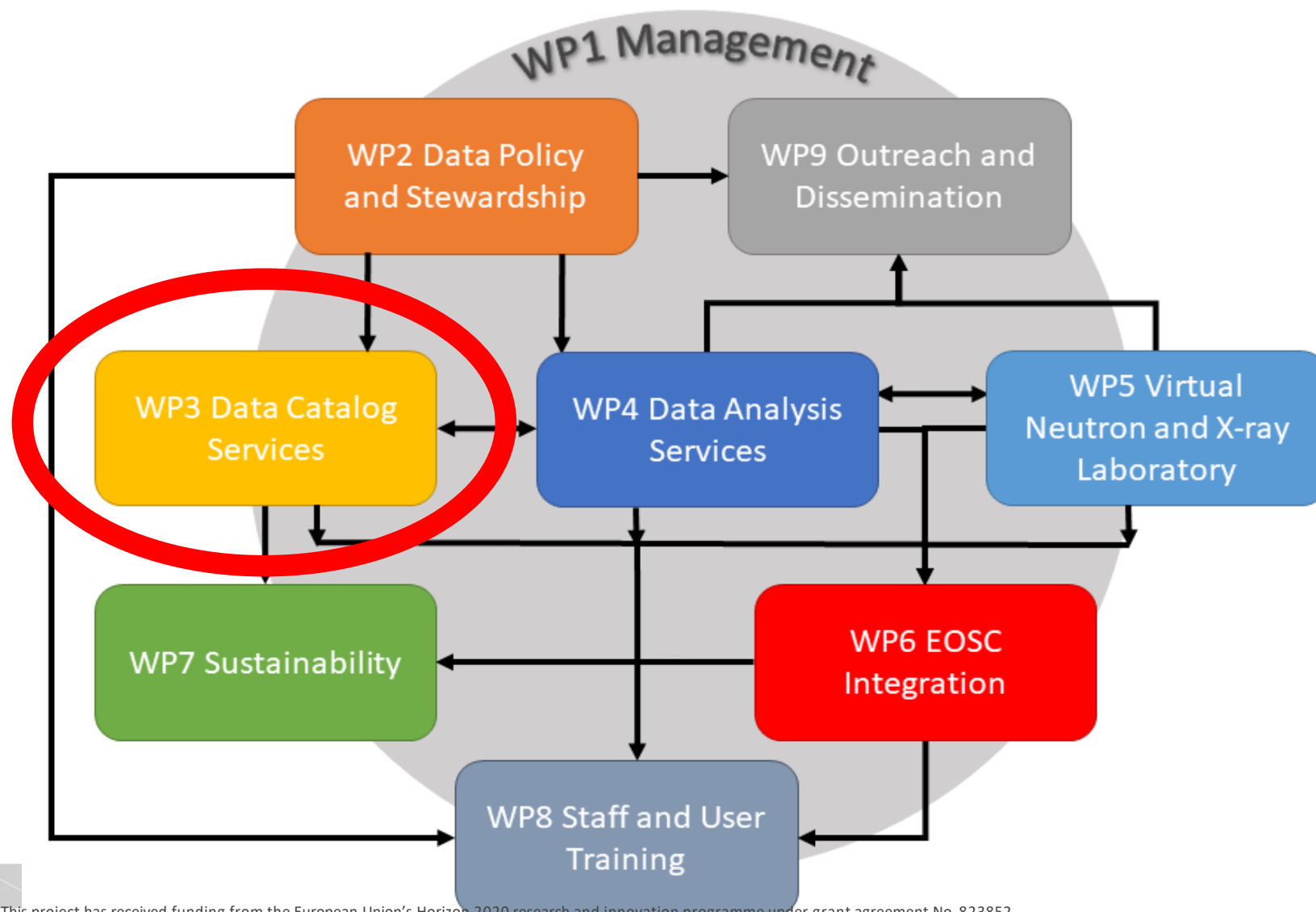# Review of Data Catalog Services (WP3)

