www.ppi4hpc.eu

**PPI4HPC**

Public Procurement of Innovations
for High Performance Computing

# Technical requirements

Open Dialogue Event
September 6, 2017 - Brussels

# Disclaimer

*For the avoidance of doubt this presentation is solely made for the purpose of informing the market and of initiating a technical dialogue with the market in order to prepare a joint procurement procedure.*

*It does not signify the beginning of a procurement procedure or constitute a commitment by the public procurers involved in the presentation to undertake such exercise at a later stage.*

*The final form of the procurement could differ from the form presented during this meeting.*

*Participation in this open dialogue event is not a precondition for responding to the planned procurement procedure.*

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

# PPI4HPC

Public Procurement of Innovations
for High Performance Computing

# Overview and organization
# of the technical specifications

Gilles Wiber, CEA

Open Dialogue Event
September 6, 2017 - Brussels

# Outline

www.ppi4hpc.eu

- Goal of the procurement

- Organization of the technical specifications

- Objectives of the next technical presentations

Open Dialogue Event
September 6, 2017

# Goal of the procurement (1/2)

The goal of the 4 public procurers involved (BSC, CINECA, JUELICH, GENCI) is, for each of them, to procure:

- An **innovative** high performance supercomputer and/or an **innovative** high performance storage system

- That will be integrated in their computing center to be used as **production** systems

- Including maintenance and support

Altogether, these systems will cover a large range of applications including traditional HPC application, HPDA and AI

Open Dialogue Event
September 6, 2017

# Goal of the procurement (2/2)

"**Innovative and production quality**" means that:

- The proposed IT equipment should implement major innovative technology and/or architecture not deployed so far in IT equipment already delivered of similar size

- Only final configuration, tuning and development/integration of tools for the site is expected to take place on site in collaboration between the supplier and the site

# Why a joint procurement ?

This joint procurement aims at fostering:

- Science and engineering applications in Europe within PRACE

- R&I on HPC architectures and technologies in Europe with strong relationship (collaboration) between the procurers and the suppliers

- A greater weight and more impact

  – on important topics of common interest

  – on the design of innovative solutions

⇨ PPI4HPC technical group leveraging the expertise of the organizations involved in the project

Open Dialogue Event
September 6, 2017

# Organization of the technical specifications

www.ppi4hpc.eu

| Technical specifications | Common | Local |
|---|---|---|
| Glossary | | |
| Definition of the system | | Sizing, target architecture, local context and constraints |
| Important topics | Important topics of common interest | Lot specific topics of interest |
| Evaluation criteria (mandatory / desirable features) | Common mandatory and desirable features | Lot specific evaluation criteria |

The PPI4HPC project has received funding from the European Union's Horizon 2020 research and innovation programme under the grant agreement № 754271.

Open Dialogue Event
September 6, 2017

8

# Organization of the technical specifications

| Technical specifications | Common | Local |
|---|---|---|
| Glossary | | |
| Definition of the system | | Sizing, target architecture, local context and constraints |
| Important topics | Important topics of common interest | Lot specific topics of interest |
| Evaluation criteria (mandatory / desirable features) | Common mandatory and desirable features | Lot specific evaluation criteria |

Common topics and features

Open Dialogue Event
September 6, 2017

# Organization of the technical specifications

| Technical specifications | Common | Local |
|---|---|---|
| Glossary | | |
| Definition of the system | | Sizing, target architecture, local context and constraints |
| Important topics | Important topics of common interest | Lot specific topics of interest |
| Evaluation criteria (mandatory / desirable features) | Common mandatory and desirable features | Lot specific evaluation criteria |

Description of lots

Open Dialogue Event
September 6, 2017

# Objectives of the next presentations

- Explain the current status of our reflection on technical specifications

    - Common

    - Lot specific

- Get feedback from the market (today, later including during one-to-one meetings)

    - State of the art

    - Objectives and roadmaps

- Investigate possible collaborations

Open Dialogue Event
September 6, 2017

# Next presentations

- Common technical specifications

  - Important topics of common interest

  - Focus on collaborations

- Description of lots

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

**PPI4HPC**

Public Procurement of Innovations
for High Performance Computing

# Important topics of common interest

Gilles Wiber, CEA

with the contribution of
Mirko Cestari, Carlo Cavazzoni, CINECA
Dorian Krause, Dirk Pleiter, JUELICH
Eric Boyer, GENCI
Javier Bartolomé, BSC

Open Dialogue Event
September 6, 2017 - Brussels

# Outline

- Motivations

- Topics of common interest

- Collaborations

# Motivations

# Motivations (1/2)

Major challenges for future generation supercomputers and storage:

- Increasing performance while keeping compute and data balanced
- Improve energy efficiency to stay within electricity costs budget
- Flexibility and ease of use despite increasing complexity of the hardware architectures
- Monitoring of system behavior and utilization
- Security
- Total Cost of Ownership optimization

