84.0	82.8	83.6		
collide() ms	196.4	157.0	156.0	157.4	164.7		
collide() GFLOPS	161.1	201.6	202.9	201.0	192.1		
	ECC N	OT Enat	oled				
propagate() ms	30.6	21.4	22.1	22.5	22.3		
propagate() GB/s	79.8	114.1	110.4	108.5	109.4		
collide() ms	192.4	151.0	150.5	151.7	158.9		
collide() GFLOPS	164.5	209.7	210.3	208.6	199.1		

ECC is crucial for the reliability of large production runs.

(人間) トイヨト イヨト

# Flow Diagram



Host-CPU runs three threads:

- T0 and T1 manage runs on GPUs
- T2 executes communication with neighbour nodes

Each GPU runs three streams:

- S0 applies propagate to the bulk
- S1 and S2 copies borders to and from memory buffers

## Results: Single GPU Performance

D2Q37 Version V1							
GPU CPU-W CPU-S							
GFLOPS	172	60	166				
R <sub>max</sub>	33%	38%	<b>48%</b>				
T/site ns	44	130	46				
MLUps	22	7.7	21.7				
D	2Q37 V	ersion V2					
D	2Q37 V GPU	ersion V2 CPU-W	CPU-S				
GFLOPS	2Q37 V GPU 188	ersion V2 CPU-W 72	CPU-S 216				
GFLOPS R <sub>max</sub>	2Q37 V GPU 188 <b>36%</b>	ersion V2 CPU-W 72 <b>45%</b>	CPU-S 216 <b>62%</b>				
GFLOPS R <sub>max</sub> T/site ns	2Q37 V GPU 188 <b>36%</b> 40	ersion V2 CPU-W 72 <b>45%</b> 110	CPU-S 216 <b>62%</b> 35				

- GPU: NVIDIA Tesla C2050 card,  $\approx$  500 Gflops DP peak-performance
- CPU-W: dual six-core Westmere system,  $\approx$  160 Gflops DP peak-performance
- CPU-S: dual eight-core Sandybridge system,  $\approx$  345 Gflops DP peak-performance

・ 同 ト ・ ヨ ト ・ ヨ ト ・

### Results: Multi-GPU Performance



• strong regime: code runs on a lattice size  $L_x \times L_y = 1024 \times 7168$ 

• weak-regime: sub-lattice size on each node is  $L_x \times L_y = 254 \times 14464$ 

### Results

performance: the single GPU-code performs

- ► a factor ≈ 2 3 better than a dual-socket system based on 6-core Westmere CPUs;
- equal to a just launched system based on 8-core Sandybridge CPUs (Launch Date Q1'12).
- efficiency: on GPUs is lower than CPUs, but still high for a production ready code  $\approx$  30 40% of peak performance.
- strong-regime scalability: is good as the size of the local lattice is large enough to hide communication and memory-copy overheads.
- weak-regime scalability: is linear as the communications has been hidden with computation of propagate phase.
- programmability: fine-tuning CPU-program has required accurate programming efforts while on GP-GPU CUDA allows to exploit data-parallelism.

イロト 不得 トイヨト イヨト 二日

# Xeon-Phi 5110P co-processor

#cores	61
frequency	1090 MHz
memory	8 GB GDDR5
L1-cache / core	32 KB
L2-cache / core	512 KB
Peak Perf. SP/DP	$\approx 2/1$ TFlops
Peak Memory Bw	320 GBytes



< 6 N

- PCIe 16x Gen2 card (8 GB/s)
- 240 threads
- 512-bit vector FPU
- L2-cache blocks are shared among the cores
- 512-bit Advanced Vector Extensions (AVX)

$$T_{exe} \ge \max\left(rac{W}{F}, rac{l}{B}
ight) = \max\left(rac{7666}{1064}, rac{592}{320}
ight) \, \mathrm{ns} = \max(7.20, 1.85) \, \mathrm{ns}$$

### **Relevant Optimization**

Xeon-Phi peak floating-point throughput:

 $P = f \times #cores \times NopPerCycle \times NflopPerOp$ 

In our case we have:

- *f* = 1.090 GHz
- #cores = 61
- NopPerCycle = 2, one fused-multiply-add per cloc-cycle
- NflopPerOp = 16 single-, 8 -double precision

3

14 TH 14

### **Relevant Optimization**

Applications running on Xeon-Phi can approach peak performance if codes exploits relevant hardware features:

#### ocre parallelism:

all cores has to be kept active and working in parallel, e.g. running different functions or working on different data-sets (MIMD/multi-task or SPMD parallelism);

#### hyper-threading:

cores have to execute 2 or more threads (up-to 4) to keep hardware pipelines busy and hide memory accesses latency;

#### vector programming:

each core has to process data-set using vector (streaming) instructions (SIMD parallelism); in the case of Xeon-Phi up-to 8 double-precision values can be processed by each vector instructions

