# Open Data for DESY and HIFIS

## A portal bundle for DESY, HIFIS, NFDI and their pilot node in EOSC Beyond

Tim Wetzel, Patrick Fuhrmann, Uwe Jandt, Paul Millar, Sophie Servan, Franz Rhee, Peter van der Reest, Regina Hinzmann, Noel Barth, Johannes Reppin, Christian Voss, Linus Pithan, Anton Barty, ...

EGI2024, Lecce, 03rd October 2024

In cooperation with

# Open and FAIR data for Photon Science

The motivation for a prototype system

**FAIR data will become the standard**

- Funding bodies and journals demand data to be open and/or FAIR

    – Public money = public data (after embargo period)

    – Supplemental data for publications
- Combatting the reproducibility crisis in science

- Reusability makes for a more sustainable (re-)use of results obtained from costly and laboriuos experiments and enables AI/ML workflows

**Starting with Photon Science**

- As one of the largest photon science laboratories in Europe, DESY will start providing a standardized way to host Open and FAIR data for her scientists

**Towards a blueprint for HIFIS, NFDI and the community**

- After successful initial operations with DESY photon science, the portal will be opened as a HIFIS service

- We also hope to create a blueprint for OpenData portals in the community that will be shared openly

# DESY Photon Science setup
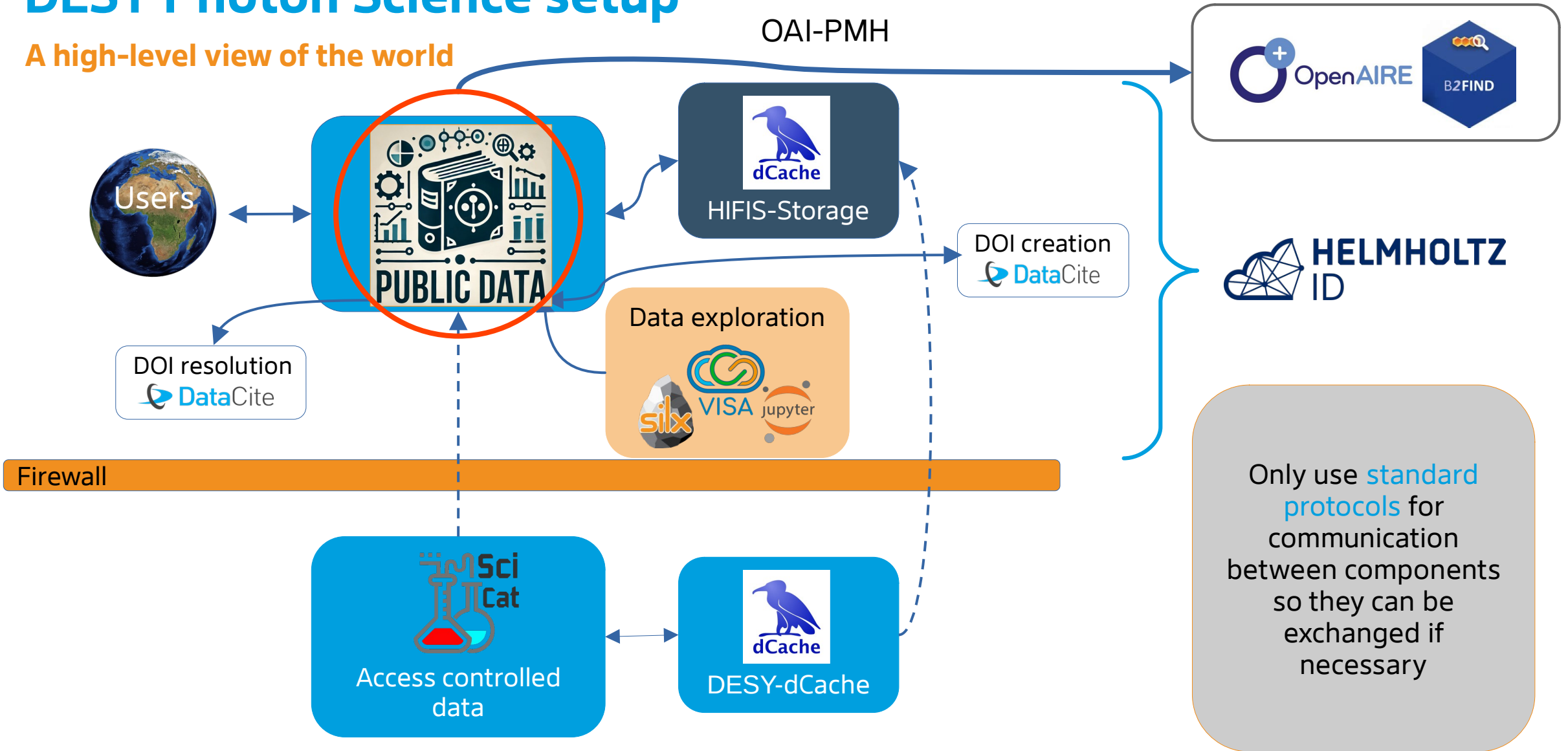
**A high-level view of the world**



*Image adapted from Anton Barty's slide*

# The minimum viable system for DESY.

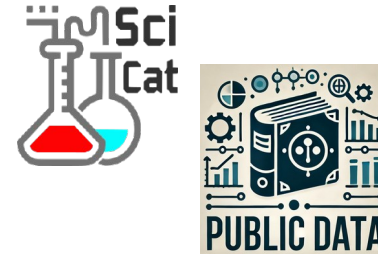## Essential components with federated access (authenticated & non-authenticated)

**Long term storage** (**dCache** via hifis-storage.desy.de**)**

- accessible via standard protocols (https, NFS, WebDAV)

**Metadata Catalogue** with

- mandatory core metadata fields

- optional domain specific metadata fields

- OAI-PMH protocol for data harvesting of core metadata by high level catalogues

**DOI Minting Service**

- In cooperation with our library

**Open Science** (**Virtual Research**) **infrastructure**

- VISA portal, currently working on it together with other synchrotron facilities in Europe under an MoU
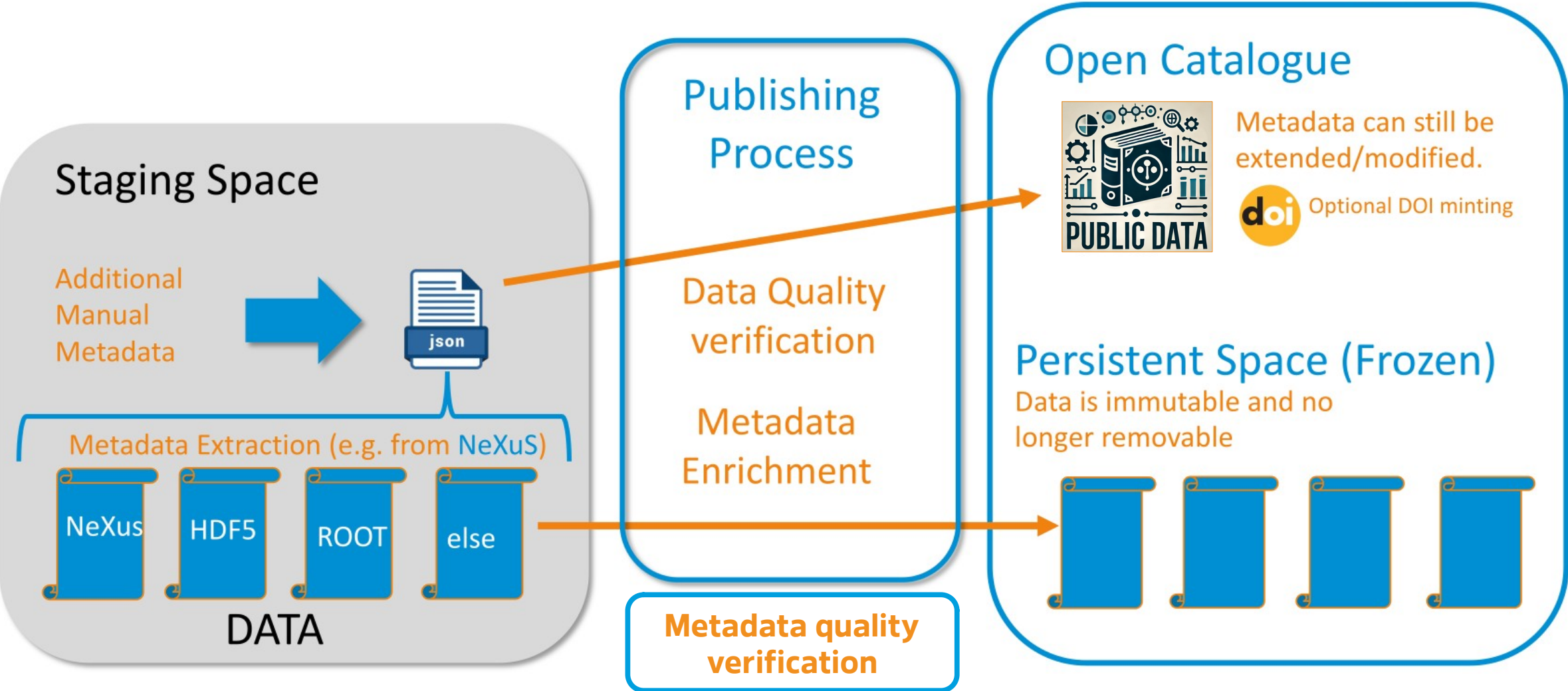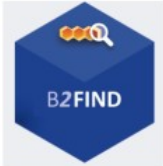
1st phase