In order to address these challenges, the PPI4HPC technical group has identified 8 important topics of common interest (part of a long term vision)

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

- **Energy efficiency and Power management**
  - Control, optimization, efficiency of energy and power are major goals of future systems

- **Data management**
  - Data storage and management is needed to keep pace and make supercomputer usable

- **Programming environment and productivity**
  - Usability and flexibility are key aspects of future computing center

- **Datacenter integration**
  - Enable links between systems and facility to optimize efficiency and operations

- **Maintenance and support**
  - Reliability, online maintenance are also a way to be more efficient

- **System and application monitoring**
  - Towards unified full-system monitoring capabilities that include production workload performance data

- **Security**
  - Security is a strategic aspect for HPC facilities, and its level will be determined according to legal regulation

- **Total Cost of Ownership**
  - Assess all the cost items related to an HPC system (OpEx & CapEx) to ensure that the best value for money is reached taking into account the future activity on the system

Open Dialogue Event
September 6, 2017

# Topics of common interest

# Nomenclature

Topics of common interest are structured according to goals, objectives, mandatory and desirable features:

## Goals

- These are long term and may not necessarily be fully achieved within the expected life-time of the system

- They come from important ideas and needs that all partners want to push forward

## Objectives

- These are measurable and attainable within the expected life-time of the system

- They might require collaborative efforts between procurers and suppliers

# Nomenclature

## Mandatory and desirable features

- Mandatory features are needed at the system installation timeframe

- All mandatory and desirable features are taken into account for the evaluation process of the different solutions

## Notice

Features may not apply to all lots and their weight may differ between lots. All partners endorse them and want to push them forward.

Open Dialogue Event
September 6, 2017

# Topics of common interest

- Energy efficiency and power management

- Data management

- Programming environment and productivity

- Datacenter integration

- Maintenance and support

- System and application monitoring

- Security

- Total cost of ownership

Open Dialogue Event
September 6, 2017

# Energy Efficiency and Power Management

## Goals

- Foster power management to dynamically cope with upper limits in power provisioning

- Improve energy efficiency of the system

## Objectives

- Correlation between power consumption and system workload

- Dynamic power capping with graceful performance degradation of the system

- Capability to optimize job execution environment for better energy efficiency

- Energy accounting mechanism

- Energy profiling of applications enabling Energy-To-Solution optimization without significant performance degradation

# Programming Environment and Productivity

**Mandatory and desirable features**

- Provide node power measurements with high level of accuracy and low impact on performance (< 1% on Linpack)

- Provide Tools/API for job energy accounting to be integrated with the resource scheduler

  - Energy values for the individual subsystem delivered in the procured system

- Provide Tools/API for analysis of the energy and correlation with performance of applications

- Provide Tools/API enabling the integration of power management framework and the resource scheduler

Open Dialogue Event
September 6, 2017

# Data Management

**Goals**

- Data access performance and capacity has to keep pace of HPC computing performance evolution

- Reduction of data access relative to computing time

- Improve actual data monitoring to avoid failures or auto-correct them (silent data corruption)

**Objectives**

- Reduce time for integrity checks and for recovery in case of failures

- Improve monitoring of data solutions healthiness

- Provide tools to obtain global overview on usage of the data infrastructure

- Reduce time spent in I/O

Open Dialogue Event
September 6, 2017

# Data Management

**Mandatory and desirable features**

- Provide an efficient automatic recovery procedure in case of failures

- Provide solutions to foresee and react ahead of future failures to keep normal functioning

- Provide efficient tools to show the usage and type of data stored in them

- Optimize and reduce time spent related to data access of the actual computing processes. Examples of Innovative solutions that may help:

  - Smart tiering or smart data placement

  - Data-aware scheduling

  - Impact reduction of data-related daemons and processes

Open Dialogue Event
September 6, 2017

# Programming Environment and Productivity

## Goals

- Provide up-to-date development software stack to European users on all PPI4HPC systems
- Support of (lightweight) virtualization mechanism enabling European users to deploy specific software stacks seamlessly on all PPI4HPC systems

## Objectives

- Provide support for recent implementations of parallel programming standards
- Provide support and system integration for customizable environment (e.g. Containers and/or Virtual machines)

Open Dialogue Event
September 6, 2017

# Programming Environment and Productivity

**Mandatory and desirable features**

- Provide Compilers C, C++ and Fortran that support:

    - Fortran 2008 or newer, C ISO/IEC 9899:2011 or newer, C++ ISO/IEC 14882:2014 or newer

- Provide a parallel development environment that supports:

    - OpenMP 4.0 or newer, MPI 3.0 or newer

- Provide an environment that supports the use of containers and/or virtual machines

Open Dialogue Event
September 6, 2017

# Data Center Integration

**Goals**

- Advance system infrastructure design (incl. system/facility interface) addressing equally:
    - Requirements for higher density, lower power consumption, better efficiency, etc.
    - Requirements for easier and safer operation, quicker installation, higher reliability, etc.

