

# Status of FZJ-IAS activities on the K computer

Wolfgang Frings

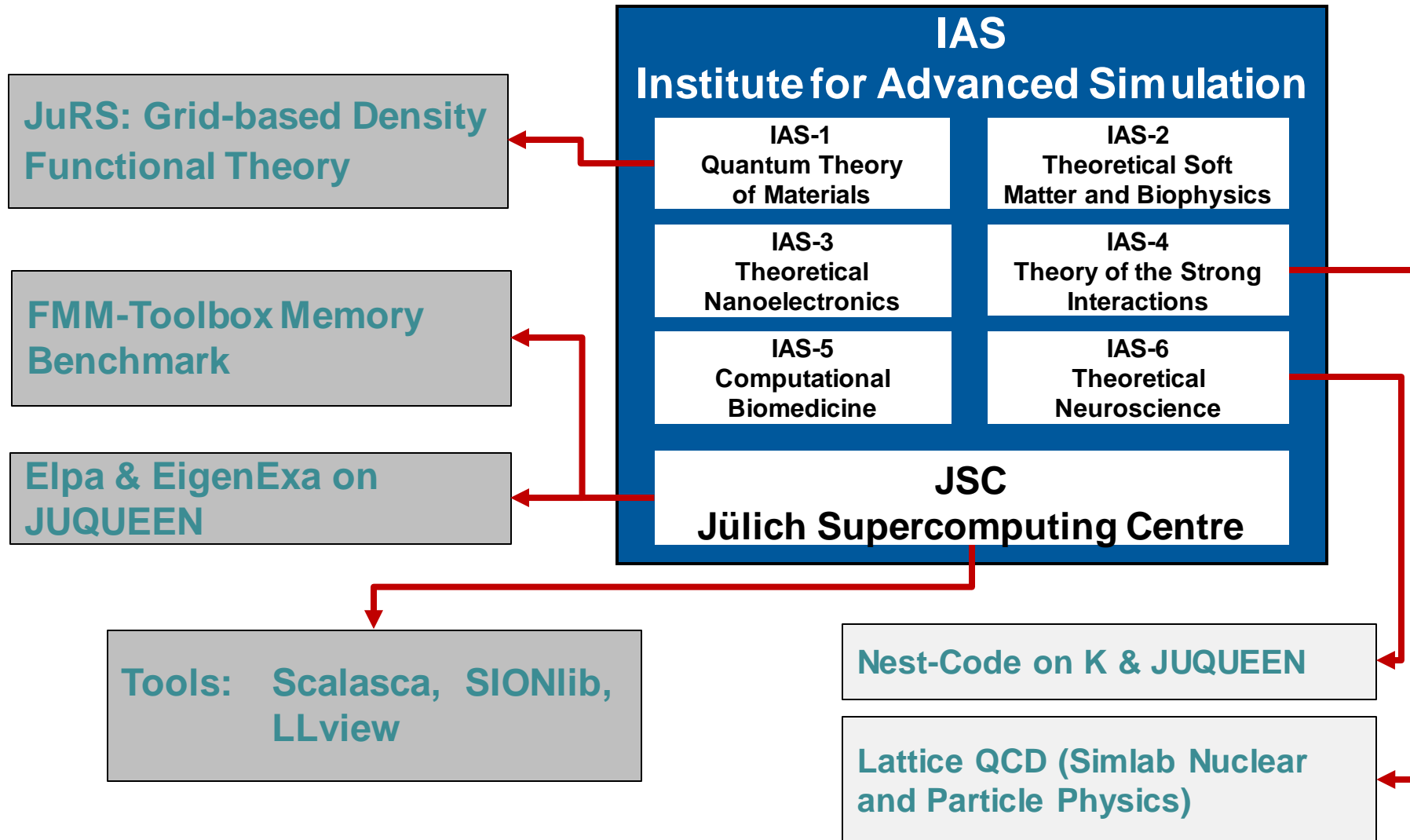
[W.Frings@fz-juelich.de](mailto:W.Frings@fz-juelich.de)

