

Metadata Practices at Synchrotron Facilities

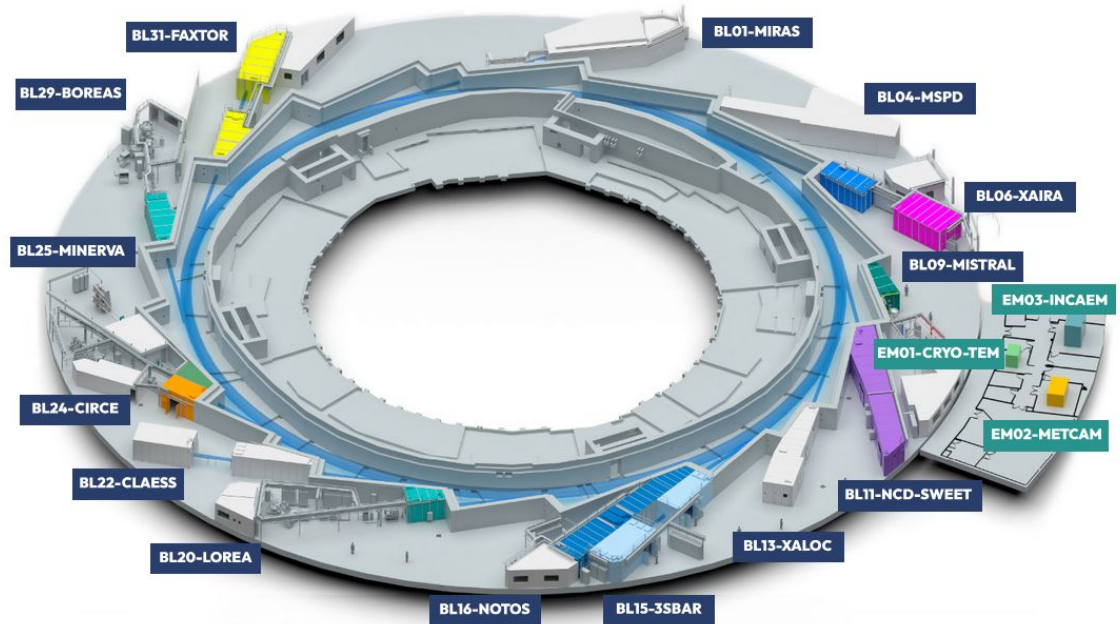
Nicolas Soler, head of SDM, ALBA synchrotron

[NAPMIX Training Workshop 2026](#)

19/05/2026

Beamlines and Electron Microscopes

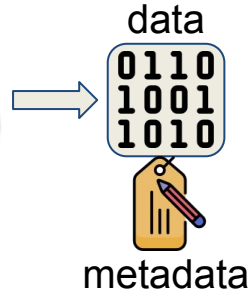
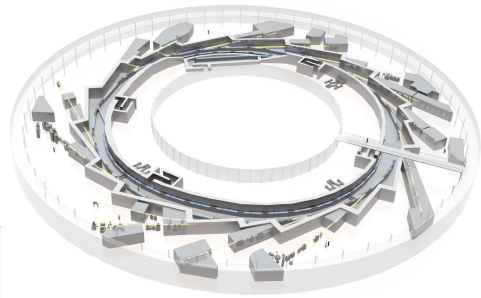
BL01 - MIRAS
BL04 - MSPD
BL06 - XAIRA
BL09 - MISTRAL
BL11 - NCD-SWEET
BL13 - XALOC
BL15 - 3Sbar
BL16 - NOTOS
BL20 - LOREA
BL22 - CLÆSS
BL24 - CIRCE
BL25 - MINERVA
BL29 - BOREAS
BL31 - FAXTOR
EM01-CRYO-TEM
EM02-METCAM
EM03-INCAEM



ALBA synchrotron at a glance: <https://www.cells.es/en/instruments/map>

The Scientific Data Management section

Our missions



Data processing
(mainly beamline support)

Programming
Optimization Java HPC
Data C/C++ Analysis
Pipelines Processing Visualization
Scientific Computing
Machine Learning Web scrapping
Matlab Python Libraries
API data reduction

Data management
In synergy with other sections in Computing

- Orchestrating FAIR metadata capture (with CT)
- Providing provenance tools
- Persistent identifiers management (with MIS)
- Participate in data catalogue implementation (with MIS)
- Data analysis platform (with MIS)
- Data management plans