**Objectives**

- Minimize need for adapting facility infrastructure to system: Enable Quicker buildup, easier operation
- Improve HW and SW installation process to shorten system setup time
- Improve resilience against short facility-related power interrupts
- Improve energy efficiency of supercomputers and facility
- Enable communication links between system and facility: Exchange of status information and alarms

Open Dialogue Event
September 6, 2017

# Data Center Integration

**Mandatory and desirable features**

- Liquid cooling loops with special/high water quality requirements shall be closed and maintained by the vendor

- Enable compute nodes to be resilient w.r.t. variations of the power source (incl. Micro power cuts, periodic oscillations, etc.) unless the power interrupt lasts longer than 300ms

- System-internal mechanisms for infrastructure failure detection and automatic system reaction

- Software installation during system bring-up in < 5 days

- Enable facility ↔ system communication of

   - Status information (power, temperatures, etc.)

   - Alarms

# Maintenance and Support

## Goals

- Keep to minimum or eliminate any maintenance downtime
- Keep or reduce job failure rate due to external causes

## Objectives

- Reduce the time needed in HPC infrastructures for maintenance
- Automatic hardware and software validation framework that helps to evaluate healthiness of all HPC components
- Predictive failure detection systems to reduce the rate of job failures to hardware issues

# Maintenance and Support

**Mandatory and desirable features**

- Improve systems designs to support rollout maintenances, reducing the disruption on a whole system due to a change/update

- Minimize the full system power on or reboot time until "ready for production" state (30 minutes or less)

- Provide tools to check and validate all hardware and software components health, which:

  - Performs auto-recovery procedures when possible

  - Perform the actions needed to avoid the usage of damaged components in future jobs/processes

Open Dialogue Event
September 6, 2017

# System and Application Monitoring

**Goals**

- Continuous lightweight profiling of production jobs in order to identify hidden optimization potentials in applications and system software

- Drive integration of monitoring capabilities

**Objectives**

- Enable collection, management and analysis of performance data for production jobs

- Enable generation of user reports
    - Desire: Use of common performance metrics (CoE POP) across European sites

- Enable anomaly detection based on collected data

- Unify system monitoring capabilities
    - Integration of different monitoring technologies (incl. external ones) for higher situational awareness for operations team

Open Dialogue Event
September 6, 2017

# System and Application Monitoring

**Requirements and Desirable Features**

- Provide lightweight continuous performance profiling capabilities
    - Min. perf (< 5%) and scalability, always applicable
    - Coverage: CPU, IPC, instr. mix; memory, cache and TLB metrics; I/O subsystem; network (RDMA); accelerators; MPI/comm. Libraries
    - Extensibility, co-existence with 3rd party profiling tools
- Scalable accumulation, compression and reduction mechanisms for data handling
- Performance report generation
    - Adjustable level of detail and focus; Integration with workload manager; API for connection with external systems (web portals)
- Anomaly detection based on data
    - Detect performance anomalies (system comp., software changes)
- Software for correlation of different metrics from different monitoring systems (including external ones)
    - Open API for near-real time access and export; graphical tools desirable

Open Dialogue Event
September 6, 2017

# Security

www.ppi4hpc.eu

## Goals

- Maintaining the security level in accordance with site policies and applicable rules/laws (including on data)
- Guaranteeing users groups isolation
- Preventing deny of service attempts

## Objectives

- Periodic security patches identification and application
- Enforce isolation of different user/group on the cluster (including data)
- Preventing deny of service attacks resulting in a significant degradation of system services by using system features or tools

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

**Mandatory and desirable features**

- Minimize the time to provide security patches solution: 48h or less after publication of a solution/workaround by the OS vendor for critical impact level breaches

- Provide live patching feature

- Improve or provide features or tools to confine resources: resources allocation feature, client and/or server side rate throttling

- Provide features for job resources usage isolation (at user / group level)

# Total Cost of Ownership

**Goal**

- Ensure that all the costs related to the system during its lifetime in a given site are identified for the selection process of the system

**Objectives**

- Select a system using a metric based upon value for money where :
    - Acquisition and operational costs are taken into account
    - Both time-to-solution and energy-to-solution matter
- Make an incentive to the vendors to propose the most efficient technology and most complete energy optimization environment.

