

JURECA: General-purpose supercomputer at Jülich Supercomputing Centre

Forschungszentrum Jülich, Jülich Supercomputing Centre *

Instrument Scientists:

- Dorian Krause, Jülich Supercomputing Centre, Forschungszentrum Jülich, phone: +49(0)2461 61 3631, email: d.krause@fz-juelich.de
- Philipp Thörnig, Jülich Supercomputing Centre, Forschungszentrum Jülich, phone: +49(0)2461 61 1472, email: p.thoernig@fz-juelich.de

Abstract: JURECA is a petaflop-scale, general-purpose supercomputer operated by Jülich Supercomputing Centre at Forschungszentrum Jülich. Utilizing a flexible cluster architecture based on T-Platforms V-Class blades and a balanced selection of best of its kind components the system supports a wide variety of high-performance computing and data analytics workloads and offers a low entrance barrier for new users.

1 Introduction

Since July 2015, the Jülich Supercomputing Centre (JSC) at the Forschungszentrum Jülich (Forschungszentrum Jülich, 2015a) operates the JURECA (Jülich Research on Exascale Cluster Architectures) cluster system as the successor of the popular JUROPA (Jülich Research on Petaflop Architectures) supercomputer. JURECA (see Figure 1) serves as a general-purpose supercomputing resource and, in accordance with Forschungszentrum Jülich's dual architecture strategy, augments the leadership-class highly scalable IBM Blue Gene/Q system JUQUEEN (Jülich Supercomputing Centre, 2015). Funding for JURECA was granted by the Helmholtz Association (Helmholtz Association, 2015) through the program "Supercomputing & Big Data". JURECA was designed by JSC together with the hardware vendor T-Platforms (T-Platforms, 2015), who won a competitive procurement initiated in 2014, to serve as a versatile scientific instrument for compute- and data-intense (simulation) science that is equally suited for capacity as for capability workloads. In the same way that JURECA's predecessor JUROPA served as a research vehicle that laid the foundation for the JURECA design, the project partners T-Platforms, the software provider ParTec and JSC are working together to incrementally improve the JURECA system during

*Cite article as: Jülich Supercomputing Centre. (2016). JURECA: General-purpose supercomputer at Jülich Supercomputing Centre. *Journal of large-scale research facilities*, 2, A62. <http://dx.doi.org/10.17815/jlsrf-2-121>

operation and to address some of the main challenges in the design and operation of clusters of the next generation.