Section created in 2021, now 6 people

- 5 software engineers
- 1 section responsible

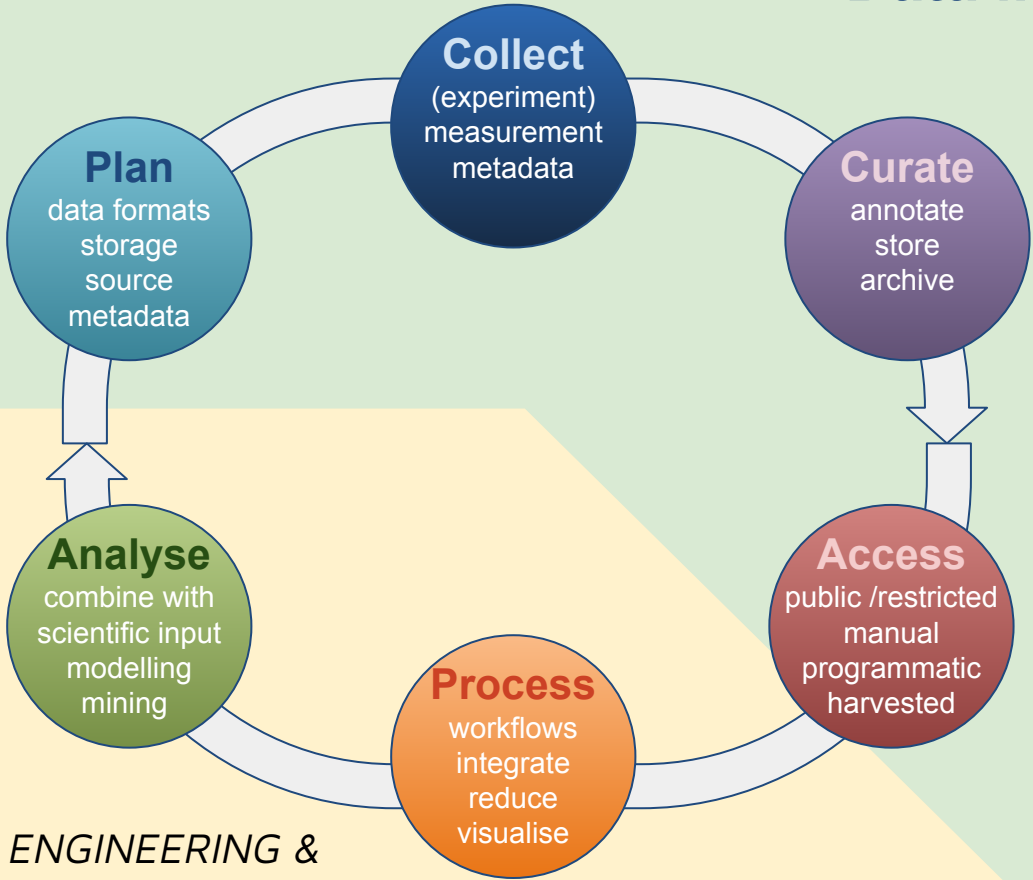


Data Management

DATA
MANAGEMENT

SCIENTIFIC
COMPUTING

SOFTWARE ENGINEERING &
ENVIRONMENT



Changes and challenges

Upgrade to [ALBA-II](#)
New instruments and detectors

More brilliance & coherence

kHz fast detectors

MX, μ CT,
Ptycho, EM, etc

European Commission, funders

Faster data collection for conventional method



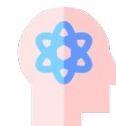
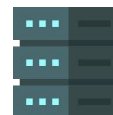
Previously marginal techniques will become routine (SSX, Ptychography)

Public research data must be browsable & **reusable** by others (FAIR)



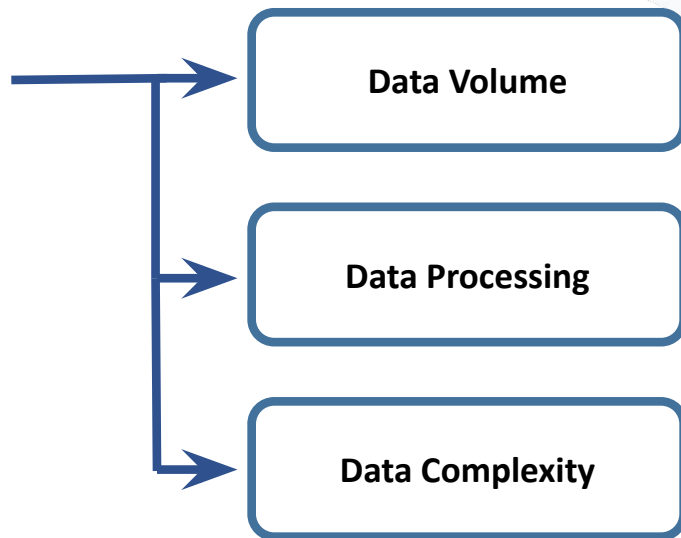
More data open, reusable

- Storage capacity
- Bandwidth
- HPC power
- Automation
- Scientific data processing expertise
- FAIR data policy
- Interfacility coordination and service federation (software and data)



Data challenges on the horizon

**Data
Challenges**



**High-quality data annotation
systems are required.**

EOSC: European Open Science Cloud



In May 2015, the [European Commission](#) proposed creating a European Open Science Cloud (EOSC) to the Competitiveness Council.

The aim was to [federate existing research data infrastructures in Europe](#) and realise [a web of FAIR data](#) and related services for science, making research data interoperable and machine actionable following the [FAIR guiding principles](#).

In the initial phase of development until [2020](#), the Commission invested around €320 million to start prototyping the EOSC through project calls in Horizon 2020 - the Commission's research and innovation funding programme.

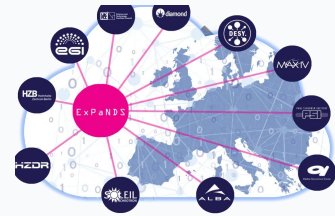
The partnership aims to deploy and consolidate by [2030](#) an open, trusted virtual environment to enable the estimated [2 million European researchers to store, share and reuse research data](#) across borders and disciplines.

The partnership will bring strategic coherence and complementary commitments at EU, national and institutional levels to bring together all advanced data infrastructures in Europe, modernise the ERA with a capability to produce "FAIR-by-design" ([Findable, Accessible, Interoperable, Reusable](#)) datasets, and populate EOSC with new FAIR data and related services, (iv) expand the FAIR data culture across Europe.



The PaNOSC and ExPaNDS projects

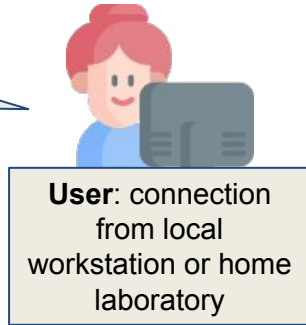
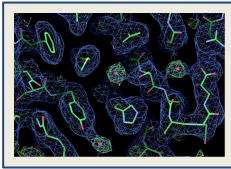
Photon and Neutron Open Science Cloud
European Open Science Cloud (EOSC) Photon and Neutron Data Service