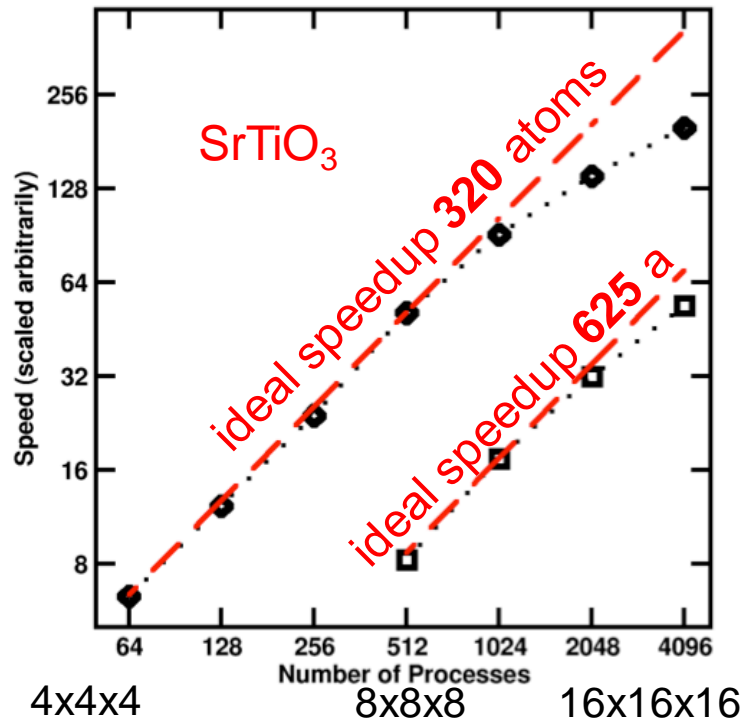
Jülich Supercomputing Centre

Riken, Japan, 8. September 2014

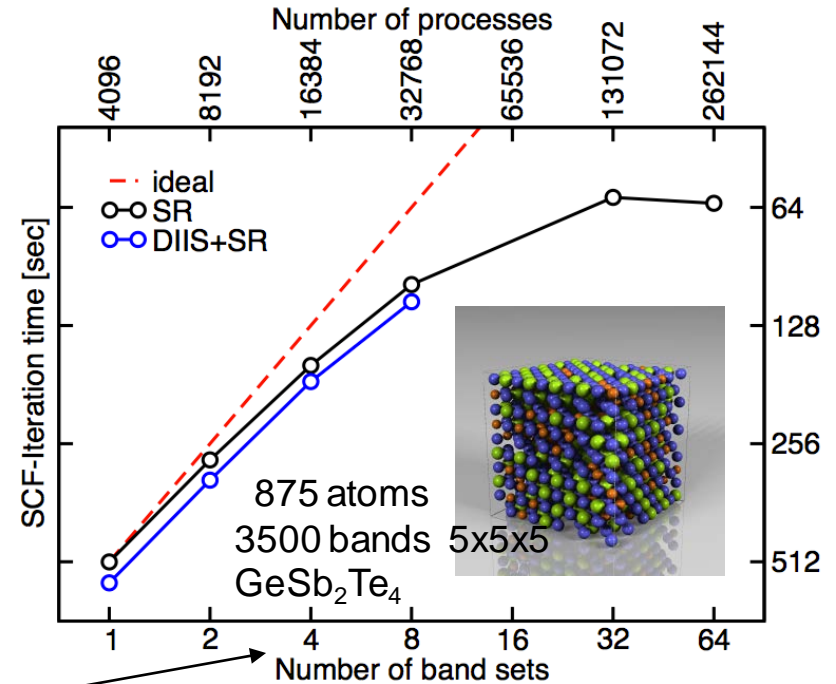


# FZJ-IAS & Activities on K and JUQUEEN





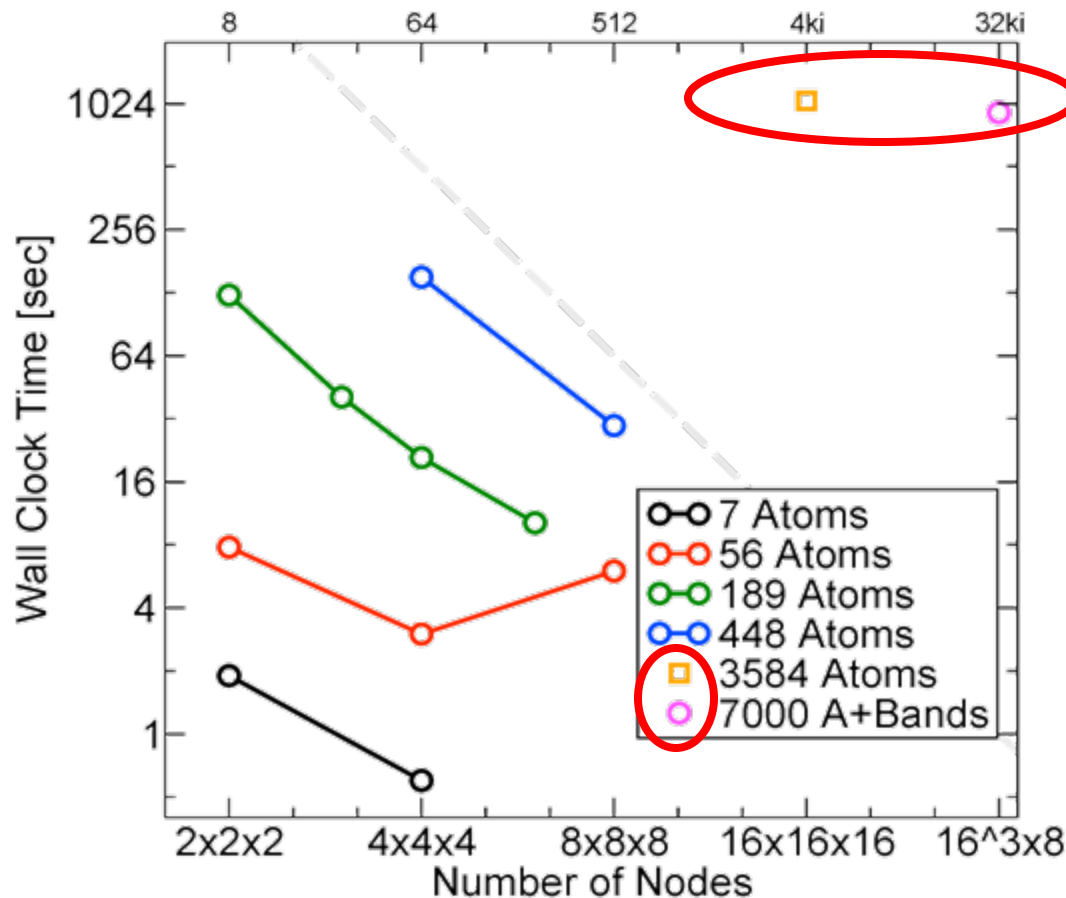
+ Band Parallelization



→ Atoms with PAW accuracy, random alloy  
10x10x10 GeSb<sub>2</sub>Te<sub>4</sub> **7000 atoms**  
(Projector Augmented Wave method)  
running on 131,072 cores

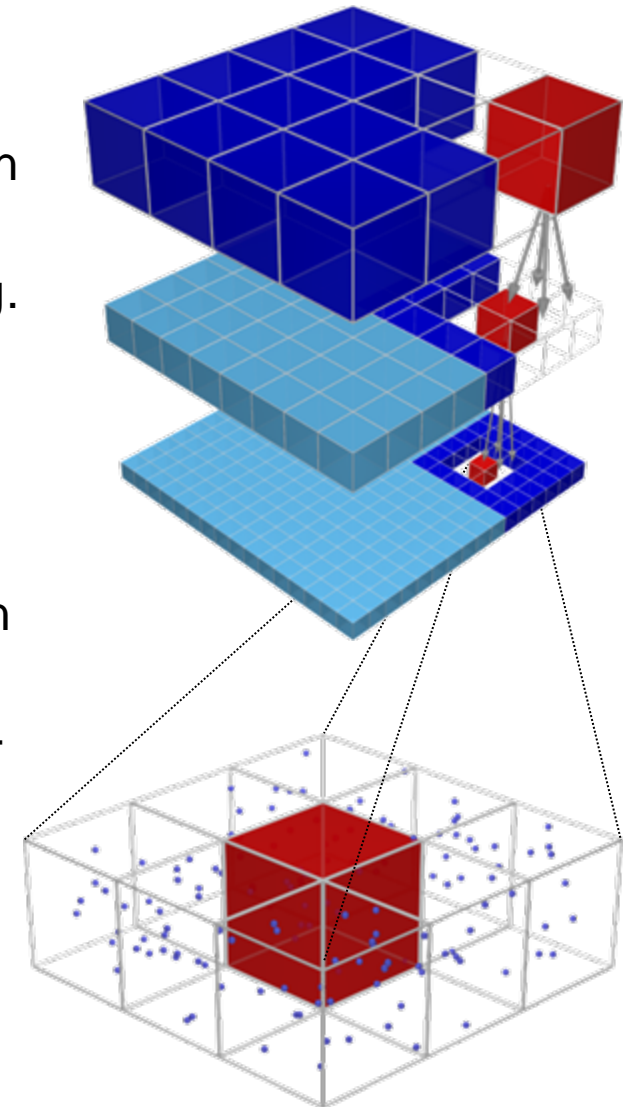
Contact: **Paul Baumeister, JSC, p.baumeister@fz-juelich.de**

- No problems porting the Fortran95 juRS-sources:  
`mpifrtpx -Ad -Kopenmp -Kfast -Koptmsg=2 -Kilfunc -Kfp_contract`
- linking with `-SSL2 -SCALAPACK -Kopenmp`



# FMM Toolbox on K: Overview

- Fast Multipole Method (**FMM**):
  - Computing Coulomb (Long-range)-interaction in  $O(N)$  instead of classical pairwise  $O(N^2)$
  - Used for molecular dynamics simulations e.g. Gromacs MD Code
- DFG-funded project *GROMEX*
  - extend Gromacs to allow FMM as backend Coulomb solver
  - extend FMM to allow for dynamic protonation within Gromacs
  - Redesign expensive compute kernels in C++
  - Develop performance portable FMM implementation for MD simulations
- First steps:
  - Test basic features of today's different HPC platforms (incl. K computer)



Contact: **Andreas Beckmann, Ivo Kabadshow, JSC, [i.kabadshow@fz-juelich.de](mailto:i.kabadshow@fz-juelich.de)**

