

## **NFDI Multi Cloud**

### **Description of the topic and possible resulting services:**

The key objective of all NFDI consortia is to systematically register and provide data in a sustainable manner such that it can be openly and unitarily accessed on national and even international level. Initially, each consortium approaches this goal independently under different disciplinary, technical, legal, and ethical conditions framed by their respective scientific community. Nonetheless, the overall goal of NFDI is to interconnect resources to enable a community-wide and cross-community access to data and applications in order to facilitate analysis, reuse, and exchange of data and establish connections. The foundation for such an overarching common infrastructure is built upon the concept of a federated Multi Cloud architecture that provides unified access by a federated Identity and Access Management (IAM) infrastructure to compute and data storage resources and allows seamless higher level integration of distributed heterogeneous services and data.

To date, it is common practice in many scientific disciplines for users to first download and locally install software repositories followed by transferring large amounts of data, e.g. from public databases. Cloud computing, in contrast, offers an economic and scalable solution by pooling compute and storage resources and providing a model in which public data is integrated or hosted by data providers so that users can perform their analyses close to where the data resides. Virtual compute environments allow, beyond that, for maximum flexibility in terms of software stacks, and portable containers enable scientists to share environments or workflows with colleagues, facilitating reproducible research. Consequently, the development of new research data infrastructures as well as the integration of services into existing cloud-based infrastructures is already a determined goal of many NFDI consortia. In order to bundle and consolidate these efforts, to include all consortia, and, above all, to promote compatibility and interoperability on multiple levels between them, a common platform is needed.

This platform is based on the ideas introduced by the Research Data Commons (RDC) concept initially supported by a number of NFDI consortia and the BMBF FAIR DS project<sup>1</sup>. By further adapting and extending the RDC according to the requirements of each consortium, a decentralized yet federated cloud computing platform, *i.e.* a Multi Cloud, allows consortia to either operate their own systems, that in turn are federated and form a cloud of clouds, or use a shared cloud system that integrates existing heterogeneous dedicated and opportunistic resources into a common federated science cloud. Thus,

---

<sup>1</sup>  Konzept "Common Infrastructures"

individual consortia can meet the specific disciplinary, ethical, and legal obligations of their communities and still offer interoperable resource access and exchange.

Following this architecture, involved NFDI centers provide the necessary capacity and expertise for (i) hosting the decentralized parts of this platform, (ii) storing the actual data and (iii) providing information about available services, computing and storage resources. Access to sensible or otherwise protected data can thus be tightly controlled and restricted to a single data center and implemented into a decentralized common IAM. Associated non critical metadata and generally non restricted data should be made publicly available in accordance to the FAIR principles either physically or virtually in the Multi Cloud storage, that interconnects and integrates the heterogeneous data landscape of the different data centers and other public resources. Data access for compute resources can be accelerated by using dynamic disk caches as provided by XCache or related implementations<sup>2</sup>.

This platform will also provide a set of basic software services for the development of scalable analysis procedures to optimally utilize the cloud resources. Such a service, *inter alia*, will also offer the required infrastructure to set up in-cloud SLURM clusters<sup>3</sup> (or HPC clusters operated by other local batch systems) as offered by the BiBiGRid<sup>4</sup> framework, deploy workflow engines like Nextflow or Snakemake, set up virtual compute environments, e.g. by using Ansible or provide a Kubernetes appliance for container orchestration and long running services. Via overlay batch systems as COBaID/TARDIS<sup>5</sup> compute resources of various centers can be combined to a single entity. Similar to the cloud SDKs from Amazon or Google, for example, these services should be accessible via a series of APIs for which language-specific client stubs are available for easy integration. In addition to a central command line interface (CLI) for essential operations, a user-friendly web frontend (WebUI) or even data portals as the PUNCH Science Data Platform foreseen by PUNCH4NFDI will be implemented for managing resources, deployed services as well as data and workflows.

