# NFDI4HPC

*Sprecher:* Prof. Dr.-Ing. André Brinkmann

Johannes Gutenberg-Universität Mainz, brinkman@uni-mainz.de

**Key questions/objectives of the consortium**

Understanding natural processes with the help of high-performance computers (HPC) has become one of the main drivers for new scientific insights. This understanding is often critically supported by pre- and post-processing experimental data, e.g., by aligning short-read data from next-generation sequencing machines or when analyzing the vast output of experiments at particle accelerators like the one at CERN. More fundamentally, though, supercomputers have become an experimental facility themselves, often referred to as the third leg of science next to experimental and theoretical science. Simulations are performed on HPC systems based on scientific models or even ab-initio and are only compared with "real life" experimental data after processing a huge evaluation space has led to results which are promising or even optimal. They are often based on the coupling of "classical" numerical HPC techniques, big data analysis, and machine learning techniques, and have successfully deployed to derive new drugs, make weather predictions, construct and simulate cars and airplanes, and to steer urban traffic.

Often overlooked, though, the output of these simulations is not restricted to the domain science alone but is often accompanied by monitoring and telemetry information collected while running a specific job on the HPC system. The information can include data about the set of compute nodes used, statistics about hardware resource utilization on each node, runtimes and scheduling times, as well as information about I/O access patterns, network traffic, or energy usage. This is further aggravated by the fact that HPC systems have become very heterogeneous, including CPU (central processing unit) and GPU (graphical processing unit) nodes, different levels of storage systems, and hierarchical interconnection networks. Additionally, beyond data from the HPC systems, it can also include data from building control systems or even weather data due to its impact on cooling systems.

Unfortunately, only part of the collected data is following standards, some of the data is highly unstructured (e.g., log data), data is originating from independent and uncoordinated systems (e.g., facility vs. system vs. application data) and data is typically not shared between facilities and often not even within the institution running the HPC system(s).

Nevertheless, initial experiments at individual facilities have shown that this data is very valuable, not only from a research, but even more so from an economic perspective when considering HPC center planning and operations. Just as one example, previous research in the scheduling domain has shown that even small changes to the configuration of HPC environments can easily change throughput by more than 25%, a huge impact if taking the annual investment into HPC at German universities and research facilities of around 100.000.000 Euro into account, e.g., 25% more science with the same investment.

However, comparison studies can only be performed if data from different HPC centers is findable, accessible, interoperable and re-usable (FAIR principles), which is currently not possible. Furthermore, many of the insights gained by collecting this information can be applied to "standard" university data centers and also be transferred into industrial settings.

**Summary of the planned research data infrastructure that is specifically intended to address the needs of research users in their respective work processes**

Currently, there are two kinds of measured data originating from HPC systems: Application measurements often use specifically instrumented versions that provide a great deal of information on how parallel applications utilize the HPC system but are limited to the context and runtime of the application. System and facility monitoring provide a continuous and more holistic view but lack the in-depth details of application measurements. A combination of both offers the most value, e.g., by enriching a trace of function calls and communications within an application with power measurements from a system-wide energy monitoring.

To provide broad coverage of HPC research data, the planned infrastructure will leverage both,

detailed application measurements and continuous monitoring. Partners of this consortium have been actively involved in defining, e.g., the Open Trace Format 2. OTF 2 enables researchers to archive application measurements in a standard format and enables interoperability for analysis and visualization with a wide range of established tools. On top of the raw storage, application traces can be indexed using automated metadata extraction. The planned infrastructure will cater to both, high-resolution metric data and unstructured log data. The storage level will exploit the time-series nature of such data by applying a multi-level aggregation. This approach enables very efficient lookups and thus makes extensive datasets accessible for analysis. By standardizing the description of components within the HPC system that act as data sources, data from the application, system, and facility level can be related more easily. This joint analysis particularly enables research on performance and resilience that considers the complex interactions between software, underlying hardware, and supporting infrastructure.

### *Planned implementation of the FAIR principles*

There are only very few and mostly outdated public archives providing research data available within the HPC community. These archives only include selected datasets and are not easily findable according to the FAIR principles. Most collected dataset are not open to the public and can only be accessed by the individual hosting sites.

Reproducibility of results and open data are nevertheless accepted to be important within the community. The IEEE/ACM Supercomputing Conference (SC), e.g., started a reproducibility initiative and it made an Artifact Description (AD) appendice mandatory for all papers submitted to the SC 2019 technical program. Also, other conferences in our domain, e.g., the ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming (PPoPP) encourage authors of accepted papers to submit their papers for an artifact evaluation.

Data and metadata repositories in our domain therefore shall comply with internationally accepted standards on research data management to ensure that reproducibility and comparability of results are not restricted to individual cases. It is therefore required to comply with the FAIR (findable, accessible, interoperable, re-usable) principles and, in particular, with metadata standards.

Some of the contributors to this document are already involved in projects building unified research data infrastructures, such as GeRDI. Leveraging this expertise, we will build as part of the NFDI activities policies and procedures on metadata standardization as well as on making data actually findable by research data search engines (e.g., Google Dataset Search or GeRDI). Minimum requirements will include the availability of core metadata properties defined by DataCite, the assignment of persistent identifiers to datasets, and an open standard interface (e.g., OAI-PMH) for metadata harvesting by search engines.