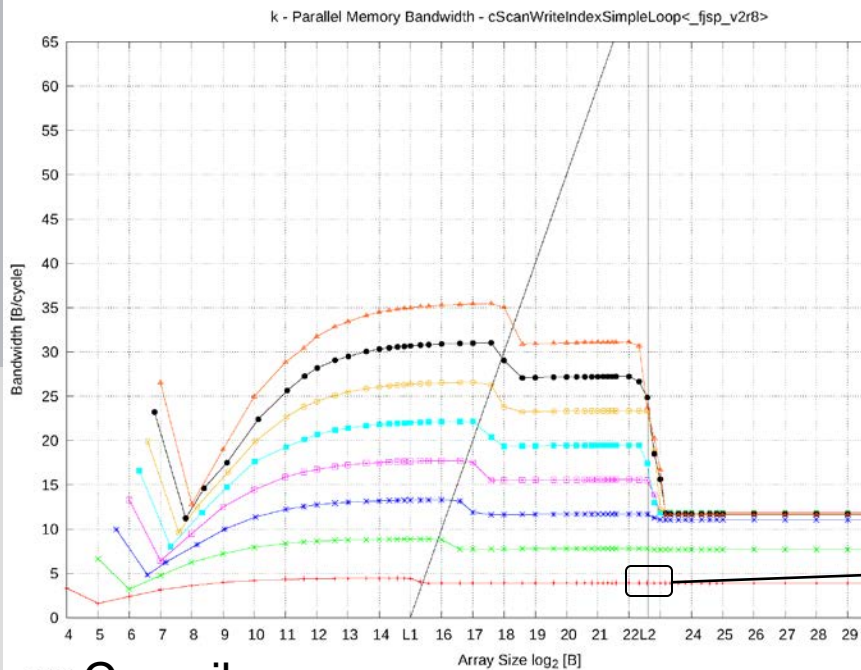
# FMM Toolbox on K: Memory Bandwidth

FMM is compute bound / memory bandwidth bound

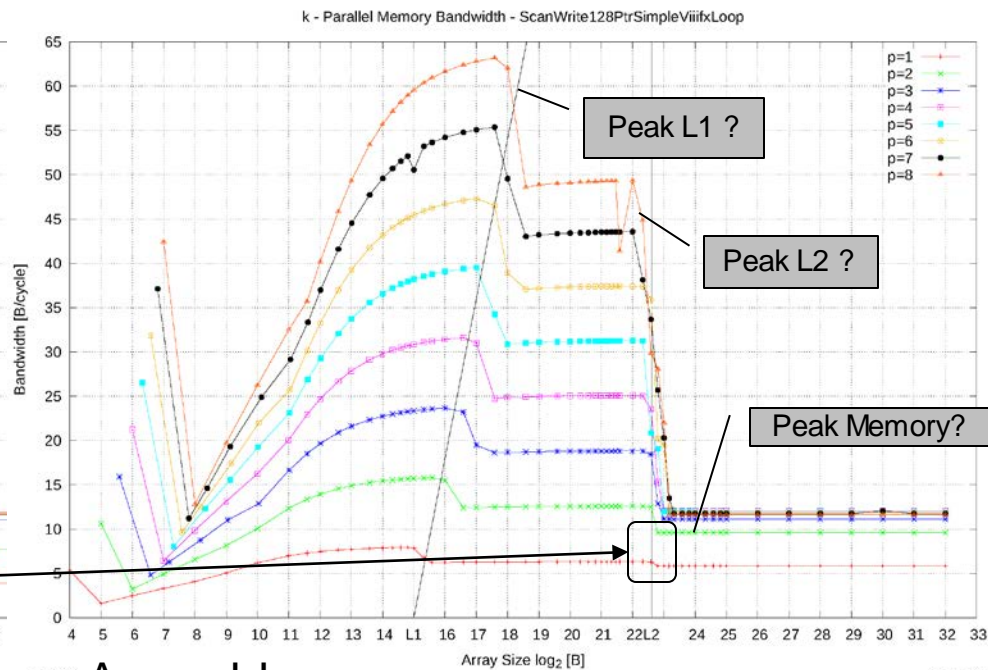
→ What is the maximum write bandwidth achievable on a single node?

Benchmark **pmbw** (parallel memory bandwidth)

- 128bit floating point stores (double x 2); linear access per thread
- separate memory per thread; pure write, no computation! single node tests only



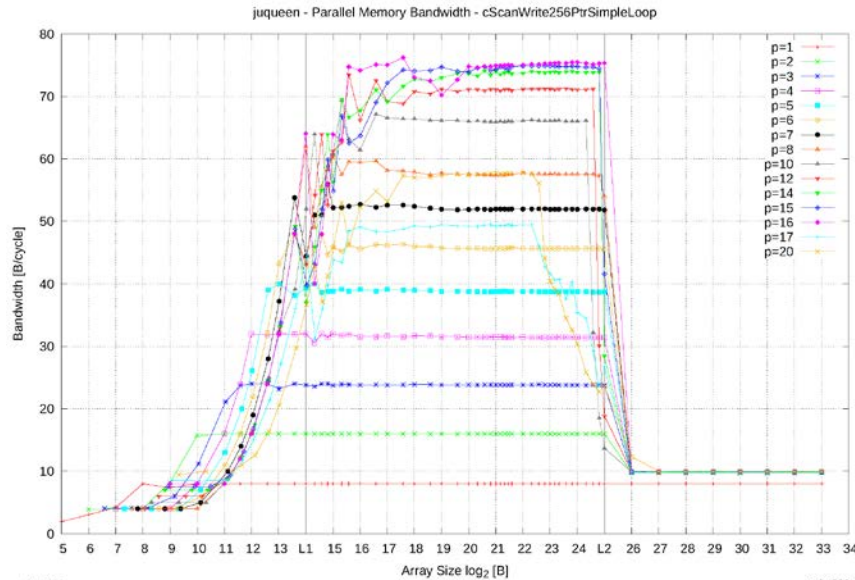
Compiler



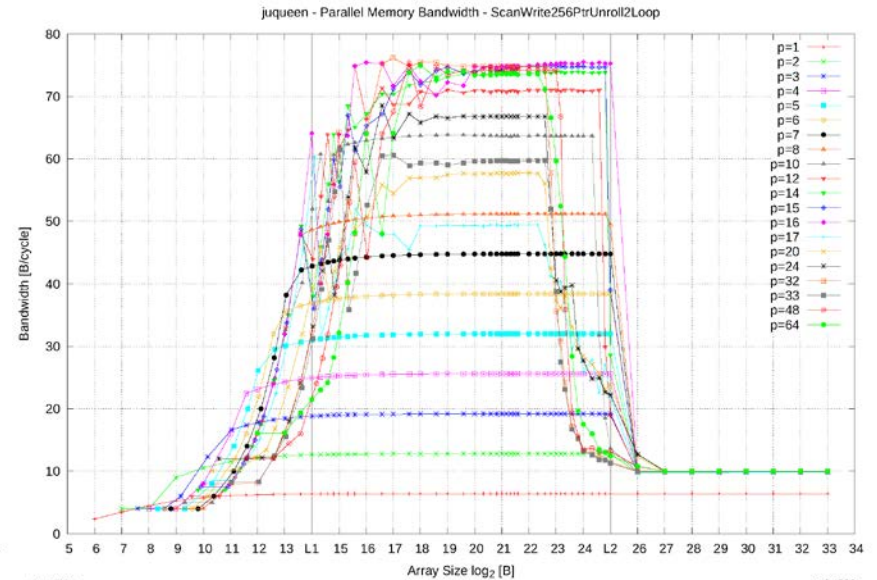
Assembler



# Benchmark pmbw on Juqueen



Compiler (bgclang, LLVM+QPX)



Assembler (Unroll2)

# Test of Elpa and EigenExa on JUQUEEN

- **ELPA** (EigensoLver for Petaflop Architectures)  
BMBF funded project (RZG Garching, Uni Wuppertal, ...)
  - Fortran 95, add-on to ScaLAPACK, OpenMP
  - ELPA 1: one step reduction and back transformation
  - ELPA 2: two step reduction to banded form, back transf. on BlueGene/Q  
→ ELPA-2 especially developed for part of the eigenspectrum
- **EigenExa Library**  
Toshiyuki Imamura, RIKEN Advanced Institute for Computational Science
  - Fortran 95, OpenMP for reduction and back transformation available
  - eigen s: one step reduction and back transformation
  - eigen sx: two step reduction to pentadiagonal form  
→ EigenExa developed to compute full eigenspectrum

