**6th HPC Café**

**06.06.2025**

# Agenda HPC Café – 06.06.2025

1. HPC Job Performance Monitoring

2. User Support through Voucher Projects

3. Questions and Answers

Scientific Computing Center (SCC)

# 2. User Support through Voucher Projects

**www.kit.edu**

# User Support Structures

1. Workshops/Course
   - Every semester: HPC introduction and advanced course

2. Wikis
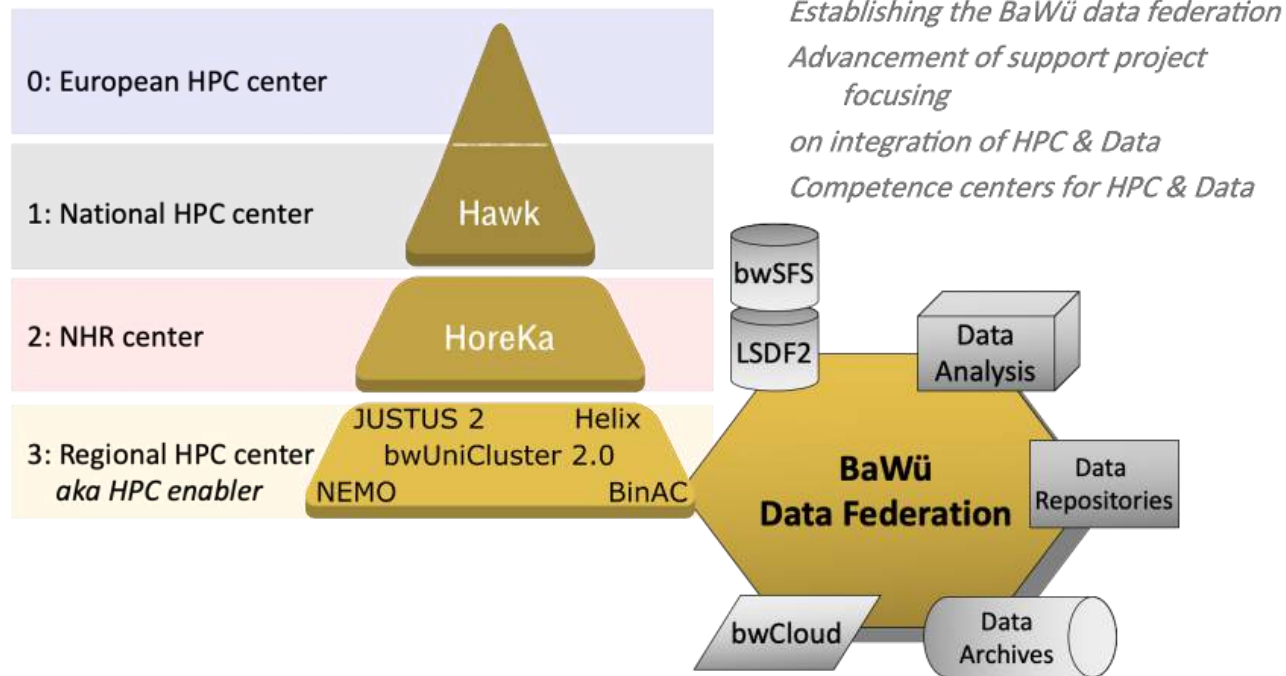   - https://www.nhr.kit.edu/userdocs
   - https://wiki.bwhpc.de

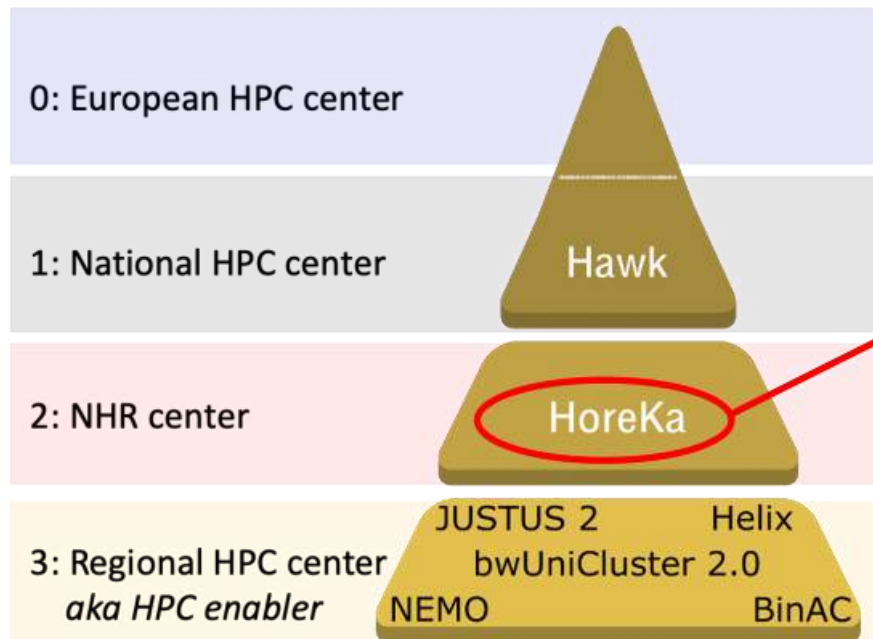3. Support projects
   - Voucher projects (-> NHR)
   - Tiger team projects (-> bwHPC)

# bwHPC:

**Baden-Württemberg's implementation strategy for HPC, Data Intensive Computing & Large Scale Scientific Data Management**



*Advancement of federated HPC@tier3*

*Establishing the BaWü data federation*

