



HPC Architectures – past ,present and emerging trends

Andrew Emerson, Cineca a.emerson@cineca.it



27/09/2016

High Performance Molecular Dynamics - HPC architectures







- □ Computational Science
- Trends in HPC technology
- □ Trends in HPC programming
 - □ Massive parallelism
 - Accelerators
 - □ The scaling problem
- Future trends
 - Memory and accelerator advances
 - □ Monitoring energy efficiency
- □ Wrap-up





Computational Science



"Computational science is concerned with constructing mathematical models and quantitative analysis techniques and using computers to analyze and solve scientific problems. In practical use, it is typically the application of computer simulation and other forms of computation from numerical analysis and theoretical computer science to *problems in various scientific disciplines."* (Wikipedia)

Computational science (with theory and experimentation),

is the "third pillar" of scientific inquiry, enabling researchers to build and test models of complex phenomena.



27/09/2016

High Performance Molecular Dynamics - HPC architectures





Computational Sciences





Computational methods allow us to study complex phenomena, giving a powerful impetus to scientific research.



The use of computers to study physical systems allows to manage phenomena

- very large

(meteo-climatology, cosmology, data mining, oil reservoir)

- very small

(drug design, silicon chip design, structural biology)

- very complex

(ffundamental physics, fluid dynamics, turbolence)

- too dangerous or expensive

(fault simulation, nuclear tests, crash analysis)







Which factors limit computer power?



we can try and increase the speed of microprocessors but ...





Memory Hierarchy





High Performance Molecular Dynamics -HPC architectures



HPC Architectures



parallelism. This can be on many levels: – Instruction level parallelism

The main factor driving performance is

- Vector processing

Computing Applications and Innovation

- Cores per processor
- Processors per node
- Processors + accelerators (for hybrid)
- Nodes in a system

Performance can also derive from device technology

- Logic switching speed and device density
- Memory capacity and access time
- Communications bandwidth and latency



Proc

Proc.

Proc.

Proc.

Proc.

Node



Proc

Proc.





HPC systems evolution in CINECA

- 1969: CDC 6600 1st system for scientific computing
- 1975: CDC 7600 1st supercomputer
- 1985: Cray X-MP / 4 8 1st vector supercomputer
- 1989: Cray Y-MP / 4 64
- 1993: Cray C-90 / 2 128
- 1994: Cray T3D 64 1st parallel supercomputer
- 1995: Cray T3D 128
- 1998: Cray T3E 256

- 2002: IBM SP4 512 2005: IBM SP5 512
- 1st MPP supercomputer 1 Teraflops
- 2006: IBM BCX 10 Teraflops
- 2009: IBM SP6
- 2012: IBM BG/Q
- 100 Teraflops 2 Petaflops
- 2016: Lennovo (Marconi) 13 Pflops

















HPC architectures/1



The are several factors that have an impact on the system architectures incl:

- 1. Power consumption has become a primary headache.
- 2. Processor speed is never enough.
- 3. Network complexity/latency is a main hindrance.
- 4. Access to memory.





HPC architectures/2



Two approaches to increasing supercomputer power, but at the same time limiting power consumption:

- 1. Massive parallelism (IBM Bluegene range).
- 2. Hybrids using accelerators (GPUs and Xeon PHIs).





IBM BG/Q



- BlueGene systems link together tens of thousands of low power cores with a fast network.
- In some respects the IBM BlueGene range represents one extreme of parallel computing

Name: Fermi (Cineca) Architecture: IBM BlueGene/Q Model: 10 racks Processor Type: IBM PowerA2, 1.6 GHz Computing Cores: 163840 Computing Nodes: 10240, 16 core each RAM: 16 GB/node, 1GB/core Internal Network: custom with 11 links -> 5D Torus Disk Space: 2.6 PB of scratch space Peak Performance: 2PFlop/s





13





Hybrid systems



- Second approach is to *"accelerate" normal processors* by adding more specialised devices to perform some of the calculations.
- The approach is not new (maths co-procs, FPGAs, video-cards etc) but became important in HPC when Nvidia launched CUDA and GPGPUs.
- Capable of more Flops/Watt compared to traditional CPUs but still relies on parallelism (many threads in the chip).



