# Site Report: Jülich Supercomputing Centre

2015-09-16  |  SLUG 2015  |

# Jülich Supercomputing Centre (JSC)

# Jülich Supercomputing Centre

## Supercomputer operation for:

- Centre – FZJ
- Region – RWTH Aachen University
- Germany – Gauss Centre for Supercomputing
  John von Neumann Institute for Computing
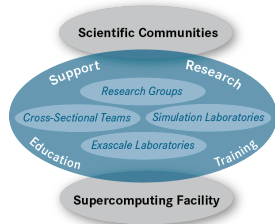- Europe – PRACE, EU projects

## Application support

- Unique support & research environment at JSC
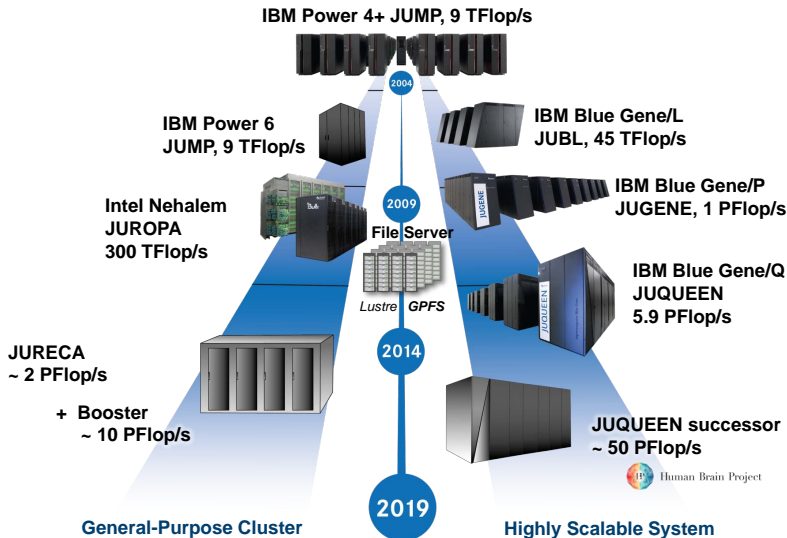- Peer review support and coordination

## R&D work

- Methods and algorithms, computational science, performance analysis and tools
- Scientific Big Data Analytics
- Computer architectures, Co-Design
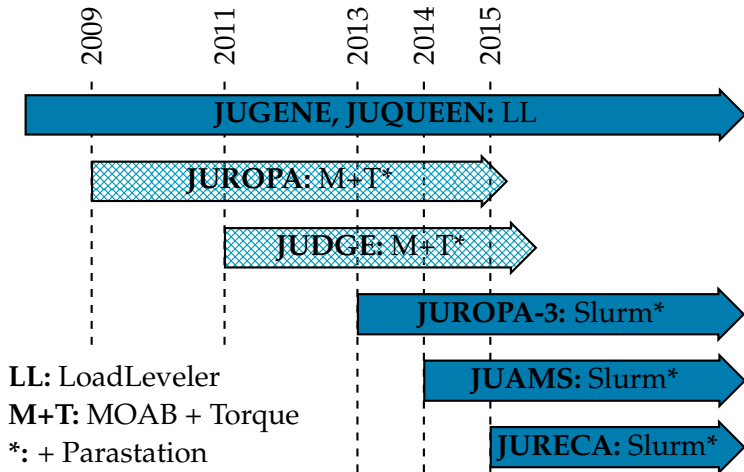  Exascale Laboratories: EIC, ECL, NVIDIA
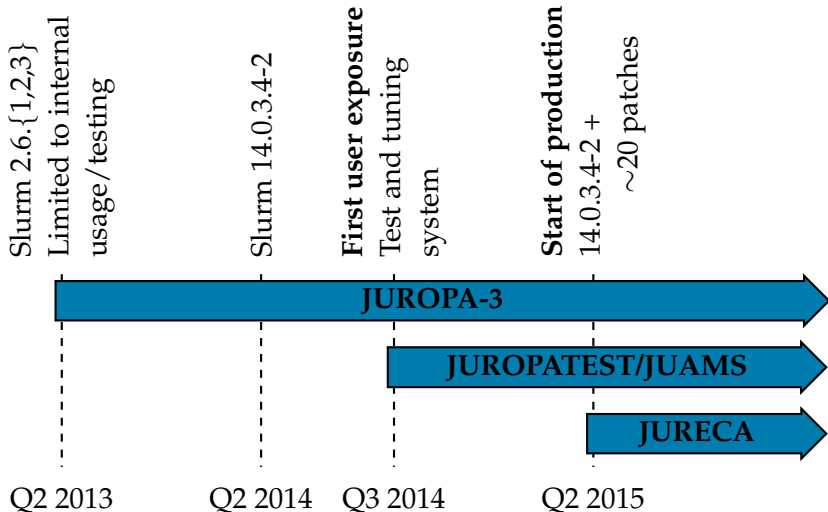
## Education and Training