**Assessment**

- A formula including CapEx and a projected OpEx under a customized set of applications is set up to evaluate the service provided

**Note**

- Formula and variables used are site depend and may differ in each local lot

Open Dialogue Event
September 6, 2017

# Collaborations

# Collaborations

## Context

- Facilitate continuous improvements during lifetime of the system
- Create impact on product development by means of a collaboration process

## Objectives

- Agree on cooperation between public procurer and contractor for improving the system (e.g. in terms of performance, energy efficiency, or functional capabilities) during the lifetime of the system

- Establish a collaboration process based on tight interaction between experts of the contractor and the public procurer

**Description of lots**

Open Dialogue Event
September 6, 2017 - Brussels

# Outline

- Partner technical projects summary

- For each partner

  - Computing center description

  - Technical project overview

  - Special topics interest

- Conclusion

Open Dialogue Event
September 6, 2017

# Partner technical projects summary

| Partner | Date | Type | Situation | Application |
|---------|------|------|-----------|-------------|
| CINECA | 2019-2020 | Supercomputer, several tens of PF | Replacement of Marconi-KNL | Large range of application including traditional HPC applications, HPDA and AI |
| BSC | 2019 | Storage up to hundreds PB, and nodes for data analytics | Complement to BSC storage | |
| GENCI/CEA | 2019 | Supercomputer, several tens of PF | Complement to CURIE2 | |
| JUELICH | 2020 | Supercomputer | Replacement of JURECA | |

www.ppi4hpc.eu

**PPI4HPC**

Public Procurement of Innovations
for High Performance Computing

# CINECA Lot

Mirko Cestari

with the contribution of Carlo Cavazzoni
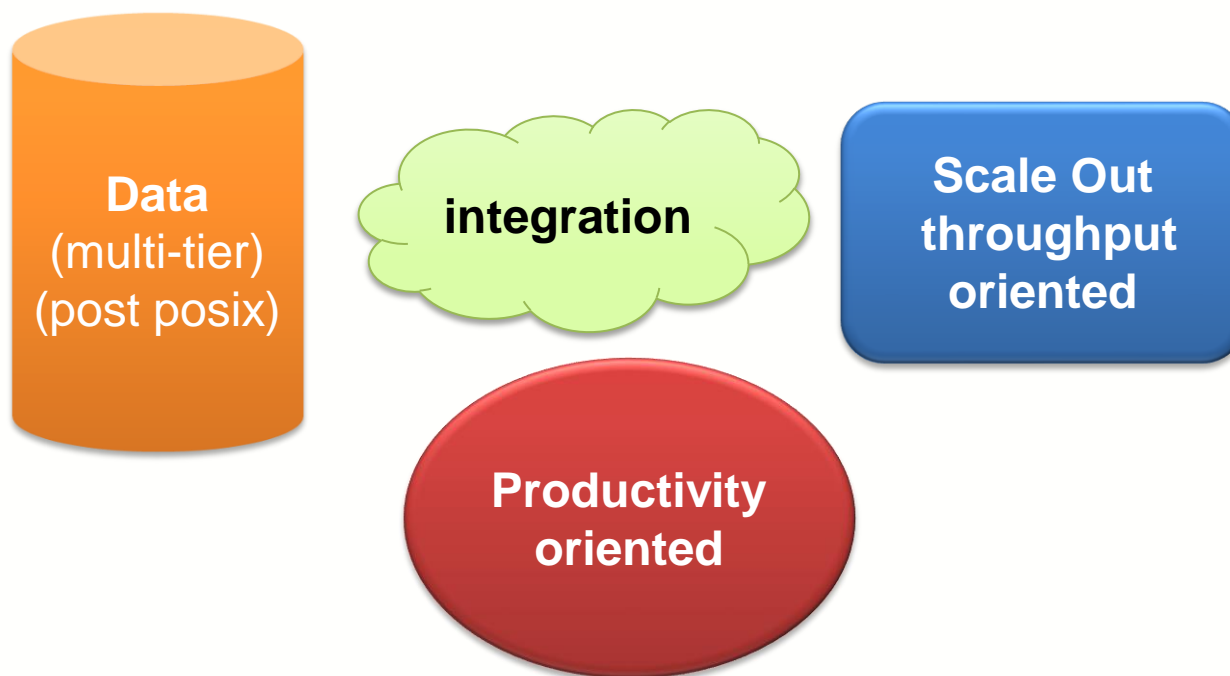
Open Dialogue Event
September 6, 2017 - Brussels

# CINECA Computing Center

**100TF**
1MW

**20x**

**2PF**
1MW

Paradigm change

**5x**

10x
(in total)

**11PF+**
9PF
3.5MW

1x
(latency cores)

**5x**

**50PF+**
10PF
~4MW

solid

2x
(latency cores)

**5x**

**>250PF+**
>20PF
~8MW

Pre-exascale

| 2009 | 2012/2013 | **2016/2017** | 2019/2020 | 2021/2022 |
|------|-----------|---------------|-----------|-----------|
| IBM SP6 Power6 | Fermi IBM BGQ PowerA2 | **Marconi Lenovo Xeon+KNL** | ?? Scalar+Vect or Acc. | ?? Scalar+Vect or Acc. |