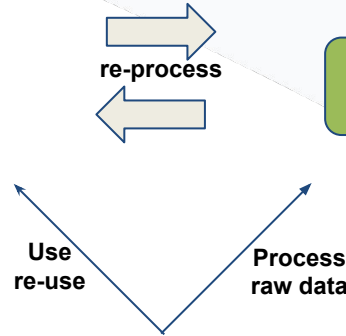
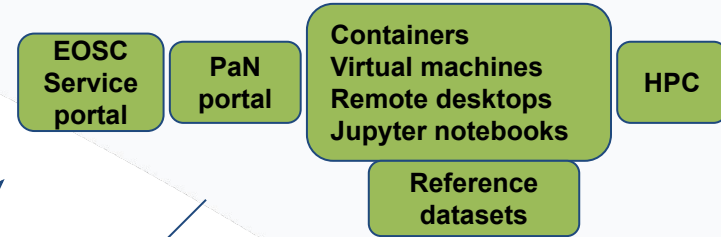
Were **4-year projects** started respectively in 2018 and 2019 to expand, accelerate and support the **data management** and **data services** provided through the EOSC for major national Photon and Neutron Research Infrastructures (PaN RIs)

- Enable EOSC services and to **provide coherent FAIR data services** to the scientific users of national Photon and Neutron sources
- Connect national PaN RIs through **a platform of data analysis as a service** (DAaaS) for users from research institutes universities, industry etc.
- Develop and maintain a **catalogue of data and analysis software** for Photon and Neutron data
- Gather feedback and **cooperate with the EOSC** governance bodies to improve the EOSC and develop standard relationships between scientific publications, Photon and Neutron scientific dataset (raw data), experimental reports, instruments and authors (via ORCID)

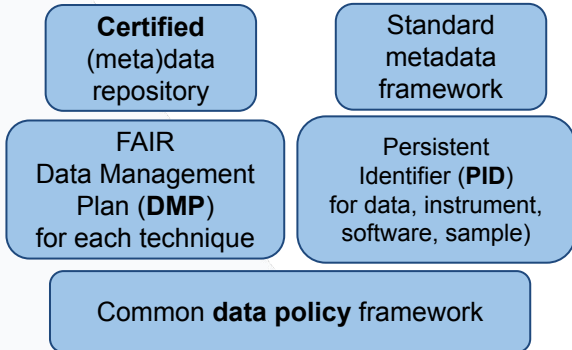
WP 2,3,4 overview



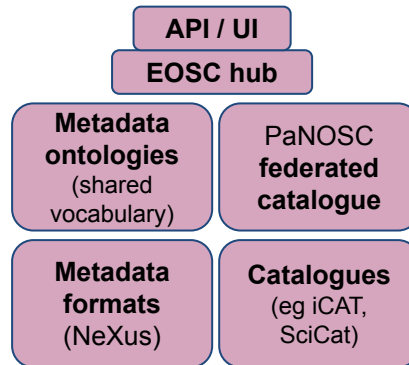
Data analysis as a service (WP4)



FAIR-ready data (WP2)



Data catalogues (WP3)

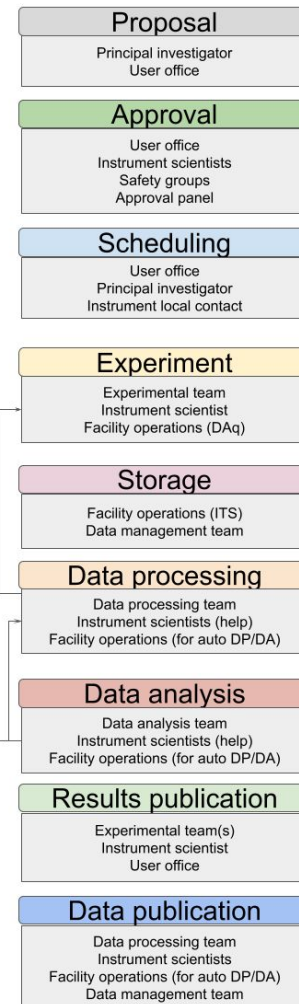


Data reduction / compression

- Metrics
- New algorithms
- Software vs hardware
- Technique-specific
- Lossy vs non lossy
- Meta-compressors (ML)

Some of the outcomes

- A common FAIR data policy framework
- Common basis for Data Management Plans (DMP)
- Metadata recipes
- Common solutions of metadata catalogues: mainly Scicat and ICAT
- A federated interface for searching data
- Ontologies (NeXus and PaNET)
- A common (meta)data format: NeXus/HDF5
- A common solution for remote data processing and visualization: VISA



PI/Main proposer	P1	FA
Co-investigators	P1	FA
Instrument requested	P1	F
Funding source	P2	F
Sample description	P1	F
Proposed experimental conditions [Safety conditions]	P1	F
Experiment description	P3	F
Prior art (related publications, proposals)	P1	F
Facility information	P2	F
Proposal identifier	P1	F
[Approval panel]	P3	/
Sample safety assessment	P2	/
Allocated day & time on instrument	P2	FA
Scheduled visiting experimental team	P2	FA
Safety Training data	P3	/
Detailed experimental planning	P2	F
Sample preparation	P2	FR
[Sample reception]	P3	/
Visiting experimental team (user id)	P1	FA
Experiment date	P1	FA
Sample information	P1	FR
Instrument information	P1	FR
Calibration information	P1	FR
xperimental planning	P2	FR
Environmental parameters	P2	FR
Laboratory notebook	P2	FR
Instrument scientist	P2	F
[Experimental report]	P3	R
Persistent Identifiers (PIDs)	P1	FA
Preservation description information	P1	AR
Dataset information	P1	F
File identifier	P2	AR
[Representation information]	P3	IR
[Instrument parameters]	P3	FR
Processing team (user ID)	P2	AIR
Original data	P1	IR
Data format (after processing)	P1	IR
Dataset information	P2	AIR
Processing information	P1	R
Software package information	P1	R
Analysis team (user id)	P2	AIR
Original data	P1	IR
Software package information	P1	IR
Dependence tracking and workflow	P2	R
Data formats (after analysis)	P1	IR
Dataset information	P1	IR
File identifier	P1	AIR
[Instrument parameters]	P3	IR
[Calibration information]	P3	IR
Authors / Coauthors (user ID)	P1	FA
Proposal information	P1	FA
Publication information	P1	F
persistent Identifier (PID)	P1	F
[Supplementary data information]	P3	F
Resource identity	P1	FI
Related resource	P2	F
Creator	P1	F
Contributor	P2	F
Title	P1	F
Publisher	P1	FI
Publication year	P1	FI
Licence	P1	IR
Release date	P1	IR