*Advancement of support project focusing*

*on integration of HPC & Data*

*Competence centers for HPC & Data*

# NHR – Nation High Performance Computing

## HPC in Baden-Württemberg

| Level | Center |
|---|---|
| 0: European HPC center | |
| 1: National HPC center | Hawk |
| 2: NHR center | HoreKa |
| 3: Regional HPC center *aka HPC enabler* | JUSTUS 2, Helix, bwUniCluster 2.0, NEMO, BinAC |

## National HPC at Tier 2

| Centers | Universities |
|---|---|
| NHR4CES@RWTH | RWTH Aachen |
| NHR4CES@TUDa | TU Darmstadt |
| NHR@FAU | Univ. Nürnberg-Erlangen |
| NHR@Göttingen | GWDG + Univ. Göttingen |
| NHR@KIT | KIT |
| NHR@TUD | TU Dresden |
| PC2 | Univ. Paderborn |
| NHR@SW | Univ. Frankfurt a.M., Mainz, Kaiserlautern-Landau, Mainz, Saarland |

Scientific Computing Center (SCC)

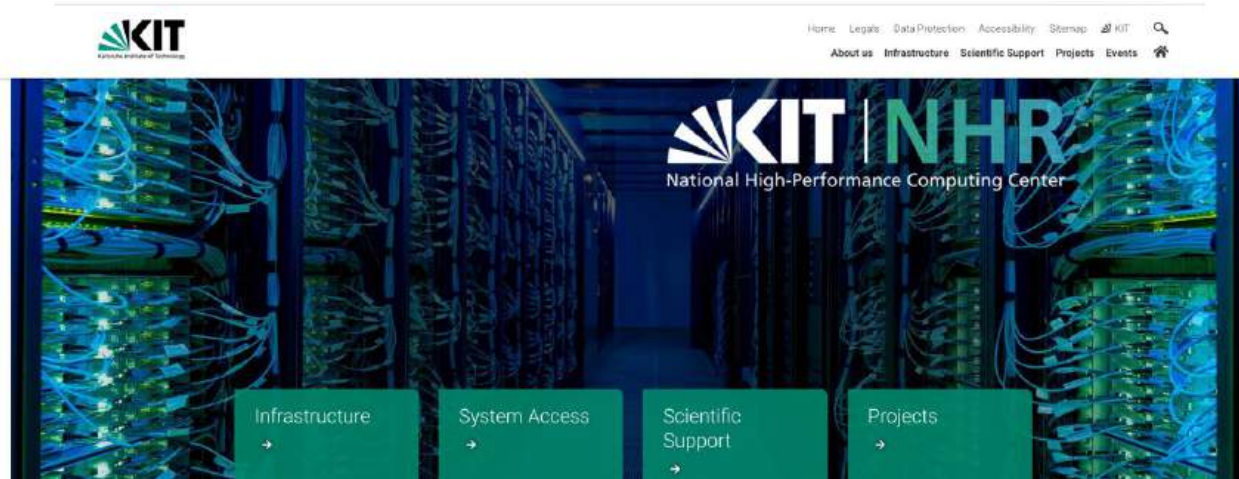# bwHPC: How to get support?