**12 May 2021**

**Authors: Andy Götz (coordinator)+ Tobias Richter (WP3 leader)**

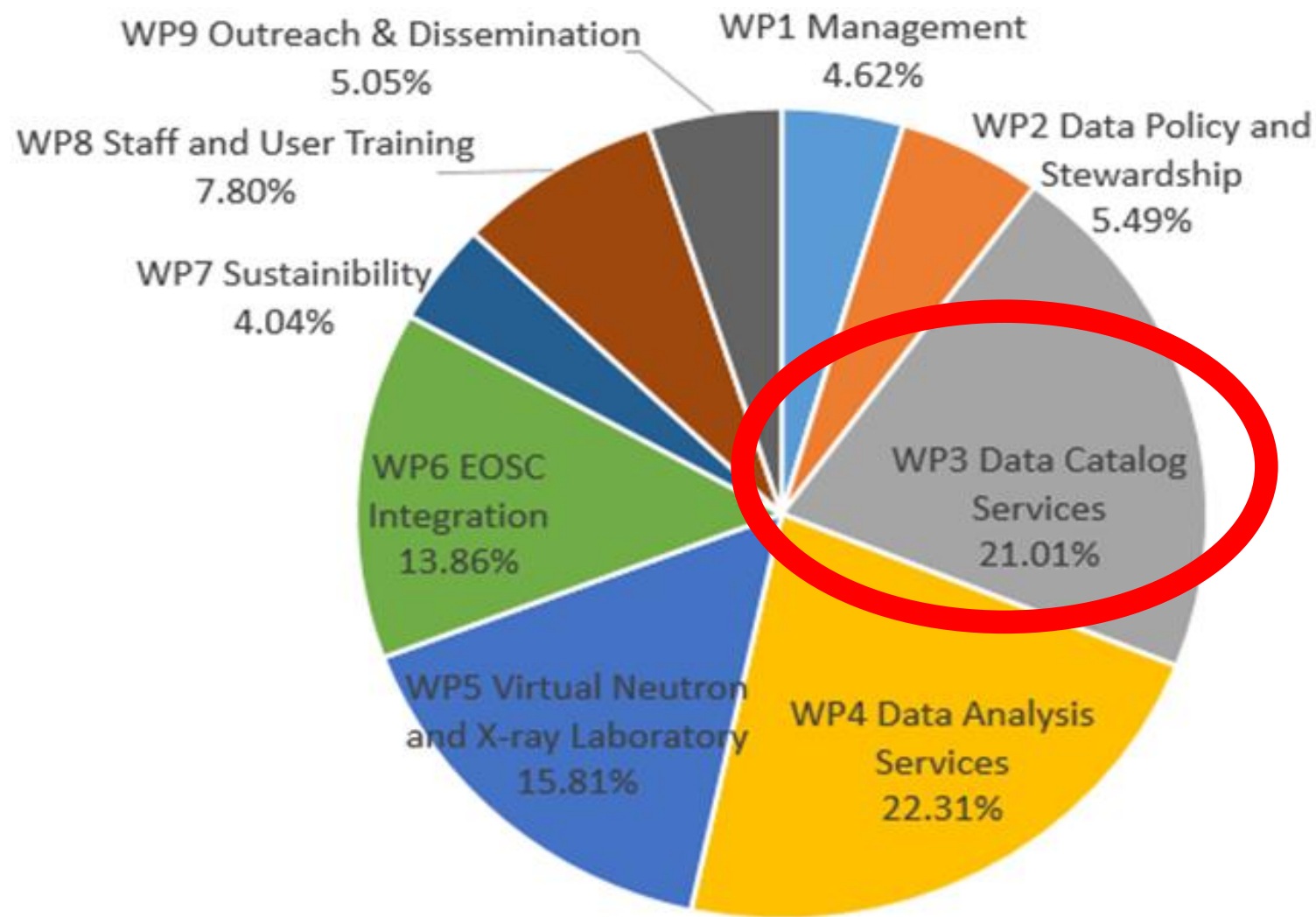**Place: