Preparation of a Petascale Supercomputing Infrastructure for European Scientists PRACE Status at Mid Year 2009

Dr. Jean-Philippe Nominé, CEA









Outline: PRACE Status at Mid Year 2009

- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009

HET: The Scientific Case

- Weather, Climatology, Earth Science •
 - degree of warming, scenarios for our future climate.
 - understand and predict ocean properties and variations
 - weather and flood events
- Astrophysics, Elementary particle physics, Plasma physics
 - systems, structures which span a large range of different length and time scales
 - quantum field theories like QCD, ITER
- Material Science, Chemistry, Nanoscience
 - understanding complex materials, complex chemistry, nanoscience
 - the determination of electronic and transport properties
- Life Science
 - system biology, chromatin dynamics, large scale protein dynamics, protein association and aggregation, supramolecular systems, medicine
- Engineering
 - complex helicopter simulation, biomedical flows, gas turbines and internal combustion engines, forest fires, green aircraft,
 - virtual power plant



















European HPC Shortcomings

- Fragmentation
 - No sustained HPC service beyond the national scale
 - No coordination of procurements
- Lack of a strong HPC industry
 - European companies mainly supply the European market (like Bull) or occupy market niches (like Numascale AS, Dolphin, ParTec)
 - Europe needs an independent access to this key technology



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009

Computational science infrastructure in



The European Roadmap for Research Infrastructures is the first comprehensive definition at the European level

Research Infrastructures are one of the crucial pillars of the European Research Area

A European HPC service – impact foreseen:

- strategic competitiveness
- attractiveness for researchers
- supporting industrial development



The ESFRI Vision for a European HPC service

- European HPC-facilities at the top of an HPC provisioning pyramid
 - Tier-0: 3-5 European Centres
 - Tier-1: National Centres
 - Tier-2: Regional/University Centres



- Shape the European HPC ecosystem involving all stakeholders
 - HPC service providers on all tiers
 - Grid Infrastructures
 - Scientific and industrial user communities
 - The European HPC hardware and software industry

Expected Structuring Effects

- Renewal of Tier-0 systems every 2-3 years required
- Investments of 200 400 Mio € every 2-3 years will create a critical mass for strengthening the European HPC industry
 - Capability to provide independent access to this key technology
 - Access to world-class HPC systems as a competitive advantage for European science and industry
- A European HPC Service as part of the European Research Area
 - Access to world-class HPC systems for the best researchers from all European countries



First Steps and Achievements





- PRACE context HPC in Europe
- PRACE objectives

PRACE organisation

- Key achievements in 2008 and 2009
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012



PRACE – Project Facts

- Workplan of the PRACE Project:
 - Prepare the contracts to establish the PRACE permanent Research Infrastructure as a single Legal Entity in 2010 including governance, funding, procurement, and usage strategies.
 - Perform the technical work to prepare operation of the Tier-0 systems in 2009/2010 including deployment and benchmarking of prototypes for Petaflop/s systems and porting, optimising, peta-scaling of applications
- Project facts:
 - Partners: 16 Legal Entities from 14 countries
 - Project duration: January 2008 December 2009
 - Project budget: 20 M € , EC funding: 10 M € 1300 PM

PRACE is funded in part by the EC under the FP7 Capacities programme grant agreement INFSO-RI-211528











Europe's current position in HPC





PRACE Work Packages (classical FP7 project...)

- WP1 Management
- WP2 Organizational concept
- WP3 Dissemination, outreach and training
- WP4 Distributed computing
- WP5 Deployment of prototype systems
- WP6 Software enabling for prototype systems
- WP7 Petaflop/s systems for 2009/2010
- WP8 Future petaflop/s technologies

Non technical

Technical



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012



Organisational Structure

- Legal Form
 - ERI European Research Infrastructure would have been optimal, but is not fully ready yet
- Still under discussion... as well as detailed governance model





Accessing the future PRACE RI (main user interface....)

Access Model

- Based on peer-review: "the best systems for the best science"
- Free-of-charge
- Technically: DEISA-like

Funding

- Mainly national funding through partner countries
- European contribution
- Access model has to respect national interests (ROI)



Guiding Principles & Flowchart from Initial Report on the Peer Review Process



How to get involved ?

If you are a national coordinator of HPC activities and your country is not yet a member:

Join the PRACE Initiative !

(scientists: influence your national coordinators...)

Port your code to the PRACE Prototypes, prepare yourselves to petascale

- Prototypes are mainly used project-internally, but...
- ... prototypes are also made available publicly for testing/porting purposes using a light-weight peer-review process
- See: http://www.prace-project.eu/prototype-access

PRACE web site

- http://www.prace-project.eu ٠
- RSS feeds, news •
- Various documents on line



PRACE newsletter 4/2008.

HPC training events 🔝

» Advanced Parallel Programming (in English), Feb. 23-25, 2009, Finland » Iterative Linear Solvers and Parallelization (in German).

