



FP7 Support Action - European Exascale Software Initiative

DG Information Society and the unit e-Infrastructures



# EESI System Software

## WG 4.2

Franck Cappello, INRIA and University of Illinois

Bernd Mohr, JSC

François Bodin, Marc Bull, Toni Cortés, Jesus Labarta,  
Jacques C. Lafoucriere, David Lecomber, Thomas Ludwig,  
Simon McIntosh-Smith, Jean-François Méhaut, Matthias Müller,  
Raymond Namyst, Olivier Richard, Pascale Rossé,  
Karl Solchenbach, Vladimir Voevodin, Felix Wolf



## WG4.2 Scope and Challenges

---

- Software Eco-system: between Applications and Hardware
  - Programming models, Compilers, Runtime, Tools, OS, I/O, System management
  - WG4.2 does not cover Framework, Numerical libraries, Workflow, Scientific data management, Visualization → other WGs
  
- Many technical challenges on all layers of the software stack
  - Scalability
  - Heterogeneity
  - Power
  - Errors, Faults, Failures
  - Reducing as much as possible Overhead, Jitter, Noise, etc.
  
- A trend toward community software

## WG4.2 TOPICS

---

- System Software
  - Operating Systems, System Management
  - Job and Resource Manager
  - Runtime Systems
  - I/O systems
- Development Environments
  - Programming Models, Compilers
  - Debuggers
  - Performance tools
  - Correctness tools
- Crosscutting Dimensions
  - Resilience
  - Power management

# WG4.2 Chairs, Leaders and Experts

## Chairs:

□ <b>Franck Cappello</b>	INRIA&UIUC,	1 FR,	<b>Resilience and FT</b>
□ Benrd Mohr,	JSC,	1 DE,	Performance tools
□ <b>Jesus Labarta</b>	BSC	1 ES	<b>Programming models</b>
□ Marc Bull	EPCC	1 UK	OpenMP / PGAS
□ François Bodin	CAPS	2/V FR	Programming/Compiler for GPUs
□ <b>Raymond Namyst</b>	LABRI Bordeaux,	3 FR	MPI&OpenMP <b>Runtime</b> , GPUs
□ Jean-François Méhaut	INRIA Grenoble	4 FR	Performance Modeling/Apps.
□ <b>Matthias Müller</b>	TU Dresden	2 DE	<b>Validation/correctness Checking</b>
□ <b>Felix Wolf</b>	GRS	3 DE	<b>Performance Tools</b>
□ <b>David Lecomber</b>	ALINEA	2/V UK	<b>Parallel debugger</b>
□ Simon McIntosh-Smith	U. Bristol	3 UK	Computer Architecture & FT
□ Vladimir Voevodin	MSU	1 RU	Performance tools
□ <b>Thomas Ludwig</b>	DKRZ	4 DE	<b>Power management</b>
□ <b>Olivier Richard</b>	INRIA	5 FR	<b>Job and Resource Manager</b>
□ <b>Jacques C. Lafoucriere</b>	CEA	6 FR	<b>I/O, File system</b>
□ Toni Cortés	BSC	2 ES	I/O, Storage
□ <b>Pascale Rossé</b>	BULL	V FR	<b>OS, System management</b>
□ Karl. Solchenbach	Intel	V BE	All

## WG4.2 EESI Roadmap content for Sys-soft eco system

- Starting point: IESP roadmap, Exascale studies, Soft. used in HPC centers
- Description of the scientific and technical perimeter
- Social benefits, societal, environmental and economical impact
- S
- A  
P  
Next slides → only a small subset of the current results  
sort of executive summaries
- E  
The roadmap with timeline, HR and cost  
is not yet completed
- S
- N  
Recommendations for EU may evolve
- P
- Existing funded projects and funding agencies
- *Timeline, needs of HR, provisional costs,*
- Needs for experimental platform: size, reconfigurable?, dedicated?
- What software in an Exascale stack and with what level of responsibility

# WG4.2 Programming model

Chair: Jesus Labarta, BSC

- Programming models
  - Clean separation and interaction between application developer and system (hardware and software)
  - Improve productivity, Significantly reduce maintenance costs
  - **Incremental parallelization will reduce development costs**
- Main issues at Exascale:
  - Concurrency, Asynchrony, Malleability, Address spaces and locality, Hierarchy, I/O, Modularity, interoperability, Productivity/portability, Incremental path,
  - **But also more decoupling between programmers and machines,**
  - Cross cutting: Power Management, Fault tolerance, Performance
- Recommendations for EU
  - Existing developments with demonstrated potential and willing to play a role in the future exascale software environment: HMPP, StarSs (OmpSs)
  - Important contributions to standards: MPI/OpenMP
  - **Programming models that provide expressiveness, incremental portability and performance** will boost performance, visibility and impact of applications having early access to them