43

# Functional View

**Data**
(multi-tier)
(post posix)

**integration**

**Scale Out throughput oriented**

**Productivity oriented**

Open Dialogue Event
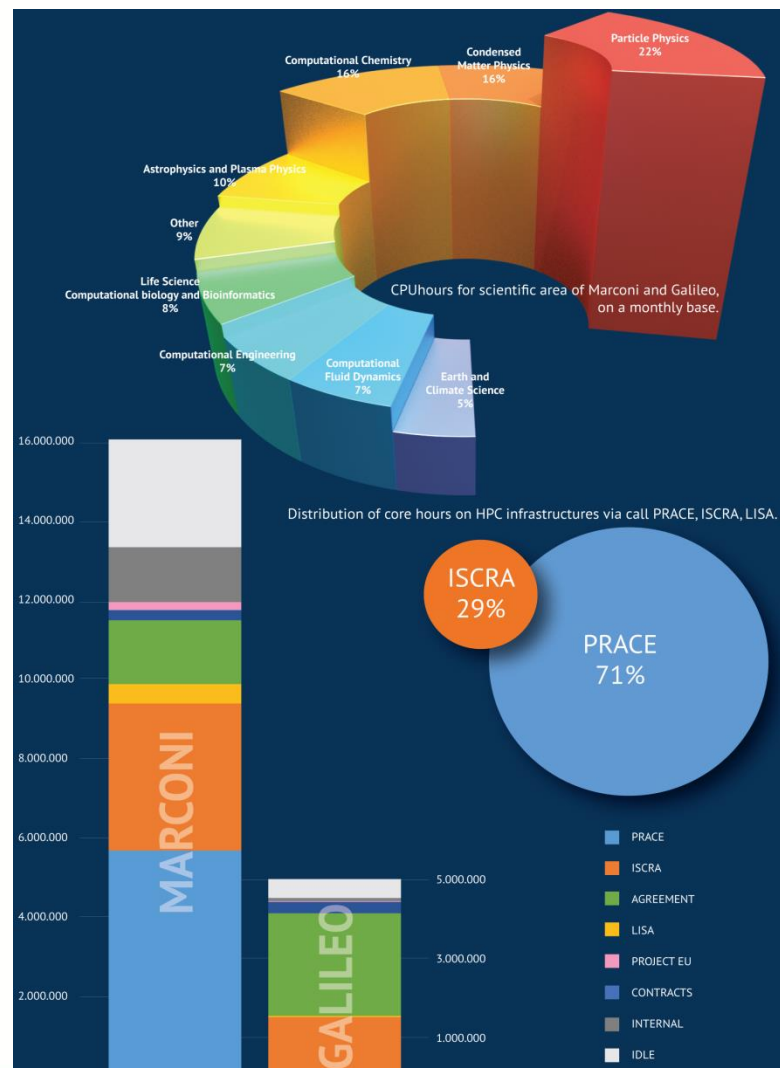September 6, 2017

# Goal


www.ppi4hpc.eu

Procure the replacement of MARCONI-KNL system:

- Largest partition of the system

  – 10 PF peak performance

  – 3600 nodes, 240K cores

  – 100 Gb network interconnection

- Scale-out and high throughput workloads

  – Maximize throughput per time unit

  – Exploit parallelism of applications

- Integrate with CINECA ecosystem

  – Connect to other partitions of Marconi (productivity, Cloud)

# Applications Workload

- **EU research**
  - PRACE Tier-0 system
  - Support for COE and FET-HPC projects

- **National research**
  - ISCRA, Italian SuperComputing Resource Allocation

Open Dialogue Event
September 6, 2017

# Procurement Target

- Hardware

  - 5x increase of system capability is achievable

  - Maximize efficiency (capability/W)

- Network

  - Next gen high speed interconnection

  - Bi-sec bdw equivalent to fat tree w/ blocking factor 2:1

- I/O

  - Working space

  - Global namespace

# CINECA Special Topics of Interest

- Power and performance

    – Special focus on efficient power monitoring and management

    – Power capping


- Architecture solution towards exascale

    – Either network, socket or nodes (or all of them)