PaNOSC Project Management Committee zoom meeting**
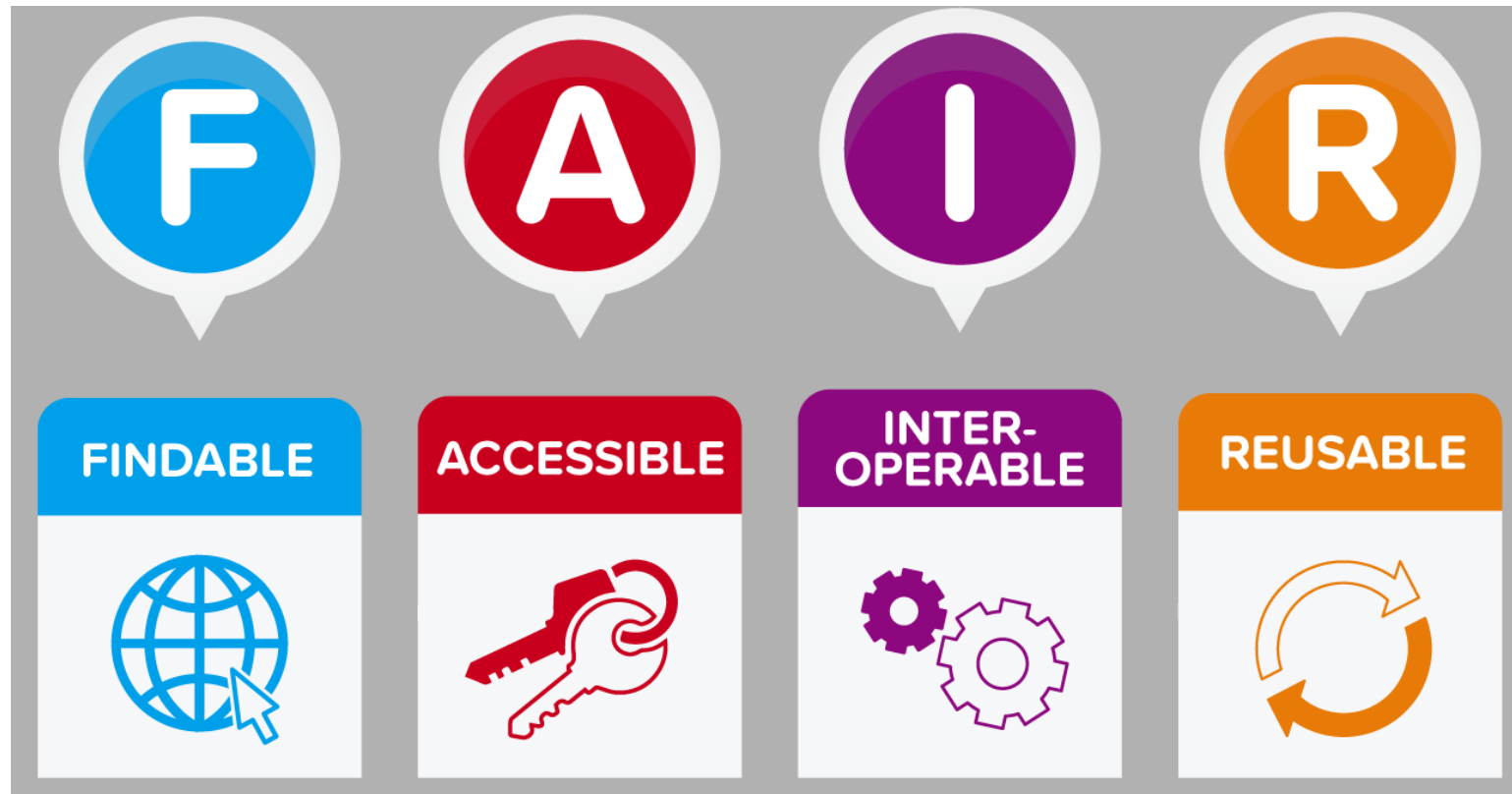
# Data Catalog Services – their role in PaNOSC