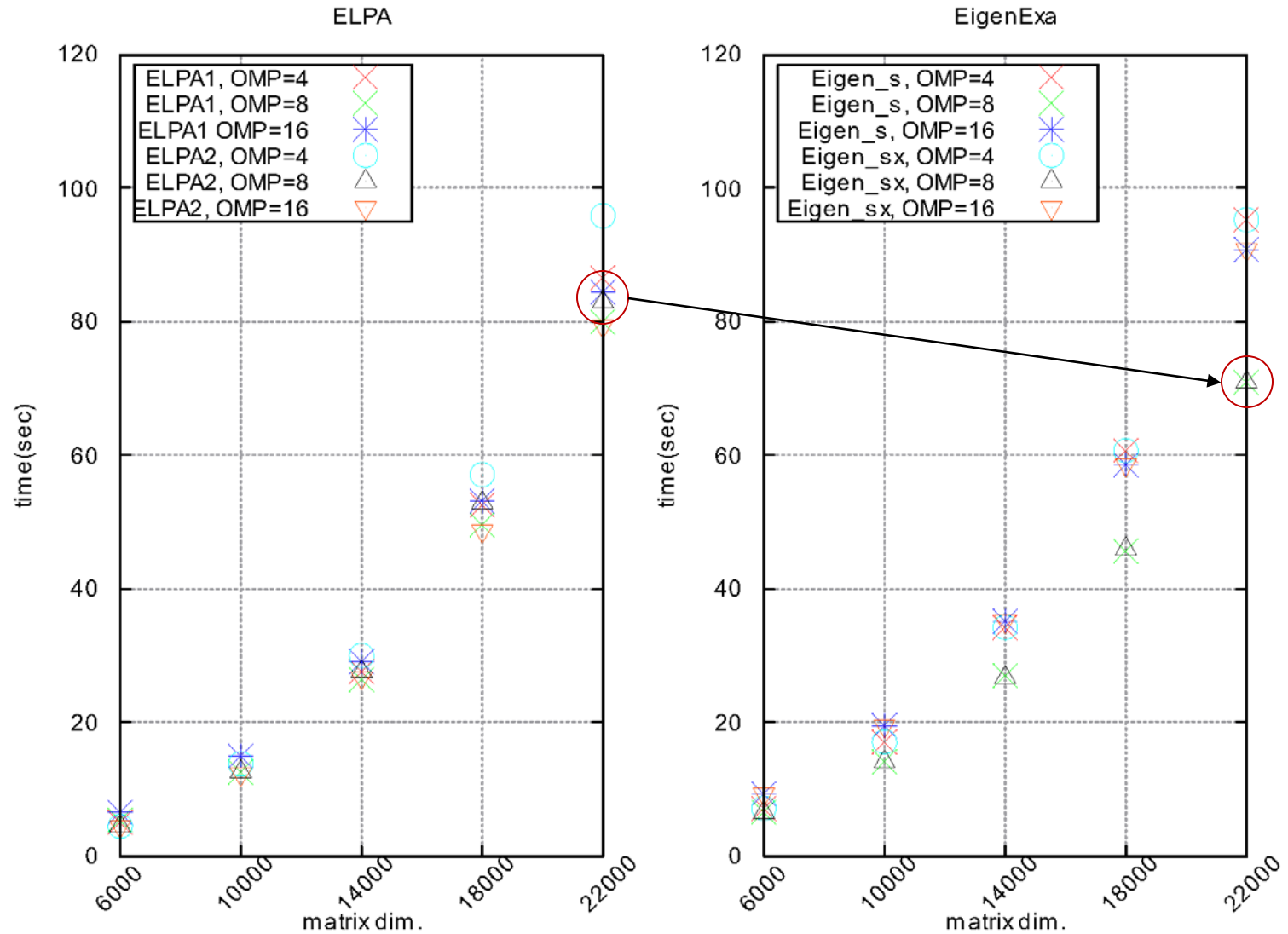
Contact: **I. Gutheil, JSC, [i.gutheil@fz-juelich.de](mailto:i.gutheil@fz-juelich.de)**



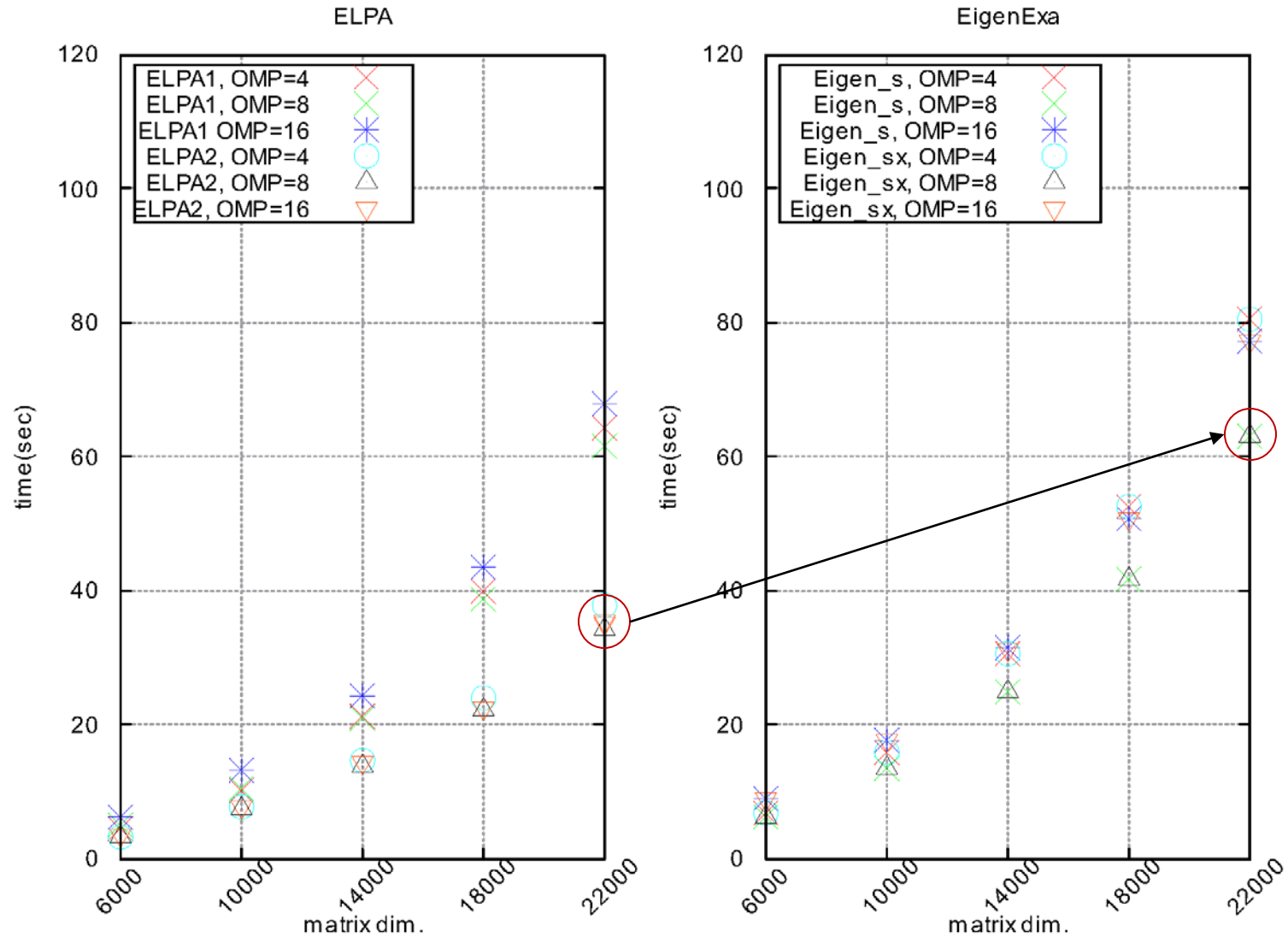
# Test Environment on JUQUEEN

- **Partition:** 1 node card, 512 cores, hybrid with 4 and 8 MPI processes per compute card
- **Matrix/Blocking:**  $N=6000$  to  $22000$  by steps of  $4000$ 
  - ELPA: block size  $nb = 16$  (ELPA 2 bandwidth)
  - EigenExa:  $nb = 48$  (reduction),  $nb = 128$  (back transformation), no accuracy improvement for eigenvalues
- **Test program ELPA:**
  - Generates diagonal matrix with random eigenvalues
  - Full matrix constructed via random Householder transformation
  - Deviation from given eigenvalues, residual, orthogonality tested
- **Test program EigenExa:**
  - Test program delivered with EigenExa, very similar to ELPA test program for pure MPI and hybrid version
  - Full eigenspectrum and 5 % of the eigenspectrum

# JUQUEEN Node card, 4 MPI processes per compute card, full eigenspectrum



# JUQUEEN Node card, 4 MPI processes per compute card, 5 % of the eigenspectrum



## Scalable performance analysis of large-scale parallel applications

- portable toolset for scalable performance measurement & analysis of MPI, OpenMP & hybrid OpenMP+MPI parallel applications
- *supporting most popular HPC computer systems*  
*Cray XT/XE/XK, IBM BlueGene L/P/Q, IBM SP & blade clusters, Fujitsu FX10 & K computer, NEC SX, SGI Altix, Linux cluster®, ...*
- Integrated instrumentation, measurement & analysis toolset
  - Customizable automatic/manual instrumentation
  - Runtime summarization (*aka* profiling)
  - Automatic event trace analysis
- Availability, Sources, documentation & publications:
  - 3-clause New BSD open-source license  
<http://www.scalasca.org>, scalasca@fz-juelich.de

Contact: **Brian Wylie, JSC**, [b.wylie@fz-juelich.de](mailto:b.wylie@fz-juelich.de)

- **Scalasca-1.4**

- ported to K computer & Fujitsu FX10 by Tomotake Nakamura (released March 2013)
- included PAPI hardware counters and 1/2/3D FJMPI topology
- 6D-Tofu network topology added March 2014

`/opt/aics/UNITE/packages/scalasca/1.4.4`

- **Scalasca-2.1**

- based on Score-P 1.3 (released August 2014)
- newly-developed instrumentation & measurement infrastructure
- now includes support for FX10 & K computer