Model: IBM PLX (iDataPlex <u>DX360M3</u>) Architecture: Linux Infiniband Cluster Nodes: 274 Processors: 2 six-cores Intel Westmere 2.40 GHz per node Cores: 12 cores/node, 3288 cores in total GPU: 2 NVIDIA Tesla M2070 per node (548 in total) RAM: 48 GB/node, 4GB/core Internal Network: Infiniband with 4x QDR switches Disk Space: 300 TB of local scratch Peak Performance: 300 TFlop/s









The Eurora supercomputer was ranked 1st in the June 2013 Green500 chart.

Hybrid Systems/2

- In the last few years Intel has introduced the Xeon PHI accelerator based on MIC (Many Integrated Core) technology.
- Aimed as an alternative to NVIDIA GPUs in HPC.

Model: Eurora prototype Architecture: Linux Infiniband Cluster Processor Type:

□Intel Xeon (Eight-Core SandyBridge) E5-2658 2.10 GHz

- □Intel Xeon (Eight-Core SandyBridge) E5-2687W 3.10 GHz
- Number of cores: 1024 (compute)

Number of accelerators: 64 nVIDIA Tesla K20 (Kepler) + 64 Intel Xeon Phi (MIC)

OS: RedHat CentOS release 6.3, 64 bit

Galileo

Model: IBM NeXtScale Architecture: Linux Infiniband Cluster Nodes: 516 Processors: 2 8-cores Intel Haswell 2.40 GHz per node Cores: 16 cores/node, 8256 cores in total Accelerator: 2 Intel Phi 7120p per node on 384 nodes (768 in total) RAM: 128 GB/node, 8 GB/core Internal Network: Infiniband with 4x QDR switches Disk Space: 2.5 Pb (Total) Peak Performance: 1 PFlop





Hybrid Systems/3 - Marconi





- New Flagship system of Cineca replaces Fermi (BG/Q).
- 3 phase introduction: phase A1 already in production, A2 has arrived and is being installed.

A1: a preliminary system going into production in July 2016, based on Intel® Xeon® processor E5-2600 v4 product family (Broadwell) with a computational power of **2Pflop/s**.

A2: by the end of 2016 a new section will be added, equipped with the next-generation of the Intel Xeon Phi product family (Knights Landing), based on a many-core architecture, enabling an overall configuration of about 250 thousand cores with expected additional computational power of approximately **11Pflop/s**.

A3: finally, in July 2017, this system is planned to reach a total computational power of about 20Pflop/s utilizing future generation Intel Xeon processors (Skylake)





27/09/2016

Top500 – November 2014



99666

17

Rank	Site	System	Cores	(TFlop/s)	(TFlop/s)	(kW)	
1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E 2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P NUDT	3, 120, 000	33, 862. 7	54,902.4	17,808	
2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Gray XK7 , Opteron 0274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560, 640	17, 590.0	27,112.5	8, 209	
3	DOE/NNSA/LLNL United States	Sequoia -BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1, 572, 864	17,173.2	20, 132. 7	7,890	
4	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect Fujitsu	705, 024	10,510.0	11,280.4	12,660	
5	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786, 432	8, 586. 6	10,066.3	3,945	BG/Q GPU
6	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Gray XC30, Xeon E5 2670 8C 2.600GHz, Aries interconnect NVIDIA K20x Cray Inc.	115,984	6,271.0	7, 788. 9	2,325	Xeon PHI
0	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi SE109 Dell	462, 462	5, 1 68. 1	8, 520. 1	4,510	
8	Forschungszentrum Juelich (FZJ) Germany	JUQUEEN BlueGene/Q, Fower BQC 16C 1.600GHz, Custom Interconnect IBM	458, 752	5,008.9	5,872.0	2,301	
9	DOE/NNSA/LLNL United States	Vulcan -BlueGene/Q, Power BQC 16C 1.600GHz, Custom Interconnect IBM	393,216	4,293.3	5,033.2	1,972	
10	Government United States	Cray XC30, Intel Xeon E5-2697v2 12C 2.7GHz, Aries interconnect Cray Inc.	225, 984	3, 143. 5	4,881.3		C