# Data Catalog Services – 2nd largest WP of PaNOSC

3

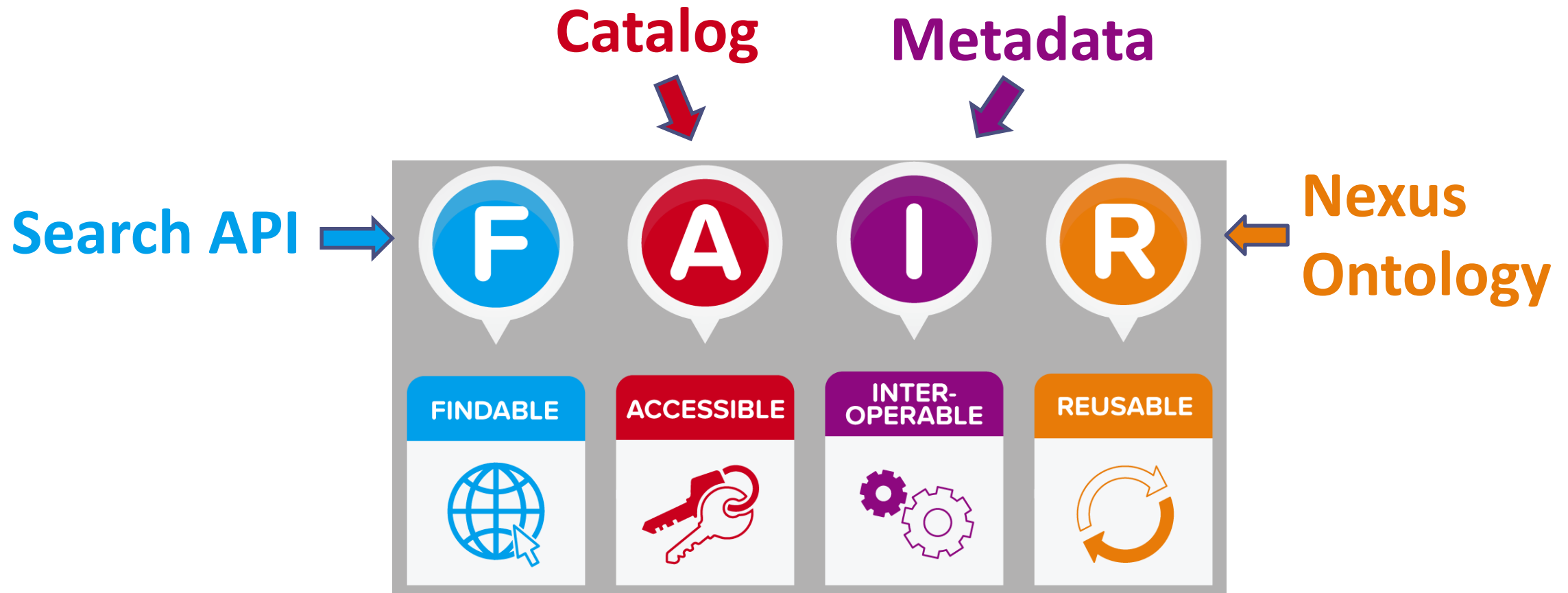# Data Catalog Services – essential for all FOUR principles of FAIR data

# WP3 present in ALL 4 FAIR principles

Catalog

Metadata

Search API

Nexus Ontology

FINDABLE

ACCESSIBLE

INTER-OPERABLE

REUSABLE

panosc

# Data Analysis Services – Effort per partner

1 - CERIC-ERIC         **88.00**

2 - ELI-DC         **78.00**

3 - ESS         **43.00** (WP leader)

4 - European XFEL         **36.00**

5 - ESRF         **25.00**

6 - ILL         **21.00**

**Total 291.00 (~ 24 years)**

| Partner | CERIC | ESS | ELI | ESRF | ILL | XFEL |
|---|---|---|---|---|---|---|
| Catalogue | VUO (online storage NOT a catalogue) | SciCat | TBD | ICAT | ILL Own | myMdC |
| URL | https://vuo.elettra.trieste.it | https://scicat.esss.se | --- | https://datahub.esrf.fr | https://data.ill.eu | https://in.xfel.eu/metadata |
| Login required | Yes | Yes | --- | Yes | Yes | Yes |
| File formats | NeXus, HDF5, ASCII and many others | NeXus | --- | EDF, SPEC, MCA, CBF, CCD, MCCD, HDF5, NeXus | NeXus and ILL Ascii | HDF5 |
| Database | Oracle | MongoDB | --- | Oracle and MongoDB | Oracle | MySQL and PostgreSQL |
| Language | Plsql, Python | Javascript | --- | JAVA and Javascript | PHP | App: Ruby(onRails), Client: Python |
| Main technologies | WebDAV, Guacamole | Angular | --- | React, NodeJS, EJB, JPA | Symfony, JQuery | Rails |
| Number of public datasets/files | 0/0 | 181/250,000 | --- | ~540K/157M | ~250K/4M | 0/0 |
| Using OAI-PMH | No | Not yet installed | --- | No | No | No |
| Minting DOIs | Yes | Yes | --- | Yes | Yes | Yes |
| Data/embargo policy | Not defined | Embargoed for 3 years | --- | Embargoed for 3 years, ESRF Data Policy | Embargoed for 3 to 5 years, ILL Data Policy | Embargoed for 3 with possible extension to 5 years, XFEL Data Policy |
| Number of instruments connected to data catalogue | None | 1 | --- | 17 | 54 | 16 |

