

National High Performance Computing for Science

HASCO Summer School 2024

Christian Köhler
christian.koehler@gwdg.de

August 5, 2024

hpc@gwdg.de
GWDG – Gesellschaft für wissenschaftliche Datenverarbeitung mbH Göttingen

The NHR alliance

HPC at GWDG

Excursion: Using Graphics Cards for Scientific Computing

Application Software

Getting access

Further resources

The NHR alliance

The NHR alliance

Why HPC?

Comparison: CPU/GPU FP Operations

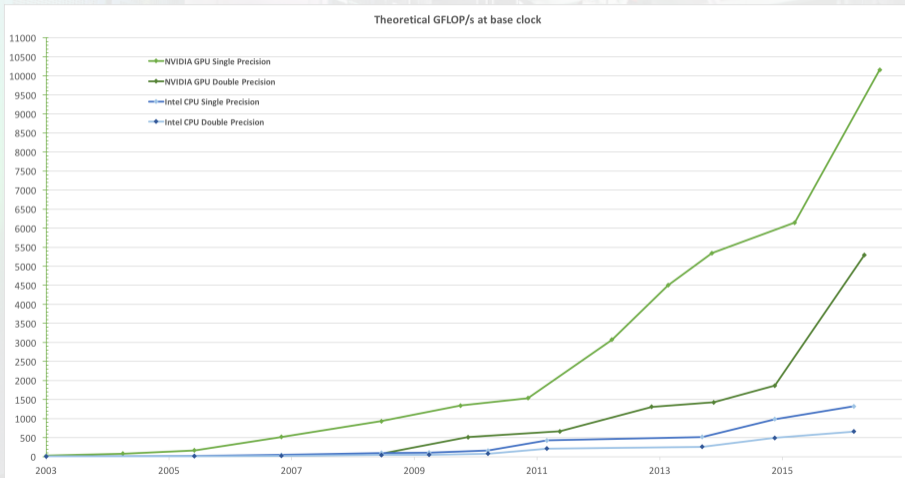


Image source: nVidia - CUDA C Programming Guide

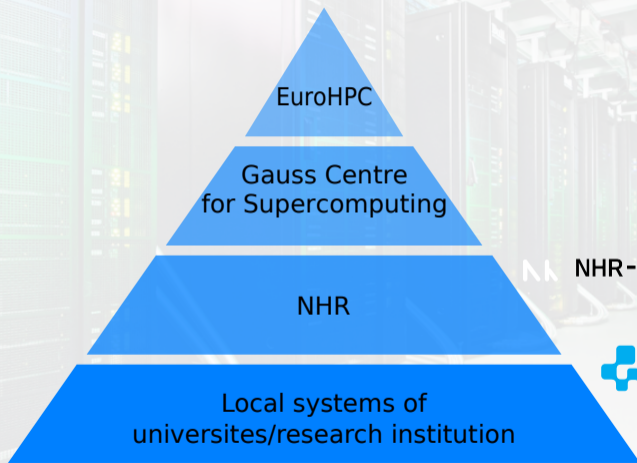
The NHR alliance

Göttingen in the HPC Landscape

HPC landscape in Germany

T0: Europe, T1/2: Germany, T3: local

NHR-NORD@GÖTTINGEN



NHR-NORD@GÖTTINGEN
"Emmy", "Grete"

GWGD
Scientific
Compute
Cluster

- HLRN (North German Supercomputing Alliance) exists since 2001
- originally joint effort of Berlin, Brandenburg, Hamburg, Mecklenburg-Vorpommern, Lower Saxony, Schleswig Holstein (Brandenburg joined in 2012) + federal funding
- HLRN-I to HLRN-III operated by RRZN Hannover (now LUIS) and Zuse Institute Berlin (ZIB)
- HLRN-IV started operation in 2018, operated by UGOE/GWDG and ZIB
- 2020: application for NHR funding by both sites
 - originally 8 centers funded in total
 - joined by NHR@SW