- Data center integration

    – More than one option could be in place

        – Current premises or new data center

    – Solutions to leverage the available options

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

# PPI4HPC

Public Procurement of Innovations
for High Performance Computing

## JUELICH Lot

Dorian Krause

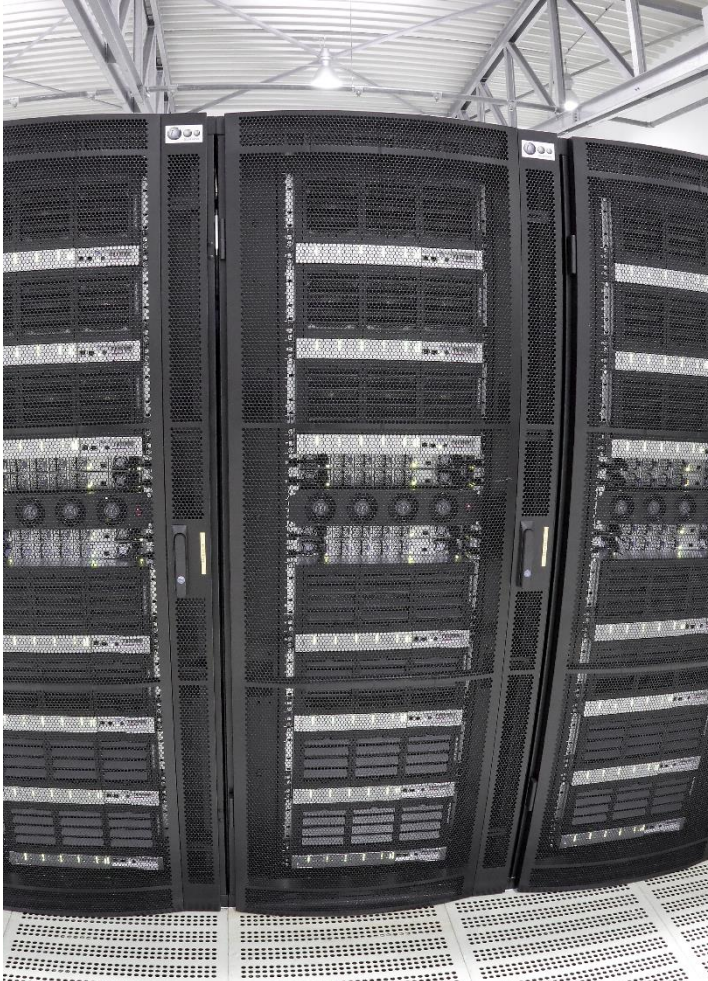Open Dialogue Event
September 6, 2017 - Brussels

