



DATA SCIENCE &
SCIENTIFIC COMPUTING



**UNIVERSITÀ
DEGLI STUDI
DI TRIESTE**

Designing and deploying a FAIR-by-design data pipeline and platform for electron microscopy laboratories

Research thesis in: Data Management

Supervisor

Dott. Federica Bazzocchi

Candidate

Nicola Perin

University of Trieste

19 settembre 2025

Outline

- 1 Electron microscopy data: what it looks like and where the problems are
- 2 FAIR principles and the NeXus (NXem) standard
- 3 The case study: LAME and the ORFEO datacenter
- 4 From lab problems to a proposal
- 5 Designing the web application
- 6 Testing with VirtualOrfeo
- 7 Walking through the app
- 8 Live demo

Electron microscopy at a glance

- ▶ Different modes: TEM, SEM, STEM → images, diffraction patterns, spectra.
- ▶ The data: large, complex, and very diverse in shape and size.
- ▶ The reality: every vendor has their own format, metadata is often incomplete.

The key question: how do we make these outputs easier to reuse and share?

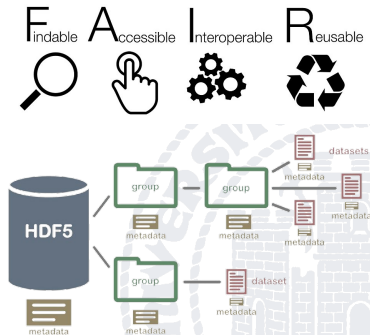
Current Problems and What's Needed

- ▶ **Fragmentation:** many formats, weak or missing metadata.
- ▶ **Friction:** manual copying, endless zip files, confusing naming.
- ▶ **Collaboration:** unclear provenance, scattered access, hard to reuse.
- ▶ **Need:** structured metadata, persistent IDs, scalable storage, simple tools.



A way forward: FAIR and NeXus (NXem)

- ▶ FAIR principles: **findable, accessible, interoperable, reusable**.
- ▶ **HDF5**: efficient format for large, structured datasets.
- ▶ **NeXus**: conventions for scientific data (NXinstrument, NXsample).
- ▶ **NXem**: application definition tailored to electron microscopy.

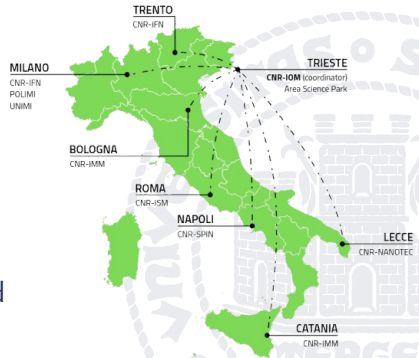


NeXus

At the national level, these FAIR practices are promoted and supported by the NFFA-DI infrastructure.

Introducing NFFA-DI

- ▶ **NFFA-DI** = Nano Foundries and Fine Analysis – Digital Infrastructure.
- ▶ Italian research initiative connecting major nanoscience centers.
- ▶ Goal: open access to advanced instrumentation, FAIR data, and computational resources.
- ▶ Acts as the national driver for FAIR data practices in nanoscience.



Source: <https://nffa-di.it/en/>

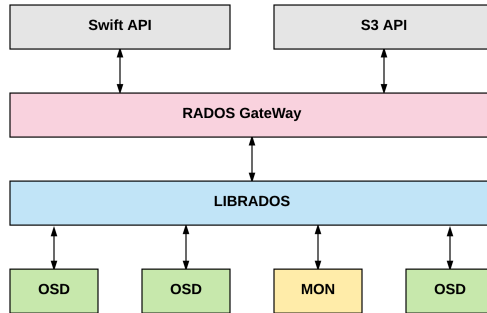
From the lab to the datacenter

- ▶ **LAME** produces multi-terabyte datasets in electron microscopy.
- ▶ As part of NFFA-DI, its work depends on sharing data with other partners.
- ▶ To support this, **ORFEO** provides the backbone: HPC resources, identity services, and S3-compatible object storage.
- ▶ The challenge: connecting LAME's lab workflows with ORFEO's infrastructure.



ORFEO's storage model

ORFEO does not use a classic file-and-folder hierarchy. Instead, data is stored as **objects** inside **buckets**, managed by Ceph RGW with the S3 protocol. Each object combines the raw data with flexible metadata, which makes it easier to describe and reuse datasets.