- User support projects:
  - bwHPC-S5 & bwRSE4HPC

- bwHPC-S5:
  - Organized by competence centers + cross-sectional support team
    - bwSupportPortal + Wiki + courses
    - Tiger team projects

- bwRSE4HPC:
  - cf. talk

# NHR: How to get support?

- Website, https://www.nhr.kit.edu
  - Resources, Documentation, Consulting, Training, Support ...



- Voucher Projects

  via. NHR support teams (e.g. SSPE) & bwRSE4HPC

# 3. Questions and Answers

# Questions and Answers

- Did you every apply for a Voucher / Tiger Team project?

- Did you every attend a course?

- Did you every use our online documentation?
  - Or used different documentation sources?

Scientific Computing Center (SCC)

# Introduction to the Job Performance Monitoring
# JobMon at NHR@KIT

Holger Obermaier | 6. June 2025

KIT

# JobMon

## History

- Started as a research project as part of J. Schmitt's Bachelor's thesis in 2021/2022
- Contributions from employees (B. Bytyqi, H. Obermaier) and student assistants (D. Schild, J. Schmitt, F. Wedler)
- Basis for D. Schild's master's thesis in 2024/2025

## Overview

- Web Service: JobMon⧉
- Access with HoreKa user account
- Same networks restrictions apply as for HoreKa (VPN)
- Visualizes performance metrics
- Aggregates performance metrics for improved clarity
- Metrics are continuously collected on cluster nodes
- Metrics are stored $\Rightarrow$ Performance changes over time can be tracked.

# Demo Application Workloads

## Categories

- *Compute Bound*
  - Performance is limited by floating point performance
  - Benchmark *DGEMM*⧉ performs a matrix matrix multiplication
- *Memory bound*
  - Performance is limited by memory bandwidth
  - Benchmark *Stream*⧉ / *BabelStream*⧉ performs vector computations
  - Benchmark *HPCG* performs the conjugate gradients method with a sparse matrix
- *Communication bound*
  - Performance is limited by interconnect bandwidth or latency
  - Benchmark *OSU Micro-Benchmarks*⧉ performs MPI point to point communication
- *IO bound*
  - Performance is limited by filesystem throughput or metadata rate
  - Benchmark *IOR*⧉ performs parallel IO operations

# Login page

# Welcome Page

# Jobs page

$\Rightarrow$ Overview of the personal HoreKa batch jobs

## Filter options

- Batch partition
- Number of nodes or GPUs
- Running or finished jobs
- Job execution time
- Exit code $\Rightarrow$ successful or non-successful job run
- Tags assigned to the job $\Rightarrow$ group jobs by your needs

# Jobs page - Filter

# Jobs page - Spider Plot

## Spider Plot

- Plot that characterizes the overall job performance
- High level view on performance limitations
- Allows categorization as:

  - IO bound
  - Memory bound

  - Compute bound
  - Communication bound

- Average and maximum values for the metrics:

  - CPU floating point operations per second
  - CPU memory bandwidth
  - GPFS IO operations
  - GPFS metadata operations

  - GPU utilization
  - GPU memory utilization
  - InfiniBand bandwidth

# Spider Plot - Compute Bound Workl. (`cpuonly`)



Id: **3015837**

User: bq0742 (hk-project-scs)