# Supercomputer Systems: Dual Architecture Strategy

JÜLICH FORSCHUNGSZENTRUM

IBM Power 4+ JUMP, 9 TFlop/s

2004

IBM Power 6
JUMP, 9 TFlop/s

IBM Blue Gene/L
JUBL, 45 TFlop/s

IBM Blue Gene/P
JUGENE, 1 PFlop/s

Intel Nehalem
JUROPA
300 TFlop/s

2009

File Server

Lustre   GPFS

IBM Blue Gene/Q
JUQUEEN
5.9 PFlop/s

JURECA
~ 2 PFlop/s

2014

+ Booster
~ 10 PFlop/s

JUQUEEN successor
~ 50 PFlop/s

Human Brain Project

2019

**General-Purpose Cluster**

**Highly Scalable System**

# (A subset of) workload managers at JSC

LL: LoadLeveler
M+T: MOAB + Torque
*: + Parastation

JUGENE, JUQUEEN: LL
JUROPA: M+T*
JUDGE: M+T*
JUROPA-3: Slurm*
JUAMS: Slurm*
JURECA: Slurm*

# Slurm at JSC

Slurm 2.6.{1,2,3}
Limited to internal usage/testing

Slurm 14.0.3.4-2

**First user exposure**
Test and tuning system

**Start of production**
14.0.3.4-2 + ~20 patches

**JUROPA-3**

**JUROPATEST/JUAMS**

**JURECA**

Q2 2013          Q2 2014          Q3 2014          Q2 2015

# Slurm on JSC clusters

```
master 1
slurmctld
slurmdbd
mariadb
```

```
GPFS or DRDB
```

```
master 2
slurmctld
(slurmdbd)
(mariadb)
```

JÜLICH
FORSCHUNGSZENTRUM

# Slurm on JSC clusters

**master 1**
`slurmctld`
`slurmdbd`
`mariadb`

GPFS or DRDB

**master 2**
`slurmctld`
`(slurmdbd)`
`(mariadb)`
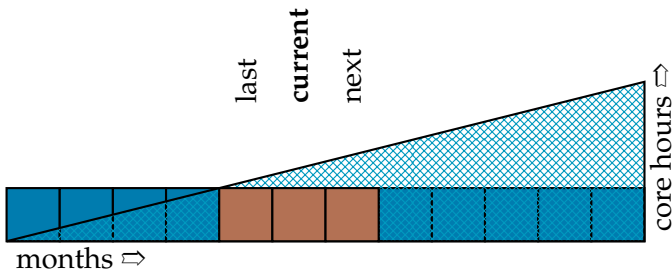
**login/batch host**
Slurm UI
99% upstream
compatible

# Slurm on JSC clusters

**master 1**
**slurmctld**
**slurmdbd**
**mariadb**

GPFS or DRDB

**master 2**
**slurmctld**
**(slurmdbd)**
**(mariadb)**

**login/batch host**
Slurm UI
99% upstream
compatible

**compute node**
**psid** + **psmunge** +
**psslurm**

# Llview: WM-centric system monitoring

# Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!



**account:** +200%
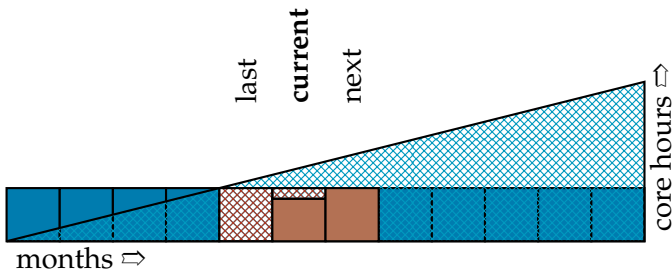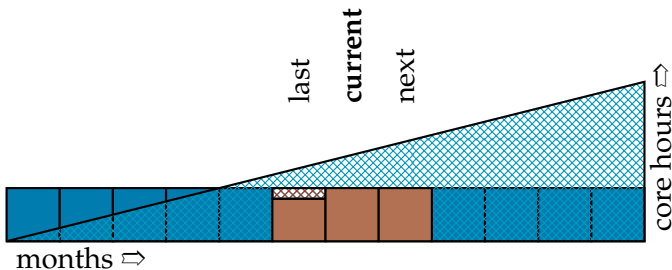**qos: normal**

Member of the Helmholtz-Association

# Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!