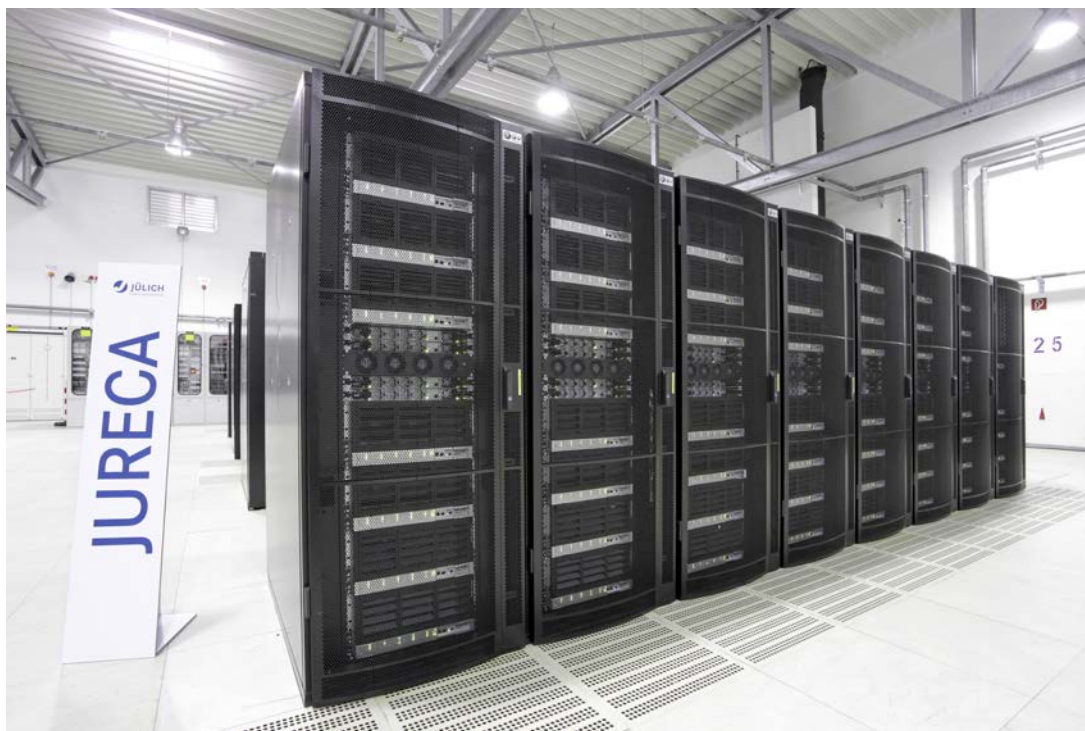


Figure 1: Jülich Research on Exascale Cluster Architectures (JURECA) at Jülich Supercomputing Centre. Copyright: Forschungszentrum Jülich.

2 JURECA system details

The JURECA architecture is an evolution of the JUROPA architecture and follows the best-of-breed approach by combining the most advanced commodity hardware and software technologies available in the industry. JURECA is a heterogeneous system offering nodes with different memory sizes (128 GiB, 256 GiB, 512 GiB as well as 1 TiB), nodes with graphics processing unit (GPU) accelerators as well as GPU-equipped nodes for visualization and other post-processing needs.

JURECA consists of 1,733 compute nodes of type T-Platforms V-Class V210S as well 75 GPU-accelerated V210F blades hosted in V5050 chassis (see Figure 2). Moreover, 64 Supermicro F618R2-RT+ twin-blade servers (512 GiB memory nodes) and 12 Supermicro 1028GR-TR visualization nodes are available.

All systems feature two Intel Xeon E5-2680 v3 12-core Haswell central processing units (Intel Corporation, 2015) (CPUs) which support up to 24 hardware threads each. Each CPU supports the AVX 2.0 instruction set architecture extension and can perform two 256-bit (i.e., four double precision floating point numbers) wide multiply-add operations per cycle. The peak performance of a (non-accelerated) JURECA node is 0.96 TFlop/s. The maximum memory bandwidth of the node is 136 GB/s. The two sockets are connected by a bi-directional 9.6 GT/s (Gigatransfers per second) Intel Quick Path Interconnect (QPI) link. In JURECA, 2133 MHz DDR4 memory technology is used.

Applications that support the use of GPU accelerators can take advantage of the additional two NVIDIA K80 graphics processing units available in 75 JURECA compute nodes. The GPUs are connected with PCI Express Generation 3.0 (16 lanes) links providing a peak of 32 GB/s bidirectional host-device bandwidth. Each K80 GPU is equipped with 2×12 GB GDDR5 memory and offers 4992 CUDA cores that

provide an additional 2.9 TFlop/s peak performance (5.8 TFlop/s per node) and 480 GB/s memory bandwidth per GPU. The 12 visualization nodes are equipped with two NVIDIA K40 GPUs intended for remote visualization usage.



(a) Back (left) and side (right) view of a T-Platforms V-Class V210S dual-socket blade server as used in JURECA. The GPU-accelerated V210F blades host two additional PCIe devices and fit in two chassis slots.



(b) Front (left) and back (right) view of the T-Platforms V5050 chassis. Each chassis can host ten V210S or, alternatively, five V210F blades.

Figure 2: T-Platforms V-Class components used in the JURECA system. Copyright: T-Platforms.

The JURECA compute nodes are connected with the industry-leading Mellanox extended data rate (EDR) InfiniBand providing 100 Gb/s (12.5 GB/s) link bandwidth and MPI latencies around one microsecond. The host channel adapters (HCA) are connected via PCI Express Generation 3.0 (16 lanes). The JURECA components are interconnected in a three-level full fat tree topology which provides full bisection bandwidth and non-blocking communication for appropriate communication patterns.