Partition: cpuonly

Name: cpu.dgemm

Nodes: 1

Start: 27.3.2025, 16:25:13

End: 27.3.2025, 17:21:28

Exit code: 0

Mean values    Max values

CPU FLOPs (per core)

CPU Memory Bandwidth

Infiniband Traffic

GPFS IOPS (HOME)

GPFS MetaOps (HOME)

# Spider Plot - Memory Bound Workl. (`cpuonly`)

Id: **3015838**

User: bq0742 (hk-project-scs)

Partition: cpuonly

Name: cpu.stream

Nodes: 1

Start: 27.3.2025, 16:25:13

End: 27.3.2025, 17:27:08

Exit code: 0

# Spider Plot - Communication Bound Workl. (`cpuonly`)

# Spider Plot - Compute Bound Workl. (accelerated)

# Spider Plot - Memory Bound Workl. (accelerated)



Id: **3013359**

User: bq0742 (hk-project-scs)

Partition: accelerated

Name: gpu.hpcg,apptainer

Nodes: 1          GPUs: 4

Start: 25.3.2025, 11:56:12

End: 25.3.2025, 13:02:09

Exit code:  0

Mean values   Max values

CPU FLOPs (per core)

CPU Memory Bandwidth

GPFS IOPS (HOME)

GPFS MetaOps (HOME)

GPU Memory Utilization

GPU Utilization

Infiniband Traffic

# Spider Plot - Comm. Bound Workl. (accelerated)

# Per job page

## Configuration options

- Select subset of nodes
- Select time range
- Select subset of metrics
- Set tags $\Rightarrow$ Mark jobs for quick access

## Toggles

- *Automatic Scaling* $\Rightarrow$ Select y-axis limits depending on metrics limits
- *Quantile View* $\Rightarrow$ Only plot quantiles
- *Changepoints* $\Rightarrow$ Identify times where performance metric behavior changes

# Per job page - Options, Toggles, . . .

# Toggle: Changepoints Off - Infiniband Packets

# Toggle: Changepoints On - Infiniband Packets

# Toggle: Changepoints - Insight

## Insight

- OSU Micro-Benchmark performs communication with increasing message size
- Message size influences performance
  $\Rightarrow$ Visible as steps in the graph
- Steps can be detected by *changepoint algorithm*
- *Changepoints* marked by vertical lines

# Per Job Page - Toggle: Quantile View

## Quantile View

- Improve clarity when too many graphs are displayed in one diagram
- Only three graphs (25% / 50% / 75% Quantile = Quartiles) are drawn
- 25% Quantile $\Rightarrow$ 25% of the measured values are below this graph
- 50% Quantile $\Rightarrow$ Median
- Difference between upper and lower Quantile $\Rightarrow$ Measure for the spread of the metrics

# Toggle: Quantile View Off - CPU Temperature

# Toggle: Quantile View On - CPU Temperature

# Quantile View - Insight

## Insight

- CPU temperatures are collected per hardware thread
  $\Rightarrow$ Diagram appears very cluttered
- Quantile View shows the distribution of the CPU temperatures much clearer

# Toggle: Quantile View Off - CPU Frequency



2025-06-06    H. Obermaier: Introduction to JobMon    Scientific Computing Center (SCC)

# Toggle: Quantile View On - CPU Frequency

# Quantile View - CPU Frequency

**Insight**

- CPU frequencies are collected per CPU core
  $\Rightarrow$ Diagram appears very cluttered
- Quantile View shows the distribution of the CPU frequencies much clearer

# Per Job Page - Roofline Plot

## Roofline Plot

- Diagram with:
  - y-Axis: Floating point operations per second (FLOP/s)
  - x-Axis: Computational intensity (FLOP/s / byte)
- Roofline shows two performance limiting factors:
  - For *low* computational intensity: Memory bandwidth
  - For *high* computational intensity: Processor peak performance
- Plot point:
  - For each CPU package
  - For each measurement interval