**account:** +140%
**qos: normal**

# Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!

**account:** +80%
**qos: normal**

## Batch model

- Associations **==** project (primary group), user (user id)
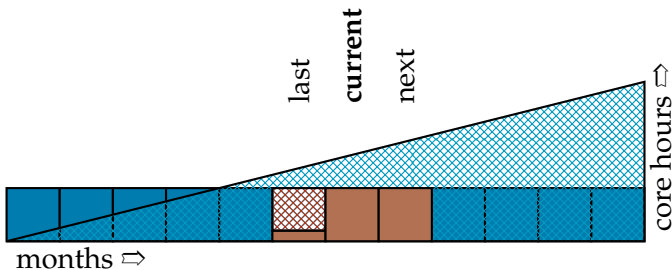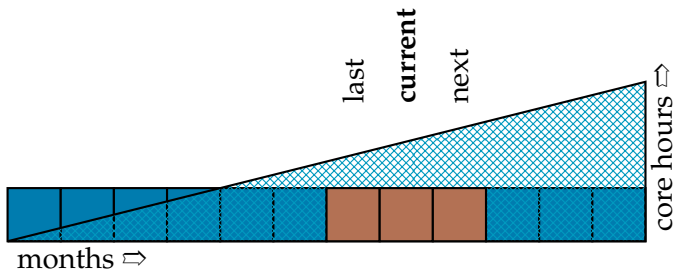- Scheduling with multifactor priorities $\Rightarrow$ QoS!



**account:** +180%
**qos: normal**

# Batch model

- Associations **==** project (primary group), user (user id)
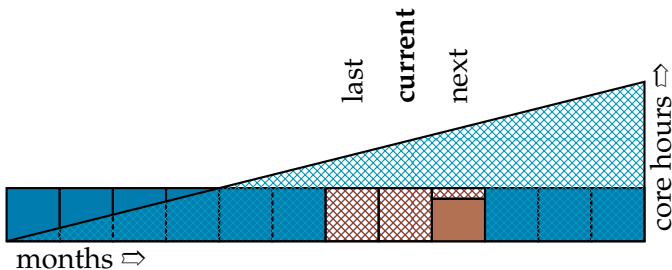- Scheduling with multifactor priorities $\Rightarrow$ QoS!
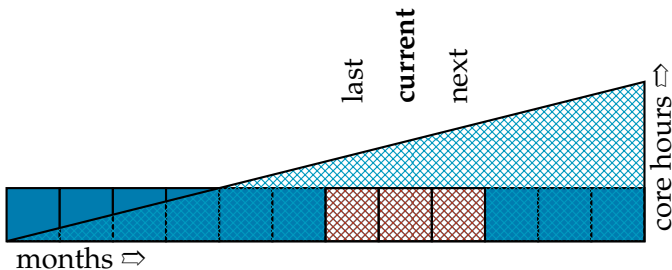


**account:** +120%
**qos: normal**

## Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!



**account:** +200%
**qos: normal**

Member of the Helmholtz-Association

## Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!



**account:** -20%
**qos: normal**

## Batch model

- Associations **==** project (primary group), user (user id)
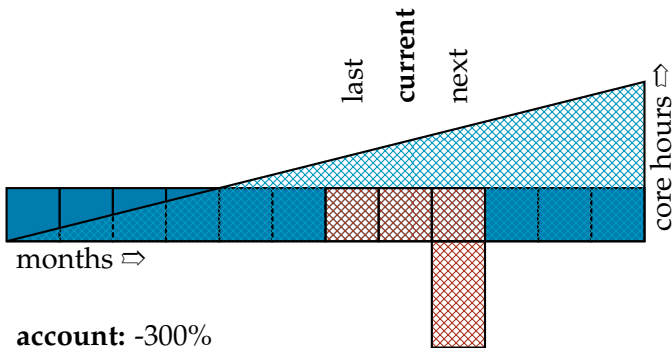- Scheduling with multifactor priorities $\Rightarrow$ QoS!