# Data catalogue survey

## WP3

panosc

# Data catalogue progress

## Fair Data API development

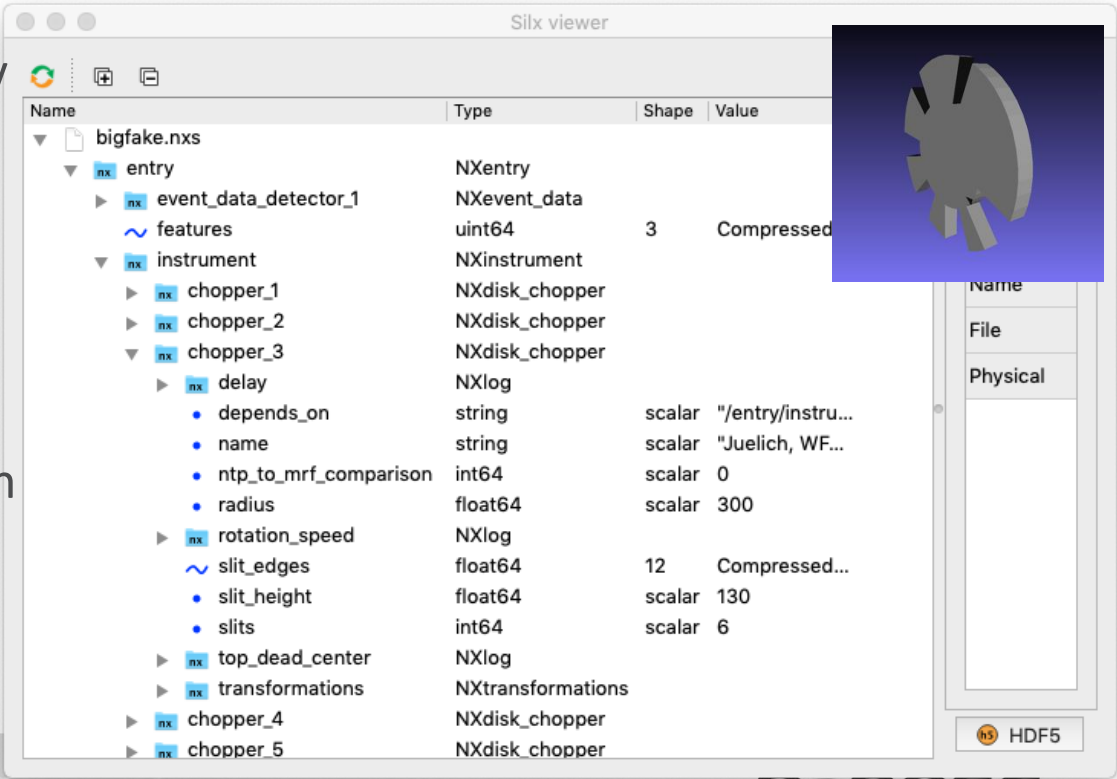Identified two sides of EOSC Integration:

- Harvesting Data by EOSC Agents and Repos
- FAIR Data Search to be federated

Existing catalogues and their capabilities as well as third-party options for implementations have been surveyed. Prototypes are being discussed at a WP3 meeting in Grenoble next week.

## NeXus Survey for Ontologies

Facility practises and plans around file formats are being shared in sessions with partners. Identified commonalities can be applied everywhere and inform the catalogue search task.



**Both Tasks link well with ExPaNDS WP3 – Alun's talk.**

panosc

# PaNOSC has 6 data catalogues with different APIs + UIs



| ESRF (icat) | CERIC (icat) | ESS (SciCat) | ILL (local) | ELI (tbd) | XFEL (MyMdc) |