- *no topologies, problem with OpenMP/C++*

`/opt/aics/UNITE/packages/scalasca/2.1` (1st release candidate)

- Used for **14th VI-HPS Tuning Workshop** (RIKEN AICS, 25-27 Mar 2014)

# Scalasca: Analysis of ABySS-P execution on K computer

Parallel genome sequence assembly using MPI (C++)

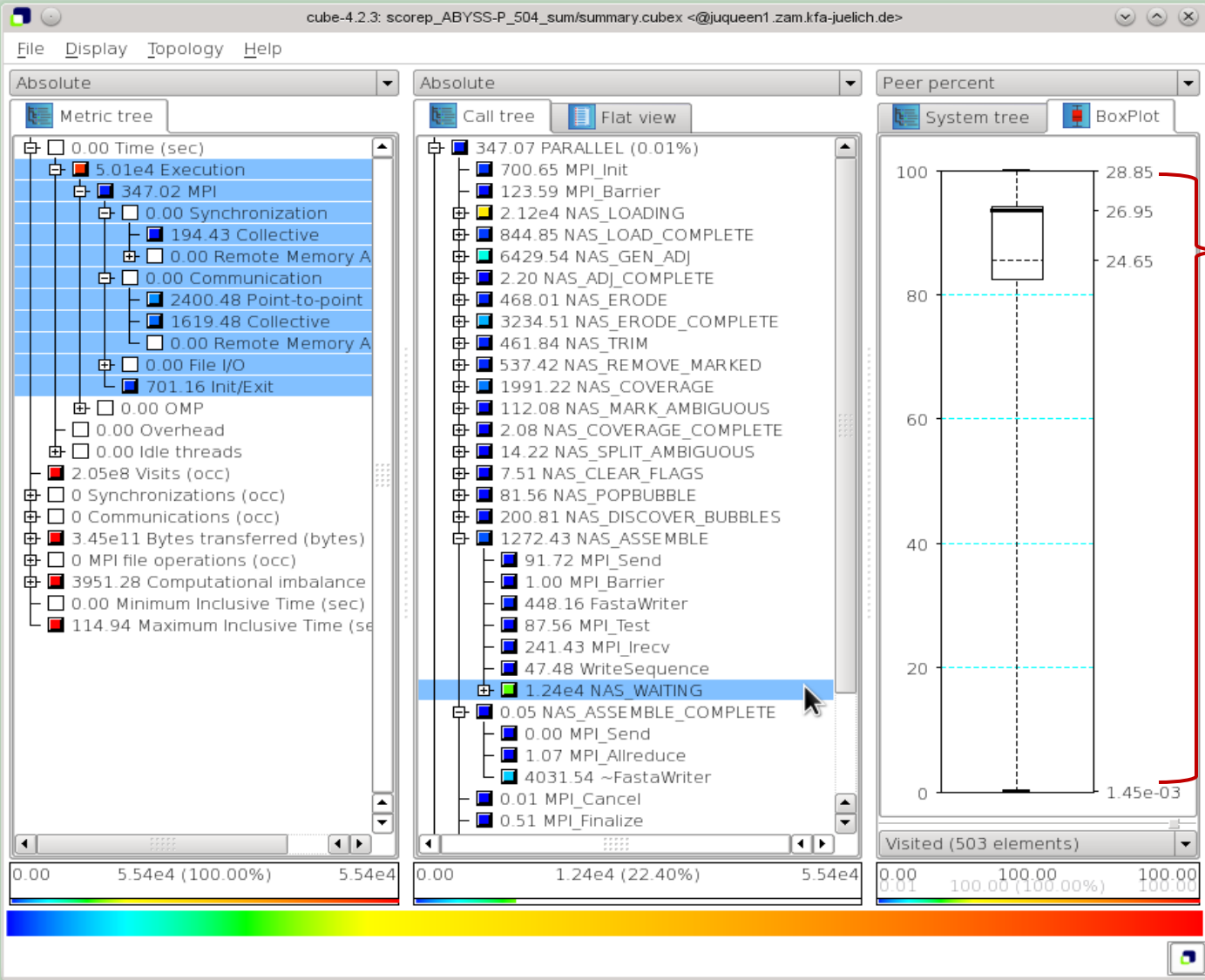
- Performance analysis on K computer by Itaru Kitayama (RIKEN)
  - *Limited scalability to around 768 compute nodes (1ppn)*

Status

- Manually annotated execution stages and file I/O operations
  - *provide context for MPI communication & synchronization*
- Summary experiments identify inefficient phases and causes
  - *NAS\_LOADING read & distribute input data*
  - *NAS\_POP/DISCOVER\_BUBBLES serialized file writing*
  - *NAS\_ASSEMBLE write sequence data load imbalance*
    - busy wait idling in NAS\_WAITING

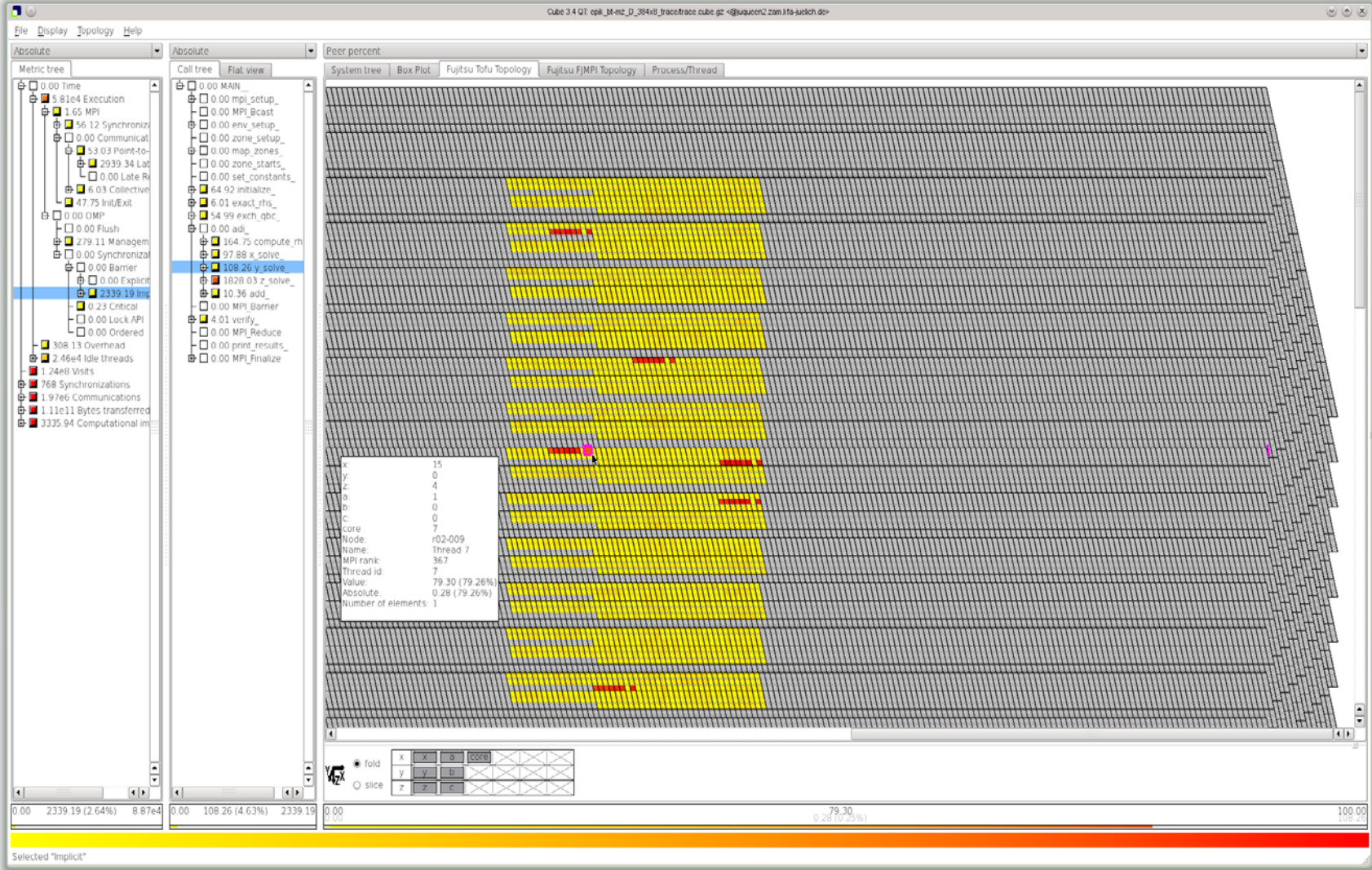


# Scalasca: Analysis of ABySS-P



load imbalance in NAS\_ASSEMBLE

# Scalasca: BT-MZ.D on K, 384x8 *y\_solve* OMP barrier imbalance



# Scalasca2 / Score-P development plans (FX10/K)

- Use of SIONlib for handling trace files
  - to avoid bottlenecks with multiple files per thread
- Record machine and application topologies
  - FJMPI, 6D-Tofu, etc. (as with Scalasca-1)
- Incorporate metrics from Tofu network counters?
- POSIX thread support?

