# Infrastructure for Research Data Management at HZB
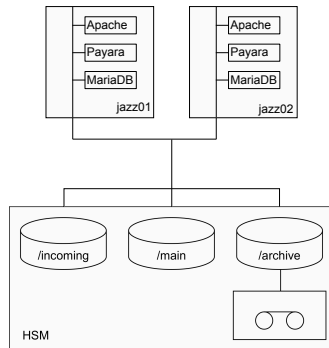
Rolf Krahl
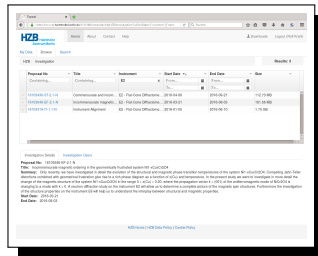
Workshop on Research Data Management, HZB, 11th June 2019

# HZB Data Policy

- HZB Data Policy regulates management of scientific data from public research at HZB's large-scale facilities.
- Distinguish raw data, results, and metadata.
- Raw data and associated metadata will be curated and stored by HZB for at least ten years.
- Raw data and associated metadata are placed in the public domain (Creative Commons CC0 Dedication).
- Access to raw data and associated metadata is restricted to their creators for an embargo period of five years. After that, they become openly accessible.
- Results and associated metadata may be stored with the raw data. They will not be curated by HZB. They may be made openly accessible upon request of their creators.
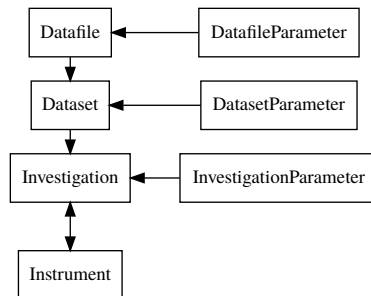
# Server Systems

- Storage systems: hierarchical storage management (HSM), the `/archive` area consists of disks and tape libraries.

- ⇒ Most data will reside on tapes most of the time.

- Two dedicated servers in active/passive configuration.

- Data volume: we calculate 2 PB/Jahr.

# ICAT Metadata Catalogue

- Access to the data is provided by the ICAT metadata catalogue.
- Search for the user's own and for public data.
- A request to download data automatically triggers the staging of that data from tape to disk. The data may be downloaded after the staging is complete.
- ICAT is developed as free software in cooperation with other Photon and Neutron sources (STFC, DLS, ISIS, ESRF, HZB).

# ICAT Schema

- Central elements for organizing the data in ICAT are:
  Investigation ← Dataset ← Datafile.
- Correspondence:
  - Investigation ≙ Proposal,
  - Dataset ≙ Measurement,
  - Datafile ≙ File.
- DatasetParameter allows storing physical metadata of the measurement in a simple keyword/value schema.

```
Datafile      ◄──  DatafileParameter

Dataset       ◄──  DatasetParameter

Investigation ◄──  InvestigationParameter

Instrument
```

# Metadata

We distinguish two classes of metadata:

## Administrative metadata

- Proposal, title, abstract, user, access rights.
- Get imported from the user office portal GATE beforehand. Analogous arrangements will be made for measurements not related to a GATE proposal.
- Relevant for the control of internal ICAT workflows.

## Physical metadata

- Parameter of the measurement, sample etc. See other talk.
- Will be collected (preferably) automatically and stored in the datafiles before ingestion into ICAT.
- A selection of the metadata may be additionally stored in ICAT as `DatasetParameter`, if they are relevant for the search of data.

# Data Life Cycle

Steps in the life cycle for scientific data at HZB:

1. User submits a proposal in GATE.
2. Proposal gets accepted.
3. Administrative metadata get imported from GATE into ICAT. (For data not related to a proposal, create an `Investigation` in ICAT instead.)
4. User comes to HZB and performs experiments. Data is collected and curated at the instrument.
5. Data are ingested from the instrument into ICAT.
6. User have exclusive access to their data.
7. User may optionally upload results from data anaylsis into ICAT.
8. After expiring of the embargo period, the raw data and associated metadata become openly accessible. The user may decide to make also the results openly accessible.

# Current Status & Implementation

Current status:

- The storage systems are available.
- ICAT operates in test production. Mostly ready for production.
- Two BER II instruments (E2 and E9) register routinely their data.
- Implementation at the first BESSY II station (NanoclusterTrap) is in preparation and will begin shortly.

Implementation:

- For each instrument, we need to hook up (the data curation and) the data ingestion to ICAT into the measurement workflow.
- This requires consideration of the prerequisites and the workflows at the instrument individually.
- We will proceed intrument by instrument.

Issues:

- We lack a proper identity management at HZB.