# Roofline Plot - Stream

# Roofline Plot - DGEMM



Roofline plot

# Roofline Plot - Insight

## Insight

- Stream (memory bound)
  - $\Rightarrow$ Low computational intensity
  - $\Rightarrow$ Limiting factor: Memory bandwidth
- DGEMM (compute bound)
  - $\Rightarrow$ High computational intensity
  - $\Rightarrow$ Limiting factor: Processor peak performance

# Per Job Page - Performance Category Energy

## Performance Category Energy

- CPU power consumption of DRAM channels and the package
- GPU power consumption
- Server system power consumption

# Performance Category Energy - Stream

# Performance Category Energy - DGEMM

# Performance Category Energy - Insight

## Insight

- Stream (memory bound)
  - $\Rightarrow$ constantly high pressure on the DRAM subsystem
  - $\Rightarrow$ Constantly high energy consumption
- DGEMM (compute bound)
  - $\Rightarrow$ Less pressure on the DRAM subsystem
  - $\Rightarrow$ Varying power consumption over time

# Performance Category Energy - HPCG

# Performance Category Energy - HPCG

# Performance Category Energy - Insight

## Insight

- GPU-HPCG benchmark
  - Preparation phase executed on the CPUs
  - Computation phase executed on the GPUs
- CPU package power consumption: Higher in preparation phase than in the compute phase
- GPU power consumption: Higher in compute phase than in preparation phase

# Per Job Page - Performance Category Filesystem

## Performance Category Filesystem

- Meta data operation
  - Number of file open / closes
  - Number of directory reads
  - Number of inode updates
  - …
- IO throughput
  - Bytes read / written
  - Number of read / writes

# Performance Category Filesystem - IOR

# Performance Category Filesystem - IOR

# Performance Category Filesystem - IOR

# Performance Category Filesystem - Insight

## Insight

- Parallel IO: Performed from multiple nodes
- Two phases:
  - Phase 1: Files are written
  - Phase 2: Files are read
- Read throughput is higher than write throughput

# Per Job Page - Performance Category Interconnect

## Performance Category Interconnect

- InfiniBand bandwidth
  - Sent
  - Received
  - Aggregated sent and received
- InfiniBand number of packages
  - Sent
  - Received
  - Aggregated sent and received

# Performance Category Interconnect - OMB

# Performance Category Interconnect - OMB

# Performance Category Interconnect - OMB



2025-06-06    H. Obermaier: Introduction to JobMon    Scientific Computing Center (SCC)

# Performance Category Interconnect - Insight

## Insight

- OSU Micro-Benchmark performs MPI point to point communication:
    - Node hkn1015 only sends data (receive bandwidth is zero)
    - Node hkn1020 only receives data (send bandwidth is zero)

# Per Job Page - Performance Category Memory

## Performance Category Memory

- Amount of memory used
  - On the system
  - On the GPU

- CPU memory bandwidth
- GPU memory
  - Utilization (in %)
  - Frequency

# Performance Category Memory - Stream

# Performance Category Memory - DGEMM

# Performance Category Memory - Insight

## Insight

- Stream (memory bound):
  $\Rightarrow$ Constantly high pressure on the memory subsystem
- DGEMM (compute bound):
  $\Rightarrow$ Less pressure on the memory subsystem
  $\Rightarrow$ Varying bandwidth over time

# Performance Category Memory - BabelStream

# Performance Category Memory - DGEMM

# Performance Category Memory - Insight

## Insight

- BabelStream (memory bound):
  $\Rightarrow$ Fully utilizes the memory subsystem of the GPU
- GPU-DGEMM (compute bound):
  $\Rightarrow$ Less pressure on the GPU memory subsystem
  $\Rightarrow$ Varying utilization over time.

