

## Introduction

The INFN community in theoretical physics is active in several scientific areas that require significant computational support; these areas stretch over a wide spectrum, requiring in some cases fairly limited computing resources – for instance in nuclear physics, high-energy physics phenomenology, spin-system simulations – while, at the other end of the spectrum huge computing power, that can only be provided at the transnational level, is needed; examples in this class are Lattice Quantum Chromodynamics (LQCD), dynamical systems and classical and ab initio simulations of bio-systems; some research groups work on areas that are acknowledged grand-challenges of High Performance Computing.

At the same time, for most groups active in these areas, it is becoming more and more difficult to independently develop their computational strategies and algorithms in a way that allows to adapt to the increasingly fast changes happening in high performance computing architectures.

Last but not least, several existing INFN projects have produced significant progress on technological developments that may be crucial building blocks for new generation HPC systems, if it can be shown that they are efficient solutions to (at least some) large scale computational problems.

SUMA plans to support all these physics goals, by providing computing resources - both in-house and through access programmes, helping develop and share the know-how needed to efficiently use new computing systems and validating and assessing performances, and at the same time aims to explore all suitable ways in which the technological developments made at INFN can be put to good use for the present and future needs of computational physics.

The SUMA project works in close collaboration with academia, research centers and computer centers in Italy, such as the Universities of Ferrara, Parma, Pisa and Rome, SISSA (Trieste) and CINECA (Bologna).

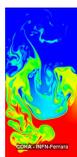
## Computational Theoretical Physics at INFN

Theoretical Physics at INFN is heavily supported by numerical simulations. In several cases INFN computational applications are recognized grand-challenges of high-performance computing.

INFN groups are active in the following projects:

### Lattice Quantum Chromodynamics (LQCD):

studying the structure and the dynamics of matter at its tiniest scales.



### Computational fluid-dynamics (CFD):

discovering the statistical properties of fluids in the turbulent regime.

### Complex systems dynamics:

understanding the behaviour of systems governed by conflicting interactions.



### Quantitative biology (QBIO):

understanding how fundamental physics interactions shape the building blocks of life.

## Work Horses

INFN researchers use a variety of HPC computer systems supporting their investigations.

INFN operates a number of Tier1 HPC clusters, and uses Tier0 facilities made available by the PRACE access program of the European Union, as well as the Blue Gene/Q system installed at CINECA in Bologna.



INFN Tier1 cluster installed in Pisa.



The Fermi BG/Q system installed at CINECA.

INFN carry on several experimental projects to optimize codes, test and assess performances on new architectures, including GPUs, MICs and FPGAs.



The FPGA-based Janus II system.



The Eurora system installed at CINECA.

## QCD on GPUs

The goal of LQCD is to compute numerically, by Monte-Carlo simulations, the theory of Quarks and Gluons, Quantum Chromodynamics, on a discretized space-time Lattice. The computational task is extraordinary, typically a  $N \times N$  sparse matrix must be inverted  $O(10^6)$  times, with  $N$  going up to  $10^8-10^9$ . For that reason, since the '80s LQCD has fostered the development of parallel HPC (e.g., APE machines, BlueGene machines).

The marriage between LQCD and GPUs was unavoidable. Pioneers entered the (video)game in 2006, at the time of OpenGL [1]. More friendly programming frameworks have then made GPUs widespread in our community.

OUR ACTIVITY: in 2009 we started the development of a production code entirely running on GPUs [2]

- ▶ Single GPU version: to cut CPU-GPU memory transfer, whole Molecular Dynamics (MD) runs on GPU.
- ▶ Mixed precision strategy: MD in single precision, Metropolis test and measurements in double.
- ▶ Performance is limited by internal memory bandwidth.  $O(100)$  sustained Gflops attained on single GPU.
- ▶ Code already in production on GPU farms (Pisa, Rome) for studying QCD matter in extreme conditions (Universe right after the Big Bang, heavy ion collision experiments) [3]
- ▶ Extension to many-GPU parallel architectures essential for studies with larger RAM requirement. Main challenge is internode communication during MD. Sinergy between communication-masking strategies at the programming level and custom interconnection architectures will be essential.