A Common approach to manage scientific data



Data policy

Regulations
Rights
Duties

Data Management Plan (DMP)

Think ahead



Data storage

Ultra fast
Disc
Tape

Data format:

Wrapping data and metadata



Data catalogue

Interface to find and browse data, logbook



Automatic pipelines

Produce processed data

During and after experiment

(Re)processing platform

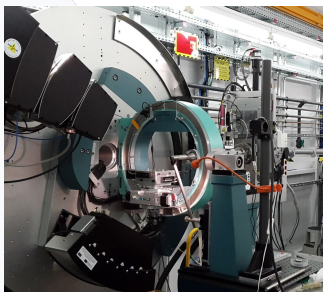
Provides remote interfaces for using ALBA's computing resources



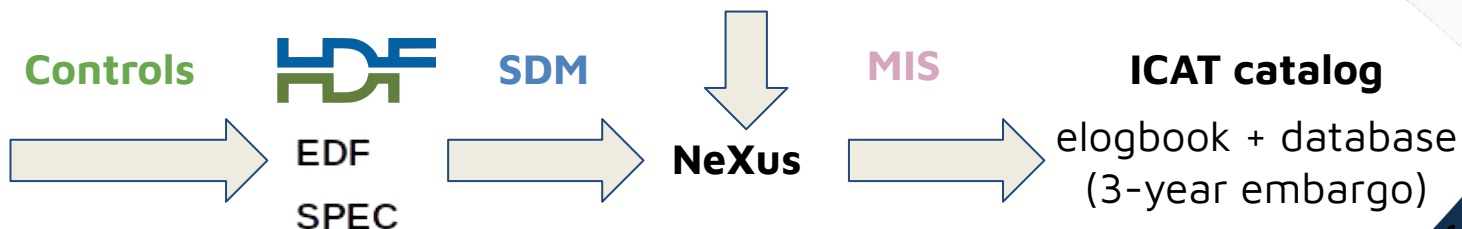
(meta)data format: NeXus/HDF5

- ALBA synchrotron has pledged to follow the [FAIR principles](#)
- [NeXus](#) taken as the standard format to store (meta)data in ALBA to comply with FAIR principles
- Since 2023, Dr. Fernán Saiz (SDM) is acting as ALBA's representative in the NeXus International Advisory Committee (NIAC)
- [NIAC](#) develops and maintains NeXus as the common format for neutron, x-ray, and muon science facilities in Europe, Asia, Australia, and North America.
- Raw data is retained for a minimum period of 5 years (archived after 1 year)

Beamline



User Office



(meta)data format: NeXus/HDF5

[nexus v2026.01 documentation](#) » [3. NeXus: Reference Documentation](#) » [3.3. NeXus Class Definitions](#) » [3.3.2. Application Definitions](#)

3.3.2. Application Definitions

A description of each NeXus application definition is given. NeXus application definitions define the *minimum* set of terms that *must* be used in an instance of that class. Application definitions also may define terms that are optional in the NeXus data file. The definition, in this case, reserves the exact term by declaring its spelling and description. Consider an application definition as a *contract* between a data provider (such as the beam line control system) and a data consumer (such as a data analysis program for a scientific technique) that describes the information is certain to be available in a data file.

Use NeXus links liberally in data files to reduce duplication of data. In application definitions involving raw data, write the raw data in the [NXinstrument](#) tree and then link to it from the location(s) defined in the relevant application definition.

Application definitions are grouped together based on the research fields where these are typically used. Definitions that address multiple research fields are listed in each category:

[Atom Probe Microscopy](#)

[Diffraction & Scattering Techniques](#)

[Electron Microscopy](#)

[Imaging Techniques](#)

[Multi-Dimensional Photoemission Spectroscopy](#)

[Optical Spectroscopy](#)

[Time-of-Flight Techniques](#)

[Complete List](#)

<https://manual.nexusformat.org/classes/applications/index.html>

(meta)data format: NeXus/HDF5

3.3.2.2. Diffraction & Scattering Techniques

Introduction

Application definitions for different diffraction and (small-angle) scattering techniques

Application Definitions

[NXioproc](#)

Application definition for any $I(Q)$ data.

[NXlauetof](#)

This is the application definition for a TOF laue diffractometer.

[NXmonopd](#)

Monochromatic Neutron and X-Ray Powder diffractometer.

[NXxbase](#)

This definition covers the common parts of all monochromatic single crystal raw data application definitions.

[NXeuler](#)

Raw data from a four-circle diffractometer with an eulerian cradle, extends [NXxbase](#).

[NXkappa](#)

Raw data from a kappa geometry (CAD4) single crystal diffractometer, extends [NXxbase](#).

[NXlaue](#)

Raw data from a single crystal laue camera, extends [NXxrot](#).

[NXlaueplate](#)

Raw data from a single crystal Laue camera, extends [NXlaue](#).

[NXmb](#)

Raw data from a single crystal diffractometer, extends [NXxbase](#).

[NXxrot](#)

Raw data from a rotation camera, extends [NXxbase](#).

[NXcanSAS](#)

Implementation of the canSAS standard to store reduced small-angle scattering data of any dimension.

[NXsas](#)

Raw, monochromatic 2-D SAS data with an area detector.

[NXsastof](#)

Raw 2-D SAS data with an area detector with a time-of-flight source.