- **NHR4CES@RWTH** - IT Center - RWTH Aachen
- **NHR4CES@TUDa** - Hochschulrechenzentrum (HRZ) - Technische Universität Darmstadt
- **ZIH** - Zentrum für Informationsdienste und Hochleistungsrechnen - Technische Universität Dresden
- **NHR@FAU** - Regionales Rechenzentrum Erlangen - Universität Erlangen-Nürnberg
- **NHR@KIT** - Steinbuch Centre for Computing (SCC) - Karlsruher Institut für Technologie
- **PC2** - Paderborn Center for Parallel Computing - Universität Paderborn

- **NHR-Nord@Göttingen** - Gesellschaft für wissenschaftliche Datenverarbeitung mbH Göttingen - Universität Göttingen
- **NHR@ZIB** - Zuse-Institut Berlin - Berlin University Alliance
- **NHR South-West**
 - Goethe-Universität Frankfurt
 - Rheinland-Pfälzische Technische Universität Kaiserslautern-Landau (RPTU)
 - Johannes Gutenberg-Universität Mainz
 - Universität des Saarlandes

Thematic focus of NHR@ZIB, NHR@Göttingen, NHR-NORD@GÖTTINGEN

- Applications
 - Life Science (RELION, BART, AMIRA, GROMACS, NAMD, ...)
 - CFD (OpenFoam, TAU, Ansys, ...)
 - Earth System Sciences (PALM, FESOM, ...)
 - Digital Humanities (Python, R, LLMs, ...)
 - Big Data and AI (Spark, TensorFlow, PyTorch, Jupyter, ...)
- Methods, Hardware, Operations
 - Data and Workflow Management
 - Virtualization and Runtime Environments
- Tools
 - ProfiT-HPC Performance monitoring
 - HPC-API for CI/CD and Workflow applications
- Teaching & Training
 - GWDG Academy, HLRN Courses
 - HPC Introduction, MPI, CUDA, HPDA, Containers in HPC, Object Storage

- Research Center for Many-core HPC IPCC in Berlin
- Berlin Big Data Center BBDC
- Berliner Zentrum für Maschinelles Lernen BMZL
- Göttinger HPC Verbund GöHPC
- Campus Institute Data Science (CIDAS) in Göttingen
- Kompetenzzentrum für Höchstleistungsrechnen Bremen BremHLR
- Computational Sciences Center CSC in Kiel
- Lothar Collatz Center in Hamburg
- → tier 3 sites in the 7 HLRN states

- Consultants (domain specific)
 - Support with system usage
 - Applying for compute projects (in particular estimating demand)
 - Bi-annual workshops: procurements (tier 2 und 3), system security, operations, projects, ...
- Local consultants (in part also active as domain spec. consultants)
 - Support with the transition tier 3 → tier 2
 - Initial account application

HPC at GWDG

HPC at GWDG

NHR System “Emmy”

- Phase 1 (Atos, 2018) → EOL
 - 448 standard nodes (2 Xeon Gold 6148, 2x20 CPU-Cores, 192 GB memory)
 - High Mem nodes: 16x 768 GB
 - Storage: 340 TiB HOME (GPFS), 8.1 PiB WORK (Lustre), 120 TiB PERM (Tape)
- Phase 2 (Atos, 2020)
 - 96 standard nodes (2 Xeon Platinum 9242, 2x48 CPU cores, 384 GB memory)
 - High Mem nodes: 16x 768 GB, 2x 1,5 TB
 - 3 GPU nodes (2 Xeon Gold 6148, 2x20 CPU cores, 768 GB memory, 4x NVIDIA V100)
- Since 2021 operated for the NHR alliance

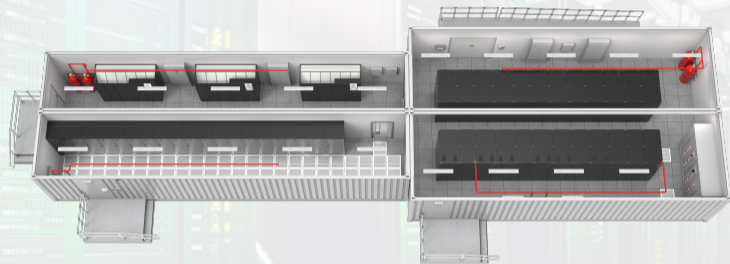


NHR System "Emmy"