Accessing ORFEO

All access to ORFEO's resources goes through a single **Ingress** endpoint. Connections are secured with **TLS certificates**, and identity is handled by **FreeIPA** and **Authentik**, providing single sign-on across services.



Welcome to authentik!

Email or Username *

Log in



From problems to a proposal

What's missing for LAME

- ▶ Data often stuck on lab machines or portable drives.
- ▶ Transfers are manual, with inconsistent folder structures.
- ▶ No ingestion standard → hard to reuse or integrate with ORFEO/NFFA-DI.

Our proposal

- ▶ **Transfer:** move data directly into ORFEO using the S3 protocol.
- ▶ **Transform:** convert outputs (e.g. TIFF) into NeXus/NXem with standardized metadata.
- ▶ **Integrate:** build on ORFEO's existing services, with a simple web interface and API.

Choosing a framework

To put our proposal into practice, we need a tool that researchers can actually use. That means building a **web application** that can:

- ▶ guide researchers through projects, samples, and experiments,
- ▶ handle uploads and metadata in a consistent way,
- ▶ connect directly to ORFEO's services (S3 storage, authentication).

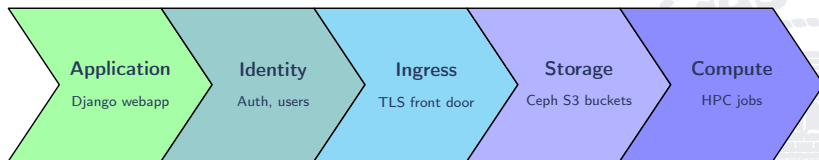
Once this was clear, the next step was to choose the right framework.

Why Django?

We needed a framework that was stable, flexible, and easy to maintain. Django fits well because it provides:

- ▶ A structured way to model projects, samples, and experiments.
- ▶ Built-in tools to create both a user interface and an API.
- ▶ Support for background tasks such as data checks and conversions.
- ▶ Strong support for authentication and long-term maintenance.

Thinking about the whole pipeline



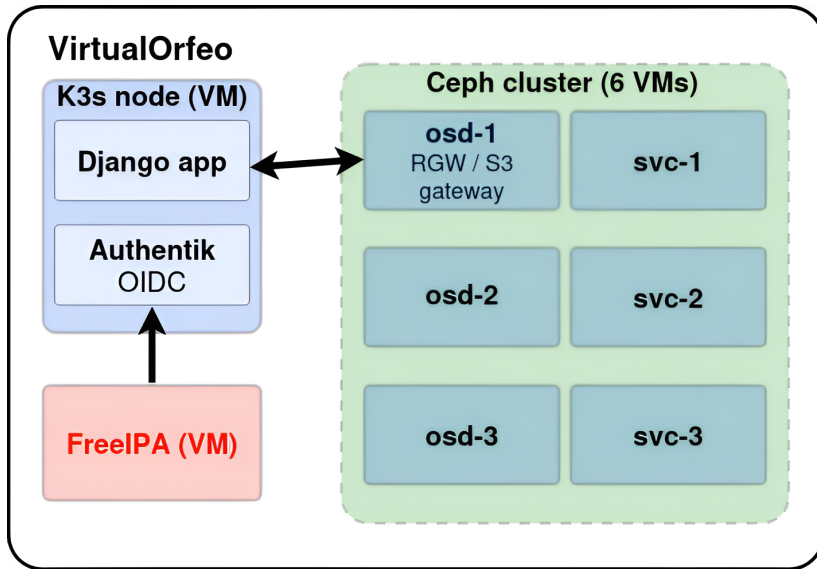
The app is one part of this chain — testing only makes sense when the whole path is reproduced. Our solution: a **digital twin** of ORFEO.

VirtualOrfeo

ORFEO is a complex infrastructure: identity, ingress, storage, and compute. Testing our Django app directly on production would be risky and slow.

- ▶ **VirtualOrfeo** is a lightweight clone of ORFEO, built on K3s.
- ▶ It uses the same Helm charts and configs as production.
- ▶ This lets us deploy the **Django app** in a realistic environment: it can authenticate through Authentik, upload to S3 buckets, and be accessed through the same ingress as in ORFEO.