Implementation of the data catalogue (ICAT)



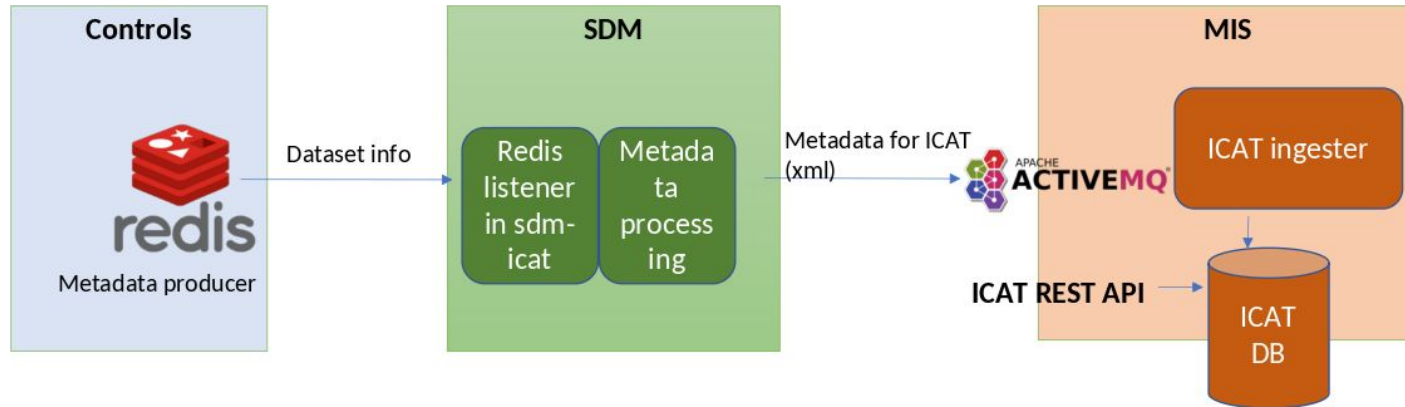
The screenshot shows the ALBA Data Portal interface. At the top, there is a navigation bar with links for Data Portal, Experiments, Publications, Logistics, Beamlines, and Manager. A search bar is present with the text "Search experiments...". Below the navigation bar, the main content area is divided into several sections:

- ALBA Data Portal:** A dark blue header with the text "Find, visualize and access data acquired at ALBA". Below this, a warning message states: "Public data is accessible to anyone with an ALBA User Office account. You need to be logged in to visualize your data when it is under embargo. See [ALBA data policy](#) for more details."
- Start searching data:** A section with a search input field containing "Experiment title, abstract, beamline, DOI..." and a "Search" button. Below it, there are tabs for "All data", "My data", "Public data", "Embargoed data", and "My calendar".
- My data:** A section displaying a list of datasets and samples. The datasets listed are: "20/04/2024 - 07/07/2009" and "07/07/2023 - 07/07/2009".
- Continue where you left off:** A section displaying a list of datasets and samples. The datasets listed are: "27/09/2024 - 27/12/2024", "27/09/2024 - 27/12/2024", "10/06/2024 - 10/06/2024", "07/07/2024 - 20/07/2024", "16/05/2024 - 16/05/2024", and "20/04/2024 - 07/07/2009".
- Logistics:** A section displaying a list of ongoing parcels. The parcels listed are: "10/06/2024 - 10/06/2024", "16/05/2024 - 16/05/2024", "16/05/2024 - 16/05/2024", and "20/04/2024 - 07/07/2009". Each parcel has a "Test" button and a "SCHEDULED" status.

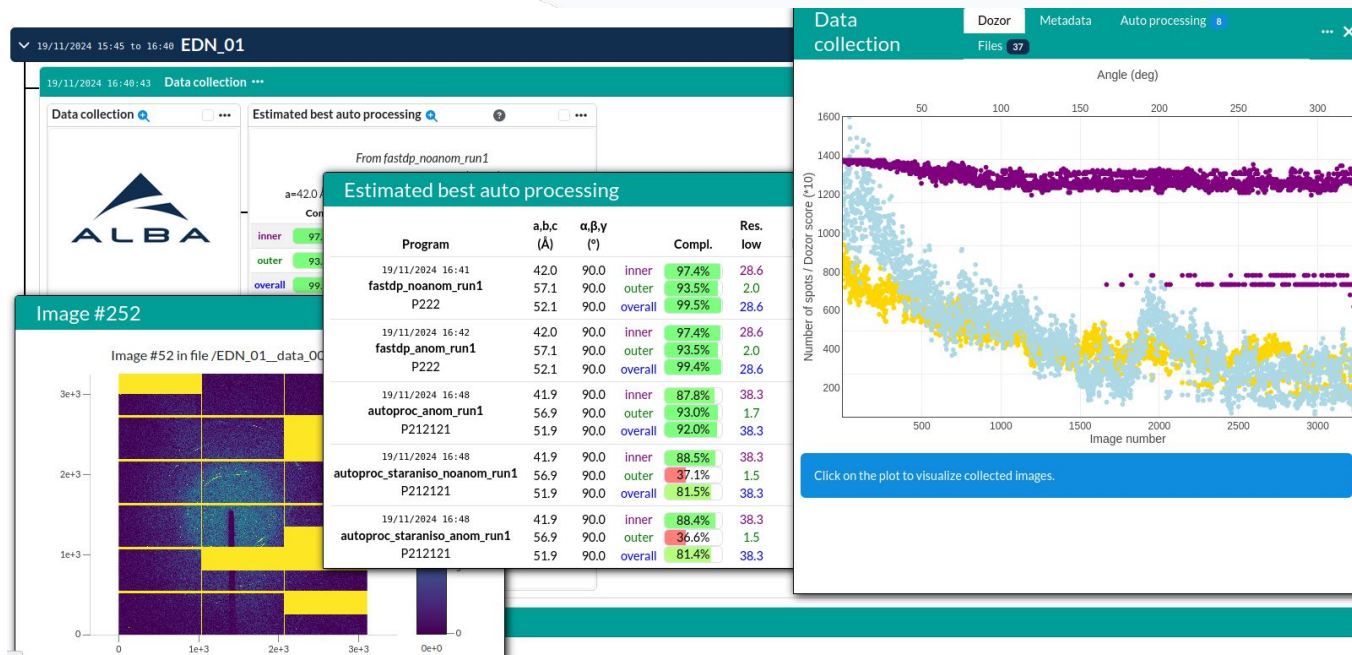
- [ICAT](#) is a data catalogue solution, resulting from a collaboration between different european facilities (ESRF, SFTC, HZB, Diamond, etc)
- It allows finding, browsing, viewing, downloading and annotating data during and after the experiment
- Now in **deployment phase** at ALBA