A particular emphasis during the design of JURECA has been put on the storage connection in order to meet the increasing data requirements of simulation sciences as well as the needs of emerging data-intensive sciences. All offered global (parallel) filesystems on JURECA are mounted from the central Jülich Storage Cluster (JUST) (Forschungszentrum Jülich, 2015c) using IBM's General Parallel Filesystem (GPFS). Users with access to several systems in the supercomputing facility at JSC work with the same filesystems on all systems so that data movement is minimized and workflows are simplified. The storage network connection is realized using InfiniBand-to-Ethernet gateways bridging JURECA's internal InfiniBand network with the facility's Terabit Ethernet backbone. This connection type was selected as it allows for >100 GB/s aggregate filesystem bandwidth as well as a high per-node filesystem performance that is hardware-wise only limited by the performance of the fourteen data rate (FDR) InfiniBand links (56 Gb/s) towards the gateways.

JURECA's software stack is largely based on open-source software. Login and compute nodes run the CentOS 7 Linux operating system with a careful setup that balances the ease of use and low entrance-barrier with the requirements, such as minimal operating system jitter, of large-scale capability clusters. JURECA uses the open-source Slurm workload manager (SchedMD LLC, 2015) in combination with the ParaStation resource management which has a proven track record in scalability, reliability and performance on several clusters operated by JSC. The ParTec Parastation ClusterSuite (ParTec Cluster Competence Center GmbH, 2015) is used for system provisioning and health monitoring.

On JURECA, the Intel and ParTec ParaStation Message Passing Interface (MPI) implementations are supported. In addition the CUDA-aware MPI implementation MVAPICH2-GDR is available for mixed MPI+CUDA applications. Different compilers, optimized mathematical libraries and pre-compiled community codes are available. We refer to the JURECA webpage (Forschungszentrum Jülich, 2015b) for more information. Monitoring of batch jobs is possible using the latest version of the LLview (Forschungszentrum Jülich, 2015d) graphical monitoring tool.

Scientists can also use UNICORE (UNICORE Forum e.V., 2015) to create, submit and monitor jobs on the JURECA system.

2.1 Hardware components

As of this writing, JURECA consists of the following hardware components. An up-to-date description of the hardware (and software) configuration of the system is maintained on the JURECA webpage (Forschungszentrum Jülich, 2015b).

- 34 racks organized in four rows
 - 1,872 compute nodes
 - * 2× Intel Xeon E5-2680 v3 Haswell CPUs per node
 - 2× 24 cores, 2.5 GHz frequency
 - Intel Hyper-Threading Technology
 - AVX 2.0 instruction set architecture extension
 - * DDR4 memory technology clocked at 2133 MHz
 - 128 GiB memory in 1,680 nodes
 - 256 GiB memory in 128 nodes
 - 512 GiB memory in 64 nodes
 - * 75 nodes equipped with 2× NVIDIA K80 GPUs each
 - 2× 4992 CUDA cores
 - 2× 24 GB GDDR5 memory
 - 12 visualization nodes
 - * 2× Intel Xeon E5-2680 v3 Haswell CPUs per node
 - * DDR4 memory technology clocked at 2133 MHz
 - Memory size of 512 GiB in 10 nodes
 - Memory size of 1 TiB in 2 nodes
 - * 2× NVIDIA K40 GPUs
 - 2× 12 GB GDDR5 memory
 - 12 login nodes
 - * 2× Intel Xeon E5-2680 v3 Haswell CPUs per node
 - * 256 GiB DDR4 memory
 - 14 service nodes for system management
- Mellanox InfiniBand EDR network organized in a three-level full-fat tree topology
 - Mellanox ConnectX-4 single port host channel adapters in nodes
 - 36-port SwitchIB-based Mellanox SB7790 leaf-level switches
 - 4× SwitchIB-based Mellanox CS7500 switches with 2× 468 and 2× 540 ports
 - 2× Mellanox SX6036G InfiniBand FDR/40 Gigabit Ethernet gateways for storage connection

2.2 Software components

- CentOS 7 enterprise-grade Linux operating system
- ParTec ParaStation ClusterSuite
- Slurm batch system with ParaStation resource management
- Intel and ParTec ParaStation Message Passing Interface implementations
- Support for OpenMP, NVIDIA CUDA, OpenCL and OpenACC programming models

2.3 Benchmark results

Using 1,764 JURECA compute nodes without accelerators a Linpack performance of 1.42 PFlop/s was measured, placing the system on spot 50 in the November 2015 Top500 list (Top500, 2015). The system consumed on average 825 kW during the Linpack run, i.e., about 1.72 GFlop/s/W. JURECA entered the Green500 list in November 2015 on place 112 (Green500, 2015). On the High Performance Conjugate Gradients (HPCG) benchmark, JURECA achieved 68.3 TFlop/s corresponding to place 18 in the November 2015 HPCG list (HPCG, 2015).