Phase 3

- Phase 3 CPU cluster (NEC, 2022, successor Emmy P1)
 - 164x 256 GB, 164x 512 GB, 12x 1 TB, 2x 2 TB
 - 4 nodes per 2U chassis, each with
 - 2x Intel "Sapphire Rapids" 8468 (48 cores) per node
 - 1x Cornelis Omni-path (100 Gbit/s) HCA
- CPU Add-on 2023 (NEC, 2023)
 - 20x 512 GB, 16x 1 TB





- Left section: air-cooled systems (max 19 Racks)
"Emmy" phase 1, GWDG SCC GPU nodes, storage
- Right section: hot water-cooled systems (max 14 racks)
"Emmy" phase 2, GWDG SCC CPU compute nodes



“Emmy” Phase 1 compute nodes

Source: GWDG

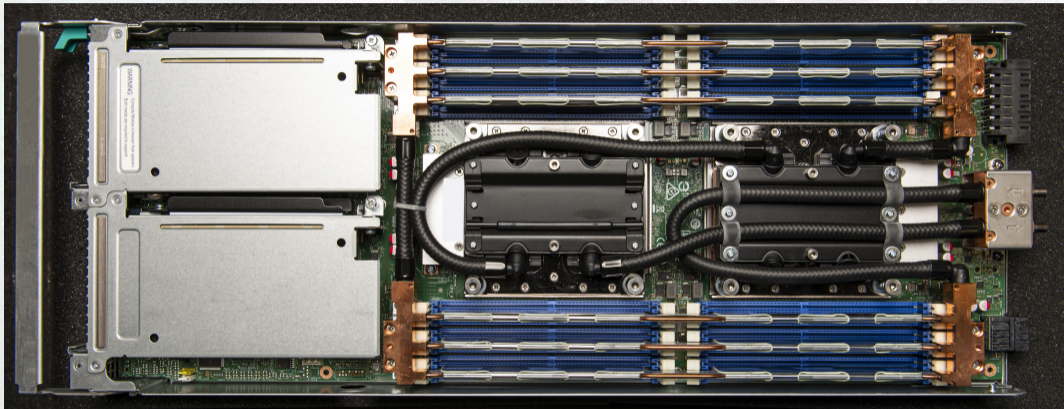


“Emmy” Phase 2 front/back

Source: GWDG

"Emmy" P2, SCC: Intel Xeon Cascade Lake AP HPC@GÖTTINGEN

Water cooling of CPUs and memory



Source: GWDG



“Emmy” Storage controllers+JBODS

Source: GWDC



MDC cooling primary loop

Source: GWDG



MDC cooling
heat exchangers

Source: GWDG



MDC
outside view

Source: GWDG

HPC at GWDG

NHR GPU System "Grete"



- Hosting various HPC systems
 - **“Grete”**, the GPU expansion for **“Emmy”**
 - **“CARO”** for the German Aerospace Center (DLR)
 - Max Planck Society (MPG) HPC systems (housing)
 - upcoming: SCC CPU+GPU expansion 2023
- IT staff+systems from GWDG, Uni Göttingen, UMG, MPG
- Expansions in construction, e.g. future CS lecture hall

GPU System “Grete”

New partition for the NHR system “Emmy”

- Installed at RZGö (MEGWARE, 2022)
- Connection to “Emmy” storage
 - Local flash storage (Atos/DDN) for the GPU cluster at RZGö → /scratch
 - Link to existing WORK (at MDC) → /scratch-emmy
 - LNet routers for connecting to Lustre (IB ↔ eth ↔ OPA)
- New login node glogin9
 - `glogin-gpu.hlrn.de`

GPU System "Grete"