**account:** -100%
**qos: nocont** $\Rightarrow$ Low priority, wallclock limitation

## Batch model

- Associations **==** project (primary group), user (user id)
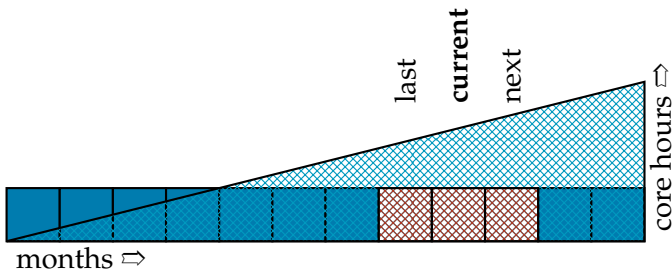- Scheduling with multifactor priorities ⇨ QoS!



**account:** -300%
**qos: nocont** ⇨ Low priority, wallclock limitation

# Batch model

- Associations **==** project (primary group), user (user id)
- Scheduling with multifactor priorities ⇨ QoS!

**account:** -100%
**qos: nocont** ⇨ Low priority, wallclock limitation

# Noteworthy/Exotic settings

- Custom X forwarding mechanism (**–forward–x**) for **srun**
- Custom CPU binding code
  - Designed for the 99% while leaving freedom for the 1%
  - Example: Do not "pack" **–exclusive** job steps
- Non-consumable gres (**mem{128,256,512,1024}**) to handle heterogeneities
  - Features are not captured in accounting database
  - **Rule of thumb:** If it is not in the accounting database we cannot use it

Member of the Helmholtz-Association

# Slurm experience

- Overall very positive experience
- Largely stable
  - Five internal tickets open against **slurmctld** (14.03.4-2):
    1. **assoc_mgr** lock contention during share updates
    2. **slurmctld** crash during assoc update
    3. Split brain occurence
    4. Invalid **protocol_version** in **_pack_cred()** during job launch
    5. Bogus projected start/end times
- Open source is a big plus
  - Increased insights
  - Ability to push features based on own needs/roadmap

Member of the Helmholtz-Association

# End of presentation

Member of the Helmholtz-Association

# JUQUEEN: Jülich's Scalable Petaflop System

IBM Blue Gene/Q JUQUEEN

- IBM PowerPC® A2 1.6 GHz,
  16 cores per node

- 28 racks, 458,752 cores

- 5,9 Petaflop/s peak
  5,0 Petaflop/s Linpack

- 448 TByte main memory

- connected to a Global Parallel File System (GPFS) with
  O(10) PByte online disk and O(50) PByte offline tape capacity

- 5D network

- Production start: Nov 5, 2012

Jun 2015:
#2 in Europe
#9 worldwide
#51 in Green500

JÜLICH
FORSCHUNGSZENTRUM

# JURECA: Jülich Research on Exascale Cluster Architectures

JURECA, an Intel-based cluster



**System integrator:
T-Platforms, Russia**

- 2 Intel Haswell 12-core processors, 2.5 GHz, SMT, 128 GB main memory

- 1,884 compute nodes or 45,216 cores, thereof
  75 nodes with 2 K80 NVIDIA graphics cards each and
  12 nodes with 512 GB main memory and 2 K40 NVIDIA graphics
  cards each for visualisation

- 2.245 Petaflop/s peak (with K80 graphics cards)
  ?? Petaflop/s Linpack

- 281 TByte memory

- Mellanox Infiniband EDR

- Connected to the GPFS file system on JUST

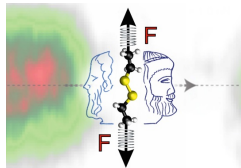# R&D and Application Support at the Jülich Supercomputing Centre
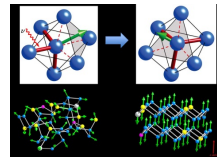
# High Impact Publications

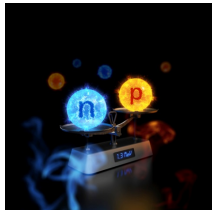Users of the facility at JSC produce about 250 publications per year



S. de Beer, M. Müser
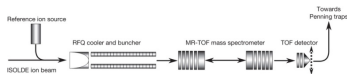Nature Communications **5**
(2014) 3781



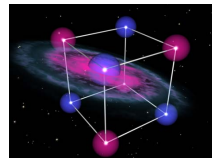D. Marx et al.,
Nature Chemistry **5**
(2013) 685



R.O. Jones et al.,
Nature Materials **10**
(2011) 129



Sz. Borsanyi et al.,
Science **347** (2015) 6229



A. Schwenk et al.,
Nature **498** (2013) 346



M. Lezaic et al.,
Nature Materials **9**
(2010) 649