# WG4.2 Runtime

Chair: Raymond Namyst, U. Bordeaux, INRIA

- Bridge the gap between underlying architecture and application requirements:
  - Scheduling, load balancing, Memory management, Communications, Sync.
  - Where accurate information is available about the actual power consumption of various hardware parts
- Main issues at Exascale:
  - Mastering heterogeneity: Unified/transparent accelerator models, Support for adaptive granularity, Fine grain parallelism, Scheduling for latency/bandwidth
  - Dealing with millions of cores/nodes: Scheduling, communication
  - **Supporting multiple programming models: MPI + threading model + accelerator**
  - Robustness: reconfigure itself when resources suddenly disappear.
- Recommendations for EU
  - Many European countries have a long-standing activity in runtime design
  - **Unified runtime system**, providing a unified API to deal with threads and lightweight tasks (together with their integration with MPI/PGAS communication systems)

# WG4.2 Validation / Correctness

Chair: Matthias Müller, TU Dresden

- Tools and methods for validation and correctness checking:
  - Validate a program in accordance to a model/specification (ex: MUST for MPI)
  - Can detect many errors, Especially: portability and non-determinism related bugs
  - Exascale will drastically increase manifestation rates!
  - **Correctness & validation helps reducing the time to solution** (better productivity)
- Main issues at Exascale:
  - Scalability, Fault tolerance, Adaption to new paradigms, Integration into new paradigms, **Integration into debugging workflow**
- Recommendations for EU
  - MUST and Marmot collaborations: Supercomputing vendors, especially those involved in DARPA's HPCS program (Cray, IBM), Software companies like Rouge Wave (Totalview), National labs (LLNL, LANL, ORNL, ANL), Universities (Utah, Houston)
  - **Foster interaction of methods/groups** within and beyond Europe
  - **Horizon. interaction** of approaches/tools: compiler + runtime + classic debugger
  - **Vertical integration** of validation and correctness at all layers of software stack

# WG4.2 Parallel Debugger

Chair: David Lecomber, ALINEA

- Any method to assist with the removal of software failures
  - Interactive tools, user tactics, Alinea DDT, Totalview Rogue Wave's, etc.
  - Primary focus on interactive tools that give user knowledge and control of program state and activity. **Limited to in single execution path.** Can be improved by static source code analysis.
  - **Some overlap with automatic Validation/Correctness tools:** Automated problem detection tools: MUST (TU-Dresden), or Umpire (LLNL).
- Main issues at Exascale:
  - Scale: responsiveness of debugger to user command and overhead), Architectural and Programming Model Unknowns (debugging support)
  - Crosscutting issues: Programmability, Resilience, I/O+Network, Reproducibility
- Recommendations for EU
  - **Continued financing/Sustainability** for Debugger tool developed in EU: DDT from Alinea, Globally only debugging product with Petascale ability. Collaborations with leading European research centres and initiatives and US
  - **Integration into debugging workflow with validation and correctness**

# WG4.2 Performance Tools

Chair: Felix Wolf, GRS

- Diagnostic programs for performance optimization of applications
  - Ideally used in combination with performance modeling
  - Huge savings in terms of time and energy, Enables new scientific discovery: by calculating larger problem sizes.
- Main issues at Exascale:
  - Million-fold concurrency, Deeper memory hierarchies, More dynamic execution, **Limited reproducibility (asynchrony, faults and power managment), Limited bandwidth to extract performance data, Multiple programming models**
  - Cross cutting: resilience (tool should be FT), power, programmability
- Recommendations for EU
  - Many of the performance tools used in production today are made in Europe, **Europe has leading position in this area**
  - OPT (Allinea), ThreadSpotter (RogueWave), Vampir (GWT-TUD), Intel trace collector and analyzer originally from Europe
  - Popular Academic tools, Paraver/Dimemas (BSC) Scalasca (Jülich/Aachen)
  - Initiative for more integration with Score-P measurement system
  - **European tool builders can make significant contribution to exascale effort, but sustainability of funding must be ensured.**

# WG4.2 I/O and File System

Chair: Jacques C. Lafoucriere, CEA

- Permanent storage: keep track of computing results for post processing and to start a new computation
  - Online (Disk, SSD, etc.) and Offline (tapes, disk based virtual tapes)
- Main issues at Exascale:
  - Extreme concurrency level (data, metadata locking), 1M disks, End2end integrity, too many transactions (contention), check usefulness of archived data
  - Storage plays a major role in Fault tolerance, NV memory → new failure scenarios
  - **Some HPC centers prefer Open Source FS** (or will be locked to a vendor)
- Recommendations for EU
  - 2 major EU companies offer storage solutions for Petascale computers
    - Xyratex: hardware and proprietary file system (Colibri)
    - Bull: hardware neutral integration of open source file systems (Lustre)
  - Significant set of small groups with good impact in the storage community around Europe. SCALUS as a 1st coordination initiative
  - Storage is a critical part of an Exascale solution
  - **EU must be a major contributor/partner to Open Source file system**