# SIONlib: Overview

- Scalable Massively Parallel I/O to Task-Local Files
- Extension of I/O-API (ANSI C or POSIX)
- C, C++ and Fortran bindings, implementation language C
- Parallel Interfaces: MPI, OpenMP, Hybrid, Generic
- Current versions: 1.5.1
- Scalability on JUQUEEN: > 100 GB/s, 1.8 mio tasks
- Open source license: <http://www.fz-juelich.de/jsc/sionlib>

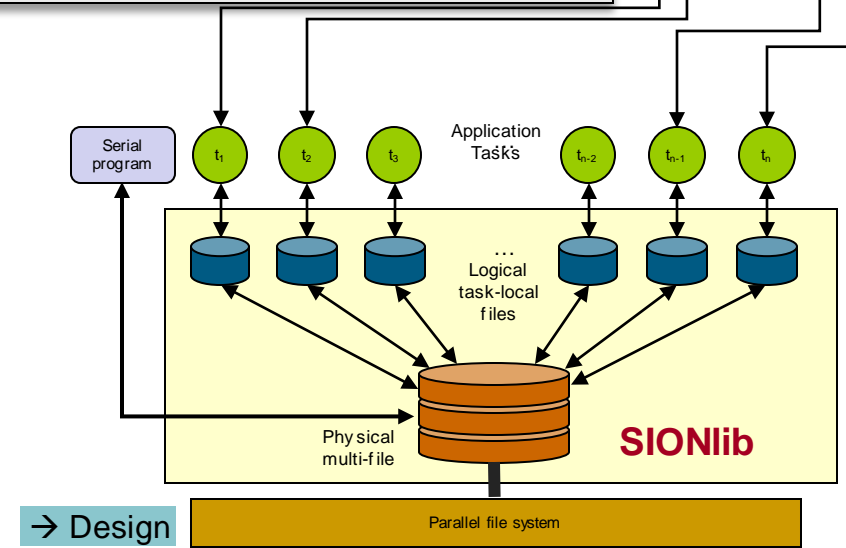
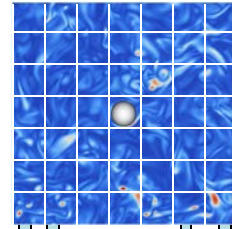
## → SIONlib-API

```

/* fopen() → */
sid=sion_paropen_mpi( filename , "bw",
                    &numfiles, &chunksizes,
                    gcom, &lcom, &fileptr, ...);

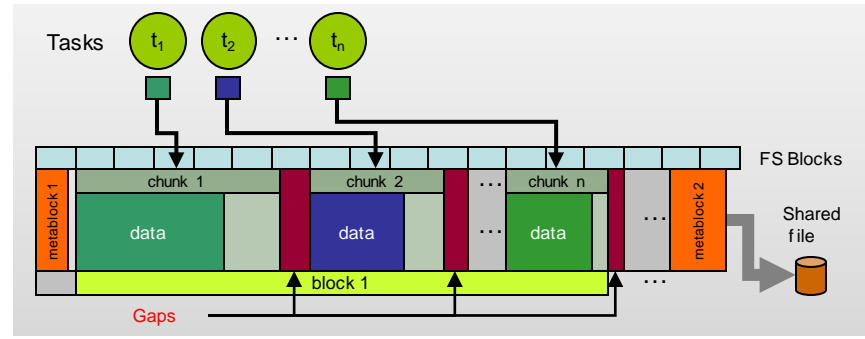
/* fwrite(bindata,1,nbytes, fileptr) → */
sion_fwrite(bindata,1,nbytes, sid);

/* fclose() → */
sion_parclose_mpi(sid)
    
```



## → Design

## → File-Format



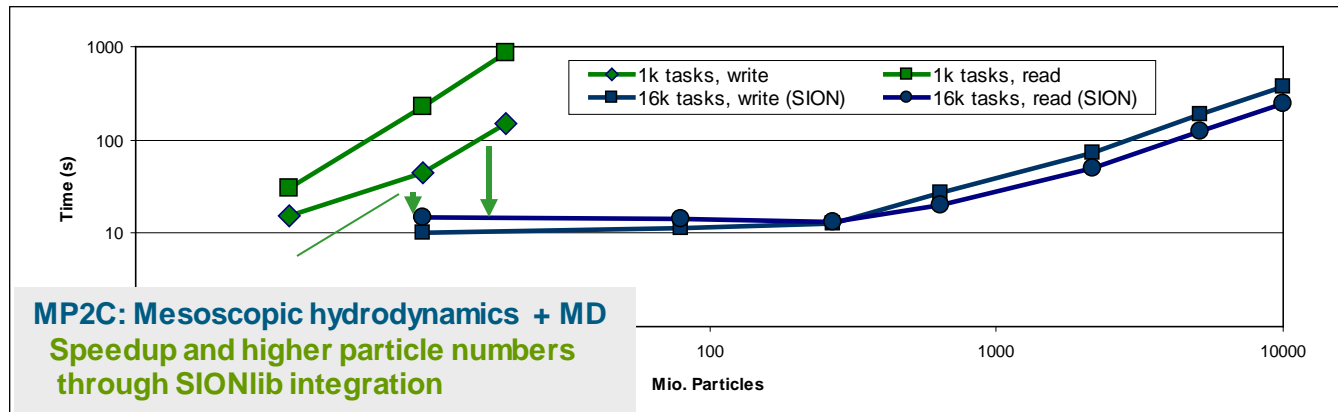
Contact: **Wolfgang Frings, JSC**  
[w.frings@fz-juelich.de](mailto:w.frings@fz-juelich.de)

# SIONlib: Applications & Use Cases

## Applications → Checkpointing, Restart Files

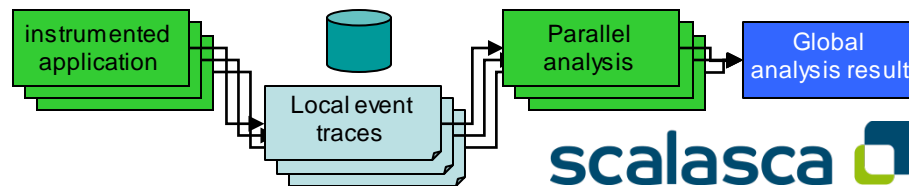
- DUNE-ISTL** (Multigrid solver, Univ. Heidelberg)
- LBM** (Fluid flow/mass transport, Univ. Marburg)
- OSIRIS** (Fully-explicit particle-in-cell code)
- Profasi**: (Protein folding and aggr. simulator)
- ITM** (Fusion-community)
- PSC** (particle-in-cell code)
- PEPC** (Pretty Efficient Parallel C. Solver)
- NEST** (Human Brain Simulation)