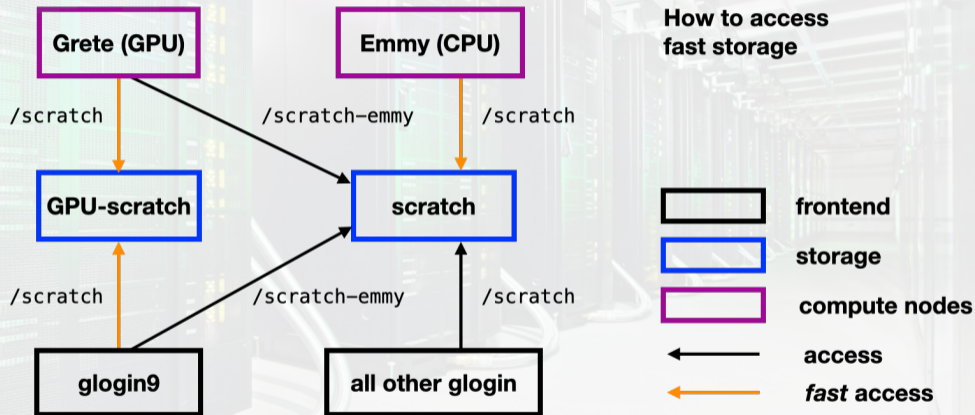
New partition for the NHR system "Emmy"

- Base system: 34 nodes (3 racks), each equipped with
 - 2x AMD Epyc 7513 CPU (32 "Milan" cores, Zen 3 microarch.)
 - 512 GB memory (DDR4, 3200 MHz)
 - 2x 1 TB NVMe SSD
 - 4x NVIDIA A100 GPU (SXM4, 40 GB HBM2 memory)
 - 2x Mellanox InfiniBand HCA (HDR)
- GPU extension 2023: 5 Knoten, each equipped with
 - 2x Intel Xeon Platinum 8468 (48 "Sapphire Rapids" cores)
 - 512 GB Speicher (DDR4, 3200 MHz)
 - 2x 1,92 TB NVMe SSD
 - 4x NVIDIA H100 GPU (SXM5, 94 GB HBM2 Speicher)
 - 2x Mellanox InfiniBand HCA (HDR)



GPU System "Grete"

Frontend Nodes



GPU System "Grete"

Specifications of the GPUs

- 136x [A100 Tensor Core GPU](#) (Ampere architecture)
 - GA100 GPU: 108 SMX → 6912 CUDA cores, 432 Tensor cores
 - 19,49/9,746/155,92 TFLOPs (SP/DP/Tensor)
 - 40 GB VRAM (HBM2e)
 - SXM4 form factor, 4 GPUs on "Redstone" baseboard
 - NVLINK full mesh, 800 GB/s bidirectional bandwidth per GPU
- GPUS can be requested individually in the same batch system (Slurm)