# JUELICH Computing Center



**IBM Power 4+**
**JUMP, 9 TFlop/s**

**IBM Power 6**
**JUMP, 9 TFlop/s**

**IBM Blue Gene/L**
**JUBL, 45 TFlop/s**

**JUROPA**
**200 TFlop/s**

**HPC-FF**
**100 TFlop/s**

**File Server**

**IBM Blue Gene/P**
**JUGENE, 1 PFlop/s**

**JURECA**
**(2015)**
**2.2 PFlop/s**

**IBM Blue Gene/Q**
**JUQUEEN**
**5.9 PFlop/s**

**JURECA**
**Booster (2017)**
**5 PFlop/s**

**Data Analytics**

**PPI4HPC**
**(2020)**

**JUST**

Open Dialogue Event
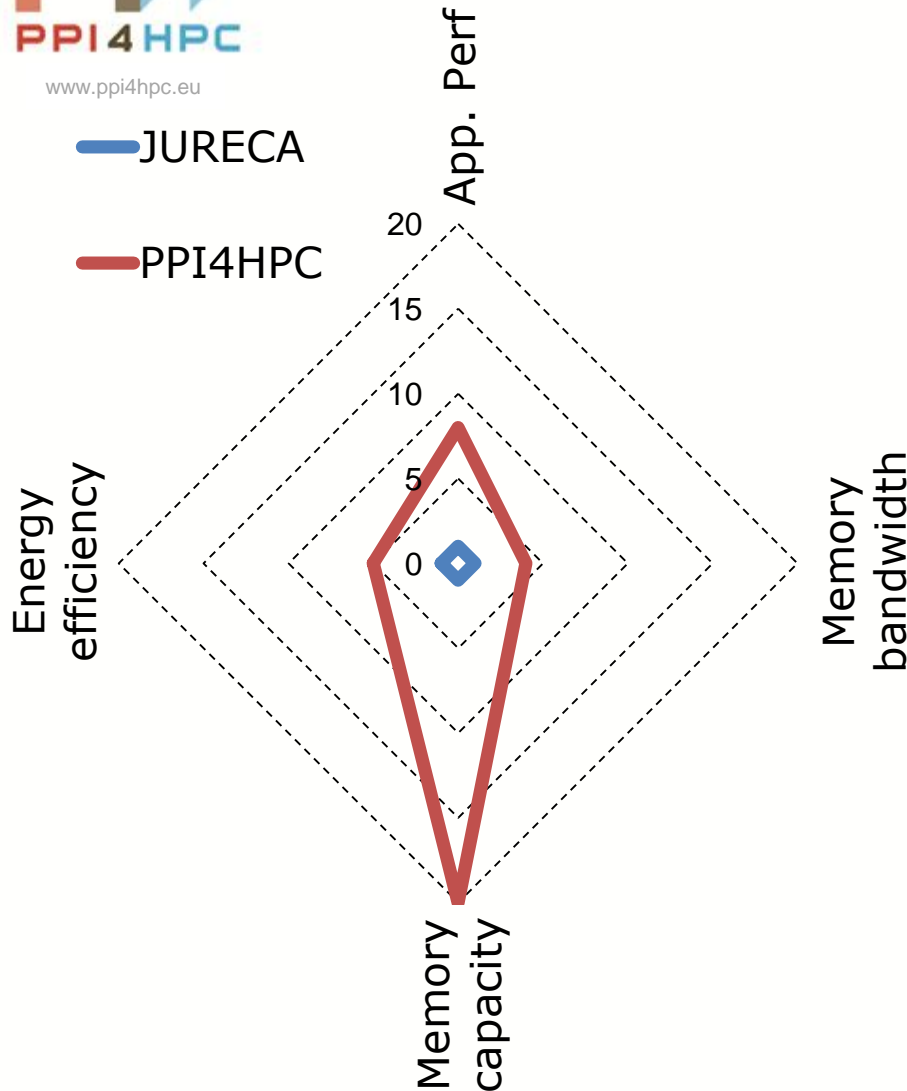September 6, 2017

www.ppi4hpc.eu

# JUELICH Technical Project

- **Target system:** Successor of the JURECA Cluster system

- JURECA characteristics:
  → 1,872 compute nodes
  → Dual socket Haswell CPUs
  → Partition with K80 GPUs
  → InfiniBand EDR full fat tree
  → Peak perf.: 2.2 PFlop/s
  → 258 TiB memory capacity
  → ca. 255 TiB/s memory bandwidth
  → Connection to JUST central storage cluster
  → High-speed network connection to Booster module
  → 34 racks, rear-door heat exchanger

Open Dialogue Event
September 6, 2017

# JUELICH Technical Project



- **Desired successor characteristics (rel. JURECA):**
  → 6-8 × application performance
  → 3-4 × memory bandwidth
  → (4 + 16) × memory capacity (volatile + non-volatile)
  → Integration with central storage
  → Same or less space
  → Free outdoor cooling capabilities

- **Time Frame:** System installation in 2020

# JUELICH Special Topics Interests

- Architecture targeting mixed capability and capacity workloads from a wide variety of communities
  → Earth systems modeling, material and life sciences, neuroscience, fundamental sciences, engineering

- **Novelty:** Features for data-intense science and data-analytics applications
  → Applicable to broader use-cases
  → Expecting innovation in hardware and software technology for NVM integration

- JURECA Booster module to be in production until 2022
  → Integration with PPI4HPC system
  → Current cross-connection bandwidth: 20 Tb/s

- Warm-water cooling enabling 65+ percentage free cooling throughout the year

www.ppi4hpc.eu

**PPI4HPC**

Public Procurement of Innovations
for High Performance Computing

# GENCI Lot

Eric Boyer

Open Dialogue Event
September 6, 2017 - Brussels

# CEA/GENCI computing center



- HPC clusters:

    – Research: General purpose cluster: CURIE2

        – => 9PF with 4.5PB local storage

    – Industrial: General purpose cluster: COBALT

        – => 1.4PF with 2.5PB local storage

- Global Storage infrastructure:

    – HSM Lustre FileSystems with 23PB on disk
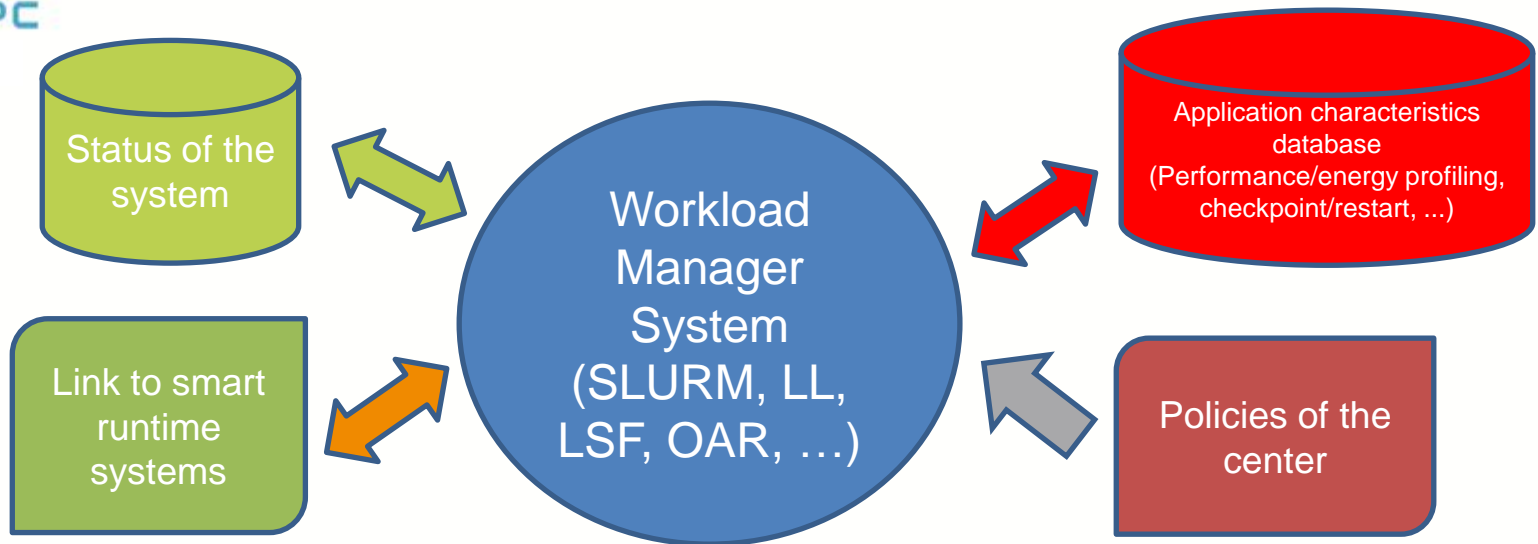      (210GB/s) and 30PB on tape

