## PRACE Brief Overview and Background to Accelerator Interest Lennart Johnsson KTH and University of Houston



# **PRACE** Vision and Mission

- Vision: enable and support European global leadership in public and private research and development.
- Mission: contribute to the advancement of European competitiveness in industry and research through the provisioning of world leading persistent High-End Computing infrastructure



## Europe's position in HPC through 2008



![](_page_3_Picture_0.jpeg)

## Supercomputing Drives Science through Simulation

![](_page_3_Figure_2.jpeg)

http://www.prace-project.eu/documents/first-industrial-seminar-presentations/2-%20Bachem%20PRACE%20Industrieseminar%20Amsterdam.pdf

![](_page_4_Picture_0.jpeg)

# **PRACE** Partners

- PRACE Initiative
  - MoU signed by14 countries April 17<sup>th</sup> 2007:
    - Austria
    - Finland
    - France
    - Germany
    - Greece
    - Italy
    - Netherlands

- Norway
- Poland
- Portugal
- Spain
- Sweden
- Switzerland
- United Kingdom

Since April 2007 Bulgaria, Cyprus, Czech Republic, Ireland, Serbia and Turkey has joined the PRACE initiative. Additional countries are in the process of joining.

# ESFRI – Estimated costs

PARTNERSHIP FOR ADVANCED COMPUTING IN EUROPE

- Unlike other European Research Infrastructures:
  - Tier-0 resources have to be renewed every 2-3 years
  - Construction cost 200 400 Mio. € every 2-3 years
  - Annual running cost 100 200 Mio. €
  - Additional effort needed for software
- A truly European challenge also in terms of funding
- PRACE The Partnership for Advance Computing in Europe
  - An Initiative created to implement the ESFRI vision of a European HPC service

![](_page_5_Figure_9.jpeg)

![](_page_6_Picture_0.jpeg)

# PRACE Project Work Packages

- 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 systems for 2009/2010
- WP8 Future petaflop technologies

![](_page_7_Picture_0.jpeg)

## Installed Petaflop Prototypes (WP5/7/8)

![](_page_7_Picture_2.jpeg)

![](_page_7_Picture_3.jpeg)

IBM BlueGene/P (FZJ) 01-2008 (MPP)

IBM Power6 (SARA) 07-2008

![](_page_7_Picture_6.jpeg)

![](_page_7_Picture_7.jpeg)

![](_page_7_Picture_8.jpeg)

![](_page_7_Picture_9.jpeg)

IBM Cell/Power (BSC)

Intel Nehalem/Xeon (CEA/FZJ): installation date April 2009

![](_page_7_Picture_12.jpeg)

![](_page_7_Picture_13.jpeg)

http://www.prace-project.eu/documents/deisa-prace-symposium-presentations/ 12-2008 AB\_Ver%20PRACE\_DEISA%20Amsterdam.pdf

![](_page_8_Picture_0.jpeg)

## FZJ

- Specific features:
  - Access to a large existing MPP system, already 1/4 PF with an architecture expandable to 1 PF
- Contribution to the PRACE project:
  - Application scaling, optimization and benchmarking including:
    - Communications
    - I/O
  - Large scale operations on selected applications
  - Assessment of electrical power usage
- Availability July 2008

FZJ	2 %
MPP IBM BG/P	
16 racks 16k nodes 64k cores (PPC 450)	
223 TF peak	

![](_page_8_Picture_13.jpeg)

## Partnership For Advanced Computing IN Europe

# Prototypes for 2009/2010 Systems

## CSC / CSCS

- Specific features:
  - Prototype installed in CSC, joint effort with CSCS
  - Funding required only for a dedicated system
  - Additional access to a larger existing system, on request – similar architecture
- Contribution to the PRACE project:
  - Access to a prototype with MPP architecture
    - AMD Barcelona, SeaStar2+ torus network
  - Early access to AMD new generation of processors (Shanghai)
  - Additional focus on hybrid MPI/OpenMP parallel programming by CSCS
  - Capability testing on the CSC existing XT system
- Availability December 2008

![](_page_9_Picture_14.jpeg)

Courtesy Francois Robin, GENCI, EC Review 2008-07-16

![](_page_9_Figure_16.jpeg)