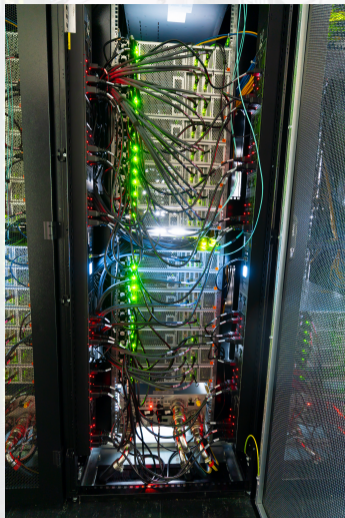
"Grete"/REACT back

Source: GWDG



"Grete"/REACT
front

Source: GWDG



“Grete”/REACT
power+fabric

Source: GWDG

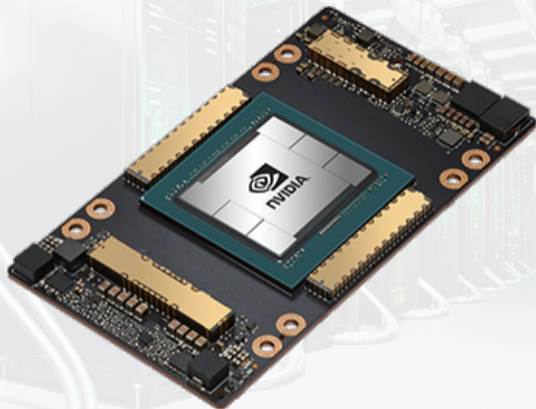


“Grete”/REACT LAN+DLC

Source: GWDG

GPU example

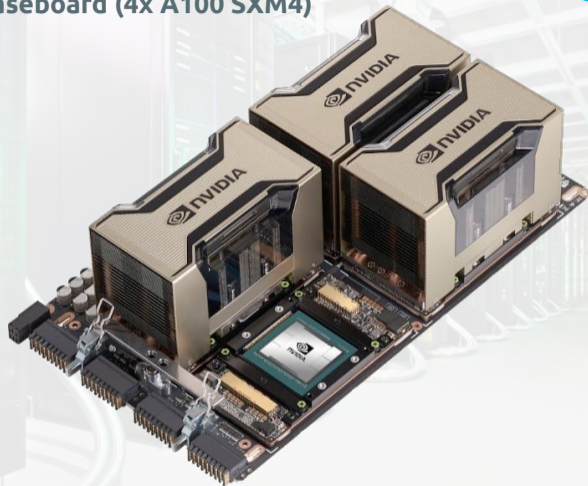
NVIDIA A100 SXM4



Source: NVIDIA

GPU example

NVIDIA "Redstone" Baseboard (4x A100 SXM4)



Source: NVIDIA

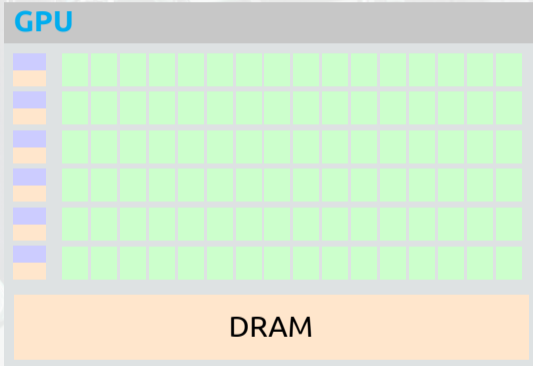
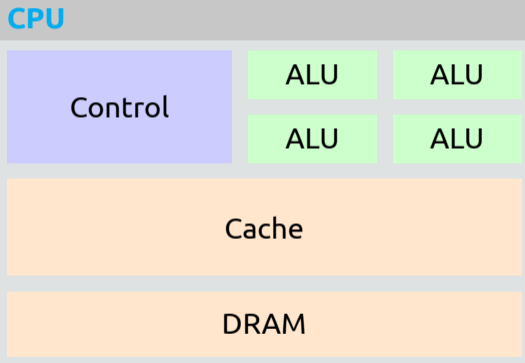
Excursion: Using Graphics Cards for Scientific Computing

Excursion: Using Graphics Cards for Scientific Computing

GPGPU

- Speedup of single-thread execution stagnating (cf. Moore's Law)
- Parallelize software at the thread level
- General Purpose Computation on Graphics Processing Unit
- Graphics Processor is used as Coprocessor

Comparison: CPU/GPU



GPGPU Programming Model

memory allocation on host

invoke mem. alloc on GPU

copy of kernel input data

invoke kernel ex. on GPU

copy of kernel output data

memory allocation on GPU

copy of kernel input data

kernel execution on GPU

copy of kernel output data

Excursion: Using Graphics Cards for Scientific Computing

Usage

AI Research

- Labs with many students doing small-scale experiments
- Mid-scale development/experiments by researchers or students
- Large-scale/scalable execution of developed pipelines

HPC Research

- Small-scale / single-node testing of an application
- Labs with many students learning to use HPC
- Scalable/Multi-node execution of existing/developed application

- **Asynchronous batch jobs**

- submitted jobs are executed later
- depending on load, individual priority (cf. FairShare)
- once running resources are granted exclusively
- most resources, main form of GPU usage

- **Interactive batch jobs**

- resources are requested in batch mode
- live usage once jobs starts
- exclusive usage, potentially long waiting time

We partition GPUs for maximum sharing (into up to 7 "devices")

- **Shared resources**

- partition with no exclusive usage (overprovisioning)
- no waiting time, but no guaranteed performance
- for testing and small-scale experiments

Excursion: Using Graphics Cards for Scientific Computing

System Configuration

- **Partitions** (-p, --partition=<partition_names>)

Partition	Purpose	Nodes	GPU
grete	huge jobs	52	4x A100
grete:shared	small jobs	62	4x-8x V100/A100
grete:interactive	testing/hyperparameters	3	V100/A100 MIG sl.
grete:preemptible	interrupt+resume	3	V100/A100 MIG sl.
gpu	small jobs	3	4x V100