High Performance Molecular Dynamics - HPC architectures



Top500 – June 2015

RANK	SITE	SYSTEM	CORES	(TFLOP/S)	(TFLOP/S)	(KW)
1	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P	3,120,000	33,862.7	54,902.4	17,808
2	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7 , Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
3	DOE/NNSA/LLNL United States	Sequoia - DiueGene/Q, Power BQC 16C 1.60 GHz, Custom ISM	1,572,864	17,173.2	20,132.7	7,890
4	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect Fujitsu	705,024	10,510.0	11,280.4	12,660
5	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786,432	8,586.6	10,066.3	3,945
6	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Chay XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect , NVIDIA K20x Cray Inc.	115,984	6,271.0	7,788.9	2,325
7	King Abdullah University of Science and Technology Saudi Arabia	Shaheen II - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries Interconnect Cray Inc.	196,608	5,537.0	7,235.2	2,834
8	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infinitend FDR, Intel Xeon Phi SE10P Sell	462,462	5,168.1	8,520.1	4,510
9	Forschungszentrum Juelich (FZJ) Germany	JUQUEEN - DiueGene/Q, Power BQC 16C 1.600GHz, Custom Interconnect IBM	458,752	5,008.9	5,872.0	2,301
10	DOE/NNSA/LLNL United States	Vulcan - BlueGene/Q, Power BQC 16C 1.600GHz, Custom Interconnect	393,216	4,293.3	5,033.2	1,972



RPEAK

POWEI

27/09/2016

High Performance Molecular Dynamics - HPC architectures 18





Top500 June 2016



77666

Rank	Site	ystem	Cores	(TFlop/s)	(TFlop/s)	(kW)
1	National Supercomputing Center in Wuxi China	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRCEP	10,649,600	93,014.6	125,435.9	15,371
2	National Super Computer Center in Guangzhou China	Tianhe-2 (MilkyWay-2) - TH-IVB- FEP Cluster, Intel Xeon 15-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P	3,120,000	33,862.7	54,902.4	17,808
3	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7, Opteron 627, 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
4	DOE/NNSA/LLNL United States	Sequoia - BlueGene/ S, P ower BQC 16C 1.60 GHz, Custom ISM	1,572,864	17,173.2	20,132.7	7,890
5	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIII 2.0GHz, Tofu interconnect Fujitsu	705,024	10,510.0	11,280.4	12,660
6	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom ISM	786,432	8,586.6	10,066.3	3,945
7	DOE/NNSA/LANL/SNL United States	Trinity - Cray XC40, Yeen 25- 2698v3 16C 2.3GHz, Aries interconnect Cray Inc.	301,056	8,100.9	11,078.9	
8	Swiss National Supercomputing Centre (CSCS) Switzerland	Piz Daint - Cray XC30, Xeon E5- 2670 8C 2.600GHz, Aries interconnect, NVIDIA K20x Cray Inc.	115,984	6,271.0	7,788.9	2,325
9	HLRS - Höchstleistungsrechenzentrum Stuttgart Germany	Hazel Hen - Cray XC40, Xeon E5- 2680v3 12C 2.5GHz, Aries interconnect Cray Inc.	185,088	5,640.2	7,403.5	
10	King Abdullah University of Science	Shaheen II - Cray XC40, Xeon E5-	196,608	5,537.0	7,235.2	2,834



27/09/2016

High Performance Molecular Dynamics - HPC architectures 19





Roadmap to Exascale (architectural trends)

Systems	2009	2011	2015	2018
System Peak Flops/'s	2 Peta	20 Peta	100-200 Peta	1 Exa
System Memory	0.3 PB	1 PB	5 PB	10 PB
Node Performance	125 GF	200 GF	400 GF	1-10 TF
Node Memory BW	25 GB/s	40 GB/s	100 GB/s	200-400 GB/s
Node Concurrency	12	32	0(100)	0(1000)
Interconnect BW	1.5 GB/s	10 GB/s	25 GB/s	50 GB/s
System Size (Nodes)	18,700	100,000	500,000	O(Million)
Total Concurrency	225,000	3 Million	50 Million	O(Billion)
Storage	15 PB	30 PB	150 PB	300 PB
I/0	0.2 TB/s	2 TB/s	10 TB/s	20 TB/s
MTTI	Days	Days	Days	O(1Day)
Power	6 MW	~10 MW	~10 MW	~20 MW



27/09/2016 High Performance Molecular Dynamics - HPC architectures 20



Parallel Software Models

CUDA C

OpenCL



- How do we program for supercomputers?
- C/C++ or FORTRAN, together with one or more of
 - Message Passing Interface (MPI)
 - OpenMP, pthreads, hybrid MPI/OpenMP
 - CUDA, OpenCL, OpenACC, compiler directives
- Higher Level languages and libraries
 - Co-array FORTRAN, Unified Parallel C (UPC), Global Arrays
 - Domain specific languages and data models
 - Python or other scripting languages







27/09/2016 High Performance Molecular Dynamics - HPC architectures 21

OpenACC



Message Passing: MPI

Main Characteristics

- Implemented as libraries
- Coarse grain
- Inter-node parallelization (few real alternatives)
- Domain partition
- Distributed Memory
- Long history and almost all HPC parallel applications use it.