## Partnership For Advanced Computing IN Europe

# Prototypes for 2009/2010 Systems

## CEA / FZJ

- Specific features:
  - Prototype distributed over 2 sites
- Contribution to the PRACE project
  - Early access to a prototype of a new product designed by BULL (GA 2H09)
    - High density blade based system with new Intel Nehalem processors
    - · Optimized for HPC
    - · Scalable to Petaflop/s
    - · Water cooling
  - Scalability testing on the FZJ part of the prototype
- Availability : March 2009

![](_page_10_Picture_13.jpeg)

![](_page_10_Picture_14.jpeg)

![](_page_10_Picture_15.jpeg)

![](_page_10_Picture_16.jpeg)

![](_page_10_Figure_17.jpeg)

![](_page_11_Picture_0.jpeg)

## NCF

- Specific features:
  - Access to a complete system with new Power 6 processors and IBM Power Cluster fat node architecture (large memory)
- Contribution to the PRACE project:
  - Focus on HPC software from US DARPA research early access
  - Specific test nodes for aggressive experimentations
  - Larger scale "production" runs on the full system
- Availability October 2008

![](_page_11_Figure_10.jpeg)

![](_page_11_Picture_11.jpeg)

![](_page_12_Picture_0.jpeg)

## BSC

- Specific features:
  - Dedicated "fine grain" hybrid system
- Contribution to the PRACE project:
  - New Power 6 + Cell processor integration
    - Analog to US PF RoadRunner but different CPUs
  - Programming techniques and tools for CPU+accelerators
  - Operation of an hybrid system: queuing system, file system, accounting, system administration, ..
  - Assessment of electrical power usage
- Availability December 2008

![](_page_12_Picture_12.jpeg)

Courtesy Francois Robin, GENCI, EC Review 2008-07-16

#### Hybrid IBM Cell + Power6

1 M€ 48.3%

12 P6 + 72 Cell blades 48 + 1296 CPUs

14 TF

![](_page_13_Picture_0.jpeg)

## HLRS

- Specific features
  - Hybrid vector + scalar architecture assessment
    - New NEC SX-9 vector processor
    - Intel x86-64
    - · Shared filesystem
  - "Coarse grain" coupled simulations
- Contribution to the PRACE project:
  - Experiment vector components part of future Japan more integrated multi-Petaflop/s systems
  - Specific I/O issues
  - Availability March 2009

![](_page_13_Picture_13.jpeg)

٠

![](_page_13_Picture_14.jpeg)

![](_page_13_Picture_15.jpeg)

![](_page_14_Picture_0.jpeg)

Summary	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 DEISA	Courtes	/ Francois Robin	GENCI, EC Rev	iew 2008-07-16	CEA new DEISA associate partner	

![](_page_15_Picture_0.jpeg)

## Prototype 2009/2010 Systems Access

- Proposals from academia and industry are eligible, as long as the organisation of the project leader is homed in Europe (or a PRACE initiative country).
- Depending on the system vendor, researchers from some countries cannot be granted access due to export regulations.

![](_page_16_Picture_0.jpeg)

# Prototype 2009/2010 System

 More details on the PRACE prototypes can be found at <u>http://www.prace-project.eu/</u> <u>documents/PRACE-Prototypes.pdf</u>

![](_page_17_Picture_0.jpeg)

# **Technology Prototypes**

![](_page_18_Picture_0.jpeg)

# Technology Prototype

Purpose: Assess technologies and architectures in regards to

- Peak performance/ efficiency
- Programmability
- Energy efficiency
- Density
- Cooling
- Cost

![](_page_19_Picture_0.jpeg)

IN HPC speed is never high enough.

Until  $\approx$  2005 partly solved by increasing clock frequency.

As energy consumption  $E = NCV^2 f$  (N = # of devices/surface unit C = capacitance V = voltage f = frequency) Choice was made to increase N, lower V, keep f as low as possible.  $\rightarrow$  More cores on chip with acceptable energy consumption thanks to possible technology shrink: 2005—2006 dual core, 2007—2009 quad-core & 6 core CPUs,

2010-2012 8 to many-core.

![](_page_20_Picture_0.jpeg)

#### Combatting lack of speed

Vector units SSE*x*, AVX, etc., have been introduced.