<https://data.cells.es>

ICAT communication setup



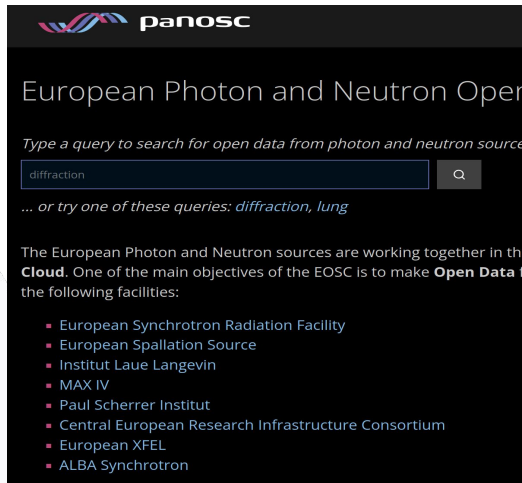
ICAT-DRAC (for XAIRA, developed by ESRF)



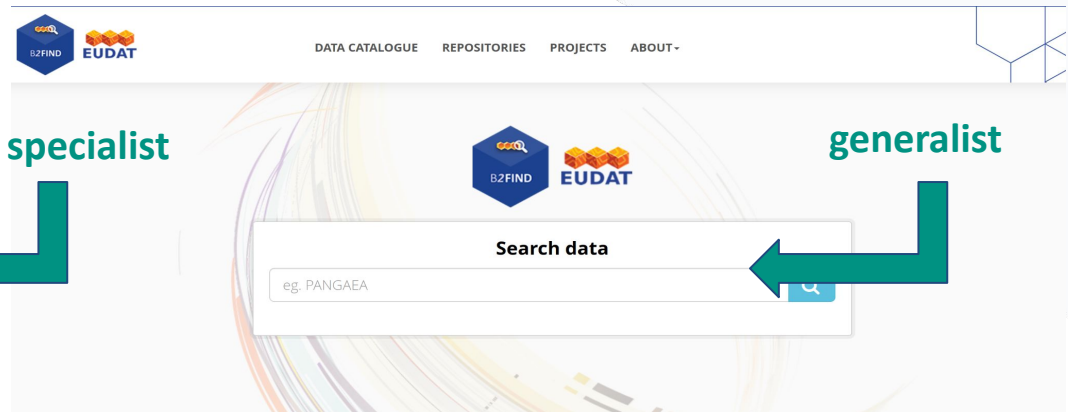
An integrated Laboratory Information Management System (LIMS)

Harvesting PaN Data catalogues

- We aim to participate into the **creation of a European PaN Data Space** and to **facilitate Spanish users** access to it.
 - Our data catalogue seamlessly **integrates with the PanOSC data catalogue**, allowing for easy discovery through general search engines.



<https://data.panosc.eu/>



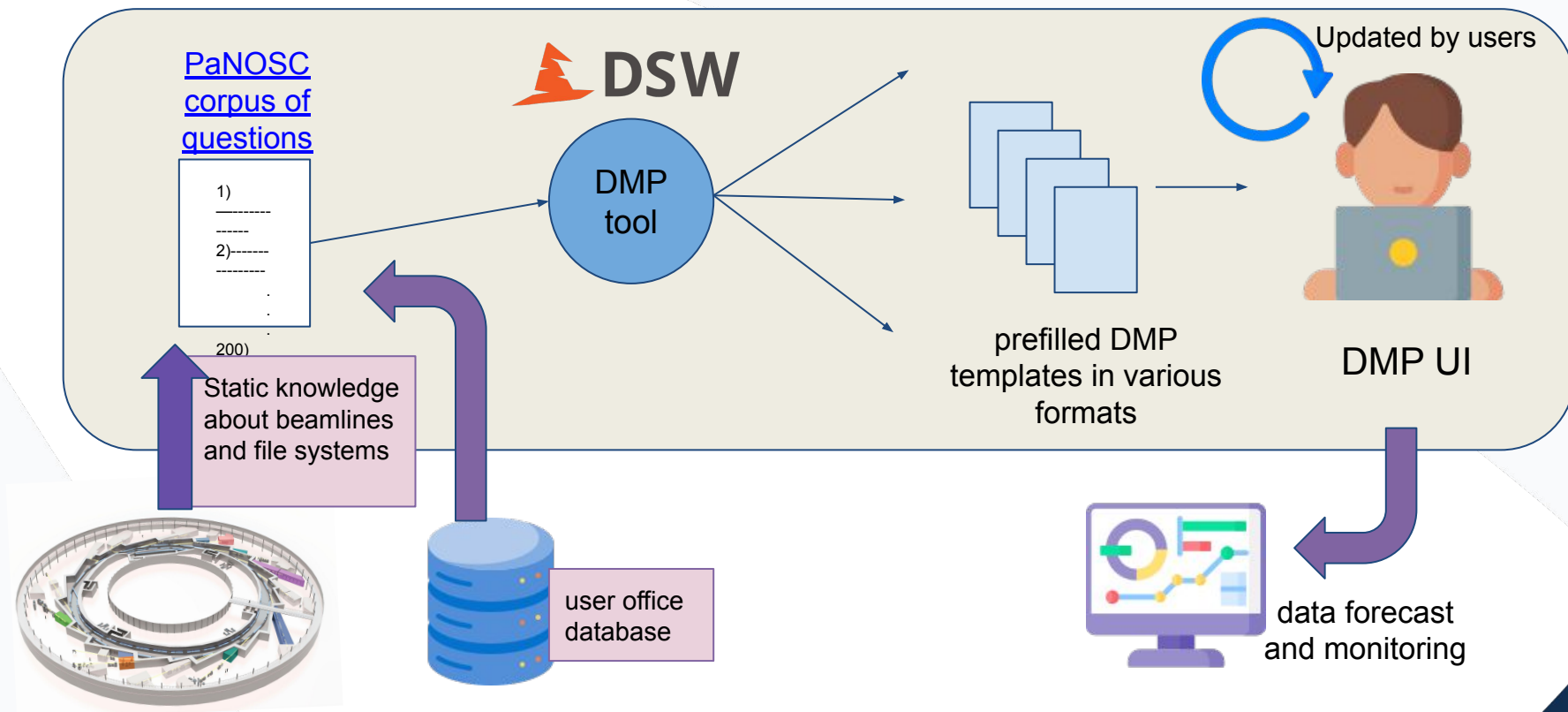
PaN specialist

generalist

<https://b2find.eudat.eu>



Data Management Plans (DMP)