3 Access to JURECA

Scientists and engineers interested in using the capacities and capabilities of JURECA for their research have to apply for JURECA compute time resources by submitting an adequate proposal in answer to corresponding compute time calls published end of January and end of July every year. Submitted proposals are evaluated scientifically through a competitive peer-review process. Additionally, the review process includes a technical assessment of the applicant's ability to efficiently perform parallel computations utilizing a larger number of compute cores on JURECA.

Basically, there are two calls available twice every year: One is conducted jointly by peers in computational science and engineering at Forschungszentrum Jülich and RWTH Aachen University, accepting proposals from the two institutions only (so-called JARA-HPC/VSR Call) (Jülich-Aachen Research Alliance, 2015). The other one (NIC Call) is performed by the John von Neumann Institute for Computing (John von Neumann Institute for Computing, 2015) (NIC), a joint organization of the three Helmholtz centers Forschungszentrum Jülich, Deutsches Elektronen-Synchrotron (Deutsches Elektronen-Synchrotron, 2015) DESY and the GSI Helmholtzzentrum für Schwerionenforschung (GSI Helmholtzzentrum für Schwerionenforschung, 2015), accepting proposals from all other German universities and research institutions. Applicants have to demonstrate that they are qualified in their respective field and that they have an appropriate knowledge in high-performance computing.

Scientists with challenging compute- or data-intense scientific problems that require access to JURECA in order to lay the necessary software foundation for the preparation of a successful proposal can obtain a limited compute time budget on JURECA along with expert support by a JSC simulation lab (Forschungszentrum Jülich, 2015f) by answering the bi-annual call for preparatory access and support resources (Forschungszentrum Jülich, 2015e).

4 Application fields

In order to exemplify the versatility of JURECA, Figure 3 shows the allocated computing time and number of projects working on JURECA in the time frame November 2015 to April 2016. In this time frame, JURECA will serve about 160 million core hours to nearly 200 projects working on advanced research in different disciplines ranging from the basic sciences to engineering applications. The employed parallel applications in different projects and scientific communities range from highly-optimized applications, which are particularly tuned for JURECA's CPU architecture, to complicated multi-program simulation



frameworks that are not uncommonly driven by dynamic scripting languages. The demands for hardware features (such as modern accelerators), main memory sizes, pre- and post-processing capabilities as well as input/output (I/O) performance vary similarly. JURECA's hardware design and software stack ensure that these requirements are met across the whole application spectrum.