# GENCI technical project

- Will complement the CURIE2 Tier0 at TGCC
  - Must be a production « usable » system – High level of stability expected
  - Balanced architecture, not only FLOPS → Compute / network performance / memory footprint / IO performance
  - Novel architectures welcome if compliant with standard programming models (aka OpenMP 4.x)
  - Will address both national and European (PRACE) research needs from academia & industry
  - Could complement other PRACE Tier0 architectures ?

- Context of convergence between HPC / Data Intensive / AI
  - Addressing needs of new communities
  - Integration of (heterogeneous) software stacks ?
  - Availability of flexible isolation / allocation resource environments

# GENCI special topics interests

- High performance computing & data analytics solution relying on high energy efficiency architectures, embedding features of power management, energy aware scheduling, energy accounting and fine grain energy application profiling and tuning tools.

- Integration of upcoming data technologies and architectures (object storage, tiering, new storage levels and technologies integration).

- Extended monitoring features and permanent information collection features targeting application resource usage and behavior.

- Strong focus on security for low delay breaches correction, isolation of resources and information access level management and prevention of deny-of-service.

- Strong focus on smart resources manager (see next slide)

- Focus on flexibility to use the machine (containers, virtualization)

- Total Cost of Ownership as a main criterion for machine selection

Open Dialogue Event
September 6, 2017

# Focus on Smart resource manager expectations

Status of the system

Link to smart runtime systems

Workload Manager System (SLURM, LL, LSF, OAR, …)

Application characteristics database (Performance/energy profiling, checkpoint/restart, ...)

Policies of the center

WMS is interconnected to all the main databases and system components
Thus it's able to:
- Automatically detect, perform first check/fix and isolate defaulting resources
- By knowing apps characteristics, (future)load of the system and policies of the center, allocate/place at the right level of resources using smart runtimes and ML techniques
- Smart application based checkpoint/restart for urgent computing, too long workloads
- Provide reports at the user or system admin level through continuous monitoring, …

www.ppi4hpc.eu

PPI4HPC

Public Procurement of Innovations
for High Performance Computing

# BSC Lot

Javier Bartolomé

Open Dialogue Event
September 6, 2017 - Brussels

# BSC computing center



- MareNostrum4

  – PRACE Tier-0 computing infrastructure

  – General purpose cluster: 11.15 PF

  – Emerging technologies clusters: 2.5 PF (NVIDIA,Power,ARM,KN*)

- Storage infrastructure

  – 14 PB GPFS accessed by HPC clusters

  – 4.7 PB GPFS archive of data on disk

  – 6 PB tape storage for backup and archive

# BSC technical project

- Creation of a data infrastructure with hundreds of petabytes capacity

  - Tiered storage solution expected

  - Provide archival storage to BSC HPC clusters

- High performance data analytics nodes

  - Attached to the previous data infrastructure

  - Perform pre and post process of data simulation

  - Include new storage/memory technologies

  - Will permit new data analytics paradigms

Open Dialogue Event
September 6, 2017

# BSC special topics interests

- Storage innovative technologies will be expected related to:

    - Smart tiering technologies to mitigate latency/bandwidth access to slower storage technologies

    - Increase efficiency of data movement between filesystems

    - Integration and usage of new memory/storage technologies to ingest and analyze huge quantities of data

Open Dialogue Event
September 6, 2017

# Conclusion

- The technical specifications are structured in common and local specifications

- The important topics of common interest are based on ideas and needs that all partners want to push forward

- We are looking for feedback from the market (today, later including during one-to-one meetings)

    – State of the art

    – Objectives and roadmaps

- We are interested to investigate possible collaborations

Open Dialogue Event
September 6, 2017

www.ppi4hpc.eu

**PPI4HPC**

Public Procurement of Innovations
for High Performance Computing

# Technical requirements – Q&A session

Moderator: Dirk Pleiter, JUELICH

Open Dialogue Event
September 6, 2017 - Brussels

The PPI4HPC project has received funding from the European Union's
Horizon 2020 research and innovation programme under the grant
agreement Nº 754271