### MP2C:



## Tools/Projects

**Scalasca:** Performance Analysis

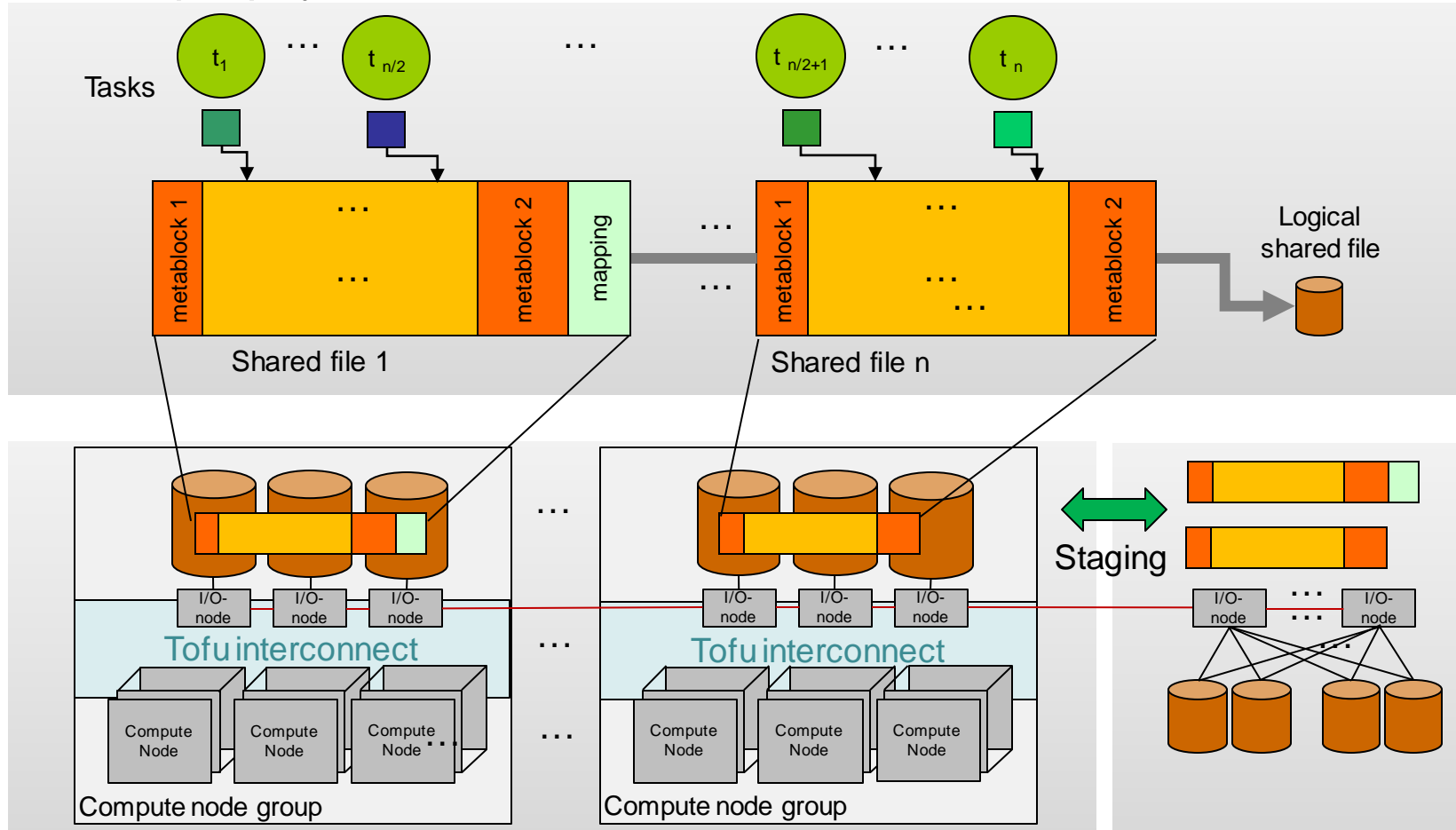


**Score-P:** Scalable Performance Measurement Infrastructure for Parallel Codes

**DEEP-ER:** Adaption to new platform and parallelization paradigm

# SIONlib: Strategy for K

- Multiple physical files for SIONlib shared file

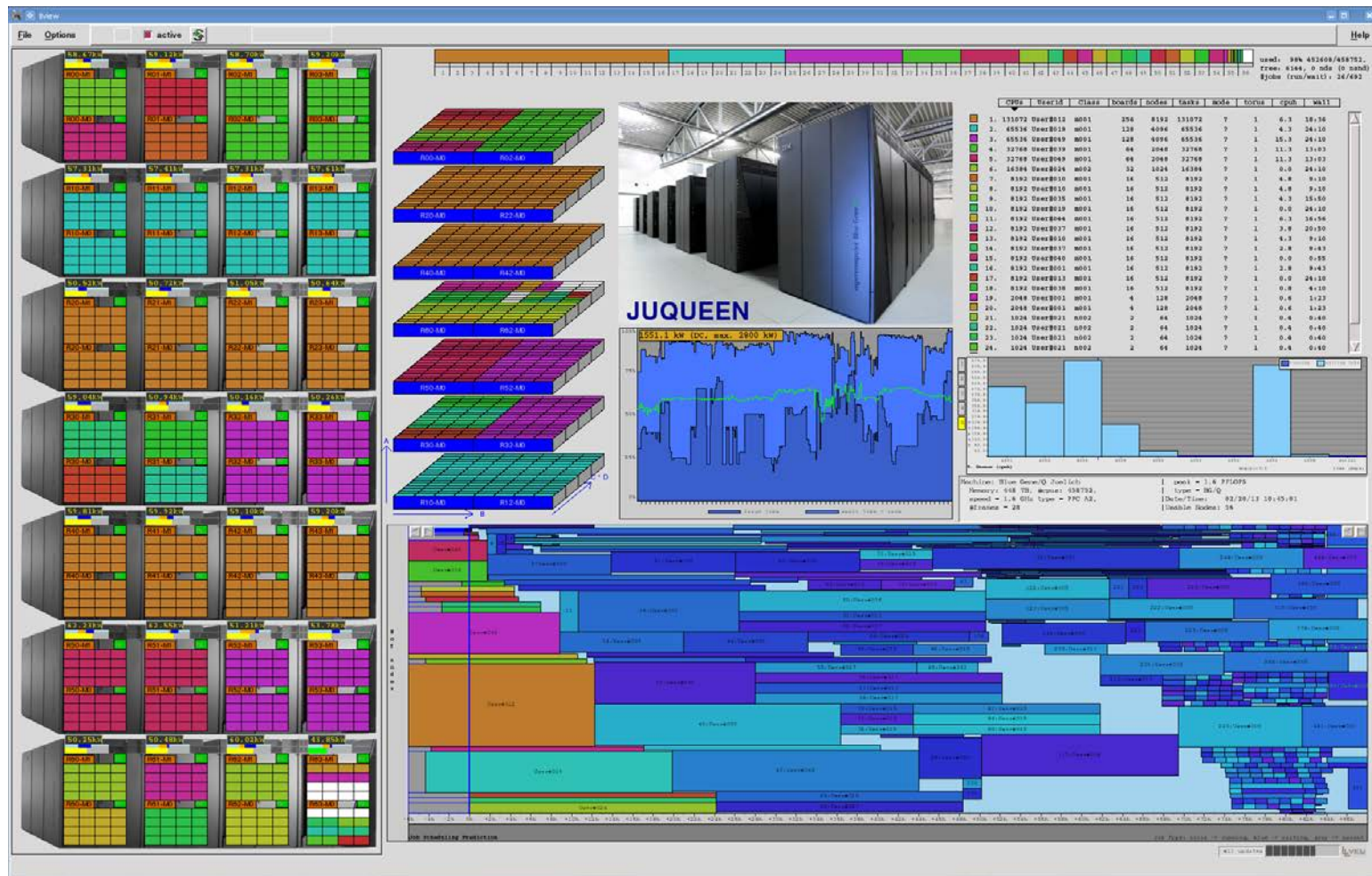


- Files on K Local Storage

K Global Storage



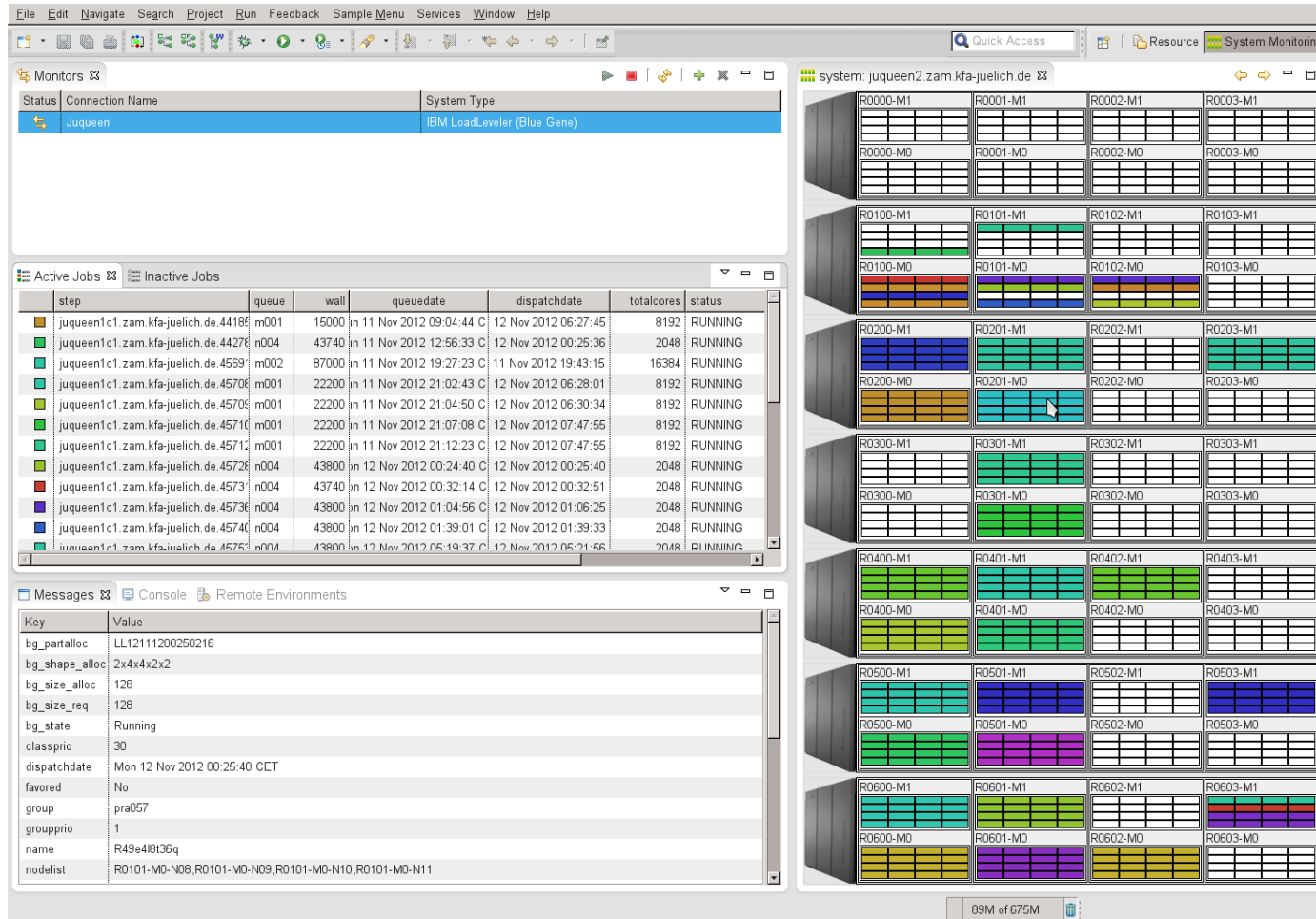
# Batch System Monitoring: LLview



Contact: [llview.jsc@fz-juelich.de](mailto:llview.jsc@fz-juelich.de)  
 WWW: <http://www.fz-juelich.de/jsc/llview>

# Batch System Monitoring: Eclipse/PTP

## Main components of LLview in PTP's monitoring perspective



The screenshot displays the Eclipse IDE with the PTP System Monitoring perspective. The main components are:

- Monitors:** A table showing the status of the batch system.
 

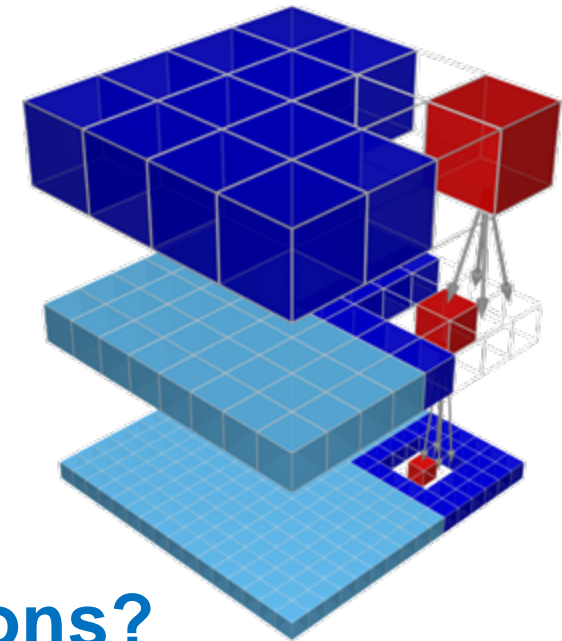
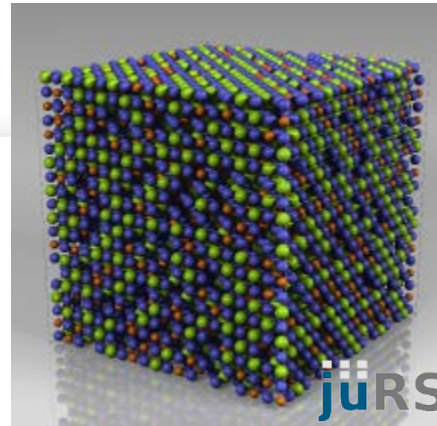
Status	Connection Name	System Type
Running	Juqueen	IBM LoadLeveler (Blue Gene)
- Active Jobs:** A table listing active jobs with columns for step, queue, wall, queuedate, dispatchdate, totalcores, and status.
 

step	queue	wall	queuedate	dispatchdate	totalcores	status
juqueen1c1.zam.kfa-juelich.de.4418	m001	15000	in 11 Nov 2012 09:04:44 C	12 Nov 2012 06:27:45	8192	RUNNING
juqueen1c1.zam.kfa-juelich.de.4427	n004	43740	in 11 Nov 2012 12:56:33 C	12 Nov 2012 00:25:36	2048	RUNNING
juqueen1c1.zam.kfa-juelich.de.4569	m002	87000	in 11 Nov 2012 19:27:23 C	11 Nov 2012 19:43:15	16384	RUNNING