2nd phase

# Our envisioned data ingestion process

**For Open Data**

# Our envisioned data ingestion process

## (Meta)data ingestion and quality verification – sisyphos.desy.de

- First installment with LinkML and Streamlit
  -
    - Metadata schema description via YAML documents setting standards that metadata has to conform to
    - Data description in terms of "classes" and "slots", allowing inheritance and mixins for creating custom types
    - 60+ different open-source tools to work with schemata for introspection, validation, format conversion, …

- Starting for the X-Ray reflectivity community within DAPHNE4NFDI

- If you are interested in details:
    - https://gitlab.desy.de/ric/opendata-metadata/
    - sisyphos.desy.de
    - Let me know so I can get you into contact with my colleagues



**PaN Reflectivity Database - Upload tool**

Welcome to our open data community. This tool will guide you through uploading your x-ray and neutron reflectometry data to the PaN Reflectivity Database. Please enter the metadata for your data set below and then upload your reflectometry curve. By submitting your data, you agree to make the data available in accordance with the Creative Commons Attribution (CC-BY) Licence. After submission, the data and metadata are written to the ORSO file format. Following curation, the data set will be published in the DESY public data catalogue.

**Administrative Data** 🔗

owner*

contactEmail*

datasetName*

principalInvestigator*

creationLocation*

**Experiment**

title*

instrument*

# hifis-storage.desy.de

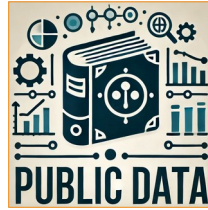## The "drop box" and and final storage space for Open Data



**Write access granted by Helmholtz VO membership.**

# public-data.desy.de

**The metadata catalog!**

PUBLIC DATA

## Search

Clear

PID

Text Search

Location

Group

Type

Keywords

Start Date — End Date

+ Add Condition

### Name | Source Folder

Reflectometry curves (XRR and NR) and corresponding fits for machine learning | ...do.6497438

spain | .../nfs

### General Information

| | |
|---|---|
| Name | Reflectometry curves (XRR and NR) and corresponding fits for machine learning |
| Description | This is a compiled dataset of raw X-ray reflectivity (XRR, reflectometry) measurements together with corresponding fit parameters, intentionally published to use as training or test data for machine learning models. (The authors aim to include NR data in further versions of this dataset and plan to include other substrates and materials for XRR. Contributions welcome!) |
| PID | undefined/10242df2-3868-42cb-bcb2-81c2c44533ec |
| Type | raw |
| Creation Time | 2024-01-25 18:34 |
| Keywords | |