- **GPUs** (-G, --gpus=[type:]<number>)

GPU/slice	Per node	Description
A100	max. num.	full GPU
2g.10gb	8	2x Compute Instance, 10 GB memory (2 per GPU)
3g.20gb	4	3x Compute Instance, 20 GB memory (1 per GPU)

- Full docs: <https://www.hlrn.de/doc/display/PUB/GPU+Usage>

- NVIDIA HPC SDK 23.3: `nvhpc/23.3`
 - for own MPI implementation: `nvhpc-nompi/23.3`
 - ...and also own compiler: `nvhpc-byo-compilers/23.3`
- OpenMPI with CUDA and InfiniBand support
 - HPC-X: `nvhpc-hpcx/23.3`
 - NVIDIA/Mellanox OFED stack: `openmpi-mofed/4.1.5a1`

Application Software

Application Software

Environment Modules

- Application domains & software
 - Quantum chemistry, Molecular dynamics
 - Bioinformatics, Genomics, Evolutionary Biology
 - Astrophysics, Cosmology, Numerical Fluid Dynamics
 - Numerical Software: MATLAB, Maple, Mathematica
 - Machine Learning
 - Medical Imaging (e.g. MRI)
 - Data Analytics: Python, R, Spark, ...
- Module system `Lmod` → software catalogue `module avail`

```
Terminal
Datei Bearbeiten Ansicht Terminal Reiter Hilfe
----- BIOINFORMATICS -----
abyss/1.5.2          fitchip/0.2.0      picard/2.10.5
abyss/1.9.0          fithichip/6.0     picard/2.20.2      (D)
abyss/2.1.0          flappie/1.1.0     pilon/1.22
allpathslg/52488    flash/1.2.11      platanus/1.2.1
amos/3.1.0          flye/2.5          plink/1.0.7
augustus/3.2.2      freebayes/1.2.0   plink/1.90        (D)
augustus/3.3.2      funannotate/1.7.4 pomoxis/0.2.2
augustus/3.3.3      gapfiller/1.10    popbam/0.3
bambam/1.4          genomethreader/1.7.1 porechop/0.2.3
bamtools/2.4.0      genomertools/1.5.6 poretools/0.6.0
bamtools/2.4.1      gmap/20190610     prank/170427
bamtools/2.5.1      grace/5.1.24      profphd/1.0.42
bbmap/37.66         grass/7.4.2       pyrad/3.0.66
bbmap/38.68         great/1.5          qqis/2.14
bcftools/1.8        guppy/CPU/3.1.5   quast/5.0.1
bcftools/1.9        guppy/GPU/3.1.5   racon/0.5.0
bcl2fastq/2.20     hail/0.2_novenv   racon/1.4.3        (D)
beagle/4.1          hail/0.2          raxml-ng/0.9.0
beast/1.8.2         hic-pro/2.11.1    raxml/8.2.4
beast/2.4.7         hic-pro/3.0.0     (D) raxml/8.2.12      (D)
beast/2.6.1        (D) hichipper/0.7.5
bedtools/2.18.0    hisat2/2.1.0     recon/1.08
                                     repeatmasker/4.0.6
lines 56-78
```

Application Software

Data Analytics Applications

- Interactive **Jupyter** notebooks can use HPC as a backend
- **TensorFlow** is a popular AI tool that can use CPU and GPU dynamically
- **Apache Spark** clusters can be spawned as multi-node HPC jobs