Originally for graphics, but also of help for simple regular arithmetic operations:

![](_page_20_Figure_5.jpeg)

![](_page_21_Picture_0.jpeg)

#### Combatting lack of speed

External accelerators are a logical extension.

Presently four types are available:

Cell-based systems (IBM)

BSC, FZJ prototypes

FPGA-based (CPUTech, Mitrion, Convey, Kuberre, SRC,...)

EPCC prototype

GPU-based (NVIDIA, ATI/AMD)

- GENCI prototype

ClearSpeed CSX700-based (PetaPath)

- CINES, LRZ, NCF prototypes

![](_page_22_Picture_0.jpeg)

## Cell - Roadrunner

![](_page_22_Figure_2.jpeg)

#### What are the trends and why?

#### Combatting lack of speed

GPU-based accelerators rely on: Many parallel processor streams (240/GPU in the Tesla S1070) At a moderate clock frequency (1.296—1.44 GHz).

8-byte precision arithmetic is much slower than 4-byte precision (< 10%) Data transport to/from host is a main issue.

![](_page_23_Figure_5.jpeg)

![](_page_23_Figure_6.jpeg)

Rank	Computer	Mflop/W	Туре	Comment
1-3	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D-Torus	722.98	Accelerator	PRACE supported prototype
4-5	BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, IB	458.33	Accelerator	Moderate size systems, no 26 and 79
6	BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, IB	444.25	Accelerator	Large system. No 2
7	GRAPE-DR accelerator Cluster, Infiniband	428.91	FPGA	
8	NUDT TH-1 Cluster, Xeon E5540/E5450, ATI Radeon HD 4870 2, IB	379.24	GPU	Large system. No 5
9 -20	Blue Gene/P Solution	378.76	Proprietary	
21	Blue Gene/P Solution	366.58	Proprietary	Large system, No 1
22-23	Blue Gene/P Solution	353.98	Proprietary	Large systems, no 4 and no 8.
24	Sun Blade X6440, Opteron 2.5 Ghz, Infiniband QDR	341.42		Top Green500 commodity system
25	iDataPlex, Xeon E55xx QC 2.26 GHz, Infiniband, Windows HPC2008 R2	299.52		
500	Cluster Platform 4000 BL465c, Opteron DC 2.4GHz, GigEthernet	13.03		

![](_page_25_Picture_0.jpeg)

#### Power issues in HPC

For Multi-Petaflop/s to Exaflop/s systems <u>real disruption in technology is</u> <u>needed</u>. With current type of technology (even with progress factored in) we would have the following situation:

2012	2018
10 Pflop/s	1 Eflop/s
9.3 MW for IT	154 MW for IT
14 MW total	213 MW total
12.3 M€/y	187 M€/y

Assuming: PUE = 1.5, € 0.10/Kwhr.

![](_page_26_Picture_0.jpeg)

# FZJ

## eQPACE

- Extend communication capabilities of eQPACE to make it suitable for a wider range of applications. Reach a top position in the Green500 list (FZJ).
- Hardware: PowerXCell8i processor nodes with custom 3D-torus interconnect.
- Benchmarks: HPL, Euroben kernels, torus network benchmark, applications & iterative solvers.
- Programming environments: Cell SDK & CellSs

![](_page_26_Picture_7.jpeg)

## FZJ

![](_page_27_Figure_2.jpeg)

![](_page_27_Picture_3.jpeg)

![](_page_27_Figure_4.jpeg)

# LRZ - CINES

#### SGI hosting platform:

- 32 Blades XE
  - 64 Processors Intel Nehalem-EP
  - 256 Cores
  - 4 GB per core
- Infiniband QDR
- Estimated peak performance is 3 TFlop/s.
- <u>ClearSpeed-Petapath</u> accelerators:
- 32 × e710 cards
  - One per ICE blade
  - 4 GFlop/s (DP) / Watt
- Software development toolkit
- Estimated peak performance is 3 TFlop/s DP.
- The total accumulated peak performance is 6 TFlop/s DP.