### Creator Information

| | |
|---|---|
| Owner | Linus Pithan |
| Principal Investigator | linus.pithan@desy.de |
| Contact Email | linus.pithan@desy.de |
| Owner Group | fsec |
| Access Groups | |

### File Information

| | |
|---|---|
| Source Folder | /desy/public-data/upload/daphne4nfdi/10.5281_zenodo.6497438 |

### Scientific Metadata

Search

▼ DIP_1

| | |
|---|---|
| Experimentalists | Kowark, Stefan |
| Layer_CAS | 188-94-3 |
| Layer_formula | C32H16 |
| Layer_material | Diindenoperylene |
| Substrate_temperature | 303 (K) |
| instrument | ESRF, ID10b |
| q_max_fit | 0.15 (1/Ang) |
| year_experiment | 2005 |

▶ DIP_2

### Path | Size

| Path | Size |
|---|---|
| calc_xrr.py | 2 KB |
| conda_env.yml | 7 KB |
| prepare_plot.py | 4 KB |
| README.html | 6 MB |
| README.ipynb | 9 MB |
| requirements.txt | 76 B |
| xrr_dataset.h5 | 254 KB |

# visa.desy.de

## Select a dataset to spawn a virtual machine



VISA database currently populated with example datasets.

Open Data to be integrated during 2024 via automated data export from public-data.desy.de

# visa.desy.de

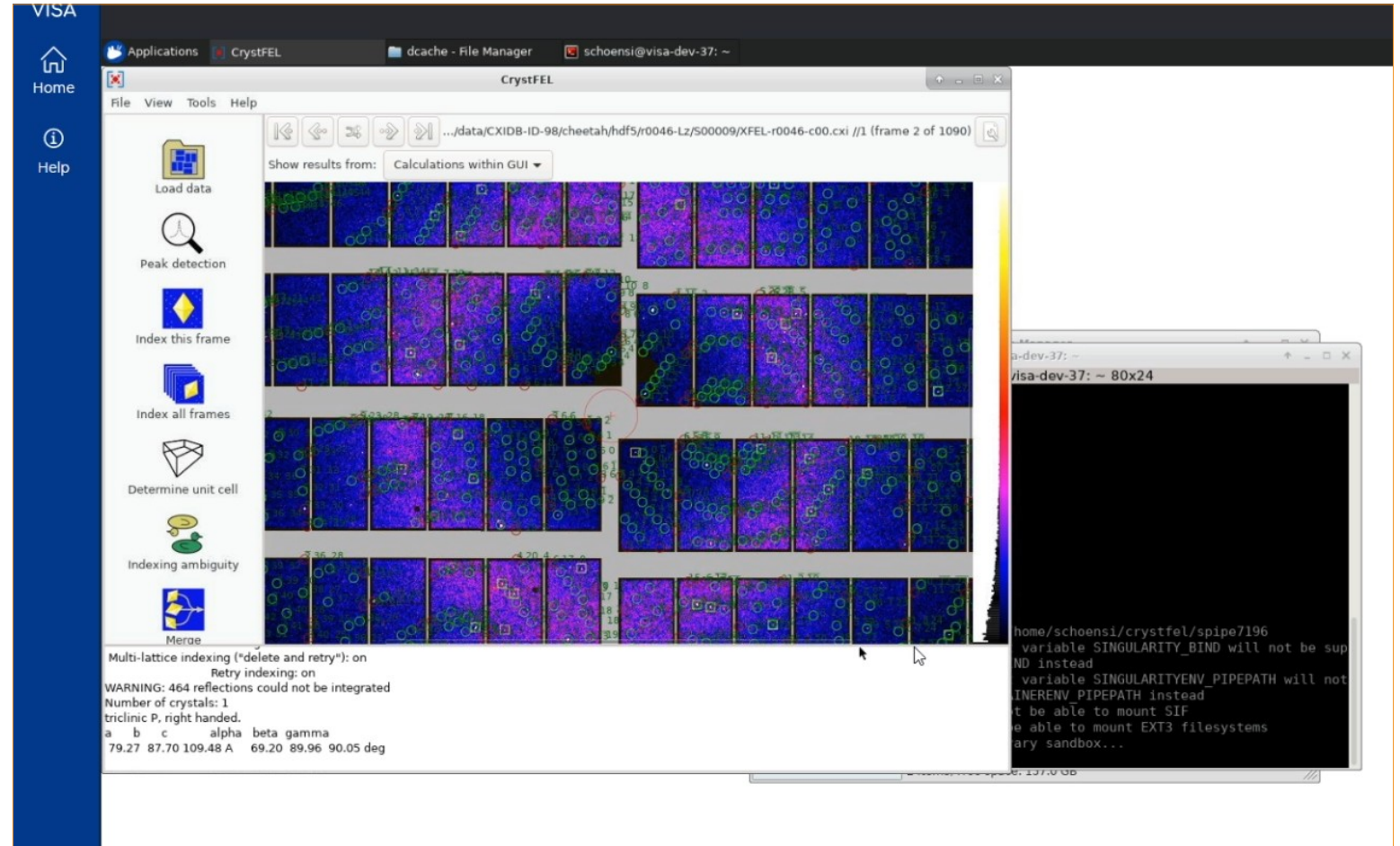## Work via remote desktop connection with graphical interfaces