juqueen1c1.zam.kfa-juelich.de.4570	m001	22200	in 11 Nov 2012 21:02:43 C	12 Nov 2012 06:28:01	8192	RUNNING
juqueen1c1.zam.kfa-juelich.de.4570	m001	22200	in 11 Nov 2012 21:04:50 C	12 Nov 2012 06:30:34	8192	RUNNING
juqueen1c1.zam.kfa-juelich.de.4571	m001	22200	in 11 Nov 2012 21:07:08 C	12 Nov 2012 07:47:55	8192	RUNNING
juqueen1c1.zam.kfa-juelich.de.4571	m001	22200	in 11 Nov 2012 21:12:23 C	12 Nov 2012 07:47:55	8192	RUNNING
juqueen1c1.zam.kfa-juelich.de.4572	n004	43800	in 12 Nov 2012 00:24:40 C	12 Nov 2012 00:25:40	2048	RUNNING
juqueen1c1.zam.kfa-juelich.de.4573	n004	43740	in 12 Nov 2012 00:32:14 C	12 Nov 2012 00:32:51	2048	RUNNING
juqueen1c1.zam.kfa-juelich.de.4573	n004	43800	in 12 Nov 2012 01:04:56 C	12 Nov 2012 01:06:25	2048	RUNNING
juqueen1c1.zam.kfa-juelich.de.4574	n004	43800	in 12 Nov 2012 01:39:01 C	12 Nov 2012 01:39:33	2048	RUNNING
juqueen1c1.zam.kfa-juelich.de.4575	n004	43800	in 12 Nov 2012 05:19:37 C	12 Nov 2012 05:21:55	2048	RUNNING
- Messages:** A console window showing system messages with keys and values.
 

Key	Value
bg_pantalloc	LL12111200250216
bg_shape_alloc	2x4x4x2x2
bg_size_alloc	128
bg_size_req	128
bg_state	Running
classprio	30
dispatchdate	Mon 12 Nov 2012 00:25:40 CET
favored	No
group	pra057
groupprio	1
name	R49e4f8136q
odelist	R0101-M0-N08,R0101-M0-N09,R0101-M0-N10,R0101-M0-N11
- Resource Monitors:** A grid of 24 resource monitors (R0000-M1 to R0603-M1) showing the status of individual nodes. Each monitor is a small grid with colored bars representing resource usage.

WWW: <https://www.eclipse.org/ptp>

# Thank You!



## Questions?