![](_page_28_Figure_16.jpeg)

## BAdW-LRZ

- SGI Altix ICE System with 48 compute blades each containing two Intel Xeon Nehalem-EP 2.53 GHz 4-core processors and 24 GB main memory;
- 4-Socket Nehalem-EX WhiteBox from Intel (α-Prototype: "Sunrise Ridge");
- SGI "UltraViolet" prototype with 16 "P1-Level" compute blades each containing two 2.0 GHz Intel Xeon Nehalem-EX 8-core processors and 32 GB main memory;
- 4× DDR Infiniband switch and cables (44 IB ports plus 2 x 10GigE ports);
- SGI UV extension unit with
  - 1 Intel Larrabee accelerator;
  - 4 ClearSpeed-Petapath Advance e710 accelerator boards;
- Two racks with system administration infrastructure (head node);
- SLES10 SP2 operating system and SGI "Tempo" and SGI "ProPack" management software;
- Storage: The hybrid system prototype is connected with the Lustre file systems of the BAdW-LRZ linux cluster (130 TB /scratch and 50 TB /project).

![](_page_29_Figure_12.jpeg)

## Partnership For Advanced Computing IN Europe

# CEA

- The prototype is a cluster of 5 nodes. The first node is a login node featuring four Nehalem-EP processors. The last four nodes are BULL R422 servers featuring 2 Haperton CPUs with 16 GB of RAM each. The BULL servers are connected through a PCI-Express 16× to two NVIDIA TESLA servers. Infiniband DDR is used to interconnect the nodes.
- A TESLA server has 2 C1060 graphic boards. A C1060 features 30 multiprocessors at 1,3 GHz and has 4 GB of memory. One multiprocessor has 8 single precision (SP) units and 1 double precision (DP) unit. Therefore, a C1060 can have 240 SP (30 DP) units running in parallel. The peak performance of a C1060 is 78 GFlop/s DP and 624 GFlop/s SP.

![](_page_31_Picture_0.jpeg)

EPCC

## Maxwell FPGA

Evaluate the performance and usability of the HARWEST Compiling Environment (EPCC).

- Hardware: FPGA prototype "Maxwell" (32 FPGAs) from both Alpha Data Ltd and Nallatech Ltd using Virtex-4 FPGAs supplied by Xilinx Corp.
- Benchmarks: 4 Euroben kernels
- Languages:
  - VHDL
  - HCE

![](_page_31_Figure_9.jpeg)

# NCF

- "Feynman e740/e780" contain 4 or 8 ClearSpeed CSX700 processors for a total of 114. The dark blue HP DL170 host nodes each connect to two Feynman boxes by PCI Express Gen. 2 16× (8 GB/s).
- The HP host nodes are connected by 4× DDR Infiniband in order to employ them, and the attached accelerator boxes.
- Internally a Feynman e780 connects to its 8 processors through a PLX 8648 PCIe switch. The bandwidth to each individual processor is 1 GB/s. The e740 boxes contain only 4 CSX700 processors but the bandwidth to/from the individual processors is doubled to 2 GB/s. This should help for bandwidth-hungry applications. The peak performances are 768 GFlop/s and 384 GFlop/s in 8-Byte precision for a Feynman e780 and e740, respectively. The energy efficiency is very high: ≈ 4 GFlop/s/Watt.
- The HP DL170h G6 host nodes contain 2 Intel X5550 Nehalem processors at 2,67 GHz and 24 GB of DDR3 memory/node. Furthermore, each host node contains dual disks of 250 and 500 GB, for system and user data, respectively. Access to the system is provided through the HP DL160 G6 head node. The system is divided in a development part (1 DL170 + 2 Feynman e740s) and a production part (7 DL170s + 2 Feynman e740s + 12 Feynman e780s).

Rack 1	Rack 2
GigE Switch	
Feynman e740	Feynman e780
HP DL170	HP DL170
HP DL170	HP DL170
Head Node (HP DL160)	KVM
IB Switch	Monitor/Keyboard
HP DL170	HP DL170
HP DL170	HP DL170
Feynman e780	Feynman e780

# SNIC-KTH

- 4-socket blade with 6core 2.1 GHz AMD Istanbul HE CPUs, 32GB/node
- 10-blade in a chassis with 36-port QDR IB switch
- 180 nodes, 4320 cores, full bisection IB interconnect
- 2TF/chassis, 12 TF/rack

![](_page_33_Figure_6.jpeg)

![](_page_33_Picture_7.jpeg)

![](_page_34_Picture_0.jpeg)