Partnership for Advanced Computing in Europe

Welcome to PRACE

The Partnership for Advanced Computing in Europe prepares the creation of a persistent pan-European HPC service, consisting of several tier-O centres providing European researchers with access to capability computers and forming the top level of the European HPC ecosystem. PRACE is a project funded in part by the EU's 7th Framework Programme.

Supercomputers are indispensable tools for solving

the most challenging and complex scientific and technological problems through simulations. To remain internationally competitive, European scientists and engineers must be provided with leadership-class supercomputer systems. PRACE, the Partnership for Advanced Computing in Europe will create a persistent pan-European high performance computing (HPC) service and infrastructure. This infrastructure will be managed as a single European entity. European scientists and technologists will be provided world-class leadership supercomputers with capabilities equal to or better than those available in the USA and Japan. The service will comprise three to five superior HPC centers strengthened by regional and national supercomputing centers working in tight collaboration through grid technologies. In other words, the partnership will become a unique entity of the pan-European HPC ecosystem.





» PRACE held All Hands Meeting in Jülich, Germany, February 12-13, 2009 2009-02-16 » Serbia joins the PRACE initiative 2009-02-12 » PRACE Part of Zero-In Magazine - Call for Papers Open 2009-02-03 » Reminder: PRACE Award 2009 2009-01-15 » PRACE Award 2009: Call for Papers has Started 2008-12-18 more...

Administrator log in

Events 🔝

» OGF25 / EGEE User Forum, 2-6 March, Catania, Italy » 24th Forum ORAP, 26 March, Lille, France » DEISA PRACE Symposium 2009: HPC Infrastructures for Petascale Applications, May 11-13, Amsterdam, the Netherlands » ISC 2009, June 23-26, Hamburg, Germany more

Partner vacancies 🔝

» Research Programmer / Systems Administrator, ICHEC, Ireland





- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and user access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012



Training & education

- Survey of HPC education and training needs
 - Showed a strong need for parallel programming training
 - cf PRACE web site <u>http://www.prace-project.eu</u>
- PRACE Petascale Summer School, August 26-29, 2008, KTH, Stockholm
 - 30+ participants
- PRACE Winter School, February 9-13, 2009, Athens, Greece
 - Almost 50 participants
- GPGPU CUDA et al. training, France, CEA+GENCI, April 2009
 - 12 participants
- Usage of PRACE prototypes for hands-on



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012



Representative Benchmark Suite

- Key goal of WP6 is a set of applications benchmarks
 - To be used in the procurement process for Petaflop/s systems
- Survey data helped us select highly-used applications which span the breadth of both scientific area and algorithmic 'dwarf'
- "Living list"
 - QCD, NAMD, CPMD, Code_Saturne, GADGET, EUTERPE, NEMO, CP2K, GROMACS, NS3D, HELIUM, AVBP, TRIPOLI-4, PEPC, GPAW, ALYA, BSIT
 - WRF, Quantum_Espresso, Octopus, SPECFEM3, ELMER
- Each application will be ported to appropriate subset of prototypes



Performance Analysis Tools and Benchmarks

- Documented:
 - Reviews of tools
 - 9 different tool suites
 - Synthetic benchmarks
 - Computation, mixed-mode, IO, bandwidth, OS, communication
 - Application benchmarks
 - Selection and porting
- Applications and Synthetic benchmarks integrated into JuBE
 - Juelich Benchmark Environment



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012





PRACE prototypes approach

- Test & try before buying
 - Assessment of technology and architectures
- Share experience between partners
- Prepare benchmarks
- Foresee technology evolutions
- Foster collaborations between providers and users
- => A selection of systems and components, with existing, near-existing or emerging technology



PRACE prototypes TWO-FOLD approach

- (Near) existing technologies for 2009-2010 (WP7)
 - Full-featured systems, either:
 - % of large existing production systems (*scaling*)
 - or extensions of existing production systems (*techno. updates*)
 - or dedicated (smaller) systems (*new technology*)
- Emerging technologies for 2011 and beyond (WP8)
 - Mostly components/subsystems
 - As of today, strong focus on application specific, attached processors (accelerators)



Selected prototypes (coordinated by WP7)

Site	Architecture Vendor/Technology	Point of contact	
FZJ	MPP	Michael Stephan	
Germany	IBM BlueGene/P	<u>m.stephan@fz-juelich.de</u>	
CSC-CSCS	MPP	Janne Ignatius janne.ignatius@csc.fi	
Finland+Switzerland	Cray XT5/XTn - AMD Opteron	Peter Kunszt <u>peter.kunszt@cscs.ch</u>	
CEA-FZJ	SMP-TN	Gilles Wiber <u>gilles.wiber@cea.fr</u>	
France+Germany	Bull et al. Intel Xeon Nehalem	Norbert Eicker <u>n.eicker@fz-juelich.de</u>	
NCF/SARA	SMP-FN	Axel Berg <u>axel@sara.nl</u>	
Netherlands	IBM Power6	Peter Michielse <u>michielse@nwo.nl</u>	
BSC	Hybrid — fine grain	Sergi Girona	
Spain	IBM Cell + Power6	<u>sergi.girona@bsc.es</u>	
HLRS	Hybrid – coarse grain	Stefan Wesner	
Germany	NEC Vector SX/9 + x86	<u>wesner@hlrs.de</u>	