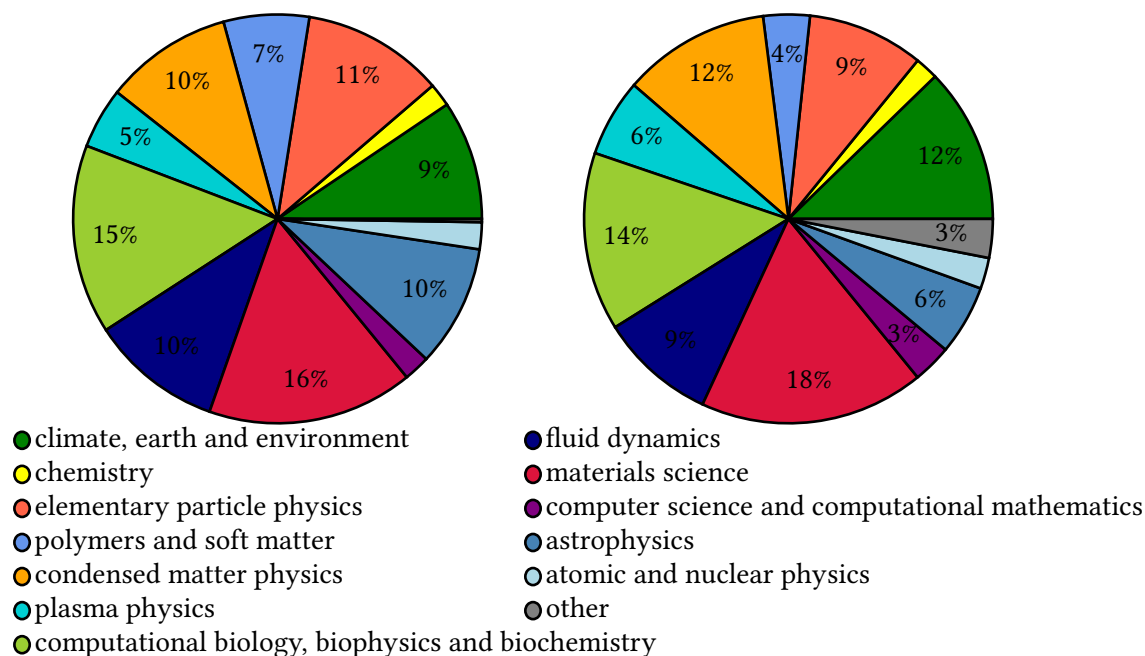


Figure 3: Allocated compute time (left) and number of projects (right) on JURECA by scientific field in the computing time period from the 1st of November 2015 to the 30th April 2016. Percentages are shown for shares above 3 %.

References

- Deutsches Elektronen Synchrotron. (2015). *Deutsches Elektronen-Synchrotron (DESY)*. Retrieved from <http://www.desy.de>
- Forschungszentrum Jülich. (2015a). *Forschungszentrum Jülich GmbH*. Retrieved from <http://www.fz-juelich.de>
- Forschungszentrum Jülich. (2015b). *JURECA*. Retrieved from <http://www.fz-juelich.de/ias/jsc/jureca>
- Forschungszentrum Jülich. (2015c). *JUST*. Retrieved from <http://www.fz-juelich.de/ias/jsc/just>
- Forschungszentrum Jülich. (2015d). *LLview*. Retrieved from <http://www.fz-juelich.de/jsc/llview>
- Forschungszentrum Jülich. (2015e). *Preparatory Access to Computing and Support Resources*. Retrieved from <http://www.fz-juelich.de/ias/jsc/prep-access.html>
- Forschungszentrum Jülich. (2015f). *Simulation Laboratories at Jülich Supercomputing Centre*. Retrieved from <http://www.fz-juelich.de/ias/jsc/simlabs>
- Green500. (2015). *Green500 November 2015 list*. Retrieved from <http://www.green500.org/lists/green201511>

GSI Helmholtzzentrum für Schwerionenforschung. (2015). *GSI Helmholtzzentrum für Schwerionenforschung*. Retrieved from <http://www.gsi.de>

Helmholtz Association. (2015). *Helmholtz-Gemeinschaft Deutscher Forschungszentren e.V. (HGF)*. Retrieved from <http://www.helmholtz.de>

HPCG. (2015). *HPCG November 2015 list*. Retrieved from <http://www.hpcg-benchmark.org>

Intel Corporation. (2015). *Intel Xeon Processor E5-2680 v3*. Retrieved from http://ark.intel.com/products/81908/Intel-Xeon-Processor-E5-2680-v3-30M-Cache-2_50-GHz

John von Neumann Institute for Computing. (2015). *John von Neumann Institute for Computing (NIC)*. Retrieved from <http://www.john-von-neumann-institut.de>

Jülich-Aachen Research Alliance. (2015). *Jülich-Aachen Research Alliance – High-Performance Computing (JARA-HPC)*. Retrieved from <http://www.jara.org/de/research/jara-hpc>

Jülich Supercomputing Centre. (2015). JUQUEEN: IBM Blue Gene/Q Supercomputer System at the Jülich Supercomputing Centre. *Journal of large-scale research facilities*, 1, A1. <http://dx.doi.org/10.17815/jlsrf-1-18>

ParTec Cluster Competence Center GmbH. (2015). *ParTec*. Retrieved from <http://www.par-tec.com>

SchedMD LLC. (2015). *Slurm Workload Manager webpage*. Retrieved from <http://slurm.schedmd.com>

Top500. (2015). *Top500 November 2015 list*. Retrieved from <http://www.top500.org/lists/2015/11>

T-Platforms. (2015). *T-Platforms*. Retrieved from <http://www.t-platforms.com>

UNICORE Forum e.V. (2015). *Uniform Interface to Computing Resources (UNICORE) webpage*. Retrieved from <http://www.unicore.eu>