# Facility Update DESY

## DESY specifics in metadata management

Regina Hinzmann, IT-InFa, for DESY.
Input from IT-RIC, IT-InFa, Library, FS-EC, FS-SC
2025-06-24

**SciCatCON 2025**
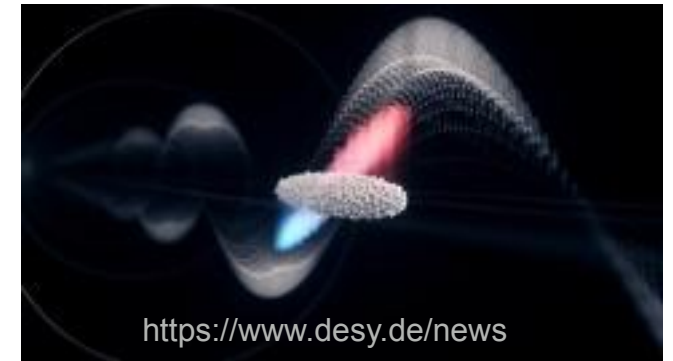**DTU Denmark**

# Metadata management at DESY FS

## General

DESY serves with its two light sources PETRA III and FLASH a broad scientific community collaborating in many national and international projects and and industrial cooperations.

- More than 3000 researchers visit DESY yearly.

- It operates 23+2 PETRA* III beamlines and 4 for FLASH**.

  - FLASH has a DAQ system to be able to record each high frequency photon pulse while at PETRA beamlines no direct DAQ is needed ⇒ different metadata capturing. Vital for FLASH, as they cannot use their data w/ metadata.

  - PETRA III data : on tape (dCache): a bit more than 27 PB (incl. 1 copy); on disk: 13 PB

  - FLASH data on dCache about 0.4 PB.

- Accelerator research and PETRA IV project are in front of DESY's door.

  (100x higher brilliance)

*Laser-plasma accelerated electron bunch show record beam quality.*



https://www.desy.de/news

# Metadata management at DESY FS

## Brief history and status

When I joined DESY-IT in August 2023, SciCat was already 2 years around (migration to new backend).

- Main setups at dedicated **demonstrator beamlines**: one at a PETRA and one at a FLASH beamline.

  - At FLASH they used an inhouse metadata system which ceased when the person left, they integrated SciCat as core element - which was not sustainable.

  - At P08 they dedicated beamline staff to write an ingestor and actively persue until today metadata ingestion into SciCat.

Today, 2 years later

- FLASH had contributed to improvments of SciCat but decided to wait for any metadata catalogue will be provided by DESY. Until then, they are happy to "survive beamtimes".

- More PETRA beamlines (P08, P05 and P10) work now with SciCat.

**(Immediate) usefulness of SciCat@DESY**

- Primary: Catalogue functionality (search and find, filter, select), also very important: issue DOIs.
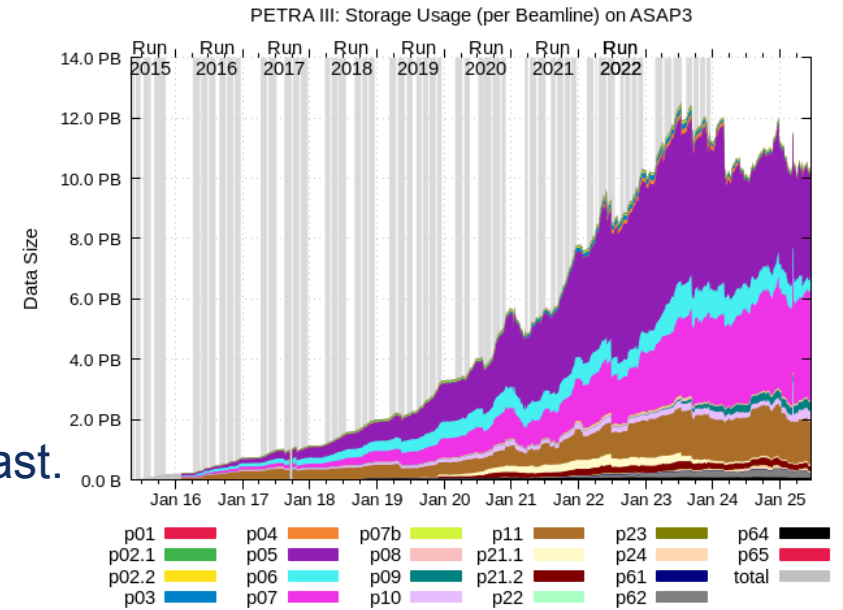
# Usefulness of SciCat @DESY

## What are our questions to answer?

**Increased usefulness of SciCat @ DESY**

➢ Service for systematic FS data access.