Using CrystFEL Docker Images to run Singularity Container and work with Crystfel 10 Graphical Interface.



*Example provided by Silvan Schön (DESY/FS-SC)*

CrystFEL

# Thank you!

## Questions?

**Contact**

**DESY.** Deutsches
Elektronen-Synchrotron

www.desy.de

Tim Wetzel, Patrick Fuhrmann
IT-RIC (Research & Innovation in Scientific Computing)
tim.wetzel@desy.de, patrick.fuhrmann@desy.de

# Backup slides

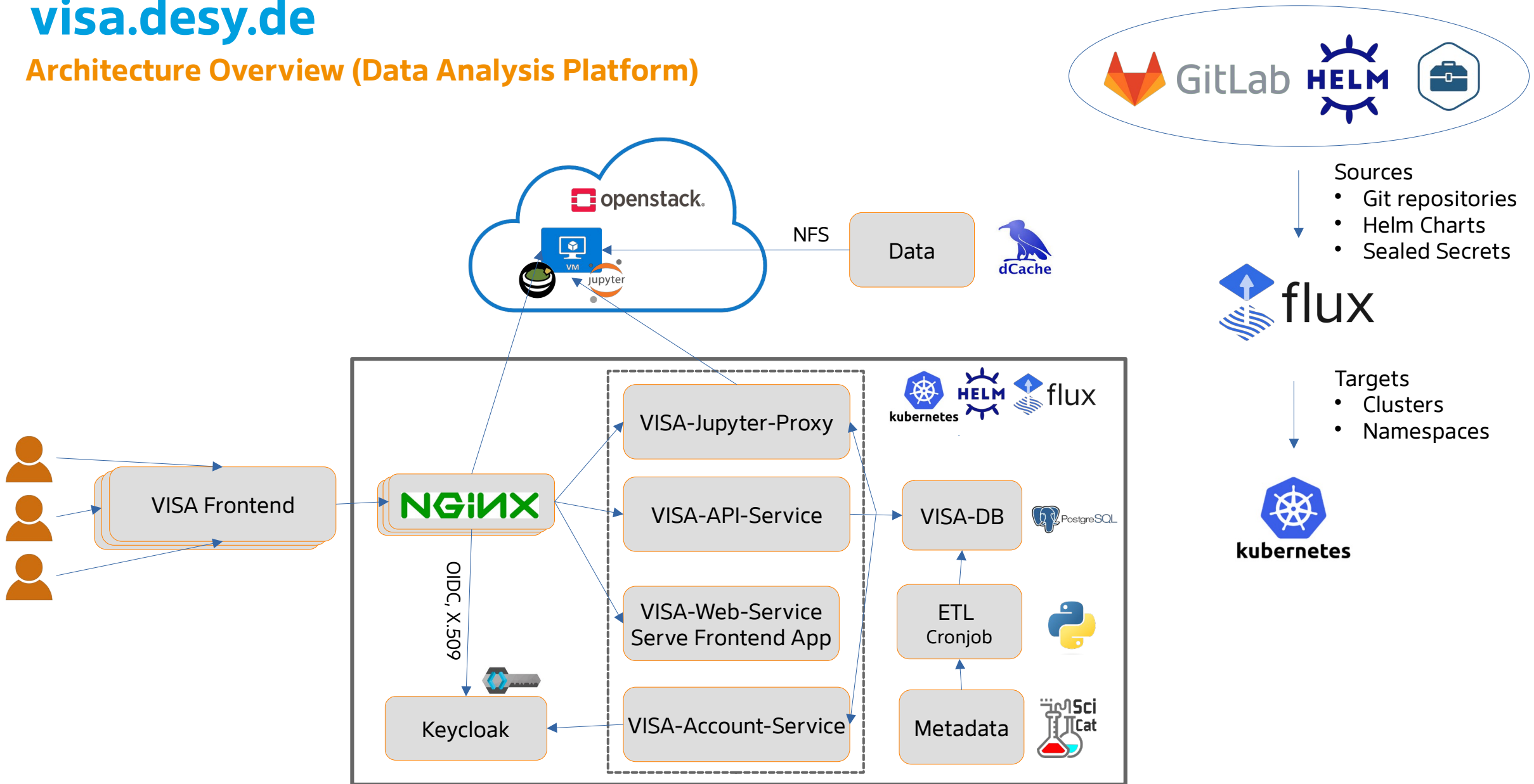# Importance of proper metadata definitions