## New Programming Approaches

Code refactoring to get extreme performances for any new architecture is not always the only choice.

We are testing the effectiveness of simpler solutions provided by modern compilers, e.g.

- ▶ Intel ICC makes array notation available, to help express vector parallelism within codes:

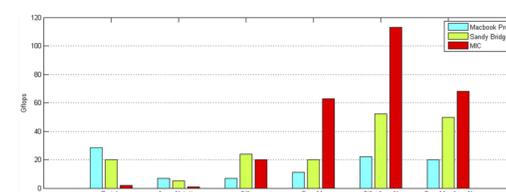
$$A[0:N] = B[0:N] + C[0:N]$$

- ▶ OpenMP is well established for multithreading, `cilk` is also available in icc.

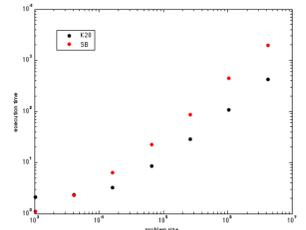
OpenACC provides support for GPU multithreading in much the same style.

All this open the way to easy, portable codes which can be as simple as in the following example:

```
# ifdef MYOpenMP
# pragma omp parallel for
# endif
# ifdef MYOpenACC
# pragma acc parallel loop present ( latti )
# endif
# ifdef MYcilk
cilk_for ( int i=0; i < setup.sz2 ; i++) {
# else
for ( int i=0; i < setup.sz2 ; i++) f
# endif
latti[y].theta1 = latti[i].theta + dt*latti[i].pi ;
}
```



(Big) Matrix multiplications on different architectures, with different combination of array notation and OpenMP or cilk

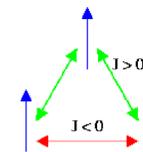
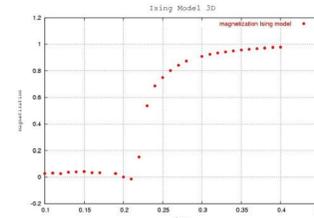


Hybrid MonteCarlo simulations of spin systems: the code is much the same for the nVidia K20 GPU (OpenACC pragmas) and for the Intel Xeon SB (OpenMP pragmas)

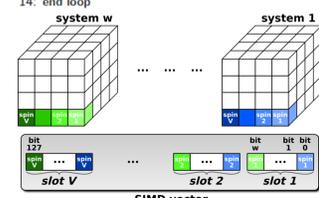
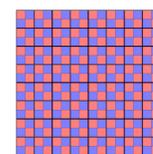
## Spin-Glass on Many-Core and GPUs

The Spin-glass is a statistic model to study behaviours of complex macroscopic systems like

disordered magnetic materials



```
Require: set of {S} and {J}
1: loop {loop on Monte Carlo steps}
2: for all  $s_i \in \{S\}$  do
3:  $s'_i = (s_i == 1) ? -1 : 1$  {flip tentatively value of  $s_i$ }
4:  $\Delta E = \sum_{ij} (J_{ij} \cdot s'_i \cdot s_j) - (J_{ij} \cdot s_i \cdot s_j)$  {compute energy change}
5: if  $\Delta E \leq 0$  then
6:  $s_i = s'_i$  {accept new value of  $s_i$ }
7: else
8:  $\rho = \text{rnd}()$  {compute a random number  $0 \leq \rho \leq 1$ ,  $\rho \in \mathbb{Q}$ }
9: if  $\rho < e^{-\beta \Delta E}$  then  $\{\beta = 1/T, T = \text{Temperature}\}$ 
10:  $s_i = s'_i$  {accept new value of  $s_i$ }
11: end if
12: end if
13: end for
14: end loop
```