VirtualOrfeo topology



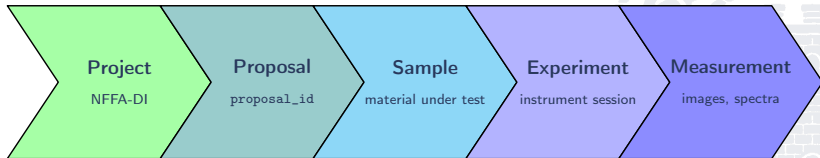
Application overview

The webapp ties together three main layers:

- ▶ **Domain model** Projects contain proposals, samples, experiments, and their measurements.
- ▶ **Data flow** Files go directly into S3 buckets (via presigned URLs), while metadata lives in Postgres.
- ▶ **Background workers** Dedicated pods run checksums, extract metadata, and generate NeXus files.

Research data model

The app organizes research work in a structured chain:



This keeps context together with raw data, making results easier to track and reuse.

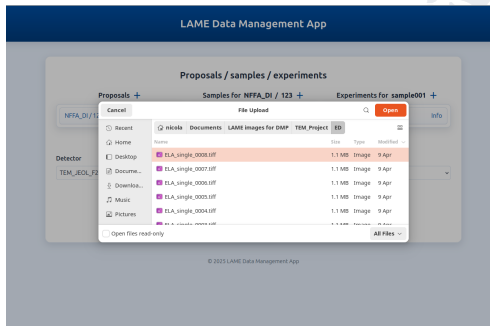
Managing research data in practice

- ▶ Three-pane board to browse and link projects, samples, and experiments.
- ▶ Metadata (context) stored alongside raw data (README.txt).
- ▶ Same operations available via REST API for automation.

The screenshot displays the LAME Data Management App interface. At the top is a dark blue header with the text "LAME Data Management App". Below this is a light blue background containing a white card with the title "Proposals / samples / experiments". The card features three tabs: "Proposals +", "Samples for NFFA_DI / 123 +", and "Experiments for sample001 +". The "Proposals" tab is active, showing a list with two entries: "NFFA_DI / 123" and "Info". The "Samples" tab is also visible, showing "sample001" and "Info". The "Experiments" tab shows "exp001 — some description" and "Info". Below the tabs, there is a section titled "Add measurement to exp001 (sample sample001)". It includes a "Detector" dropdown menu with the selected value "TEM_JEOL_F200 - TVIPS_camera". At the bottom of this section are three buttons: "Choose file(s)", "Choose folder", and "Upload & register". A "Homepage" button is located below these. At the very bottom of the app window, there is a copyright notice: "© 2025 LAME Data Management App".

From upload to storage

- ▶ The app issues a one-time presigned URL.
- ▶ Browser streams data **directly to S3**, bypassing the webserver.
- ▶ Uploads automatically trigger a background job: checksum → metadata extraction → NeXus build.

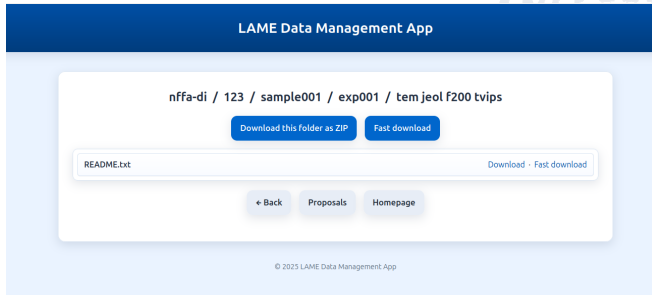


Background workers: from raw files to NeXus

- ▶ A **worker pod** runs in the cluster, always listening to Redis (a shared **to-do list**).
- ▶ When a file is uploaded, the app enqueues jobs such as:
 - Integrity check (checksum).
 - Metadata extraction (parse TIFF headers or JSON).
 - Normalization into NXem fields.
 - NeXus generation: structured `.nxs` file with metadata.
- ▶ Jobs are picked up one by one, retried automatically if they fail.

Browsing and sharing data

- ▶ Browse buckets and datasets directly from the interface.
- ▶ Download with short-lived presigned links.
- ▶ Create on-the-fly ZIP archives for folders.
- ▶ Derived data stored in mirrored namespaces for clarity.



Live Demo

Let's see the workflow in action.



Thank you for your attention!