*Which data* should be accessible?

- Operation of PETRA has increasingly produced more data over the past.

*Who* should access the data?

- Depends on *when?* Data under embargo? Hard to say w/o data policy. Clear: it's *not immediate*, but what would be the workflow once the embargo period ends?

- Journals require already *now* a DOI; *reviewer* may need to access it, *other scientists* as well when paper is published they'd like to access the data of that specific DOI.



PETRA III: Storage Usage (per Beamline) on ASAP3

# SciCat: The most promising tool for DESY?

## One goal, two SciCat instances.

Address with two setups of SciCat that run at different timescales for public data.

- IT-RIC deals with *open data*:

    **public-data.desy.de**

    provides access to actual data in HiFis dCache

- IT-InFa deals with *metadata under embargo*:

    **scicat.desy.de**

    provide catalogue functionality, but also pathway          to data publication, first manually, but eventually     automatically?

**Workflow from embargo to publicly accessible data – independent of past, current and future data taking?**

---

public-data.desy.de
no login required

public-doi.desy.de
detailed landing page with link to DataCite

---

scicat.desy.de
can only see either published or one own's private dataset records
currently one DB for all beamlines ... (P08, P05, FLASH have 114,000)
doi.desy.de
should only see registered records (with DOI).

---

# SciCat for DESY's public-data portal

**Beyond DESY scope.**

**Advanced pilot features**

- **Manual DOI minting** is set up, as final automated DOI minting as agreed with L is still missing.

  - Uses already productive DataCite prefix.

- Extended landing page by information of locality of data by which enables download by DOI. Successful tests with **download tools** aria2 and DataHugger.

**Emphasis on data curation**

- Advancement planned of ingest form generator sisyphos (community driven GUI for PaN reflectometry).

- Plan to improve data curation and more user friendly functionalities.



public-data.desy.de

20 curated datasets

no login required

public-doi.desy.de

detailed landing page with DataCite entry

# Usefulness of SciCat @ DESY

## DESY FS scope: What are our questions to answer?

**Data taking time axis and correlation with storage: systems GPFS (fast access, limited disk) and tape.**

$t_0$:  Moment beamtime is started, fix ressources on GPFS: proposal metadata, ownership, access rights, etc.

$t_1$:  Moment beamtime is stopped. Data is on GPFS until *dwell timer* expires and not activly accessed ($t_2$).

$t_{1. archive}$: First copy to tape is triggered and data is removed from GPFS.

$t_{freeze}$: Only a filelist of all data in that beamtime is left in directory on GPFS.

$t_3$:  User requests re-stage from tape to GPFS.

$t_4$:  User publishes and a DOI is minted for the data. Should data be accessible? To whom, for how long?

# Systems setup of SciCat @ DESY

## scicat.desy.de

- Run (both) SciCat in Kubernetes cluster in using argoCD.

- Own Helm Chart, published to OCI registry oci://tollerort.desy.de/scicat/scicat

- Use official images for prod. Key facts for prod: 3 APIs, 3 MongoDB, 3 traefik ingress controller, 1 frontend.

- Plan more detailed performance studies, see Keynote talk by Igor K.

- Re-design to have separation of embargo-ed and publicly accessible data.

# Outlook

## Where will we be in 2 or 5 years?

- Continue with the strategy demonstrator beamlines: gain happy users by providing nice search functionality.

- Work out how to handle past FS data. Investigate "1 BT = 1DS = 1 DOI"?

- Work out how to handle current FS data.

- By then we can have sufficient experience to handle future FS data.

- Thanks to DAPHNE project many important issues could be and are being addressed (proposal hierarchy).

- Dream or FS-IT goal: Provide functional, reliable and useful catalogue and support for DESY FS users beyond project funding.