## Consensus and standards are key

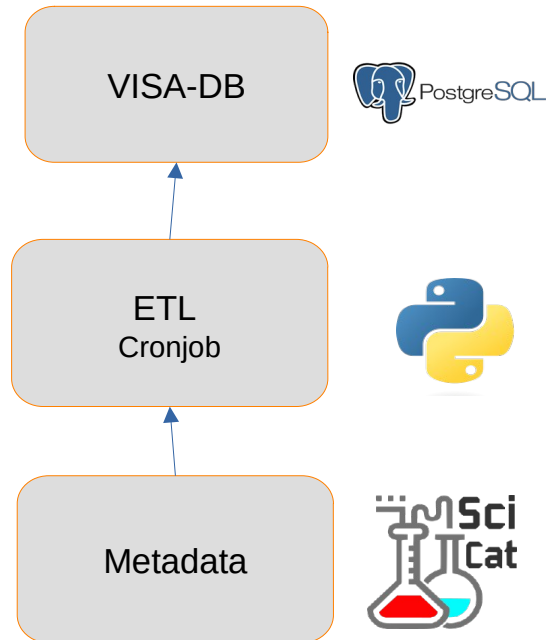| Mandatory **core metadata** fields | Defined in prior activities and by responsible reference bodies<br>e.g. DublinCore, DataCite v4.4 |
|---|---|
| Optional **domain-specific metadata** fields | To be provided by the community<br>e.g. former PaNOSC/ExPaNDS, Daphne4NFDI, Photon Science Community |
| **Additional metadata** fields | Experiment/Beamline/Facility-specific metadata might be needed |

Special challenge for open data:

**Heterogeneous origin** of data sets from different experiments with different specific metadata need to be mapped into the same catalogue
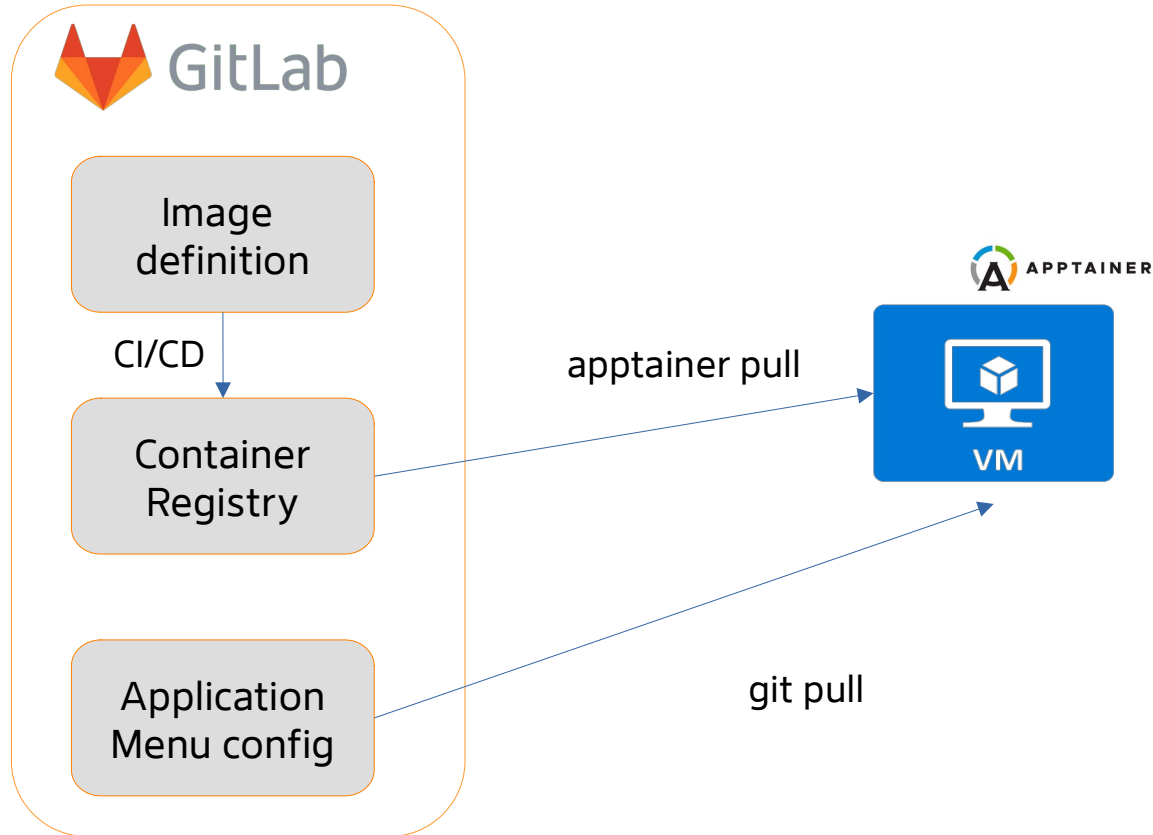
➡️ Metadata **input and verification** need to be handled properly in the publication process

# visa.desy.de
## Architecture Overview (Data Analysis Platform)

# visa.desy.de

## Metadata import via custom ETL process

VISA-DB

ETL
Cronjob

Metadata

- Python script
- Customizable depending on the metadata source (catalogue API format, authN/Z, ...)
- Can be run once for static data or as a cronjob for dynamic data
- Event-based execution would be nice to have (e.g. webhooks)

- Metadata import
  - Experimental specifications
  - Dataset status (embargoed or public)
  - User access rights
  - Storage paths

- Database backup

# visa.desy.de
## Analysis software provisioning via Apptainer images



- Software in Apptainer images
  - Many applications already available as Apptainer image from HPC workflows
- Built from .def file in CI/CD pipeline
- Image publicly available in Gitlab registry
- Pulled on application startup

- Application menu entries defined separately in git repository
- Seamless integration into the OS applications
- Menu entries updated from menu config by cronjob pulls the repository regularly
- Seamless updates to the menu by admins