Within the NFDI, metadata schemes and data management concepts will be developed for the computer science community in general and specifically the "HPC as a science" community. The general idea for serving the specific community is to add upon open standard data repository frameworks (e.g. ckan, Invenio), providing easy access to the data, and on basic metadata standards as mentioned above. Standard metadata (e.g., Dublin Core, DataCite) which are perfectly interoperable with other scientific domains will be enriched by extensions specific to our community, to be discussed in a co-development process of repository providers, computer science and HPC scientists. A key aspect in defining community metadata schemes will be ensuring sufficient richness to enable the re-use and – as far as experimentally possible – the reproduction of scientific data.

The NFDI's co-development process for metadata and repository policies, as also actual implementations of these, will be driven and accompanied by a series of workshops with the NFDI consortium, financial and research institution stakeholders as well as data-infrastructure providers and projects from Germany and the EU. Collaboration with inclusive and horizontal initiatives such as EOSC-hub, EUDAT, OpenAire and GeRDI, also in the scope of the workshops, will be launched, based on already existing contacts of the HPC centers to these infrastructures. This will make our approach interoperable across disciplines and national borders and prepare us for immersion in the EU-wide (e.g., EOSC) e-infrastructure landscape.

The outcomes of these workshops will also serve as input for the development of training

material, which will both be designed to be taught within online courses and in hands-on training sessions at partner facilities and as tutorials at conferences.

### *User Impact and Networking*

Facility wide data collection provides benefits not only for center operations, in terms of increased system efficiency due to continuous data-driven optimization, but also has a direct impact on any user of the system. Users can query performance and application metadata of their executions without having to individually include such tracking themselves, which typically leads to one-off solutions and hacks, and a homogenization of data, as proposed by this activity, allows users to have consistent metrics across system generations, systems and facilities. Access to such        data then enables feed-back driven optimizations at runtime, continuous and gap-free performance tracking and regression, as well as data-driven input into future application development and design, all which currently has to be based on data with huge gaps and inconsistent and system specific collection procedures. Mitigating this using a coordinated data collection, management and analysis approach, as proposed here, therefore has the potential to significantly increase application and application development efficiency, and hence to further contribute to improved machine utilization and with that, science per Euro achieved.

Challenges specifically for this use case that need to be addressed, beyond the general ones discussed already, include data protection and access rights, user-level access to data through both visual and analysis tools, and user-level customization, e.g., by adding application level information to the data stored. This consortium will work close with the core user communities at each center to address these needs based on concrete use cases.

### *International Perspective and Networking*

Due to the potentially large positive impact of ubiquitous and continuous monitoring, combined with a systematic storage and analysis of the resulting data as laid out above, most large centers world-wide are working on establishing large-scale monitoring systems. In particular the US Department of Energy (DOE) sites are investing significantly in these technologies as part of their Exascale Computing Project (ECP). Examples are the Sonar system at Lawrence Livermore National Laboratory, Oak Ridge National Laboratory's Fault Database or Sandia National Laboratories' Lightweight Data Monitoring System (LDMS). Similar efforts are ongoing in other centers in Europe, such as BSC with their power/energy tracking in EAR, or in Japan, such as RIKEN as part of the Post-K efforts. However, all these efforts are currently still in their early steps and offer only insular solutions, providing the right timing for an international activity bringing them together.

Members of this consortium are in close contact with these centers and are already working with several of these sites on establishing common use cases, data formats, as well as the surrounding infrastructures. For example, members of this consortium have been participating at various special sessions led by the team in Sandia at large scale conferences such as SC and SIAM CSE, and have close ties to the development efforts at Lawrence Livermore and BSC. This will allow a) a leveraging of development efforts leading to open community code, b) an international homogenization of such a data store while these efforts worldwide are still in their forming stages, and c) a coordinated and strengthened impact on the vendor community to provide the necessary information on future systems.q

**Vorgesehene Mitglieder des Konsortiums (Co-Sprecherinnen/Co-Sprecher und die weiteren, beteiligten Institutionen):**

| Co-Sprecher/in | Zugehörige Institution |
| --- | --- |
| Martin Schulz Professor / Chair for Computer Architecture and Parallel Systems | Technical University of Munich Boltzmannstr. 3 85748 Garching bei München |

| schulzm@in.tum.de | |
|---|---|
| Thomas Ludwig<br>Professor and<br>Managing Director<br>ludwig@dkrz.de | Deutsches Klimarechenzentrum (DKRZ)<br>Bundesstraße 45a<br>20146 Hamburg |

Additional members of the planned consortium:

- Prof. Dr. Nico Gauger, TU Kaiserslautern
- Prof. Dr. Wolfgang Karl, KIT
- Prof. Dr. Dr. Thomas Lippert, Forschungszentrum Jülich
- Prof. Dr. Thomas Ludwig, DKRZ & Universität Hamburg
- Prof. Dr. Wolfgang Nagel, TU Dresden
- Prof. Dr. Christian Plessl, Universität Paderborn
- Prof. Dr. Michael Resch, HLRS Stuttgart
- Prof. Dr. Olaf Spinczyk, Universität Osnabrück
- Prof. Dr. Ramin Yahyapour, Gesellschaft für wissenschaftliche Datenverarbeitung mbH

External supporting participants

- Dr. Stefan Hachinger, Leibniz Supercomputing Centre (LRZ)