## First RI call 2010

Merit based access as determined by evaluation conducted by PRACE Scientific Committee comprised of leading European scientists in disciplines in need of High-End Computing

## **Outreach and Education**

![](_page_35_Picture_2.jpeg)

#### Industry seminars:

1st Seminar Sept. 3, 2008 Amsterdam, The Netherlands

#### Summer & winter schools:

Stockholm, Athens

#### PRACE booths: ISC, ICT, SC,

![](_page_35_Picture_8.jpeg)

PRACE website

![](_page_35_Picture_10.jpeg)

![](_page_35_Picture_11.jpeg)

![](_page_35_Picture_12.jpeg)

![](_page_35_Picture_13.jpeg)

![](_page_35_Picture_14.jpeg)

PRACE Winter Schoolat the OTE academy, Athens 26-29.8.2009

![](_page_35_Picture_16.jpeg)

ICT 2008, PRACE-Booth

![](_page_35_Picture_18.jpeg)

PRACE Summer School Stockholm

http://www.prace-project.eu/documents/deisa-prace-symposium-presentations/AB\_Ver%20PRACE\_DEISA%20Amsterdam.pdf

![](_page_36_Picture_0.jpeg)

## High-End Computing Software Survey

![](_page_37_Picture_0.jpeg)

European Large Scale Applications characteristics (Linpack Equivalent Flops (LEFs) used in TF units)

Area/Dwarf	Dense linear algebra	Spectral methods	Structure d grids	Sparse linear algebra	Particle methods	Unstructu red grids	Map reduce methods
Astronomy and Cosmology	0	0.62	4.91	3.59	5.98	2.99	0
Computational Chemistry	15.35	26.09	1.80	3.45	7.49	0.53	12.98
Computational Engineering	0	0	0.53	0.53	0	0.53	2.8
CFD	0	1.70	7.37	3.05	0.32	3.00	0
Condensed Matter Physics	9.10	15.07	1.62	0.73	1.76	0.28	5.70
Earth and Climate Science	0	2.03	5.83	1.33	0	0.26	0
Life Science	0	4.72	0.94	0.13	0.94	0.28	3.46
Particle Physics	12.50	0	4.59	0.92	0.10	0	89.27
Plasma Physics	0	0	1.33	1.33	3.55	0.42	0.63
Other	0	0	0	0	0	0	0

http://www.prace-project.eu/documents/PRACE\_LJ.pdf

## Mapping Applications to Architectures

PARTNERSHIP FOR ADVANCED COMPUTING IN EUROPE

- Identified affinities and priorities
- Based on the application analysis - expressed in a condensed, qualitative way
  - Need for different "general purpose" systems
  - There are promising emerging architectures
- Will be more quantitative after benchmark runs on prototypes

Code	MPP (i.e. BlueGene L/P or CRAY XT4/5)	Thin node clusters (i.e. Bull INCA or SGI ICE)	Fat node clusters (i.e. Bull MESCA, SGI UltraViolet or IBM Power6)	Vector systems (NEC SX8-9, Cray X2	Accelerated systems (i.e. scalar or vector + GPU, FPGA or Clearspeed).	Accelerated systems – Cell based (i.e. Roadrunner, Maricell)	
NAMD				E	E	E	
CPMD						E	
VASP						E	
QCD				E	E	E	
GADGET							
Code Saturne				E	E	E	
TORB						E	
NEMO							
ECHAM5							
CP2K	E						
GROMACS					E	E	
N3D		E	E				
AVBP	E						
HELIUM							
TRIPOLI 4	E		E				
GPAW							
ALYA						E	
SIESTA						E	
BSIT						E	
PEPC				E	E	E	
Table 4 : application mapping to Petaflop/s systems archite							

![](_page_39_Picture_0.jpeg)

# Acknowledgements

## Workshop Sponsors

![](_page_39_Picture_3.jpeg)

![](_page_39_Picture_4.jpeg)

![](_page_39_Picture_5.jpeg)

![](_page_39_Picture_6.jpeg)

![](_page_39_Picture_7.jpeg)

www.prace-project.eu

 $http://www.prace-project.eu/documents/deisa-prace-symposium-presentations/AB\_Ver\%20PRACE\_DEISA\%20Amsterdam.pdf$