# WG4.2 OS and System Management

Chair: Pascale Rossé, BULL

- OS and System management:
  - OS is a key element between hardware and runtime/application
  - Developing and debugging large scale HPC systems requires experts at all layers of the software stack, including OS.
- Main issues at Exascale:
  - Many cores, Low overhead, low noise, etc.
  - Scalability: Huge amount of statistics, log data, events VS Centralized tools
- Recommendations for EU
  - **Hardware initiative in Europe** around ARM, AMD Fusion, Nvidia Denver, etc. **should be complemented by an OS initiative.**
  - Sustain and develop further OS R&D for HPC in Europe
  - **Evolution rather than revolution.** Revolution would require too much effort for runtime/apps porting
  - **Standardization of event messages across sources, Scalable tools**
  - Automatic or assisted diagnostic, root cause analysis, etc.

# WG4.2 Power Management

Chair: Thomas Ludwig, DKRZ

- Understand power consumption as a function of system usage:
  - Develop SW to control the HW mechanisms: SW in the OS, run-time system, etc.
  - Deploy next generation HW with energy saving mechanisms
  - A cost efficient and cost aware high performance computing is crucial for the competitiveness in science and engineering
- Main issues at Exascale:
  - Resilience (switching off/on components), Programmability (manual instrumentation), Performance optimization, Reproducibility (non-deterministic methods influence performance predictability and system noise)
  - **Power management API standardization**
- Recommendations for EU:
  - Trace tools: Vampir – Dresden, Scalasca – Jülich
  - Control daemons: Grid monitor – ParTec Munich, Germany, Power manager – Bull Cluster management France → Power capping, Power accounting)
  - **Leverage European leadership** in performance analysis tools, energy efficient hardware design and Know-how from embedded system

## WG4.2 Resilience

Chair: Franck Cappello, INRIA and UIUC

- Fault tolerance is a multi-facets and crosscutting issue:
  - **No guarantee that all faults will be masked by hardware or tolerated by the applications efficiently** (hardware correction will generate noise, global restart)
- Main issues at Exascale:
  - #errors, #failures, Scale, Power consumption, Performance impact
- Recommendations for EU:
  - 1) Establish a fault-error-failure model,
  - 2) Develop automatic root cause discovery and failure prediction,
  - 3) **Extend the applicability of checkpoint-restart** (API, Partial RST, NV mem.),
  - 4) Need a inter software layer communication system for consistent fault-error-failure management (detection, notification, decision, etc.),
  - 5) **New fault tolerance paradigms** using non-volatile memory technologies
  - 6) **International Coordination:**
    - Common FT API (interface and semantic) and inter layer communication system for portability (to limit applications rewriting)
    - G8-ECS (France, Germany, USA, Japan, Canada), ARN-JST FP3C (France and Japan), INRIA-Illinois Joint-laboratory (France and USA).

# WG4.2 Testbed needs in 2015 (draft)

————— < 100 Petaflops enough 
 ————— 100 Petaflops needed 
 —————

Dedicated/reconfigurable	<ul style="list-style-type: none"> <li>-<b>Prog. Models and Runtime</b> (interaction of runtime with kernel scheduler)</li> <li>-<b>Performance tools</b>: (interaction with resilience / power components)</li> <li>-<b>OS</b> (scheduling, memory management)</li> <li>-Measure noise generation (most of the syst software)?</li> <li>-<b>I/O-File system</b> (may be done at lower scale) (need root access for reconfiguration)</li> <li>-<b>Job and resource manager</b> (need root access for reconfiguration)</li> </ul>	<ul style="list-style-type: none"> <li>-<b>Performance tools</b> (system level): measure and minimize overheads induced at the full scale by low-level performance monitoring infrastructure tightly integrated with OS</li> <li>-<b>Prog. Models and Runtime</b> (interaction of runtime with kernel scheduler)?</li> <li>-...</li> </ul>
Production	<ul style="list-style-type: none"> <li>-<b>Compilers</b> (node level)</li> <li>-<b>Programming models</b> (node level API)</li> <li>-<b>Performance tools</b> (node level)</li> <li>-<b>Power management?</b></li> <li>-<b>Validation and correctness checking?</b></li> </ul>	<ul style="list-style-type: none"> <li>-<b>Resilience</b> (FT protocols, ABFT, NFTA, execution state storage)</li> <li>-<b>Parallel debuggers?</b> (scalability test)</li> <li>-<b>Performance tools</b> (scalability of data collection and analysis)</li> <li>-<b>Runtime</b> (scalability test)</li> <li>-<b>Prog. models</b> (system level scalability)</li> <li>-<b>Performance modeling at scale</b></li> </ul>