Current data policy for public experiments

A few keypoints



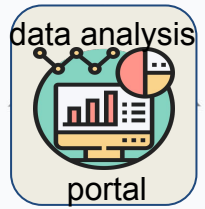
- **No need to copy data home**, ALBA custodies it
- Raw data is **retained** for a minimum period of **5 years** (archived after 1 year)
- A **3-year embargo period** is applied to raw data, afterwards it becomes **publicly available**
- High-level metadata, such as proposal title, authors and abstracts, are visible immediately after the experiment is completed.

We must update our current data policy

- Based on the template and guidelines from the PaNOSC and ExPaNDS projects
- Explicitly states our commitment to embrace the FAIR principles
- E.g. commits to rich metadata annotation, persistent identifiers, interoperable formats, etc
- Must include derived/processed data (and software)
- Encompasses our new Microscopes (InCAEM)
- Must take into account operational sustainability of the long term access and custody of the data

3) Keystones

Data Analysis as a Service: VISA



Portal for (re)processing / remote data analysis using HPC resources. (leading group: **MIS**)

- Active collaboration (MoU) between ALBA and other European list institutes
- can be run on premises or even on external backends (e.g Google Cloud)
- MIS team is exploring **new ways of deploying it** (**Kubernetes**, **Windows**)
- ALBA is candidate for the first in-person workshop in 2025

FACILITIES INVOLVED:

ALBA, DESY, ILL, ESRF, SOLEIL, ESS, EuXFEL

Why VISA?

- Some data are becoming too big to be transferred to home institution
- Users take advantage of the facility HPC
- Exhaustive software environment and provenance

The screenshot shows the VISA portal landing page. At the top center is the VISA logo, which consists of two interlocking loops in blue and green, with the word 'VISA' in blue below it. Below the logo is the text 'Data Analysis, in the cloud'. Underneath that is a paragraph: 'VISA (Virtual Infrastructure for Scientific Analysis) makes it simple to create compute instances on the data analysis infrastructure to analyse your experimental data using just your web browser'. Below this paragraph is an orange button with the text 'Sign in with your user account'. To the left of the button, there are three sections of text: 'Analyse your data' (with a sub-point: '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' (with a sub-point: 'Share your compute instance with other members of your team to collaborate together in real time'), and 'No need to install software' (with a sub-point: 'The compute instances come with pre-installed data analysis software so you can start analysing your experimental data immediately'). To the right of these sections is a screenshot of a web browser displaying a data visualization with many colored lines radiating from a central point. At the bottom of the page, there is a link for 'Questions or feedback?' and a note: 'Please send the IT service an email to mis@boell.es'.

Tools to manage scientific data



Data policy

Regulations
Rights
Duties



Data storage

Ultra fast
Disc
Tape



Data Management Plan (DMP)

Think ahead



Data format:

Wrapping data and metadata



Data catalogue

Interface to find and browse data, logbook



Automatic pipelines

Produce processed data



During and after experiment

(Re)processing platform

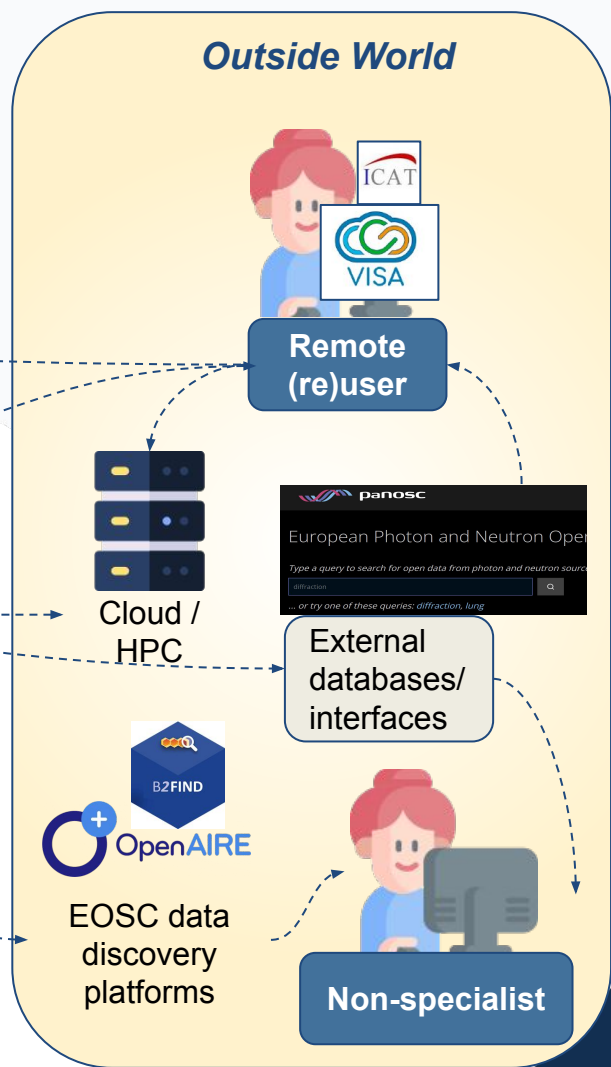
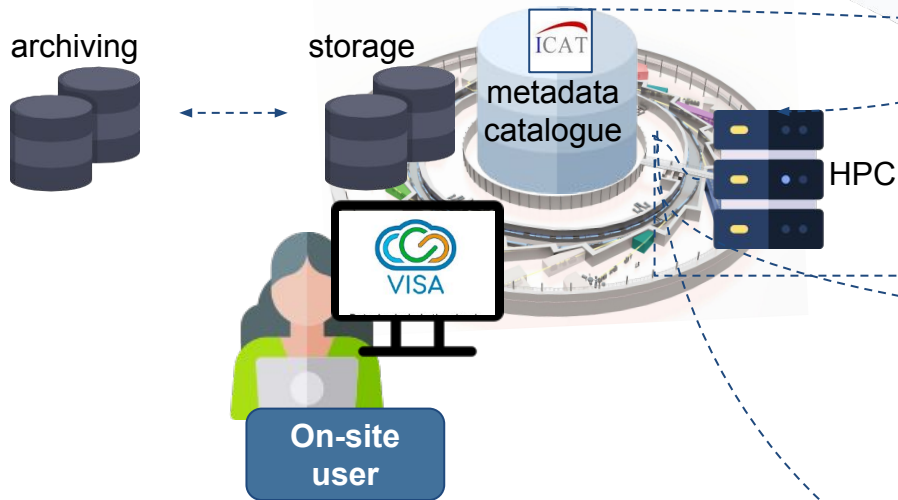
Provides remote interfaces for using ALBA's computing resources