? How do other labs handle data publication and data access?

# Thank you
## *for your attention* ❤️

# List of issues

Catalogue functionality:

- Elastic Search

- Scientific metadata search

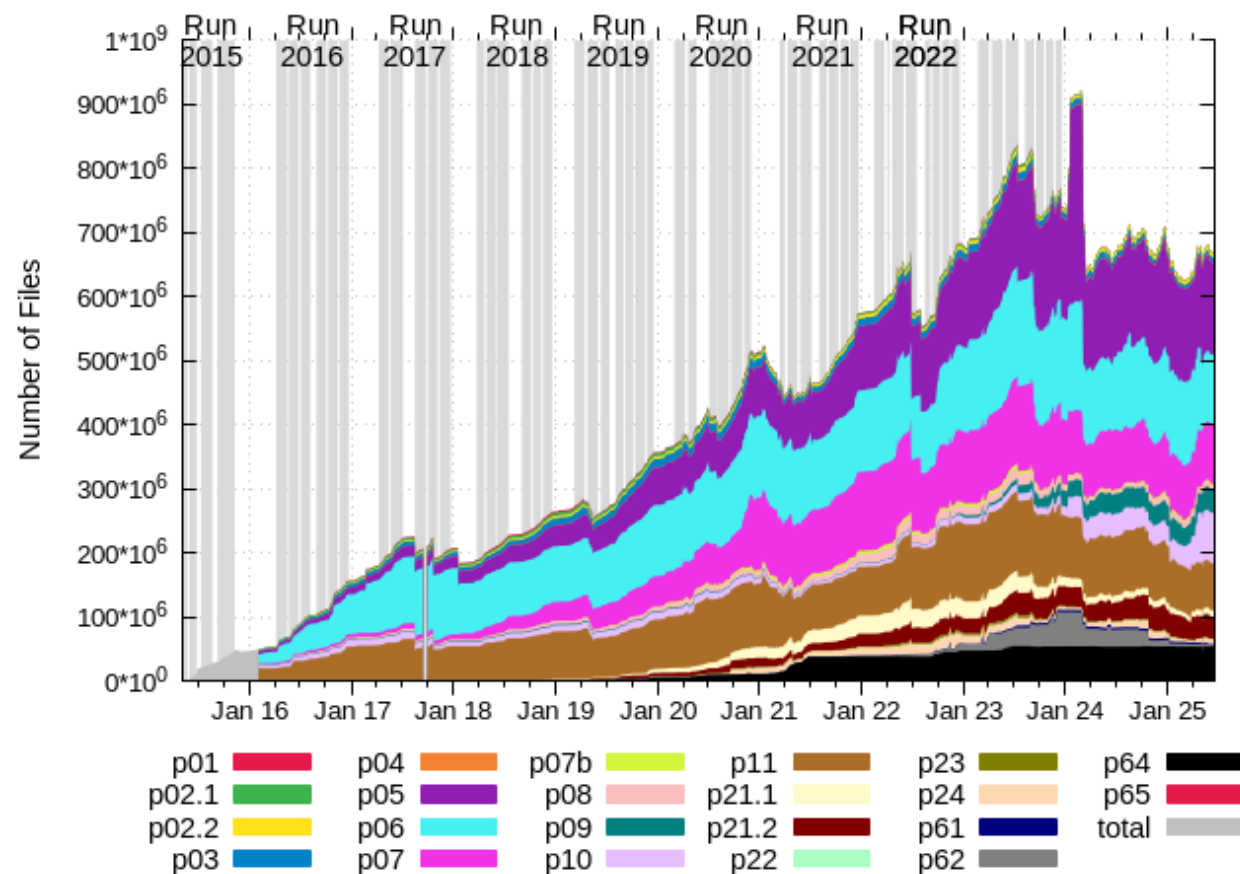# Data Taking Rate

## PETRA data on GPFS

# Data Taking Rate

## FLASH data on GPFS

# Data on Tape (dCache) incl. 1 copy

## Data on tape per beamline



Storage consumption in tapesize (per Beamline)

P05: ~10 PB

PETRA: ca. 27 PB
FLASH: ca. 0.4 PB