Installed prototypes





IBM BlueGene/P (FZJ) 01-2008



NEC SX9 vector part (HLRS) 02-2009 x86 part 04-2009 IBM Power6 (SARA) 07-2008



Cray XT5 (CSC) 11-2008



IBM Cell/Power (BSC) 12-2008

Bull Nehalem/INCA (CEA) 06-2009

	NCF	CSCS/CSC	BSC	FZJ	CEA/FZJ	HLRS
Dedicated system - makes possible unfriendly tests						
Shared large system - makes possible large runs and assessment under real production						
MPP						
Cluster with thin-nodes						
Cluster with fat nodes						
Advanced (Hybrid)						
Specific Hardware Technologies	Power6	AMD Barcelona and Shaghai	IBM Cell	Blue Gene	Intel NehalemEP	SX9
Specific Software technologies	PERCS Power7 simulator	MPI/OpenMP CAF UPC				
Full featured system with storage and IO						
Connected to the DEISA network					CEA new DEISA associate partner	
Collaboration with vendors	Software technology (IBM)	system reliability, performance, functionality (CRAY)	System design (IBM)		System design (BULL)	System design (NEC)



- PRACE context HPC in Europe
- PRACE objectives
- PRACE organisation
- Key achievements in 2008 and 2009
 - Preparation of legal entity, governance and access
 - Education and training
 - Applications benchmarking and petascaling
 - Prototypes: systems for 2010
 - Prototypes: components and technologies for 2011/2012

Prototypes selected by WP8

Sites	Hardware/Software	Porting effort
CEA " <i>GPU/CAPS</i> "	1U Tesla Server T1070 (CUDA, CAPS, DDT) Intel Harpertown nodes	"Evaluate GPU accelerators and GPGPU programming models and middleware." (e.g., <i>pollutant</i> <i>migration code</i> (ray tracing algorithm) to CUDA and HMPP)
CINES-LRZ " <i>LRB/CS</i> "	Hybrid SGI ICE2/UV/Nehalem-EP & Nehalem-EX/ClearSpeed/ Larrabee	Gadget ,SPECFEM3D_GLOBE, RaXml, Rinf, RandomAccess, ApexMap, Intel MPI BM
CSCS "UPC/CAF"	Prototype PGAS language compilers (CAF + UPC for Cray XT systems)	"The applications chosen for this analysis will include some of those already selected as benchmark codes in WP6 ."
EPCC " <i>FPGA</i> "	Maxwell – FPGA prototype (VHDL support & consultancy + software licenses (e.g., Mitrion-C))	"We wish to port several of the PRACE benchmark codes to the system. The codes will be chosen based on their suitability for execution on such a system."

Prototypes selected by WP8 (cont'd)

Sites	Hardware/Software	Porting effort
FZJ (BSC) "Cell & FPGA interconnect"	eQPACE (PowerXCell cluster with special network processor)	Extend FPGA-based interconnect beyond QCD applications.
LRZ " <i>RapidMind</i> "	RapidMind (Streaming Processing Programming Paradym) X86, GPGPU, Cell	ApexMap, Multigrid, FZJ (QCD), CINECA (linear algebra kernels involved in solvers for ordinary differential equations), SNIC
NCF "ClearSpeed"	ClearSpeed CATS 700 units	Astronomical many-body simulation, Iterative sparse solvers with preconditioning, finite element code, cryomicrotome image analysis
CINECA	I/O Subsystem (SSD, Lustre, pNFS)	-
	+ 1 additional prototype targeting energy efficiency	-



As a temporary conclusion...

- You do have access to PRACE prototypes
- We want to foster long-term relationships with scientific communities
 - PRACE governance should include a Users' Forum and a Scientific Steering Committee
- Training and education are a main concern
- We will go on watching and assessing technology and architectures, then take part in R&D and influencing vendors' roadmaps
- More to come about 2010 RI organisation
 - Keep in touch, ask us...

jean-philippe.nomine@cea.fr



Europe has aleady gone Petascale...

- Jugene 1 PFlop/s BlueGene at FZJ (TOP#3, June 2009)
- CEA will have 1 PFlop/s Bull TERA 100 Cluster mid 2010
- CEA and FZJ collaborate...



PRACE and CEA

- GENCI coordinates the French PRACE effort
 - Half of French manpower provided by CEA
 - Strong involvement in prototypes
 - Preparing TGCC new facility to host one of the first PRACE machines at CEA/DIF, Bruyères near Paris
- CPU Nehalem/Bull INCA drawers
- 1024 cores, ~10 TF





- GPU Tesla/NVIDIA
- 2 S1070 servers
- Programmings models and tools
 CUDA, HMPP



TGCC as of June 2009