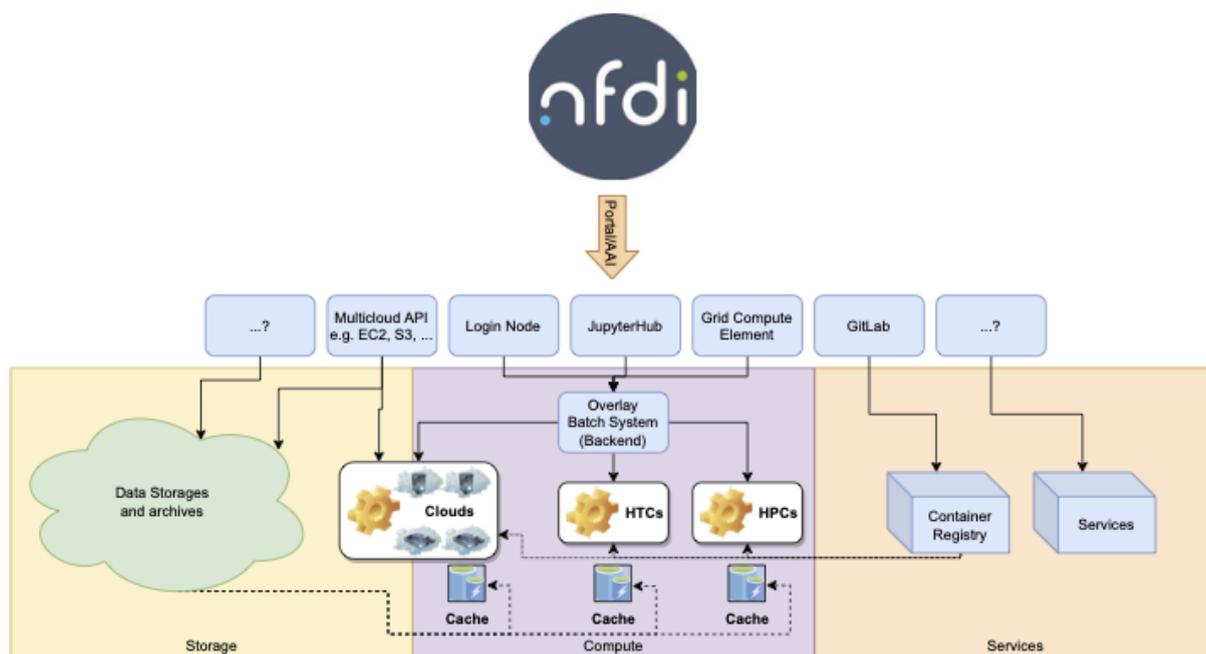
---

<sup>2</sup> <https://doi.org/10.1051/epjconf/201921404005>

<sup>3</sup> [https://en.wikipedia.org/wiki/Slurm\\_Workload\\_Manager](https://en.wikipedia.org/wiki/Slurm_Workload_Manager)

<sup>4</sup> <https://github.com/BiBiServ/bibigrid>

<sup>5</sup> [Interactive distributed computing for HEP on multi-managed cluster resources](#)



**Image 1:** Schematic overview of the interaction between the individual components of a Multi Cloud. The individual components can be deployed across multiple compute centers. The middleware components that interconnect the components can use topology information of the whole system to optimize data flow. User access can happen either via APIs, SDKs (that use the APIs) or websites depending on the audience and use-case.

**NFDI Section(s) and their Working Group(s) responsible for negotiation/development/evaluation of this topic:**

- Section Common Infrastructures, Working groups infra-Multi Cloud, Identity & Access Management, Data Integration

**Potential partners with existing expertise (list the potential partner institutions, their consorti(a) and their roles):**

NFDI consortium	Institution	multi cloud expertise
DataPLANT	Uni Tübingen	de.NBI Cloud
GHGA	EMBL Heidelberg, DKFZ, Uni Tübingen	de.NBI Cloud
NFDI4Biodiversity	Uni Bielefeld, Uni Giessen	de.NBI Cloud
NFDI4Microbiota	Uni Bielefeld, Uni Giessen	de.NBI Cloud
PUNCH4NFDI	KIT, Uni Bonn, LMU, DESY, FZJ, JGU Mainz	Overlay Batch System Compute4PUNCH
PUNCH4NFDI	GSI, KIT, Uni Freiburg	Dynamic Disk Caches
PUNCH4NFDI	AIP, GSI, KIT, DESY Zeuthen, TIB	Data Portal PUNCH SDP

## **Description of the needs addressed by this potential service on NFDI consortia:**

The goal of this topic is to design and develop an IT architecture that provides a common infrastructure for all NFDI consortia by integrating available systems into a federated and unified Multi Cloud environment. This platform will be based on the general ideas introduced by the Research Data Commons (RDC) concept that is already supported by a number of NFDI consortia and the BMBF funded FAIR DS project. It will be adapted and expanded according to the requirements of the individual consortia.