Entire setup ready: foreseen for 2026

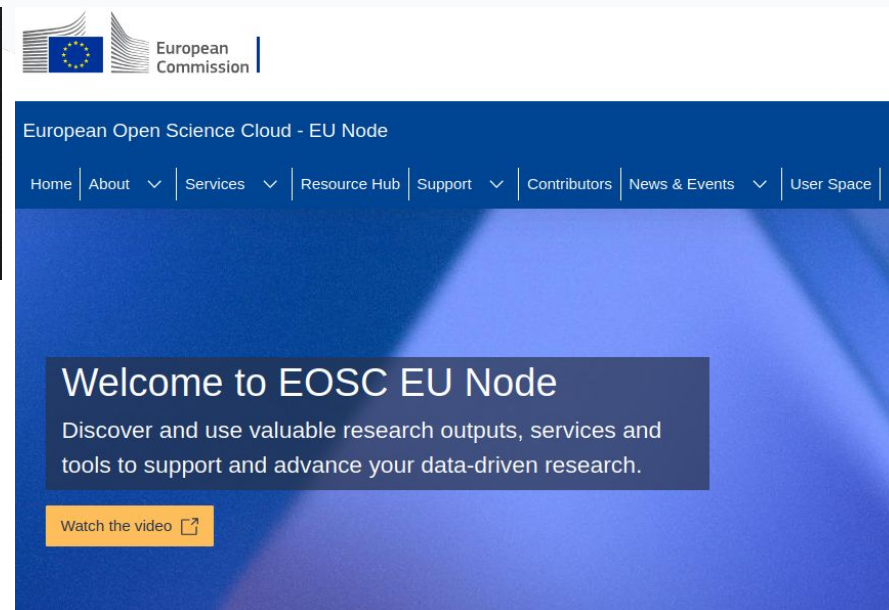
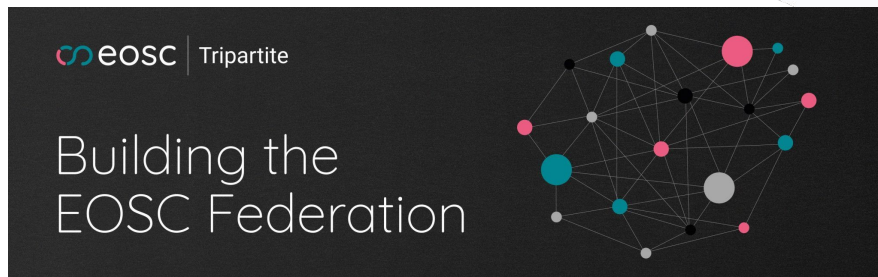


Vision of future data services for users



Finding, Accessing, Understanding, Reprocessing data

PaN EOSC node



The European Open Science Cloud (EOSC) is federating European data and services under various thematic nodes.

The first NODE (EU) was launched in October 2024!

ALBA member of the future PaNOSC node

Conclusions

- As a prerequisite for new developments in data processing & automation, ALBA needs to fully implement a FAIR-compliant data management setup. This includes:
 - An update of our **data policy**
 - Interoperable data **formats** (mainly NeXus) and standards
 - The deployment of a data **catalogue**, which includes
 - findability tools, (also via the [PaN data portal](#) and [OpenAIRE/B2find](#))
 - persistent identifiers (DOIs)
 - data landing pages
 - electronic lab. Notebooks
 - data visualization tools (e.g. NeXusviewer, ICAT-DRAC)
 - A **Data Analysis as a Service** platform (VISA), available on-site and from home institution
 - Alignment with / participation in joint PaN **Open Science** initiatives at European level

The first beneficiaries of this FAIR setup are the **users** themselves!

What's next?

- *Harmonized and exportable data processing workflows, providing rich metadata (provenance)*
- *More metadata standards, in particular for the **samples***
- *Combining data from different experiments*
- *AI to evaluate data quality*
- *Real time feedback (experiment steering)*
- *Tailored data compression*
- *AI- assisted DMP generation*
- *Software sharing among PaN facilities*
- *etc*

For all this, we need a good data management basis !

Acknowledgements

Controls

Oriol Vallcorba
Fulvio Becheri

ITS

Gemma Rosas
Sergi Pusó
Sergio Vicente

MIS

Rodrigo Cabezas
Marc Armenter
Toni Fernández

Computing Head: Óscar Matilla

SDM

Emilio Centeno, Albert Castellví, Fernán Saiz
Przemyslaw Karczmarczyk, Joaquín Gómez, Gabriel Jover

ESRF: Alex de Maria, Marjolaine
Bodin, Andy Goetz

Questions? nsoler@cells.es