Debatable whether MPI can handle millions of tasks, particularly in collective calls.

call MPI_Init(ierror) call MPI_Comm_size(MPI_Comm_World, size, ierror) call MPI_Comm_rank(MPI_Comm_World, rank,ierror) call MPI_Finalize(ierror)



Open Issues

- Latency
- OS jitter
- Scalability
 - High memory overheads (due to program replication and buffers)



Shared Memory: OpenMP



Main Characteristics

- Compiler directives
- Medium grain
- Intra-node parallelization (p-threads)
- Loop or iteration partition
- Shared memory
- For Many HPC Applications easier to program than MPI (allows incremental parallelisation)

Open Issues

- Thread creation overhead (often worse performance than equivalent MPI program)
- Memory/core affinity
- Interface with MPI



Threads communicate via variables in shared memory







Accelerator/GPGPU



Exploit massive stream processing capabilities of GPGPUs which may have thousands of cores

Sum of 1D array

```
global void GPUCode( int* input1,
int*input2, int* output, int length)
  int idx = blockDim.x * blockIdx.x +
threadIdx.x;
  if (idx < length ) {
     output[ idx ] = input1[ idx ] +
input2[ idx ];
```





NVIDIA/CUDA



Main Characteristics

- Ad-hoc compiler
- Fine grain
- offload parallelization (GPU)
- Single iteration parallelization
- Ad-hoc memory
- Few HPC Applications

Open Issues

- Memory copy (via slow PCle link)
- Standards
- Tools, debugging
- Integration with other languages





Accelerator/Xeon PHI (MIC)



The Xeon PHI co-*processor based on Intel's* Many Integrated Core (MIC) Architecture combines many cores (>50) in a single chip.

Main Characteristics

- Standard Intel compilers and MKL library functions.
- Uses C/C++ or FORTRAN code.
- Wide (512 bit) vectors
- Offload parallelization like GPU but also "native" or symmetric modes.
- Currently very few HPC Applications

ifort -mmic -o exe_mic prog.f90

Open Issues

For Knight's Corner:

- Memory copy via slow PCIe link (just like GPUs).
- Internal (ring) topology slow.
- Wide vector units need to be exploited, so code modifications probable.
- Best also with many threads





27/09/2016

Putting it all together -Hybrid parallel programming (example)



Python: Ensen	oble simulations					
	OpenMP: External loop partition					
	CUDA: assign inner loops Iteration to GPU threads					

High Performance Molecular Dynamics - HPC architectures

29



Software Crisis



Real HPC Crisis is with Software

A supercomputer application and software are usually much more long-lived than a hardware

- Hardware life typically four-five years at most.
- Fortran and C are still the main programming models

Programming is stuck

- Arguably hasn't changed so much since the 70's

Software is a major cost component of modern technologies.

- The tradition in HPC system procurement is to assume that the software is free.

It's time for a change

- Complexity is rising dramatically
- Challenges for the applications on Petaflop systems
- Improvement of existing codes will become complex and partly impossible.
- The use of O(100K) cores implies dramatic optimization effort.
- New paradigm as the support of a hundred threads in one node implies new parallelization strategies
- Implementation of new parallel programming methods in existing large applications can be painful







1975



400 Mflops

2015

2015

128 GB

128Gb



27/09/2016

CINEC/

1965

SuperComputing Applications and Innovation

8Mb

STORAGE

High Performance Molecular Dynamics - HPC architectures 31

2015

173 Gflops

(GPU)

PERFORMANCE







The problem with parallelism...

In a massively parallel context, an upper limit for the scalability of parallel applications is determined by the fraction of the overall execution time spent in non-scalable operations (Amdahl's law).



i.e. the max speedup is not dependent on N. Must minimise P if we want to many processors.





The scaling limit

20

18

16

14

0.50

0.00

0

- Most application codes do not scale up-to thousands of cores.
- Sometimes the algorithm can be improved but frequently there is a hard limit dictated by the size of the input.
- For example, in codes where parallelism is based on domain decomposition (e.g. • molecular dynamics) no. of atoms may be < no. of cores available.