# Per Job Page - Performance Category Performance

## Performance Category Performance

- Floating point operation per second (FLOP/s),
  collected per *hardware thread*,
  aggregated per *core* or per *socket*

- Instructions per cycle (IPC),
  collected per *hardware thread*,
  aggregated per *core* or per *socket*

- CPU time spend in kernel and in user space

- One minute Linux load average

- GPU utilization

- CPU and GPU frequency

# Performance Category Performance - DGEMM

# Performance Category Performance - DGEMM

# Performance Category Performance - Insight

## Insight

- Floating point operation per second (FLOP/s) are collected per hardware thread
- Examine even utilization of cores
  $\Rightarrow$ Aggregate per core
- Examine even utilization of CPU sockets
  $\Rightarrow$ Aggregate per socket
- Summed FLOP/s is the same in both diagrams

# Per Job Page - Performance Category Temperature

## Performance Category Temperature

- CPU
- GPU

# Performance Category Temperature - DGEMM



2025-06-06     H. Obermaier: Introduction to JobMon     Scientific Computing Center (SCC)

# Performance Category Temperature - Insight

## Insight

- GPU-DGEMM only utilizes one GPU
- Only this GPU gets hot
- Other GPUs maintain lower temperatures

# Performance Category Temperature - HPCG

# Performance Category Temperature - Insight

## Insight

- GPU-HPCG
    - Preparation phase executed on the CPUs
      $\Rightarrow$ GPUs not utilized
      $\Rightarrow$ Low temperature
    - Computation phase executed on the GPUs
      $\Rightarrow$ Higher temperature

# Per Job Page - Additional Features

## Additional Features

- For multi-node jobs there is a configuration option to select the per-node aggregation function used (e.g. average, sum, maximum)
- Live view of running jobs
- Download CSV file
  - All metrics
  - Use in spread sheet application or Python
- Outlook
  - Availability for Cluster uc3
  - Automatic job analyzer
    - $\Rightarrow$ Assign tags for detected characteristics

# bwRSE4HPC

*Offering Research Software Engineering (RSE) services for HPC Users in Baden-Württemberg*

Jasmin Hörter, Dominic Kempf, René Caspart, Marcel Koch, Inga Ulusoy, Andreas Baer, Glen Hunter, Thomas Isensee, Kai Riedmiller, Tim Schrader

SCC
Scientific
Computing Center

KIT
Karlsruher Institut für Technologie

SCIENTIFIC
SOFTWARE
CENTER

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

# Our Mission

*Support researchers in achieving their goals through software development.*

Why might this be of interest to you?

- You lack the man-power to realize the changes you want to see.

- You want better software but are unsure how. Your current software might be:

  - too slow

  - too unintuitive

  - too hard to use

Marcel Koch - bwRSE4HPC

# Overview

Provide software development services.

Strengthen Research Software Engineering practices.

Aimed specifically at users of bwHPC clusters.

Marcel Koch - bwRSE4HPC

# Our Services

## Short term projects

- Less than 6 months, free of charge

## Long term projects

- Longer than 6 months, requires third-party funding

Marcel Koch - bwRSE4HPC

# Short Term Projects

Initiated by filling out our request form.

- We need contact information and a short description.

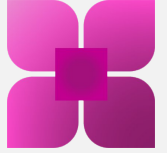The next step is a consultation with us to refine the project.

- We will create a concrete work plan based on the discussion.

After approval 1.5 RSEs provided by us will work on the project.

- We provide regular progress updates to the users.

We provide a final report to ensure the sustainability of our solution.

Marcel Koch - bwRSE4HPC
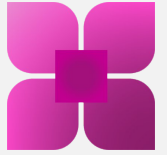
# Long Term Projects

Initiated through third-party funded projects.

You want to pursue a new research idea:

- It has an integral software component

- Your team lacks the relevant software expertise

- We can collaborate with you to cover the RSE aspects

Might naturally evolve from a short term project.

# Project Ideas

First short term project:

- Replace homegrown linear algebra backend with external HPC library.

Examples within our expertise:

- Adding GPU support
- Enable distributed computing
- …

- Performance optimization
- Prototyping

Anything missing? Come talk to us!

Marcel Koch - bwRSE4HPC

# Get in Touch

https://www.bwrse4hpc.de/

support@bwrse4hpc.de