From a general user's point of view, the Multi Cloud architecture has to address the needs of different user profiles from a broad range of scientific disciplines covered by the NFDI consortia. Application users who are mainly interested in visualizing and analyzing various data sets in a user-friendly and highly interactive manner or in doing large scale batch processing on big data sets using their expert domain knowledge do not need or even do not want to understand the details of the technical implementation. But they will benefit most from a seamless access to corresponding services and federated and extremely diverse data, especially when highly integrative analyses are desired. A typical expert data analyst on the other hand might need more direct and much more flexible access to data streams and computing and storage resources via well-defined interfaces that also support the development of automated processing workflows. Therefore, our platform will offer basic services to easily build scalable analysis pipelines based on workflow engines and we will provide easy access to suitable HPC and HTC resources. Finally, software engineers require low level access to APIs and SDKs to facilitate the development of novel algorithms and software tools. These SDKs are language specific, automatically generated implementations of the Multi Cloud APIs to avoid redundant implementations of the API and therefore reduce implementation overhead for the individual programmer. This also includes support for the standardized and simplified deployment of software solutions that are needed as components for new developments (e.g. MongoDB or Elasticsearch/Opensearch for accelerating access patterns and data aggregation). In addition to these general use-cases domain-specific services can be integrated into the Multi Cloud to reduce programming overhead for these services by using parts of the common infrastructure. This would also improve findability of the individual services across NFDI boundaries.

In any case, the availability of a unified Multi Cloud architecture will significantly decrease the complexity for handling and using highly heterogeneous data across scientific disciplines and it bears a huge potential to open new research fields in the context of data integration. At the same time, we can massively reduce the required effort for the initial skill adaptation training of data analysts and developers (i.e. quite limited human resources) when they are switching projects once they are familiar with the Multi Cloud infrastructure.

## **State of the art for this potential service:**

- The federated de.NBI Cloud can serve as an extremely successful example of a Multi Cloud infrastructure including the already well-established concept for data storage and analysis federation as well as the elaborated but lean governance model.
- Established compute and storage providers, universities and research centres that are actively engaged in the PUNCH4NFDI consortium and others contribute many years of experience in the field of grid computing, such as the world's largest scientific computing grid WLCG (Worldwide LHC Computing Grid). In this context dynamic disk caches and overlay Batch systems as COBaID/TARDIS have been developed in order to include opportunistic and heterogeneous resources into existing production environments. This deep knowledge will be essential to build integrated computing solutions within our Multi Cloud environment.

## **Describe the overall strategy for the possible service with regard to the following stages:**

### **1.) Service initialisation strategy**

- a.) Identification of cloud sites and evaluation of their technical bases
- b.) Agree on a basic set of use cases which need to be addressed by the Multi Cloud implementation
- c.) Identification of a small basic set of required services for a prototype
- d.) Agree on a set of technical standards for the implementation
- e.) Agree on monitoring, accounting and logging infrastructure
- f.) Agree on an extendable architecture for the Multi Cloud to include services from external contributors
- g.) Agree on the storage architecture and its technical foundation
- h.) Agree on a governance infrastructure
- i.) Gap analysis: which of the required services can be addressed by existing solutions and where is development work required

### **2.) Service integration strategy**

- a.) Implement storage system
- b.) Implement the identified services for a prototype
- c.) Create guidelines for long-term sustainability for internal and external services
  - i.) Quality standards
  - ii.) Support
  - iii.) Security
- d.) Define roadmap towards production environment

### **3.) Ramp-up strategy for service operation**

- a.) Start production environment based on the prototypes developed in phase 2
- b.) Identify additional use cases which need to be addressed before implementing the full production environments
- c.) Increase the number of services; Services will be implemented according to community needs and technical feasibility
- d.) Extend the number of cloud sites

**Address possible challenges and risks:**

- Very heterogeneous landscape from VM and container based cloud environments with S3 object-storage to HPC and HTC based providers with shared file systems need to be integrated
- AAI solution from IAM group required
- Large variety of potential services with only small resources requires focus
- Users need to adapt the new programming paradigms

**List other NFDI basic or subject specific topics or services with which this service will interact:**

- IAM: Will be used for Authentication & Authorization across all services
- DI: Uses the services of the Multi Cloud to implement more specific data integration services