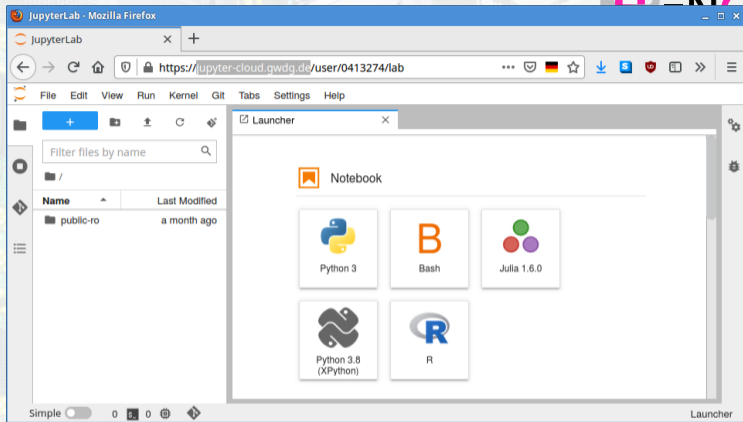


- Apache Spark
 - Spark clusters can be set up automatically
 - interaktive sessions (e.g. Scala)
 - Monitor cluster status via web interface
- TensorFlow
 - Integration in users' Python environment
 - Using GPU nodes with CUDA possible



Overview & Documentation of HPC applications:

<https://docs.hpc.gwdg.de>



Data Analytics Jupyter

- JupyterHub on HPC: https://docs.hpc.gwdg.de/getting_started/jupyterhub
- JupyterCloud instance: <https://jupyter-cloud.gwdg.de>

Spawner Options

Select a job profile:

GWDG HPC with own Container

Set your own Singularity container location (allowed characters: [a-zA-Z.~-])

\$HOME/jupyterhub-gwdg/jupyter.sif

Set the duration (in hours):

8

Set the number of cores:

10

Set the amount of memory (in GB):

32

Jupyter Notebook's Home directory

\$HOME/jupyterhub-gwdg

[Documentation](#)

Spawn

- jupyter-hpc.gwdg.de spawns Jupyter on HPC
- supports IPython Parallel
- Users can choose resources, start individual Singularity Container
- Support for Jupyter on (shared) GPU nodes



Getting access

Getting access

Account creation

- **Scientific Compute Cluster (SCC)** (Tier 3) for UGOE/MPG
 - Covering basic usage, Throughput-Computing
 - Consulting for the transition to T2, Scaling tests
 - Requirements
 - GWDG Full Account (User/affiliation UGOE/MPG)
 - Thesis projects: Apply for an account via supervisor's institute
 - Courses: temporarily usable guest accounts
- **NHR systems "Emmy"/"Grete"** (Tier 2)
for Universities and Research Institutions in the NHR alliance
 - Proposal preparation, Software setup
 - Scientific review of compute time projects

- **Scientific Compute Cluster (SCC)**

- Account activation: request via hpc@gwdg.de
- Job prioritization via FairShare

- **NHR systems “Emmy”/“Grete”**

- Apply for an account at nhr-support@gwdg.de
- Starting from 300.000 Coreh per quarter
 - Project application for further compute time
 - NHR-wide portal JARDS: <https://jards.nhr-verein.de/>
 - Review by scientific board (+techn. review by HPC sites)
 - Scientific review can be omitted , if project has been granted by BMBF/DFG/EU/GCS/NHR/... (Whitelisting)

Further resources

Where to get more information?

- **GWDG HPC docs** <https://docs.hpc.gwdg.de> shows you information about
 - Basic access
 - Hardware overview
 - How to submit jobs
 - Advice on most common applications
- **Support** mail address: hpc@gwdg.de / nhr-support@gwdg.de
- General GWDG support address: support@gwdg.de
- Rocket.Chat channel #hpc-users (<https://chat.gwdg.de>)

- We offer **courses** on SCC usage, programming with MPI, OpenMP, CUDA, Python, Research Data Management and more
→ **GWDG Academy** <https://www.gwdg.de/academy>
- Regular introductory course **“Using the Scientific Compute Cluster”**
 - will be held online via BigBlueButton
- Yearly event: **“Parallel Programming Day”**
<https://user.nhr.zib.de/wiki/spaces/PUB/pages/429900/Parallel+Programming+Day+series>

- We offer **courses** on SCC usage, programming with MPI, OpenMP, CUDA, Python, Research Data Management and more
→ **GWDG Academy** <https://www.gwdg.de/academy>
- Regular introductory course **“Using the Scientific Compute Cluster”**
 - will be held online via BigBlueButton
- Yearly event: **“Parallel Programming Day”**
<https://user.nhr.zib.de/wiki/spaces/PUB/pages/429900/Parallel+Programming+Day+series>
- **Thanks for your attention! Questions?**