System	Core 2 Duo	CBE (16 cores)	Janus	C1060	NH (8 cores)	C2050	SB (16 cores)	K20X	Xeon-Phi	Janus 2
Year	2007	2007	2008	2009	2009	2010	2012	2012	2013	2013
Power (W)	150	220	35	200	220	300	300	300	300	25
SUT (ps/flip)	1000	150	16	720	200	430	60	230	52	2
Energy/flip (nJ/flip)	150	33	0.56	144	244	129	18	69	15.6	0.05

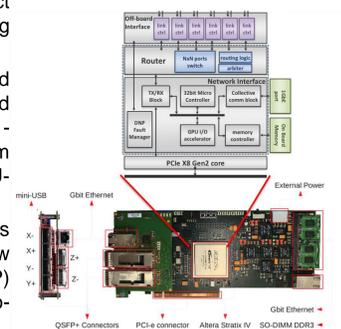
Spin-update-time (SUT) of EA simulation codes on a  $64^3$  lattice on several architectures.

## Networks for HPC

In the race towards exaFLOPS systems one of the most critical aspects is the design of smart, efficient and robust network able to interconnect the huge number of many-core high performance processors equipping modern HPC platforms.

Leveraging on past APE experiences - 3D Torus network for specialized systems, proven effective for large scale scientific computing - and taking into account the emerging of powerful many-core architectures - mainly GPU processor - we designed APENet+. APENet+ is a custom 3D torus interconnection architecture optimized for hybrid clusters CPU-GPU [4].

The APENet+ interconnect fabric is built on a FPGA-based PCI-express board with 6 bi-directional off-board links showing 34 Gbps of raw bandwidth per direction, and leverages upon peer-to-peer (P2P) capabilities of Fermi and Kepler-class NVIDIA GPUs to obtain real zero-copy, GPU-to-GPU low latency transfers.



The minimization of APENet+ transfer latency is achieved through the adoption of a simple RDMA protocol implemented in FPGA with specialized hardware blocks tightly coupled with embedded microprocessor.

In the framework of SUMA project we will refine the architecture and improve the performances of APENet+. First, the adoption of the last generation 28nm FPGA will enable us to switch to Gen3 PCIe protocol on host side and to integrate faster Torus channels, doubling the I/O bandwidth on both sides of the NIC. Furthermore, the huge amount of HW resources available in the 28nm FPGA will be used to develop specific computational task accelerators in the form of an ASIP (Application Specific IP) or as a custom hardware blocks. Lastly, we will perform customization of APENet+ IP in order to evaluate the use of torus-based unconventional interconnect topologies for computing systems dedicated to large-scale simulation in INFN non-traditional research areas (Bio-computing and Brain simulation).

## References

- [1] G. I. Egri, Z. Fodor, C. Hoelbling, S. D. Katz, D. Nogradi, K. K. Szabo, Lattice QCD as a video game, Comput. Phys. Commun. 177, 631 (2007) arXiv:hep-lat/0611022.
- [2] C. Bonati, G. Cossu, M. D'Elia and P. Incardona, QCD simulations with staggered fermions on GPUs, Comput. Phys. Commun. 183, 853 (2012) [arXiv:1106.5673 [hep-lat]].
- [3] C. Bonati, G. Cossu, M. D'Elia and F. Sanfilippo, Phys. Rev. D 83, 054505 (2011) [arXiv:1011.4515 [hep-lat]]; M. D'Elia, M. Mariti and F. Negro, Phys. Rev. Lett. 110, 082002 (2013) [arXiv:1209.0722 [hep-lat]]; C. Bonati, M. D'Elia, M. Mariti, F. Negro and F. Sanfilippo, arXiv:1307.8063 [hep-lat].
- [4] R. Ammendola et al, "APENet+: a 3D Torus network optimized for GPU-based HPC Systems", (2012) J. Phys.: Conf. Ser. 396 042059; R. Ammendola et al, "GPU peer-to-peer techniques applied to a cluster interconnect", Proceedings of CASS2013 workshop, accepted for publication (arXiv:1307.8276).