# WG4.2 Contributions to an Exa Stack (draft,



- Programming models
  - **StarSs**: Node level programming model offering programmers a natural interface to an underlying data-flow execution model and aiming allowing the “same” source code to run on “any” target architecture. In particular, **OmpSs** implementation that integrates OpenMP and the StarSs concept (SMP, GPUs,...) and nicely integrates with MPI. <http://www.bsc.es/smpsuperscalar> and <http://www.bsc.es/OmpSs>, BSC
  - **HMPP**: Based on a set of directives, programming model designed to handle hardware accelerators, [www.caps-entreprise.com/hmpp.html](http://www.caps-entreprise.com/hmpp.html) CAPS
  - **OpenHMPP** will be announced at SC11. Members of the consortium will be major HPC users, compiler companies such as Pathscale, INRIA
  - **T-Platforms Exascale API environments**: a new scalable version of MPI and a PGAS library with a new set of compilers and tools, Russia
  - European PGAS development called GPI (from Fraunhofer), [www.gpi-site.org](http://www.gpi-site.org)
- Runtime
  - **StarPU runtime system**, StarPU is a task scheduler exploiting multi-GPU multi-core platform efficiently, using adaptive performance models to carefully dispatch tasks over the available heterogeneous processing units. <http://runtime.bordeaux.inria.fr/StarPU/>
- Parallel Debugger
  - **DDT**: scalable parallel debugger for debugging MPI and multi-threaded codes, Alinea
- Validation/correctness
  - **MUST**: scalable Runtime Error Detection in MPI Programms, TU Dresden

# WG4.2 Contributions to an Exa Stack (draft,



## □ Performance tools

- **Scalasca**: Performance-analysis tool German Research School for Simulation Sciences and Juelich Supercomputing Centre
- **Score-P**: Generic measurement infrastructure for several performance tools including Periscope, Scalasca, TAU, Vampir, Virtual Institute - High Productivity Supercomputing <http://www.vi-hps.org>
- **BSC Integrated tools environment**: instrumentation (**Extrae**), visualization (**Paraver**), Simulation (**Dimemas**) and other analysis modules such as clustering, combined sampling and instrumentation, **Par** <http://www.bsc.es> , BSC
- Persicope (Technical University of Munich) <http://www.lrr.in.tum.de/~petkovve/psc/>
- ThreadSpotter (RogueWave) <http://www.roguewave.com/products/threadspotter.aspx>
- Vampir (Technical University of Dresden) <http://www.vampir.eu/>

## □ OS and system management

- **T-Platforms Clusterx** Exascale OS and system management: a highly scalable exascale-oriented OS and system software stack (highly scalable real-time monitoring and management system), Russia
- **Clustershell**: start commands on the cluster nodes (in parallel, with a high scalability based on tree) and to get a usable output (summary, differences, ...) for human [sourceforge.net/projects/clustershell/](http://sourceforge.net/projects/clustershell/), CEA
- **Ganesha**: "protect" the computing centre NFS servers from the load of the NFS clients [sourceforge.net/projects/nfs-ganesha/](http://sourceforge.net/projects/nfs-ganesha/), CEA
- **Robinhood** : a FS monitor/space manager optimized for huge FS (PB today, EB in the future) [sourceforge.net/projects/robinhood](http://sourceforge.net/projects/robinhood/), CEA
- **INRIA OAR**, *res. Job manager* , **Kadeploy**, fast and scalable deployment system, **Taktut**: parallel and scalable remote execution tool for cluster. <http://mescal.imag.fr/software.php>

## WG4.2 Roadmap to be completed by June

---

- Description of the scientific and technical perimeter
- Social benefits, societal, environmental and economical impact
- Scientific and technical hurdles
- Address cross cutting issues : Resilience, Power Mngt, Programmability, Perf optimization and Reproducibility of the results,
- European strengths and weaknesses in the worldwide competition
- Sources of competitiveness for Europe
- Needs of education and training
- Potential collaborations outside Europe
- Existing funded projects and funding agencies
- **Timeline, needs of HR, provisional costs,**
- Needs for experimental platform: size, reconfigurable?, dedicated?
- What software in an Exascale stack and with what level of responsibility