5000

#cores



64

48

120



27/09/2016

High Performance Molecular Dynamics - HPC architectures 33

10000



Parallel Scaling



The parallel scaling is important because funding bodies insist on a minimum level of parallelism.

Computer System	Minimum Parallel Scaling	Max memory/core (Gb)
Curie	Fat Nodes 128 Thin Nodes 512 Hybrid 32	4 4 3
Fermi	2048 (but typically >=4096)	1
SuperMUC	512 (typically >=2048)	*
Hornet	2048	*
Mare Nostrum	1024	2Gb

* should use a substantial fraction of available memory

Minimum scaling requirements for PRACE Tier-0 computers for calls in 2013

27/09/2016 High Performance Molecular Dynamics - HPC architectures 34





Other software difficulties



- Legacy applications (includes most scientific applications) not designed with good software engineering principles. Difficult to parallelise programs with many global variables, for example.
- Memory/core decreasing.
- I/O heavy impact on performance, esp. for BlueGene where I/O is handled by dedicated nodes.
- Checkpointing and resilience.
- Fault tolerance over potentially many thousands of threads.
 - In MPI, if one task fails all tasks are brought down.





Memory and accelerator advances – things to look out for

Memory

Computing Applications and Innovation

- In HPC memory is generally either fast, small cache (SRAM) close to the CPU or larger, slower, main memory (DRAM).But memory technologies and ways of accessing it are evolving.
 - Non-volatile RAM (NVRAM). Retains information when power switched off. Includes flash and PCM (Phase Change Memory).
 - 3D Memory. DRAM chips assembled in "stacks" to provide a denser memory packing (e.g. Intel, GPU).
- **NVIDIA GPU**
 - NVLINK, high-speed link (80 Gb/s) to replace PCI-E (16 Gb/s).
 - Unified Memory between CPU and GPU to avoid separate _ memory allocations.
 - GPU + IBM Power8 for new hybrid supercomputer (OpenPower).
- Intel Xeon PHI (Knights Landing)
 - Upgrade to Knights Corner. More memory and cores, faster _ internal network and possibility to boot as standalone host.

Word Line Word Line













Energy Efficiency



- Hardware sensors can be integrated into batch systems to report the energy consumption of a batch job.
- Could be used to charge users according to energy consumed instead of resources reserved.

PowerDAM commands

Measures directly the energy in kWh (=3600 kJ). Current implementation still very experimental.

ets --system=Eurora --job=429942.node129

- EtS is: 0.173056 kWh
- Computation: 99 %
- Networking: 0 %
- Cooling: 0 %
- Infrastructure: 0 %





Energy Efficiency



Energy consumption of GROMACS on Eurora.

PBS Job id	nodes	Clock freq (GHz)	#gpus	Walltime (s)	Energy (kWh)	Perf (ns/day)	Perf- Energy (ns/kJ)
429942	1	2	0	1113	0.17306	10.9	69.54724
430337	2	2	0	648	0.29583	18.6	62.87395
430370	1	3	0	711	0.50593	17.00	33.60182
431090	1	3	2	389	0.42944	31.10	72.42023

Exercises:

- compare clock freq 2 Ghz with 3 Ghz
- clock freq 3 Ghz with and without GPU





"approximate computing"



- Energy efficiency is a big deal next-generation exaflop"machines, which are capable of 10¹⁸ operations a second, could consume as much as 100 megawatts, the output of a small power station.
- (In terms of flops, the human brain is 10,000 times more efficient.)
- Solution? Reduce the accuracy (precision) of calculations by lowering the voltage supplied to least significant bits.
- Already being used in audio-visual applications. Could become important as well in HPC, e.g. weather modelling.











- HPC is only possible via parallelism and this must increase to maintain performance gains.
- Parallelism can be achieved at many levels but because of limited code scalability with traditional cores increasing role for accelerators (e.g. GPUs, MICs). The Top500 is becoming now becoming dominated by hybrid systems.
- Hardware trends forcing code re-writes with OpenMP, OpenCL, CUDA, OpenACC, etc in order to exploit large numbers of threads.
- Unfortunately, for many applications the parallelism is determined by problem size and not application code.
- Energy efficiency (Flops/Watt) is a crucial issue. Some batch schedulers already report energy consumed and in the near future your job priority may depend on predicted energy consumption.