3

イロト 不得 トイヨト イヨト

### Core Parallelism: Threads Management

```
for(t = 0; t < NTHREAD; t++) {
    pthread_create(&threads[t], NULL, threadFunc, (void *) &tData[t]);
}
for(t = 0; t < NTHREAD; t++) {
    pthread_join(threads[t], NULL);
}</pre>
```

```
#pragma omp parallel private(tid)
{
   tid = omp_get_thread_num();
   theadFunc((void *) &targv[tid]);
}
```

June 2-8, 2013 90 / 101

3

イロト イポト イヨト イヨト

### **Vector Parallelism**



June 2-8, 2013 91 / 101

3

イロト 不得下 イヨト イヨト

### Lattice Boltzmann Propagate Benchmark



June 2-8, 2013 92 / 101

< 6 N

### Lattice Boltzmann Collide Benchmark



S. F. Schifano (Univ. and INFN of Ferrara) Multi- and many-core computing for Physics

### Lattice Boltzmann Summary

	C2050	2-WS	2-SB	Xeon-Phi	K20
propagate GB/s	84	17.5	60	94	160
$\epsilon$	58%	29%	70%	29%	64%
collide GF/s	205.4	88	220	394	506
$\epsilon$	41%	55%	63%	37%	38%
$\xi$ (collide)	_	1.19	1.27	0.76	_

$$\xi = \frac{P}{N_c \times \mathbf{v} \times f}$$

June 2-8, 2013 94 / 101

イロト イヨト イヨト イヨト

# PC-based L0 Trigger

#### L0TP = FPGA + CPU

- In primitives ⇒ to memory of the CPU
- OPU elaborates primitives
- OPU sends back *L0 trigger signals*



### **First Exercise**

measure the distribution (and maximum) Round-Trip Time (RTT) FPGA-CPU-FPGA (i.e. estimation of t(L0)).

## Test Setup: PC

- Intel Core i7 930 2.80 GHz
- 8 MB L3-cache
- 4 core
- hyper threading disabled
- northbridge X58 (rev. 13)
- QPI @ 4.8 GT/s (≈ 22.5 GB/s)
- memory 3 x 2 GB DDR3 1.067 GHz
- OS Fedora core 12, kernel 2.6.32.26-175.fc12.x86\_64



4 **A A A A A A** 

### Test Setup: FPGA

- devkit altera
- Stratix IV GX 230 (EP4SGX230N) -C2 (fast) speed-grade
- I PCI Express 8X IP core



< 🗇 🕨

-

### Max RTT Distribution: 20 sec. runs



June 2-8, 2013 98 / 101

< 31

< (17) >

# Max Round-Trip Time (MRTT)

Average/Maximum MRTT over 100 runs of 20 seconds each.

CB	NB	14.7 MHz	12.5 MHz	11.9 MHz	10.0 MHz	9.3 MHz
8192	32	16 / 313	17 / 380	11 / 118	11 / 74	9/12
	64	12 / 39	15 / 131	12 / 88	13 / 75	15 / 44
	128	17 / 45	17 / 82	16 / 80	18 / 63	18 / 42
	256	25 / 193	26 / 54	27 / 52	32 / 60	33 / 96
16384	32	13 / 365	18 / 232	12 / 182	12 / 76	12 / 71
	64	12 / 38	15 / 97	14 / 67	15 / 73	13 / 38
	128	22 / 199	19 / 92	16 / 42	18 / 84	21 / 47
	256	26 / 53	28 / 91	29 / 54	33 / 81	31 / 45
32768	64	20 / 275	17 / 150	14 / 78	19 / 74	18 / 76
	128	22 / 209	17 / 54	22 / 79	23 / 81	23 / 85
	256	29 / 209	28 / 54	31 / 90	34 / 93	38 / 90
65536	128	25 / 224	27 / 92	25 / 84	24 / 83	26 / 82
	256	33 / 200	33 / 89	35 / 91	35 / 93	38 / 92

NB and CB in unit of packets (64 Bytes), data in microsec.

3

## Conclusions

Multi-core architectures have a big inpact on programming.

- Efficient programming requires to exploit all features of hardware systems:
  - core parallelism
  - data parallelism
  - cache optimizations
  - NUMA (Non Uniform Memory Architecture) system
- Accelerators are not a panacea:
  - good for desktop-applications
  - hard to scale on large clusters

### the one million dollar question

So ... which is the best computing system to use ?

### Acknowledgments: I would like to thank

- Luca Biferale, Mauro Sbragaglia, Patrizio Ripesi University of Tor Vergata and INFN Roma, Italy
- Andrea Scagliarini, University of Barcelona, Spain
- Filippo Mantovani, University of Regensburg, Germany
- Marcello Pivanti, Sebastiano Fabio Schifano, Raffaele Tripiccione University and INFN of Ferrara, Italy
- Federico Toschi Eindhoven University of Technology The Netherlands, and CNR-IAC, Roma Italy
- Fabio Pozzati, Alessio Bertazzo, Gianluca Crimi University of Ferrara

Results presented here are developed in the framework of the INFN COKA and SUMA projects.

I would like to thank CINECA, INFN-CNAF and JSC for access to their systems.

June 2-8, 2013 101 / 101

э.

イロト 不得 トイヨト イヨト