# PaNOSC common API across all sites



Federated Search API

# Data Catalog Services status in 5/2021

1. Recruitment of staff
2. **Define and develop needs for Search API (ALL)**
3. **Implement Search API locally (ESS, ELI, ILL, EuXFEL, ESRF)**
4. **Setup Federated Search demonstrator (ELI-ALPS )**
5. **Extend Nexus standard (EuXFEL, ELI)**
6. **Local deployment of metadata catalogs (ALL)**
7. **Collaborate with other Work Packages and ExPaNDS**
8. **Extend catalogs and e-logbooks (ESRF, ELI, EuXFEL, ESS)**
9. **Contribute to maintenance of HDF5 (h5py, EU HDF meeting)**
10. **Federated search demonstrator integrated in EOSC (ESS+ELI)**

panosc

# Data Catalog Services issues in 5/2021

1. **Deployment of metadata catalogs (CERIC + ELI)**

2. **Deployment of federated search demonstrator (ELI/ESS)**

3. **Integrating search API in PaN portal (ELI, ESS, ILL)**

4. **Extending Search API to implement authentication + closed data (ILL)**

# Data Catalog + Analysis services – big picture



VISA

## Data Analysis, in the cloud

VISA (Virtual Infrastructure for Scientific Analysis) makes it simple to create compute instances infrastructure to analyse your experimental data using just your web browser

🔒 Sign in with your ILL account

panosc
photon and neutron
open science cloud

# Dataset Search Beta

Search for Datasets 🔍

Common API to search across all PaNOSC catalogues

ESRF (icat)  CERIC (icat)  ESS (SciCat)  ILL (local)  ELI (tbd)  XFEL (MyMdc)

Analyse your data

Create a new compute instance and use your web browser to access a Remote Desktop or JupyterLab to start analysing your experimental data

Collaborate with your team

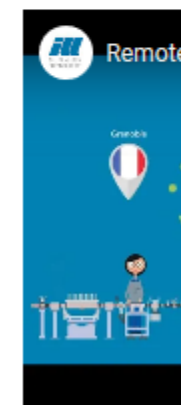Share your compute instance with other members of your team to collaborate together in real time

Remotely control your experiment

Use NOMAD remote, the ILL remote instrument control software, to remotely control your experiment

No need to install software

The compute instances come with pre-installed data analysis software so you can start analysing your experimental data immediately

panosc

# PaNOSC KPIs

| | ILL | ESRF | CERIC | XFEL | ELI | ESS |
|---|---|---|---|---|---|---|
| Open Data 2018 | 100s | 2 | 0 | 10s | 0 | 0 |
| Open Data 2023 | 1000s | 1000s | 100s | 1000s | 100s | 10s |
| Data Services 2018 | Pilot | In progress | Remote | In progress | ? | In progress |
| Data Services 2023 | Desktop Jupyter | Jupyter Desktop | Jupyter Desktop | Jupyter Desktop | Desktop Jupyter | Jupyter Desktop |
| Common data API 2018 | No | No | No | No | No | No |
| Common data API 2023 | Yes | Yes | Yes | Yes | Yes | Yes |
| User training 2018 | No | No | No | No | No | No |
| User training 2023 | Yes | Yes | Yes | Yes | Yes | Yes |

panosc

PaNOSC Milestones

# PaNOSC WP3 Deliverables

**D3.1 API definition (M18, R, PU, ESS)**

**D3.2 Demonstrator implementation (M28, Other, PU, ESS)**

**D3.3 Catalog service (M40, DEC, PU, ESS)**

**D3.4 Implementation Report from Facilities (M44, R, PU, ESS)**

**D3.5 NeXus Metadata Mapping Schema and Proposed New Definitions (M42, R, PU, ESS)**

# Data Catalog Services

# Dashboard

## WP3 : Data Catalog Services

WP Leader: Tobias Richter

| Partner | Correspondents | API defined | API end point | API serving live data | API Authentication | API in PaNOSC Portal | API registered in EOSC | OAI-PMH | OAI-PMH serving live data |
|---------|----------------|-------------|---------------|-----------------------|--------------------|----------------------|------------------------|---------|---------------------------|
| ESRF | A de Maria | ✓ | Yes, ✓ or NO | | | | | | |
| ILL | S Caunt | ✓ | | | | | | | |
| ESS | T Richter | ✓ | ✓ | ✓ | | | | ✓ | ✓ |
| XFEL | ? | ✓ | | | | | | | |
| CERIC | ? | ✓ | | | | | | | |
| ELI | L Schrettner | ✓ | | | | | | | |
| EGI | ? | ✓ | | | | | | | |

Upcoming deliverables & milestones:

- D3.2 Demonstrator implementation M28 (March 2021) ✓
- MS3.3 Catalog Integration Best Practices Meeting M30 (May 2021)
- D3.3 Catalog Service M40 (March 2022)