# Scientific Data Management
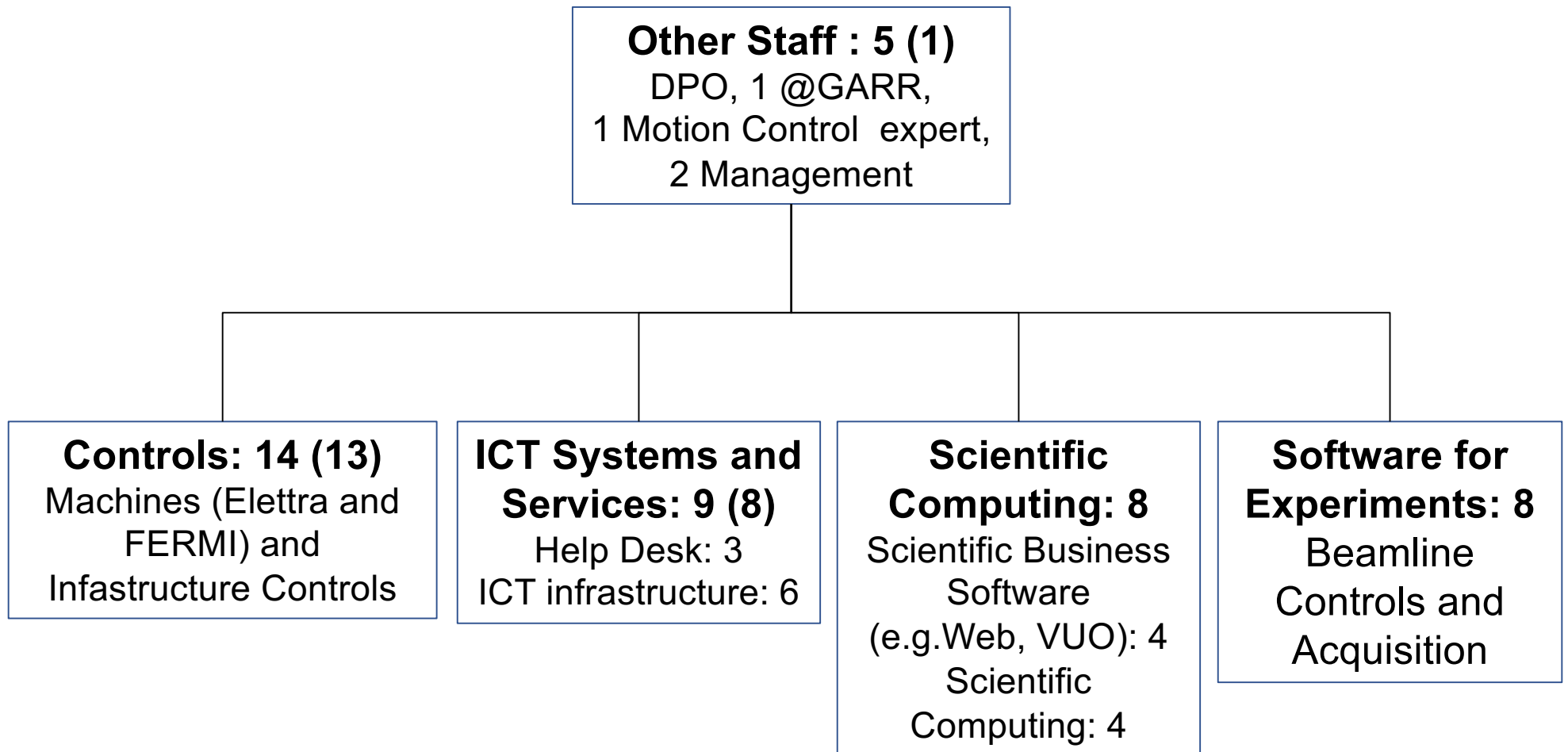
## Roberto Pugliese

*roberto.pugliese@elettra.eu*
ICT Group

✓Elettra ICT Group

✓ICT Infrastructure

✓FAIR Data

- Sientific Data Policy
- DOIs and Scientific Data Storage implementation
- KPIs

# Elettra ICT Group: 44 (38)

**Other Staff : 5 (1)**
DPO, 1 @GARR,
1 Motion Control expert,
2 Management

**Controls: 14 (13)**
Machines (Elettra and FERMI) and Infastructure Controls

**ICT Systems and Services: 9 (8)**
Help Desk: 3
ICT infrastructure: 6

**Scientific Computing: 8**
Scientific Business Software (e.g.Web, VUO): 4
Scientific Computing: 4

**Software for Experiments: 8**
Beamline Controls and Acquisition

**FERMI FL:** networking, HCC controls, offline Tape Lib

**General T2:** networking, HCC general, HPC cluster, scratch, online

**Elettra SB:** networking, HCC controls

**HPC cluster kalculus** 1296 core (2592 thread) CPU, 20 TB RAM, 27 blades; 20736 CUDA core, 1296 Tensor core, 120GB RAM, GPU HBM2 3 blades; blade link at 25 Gb/s; LAN connection to other virtualisation clusters and storage at 100 Gb/s (+10 CPU blades to be added soon)
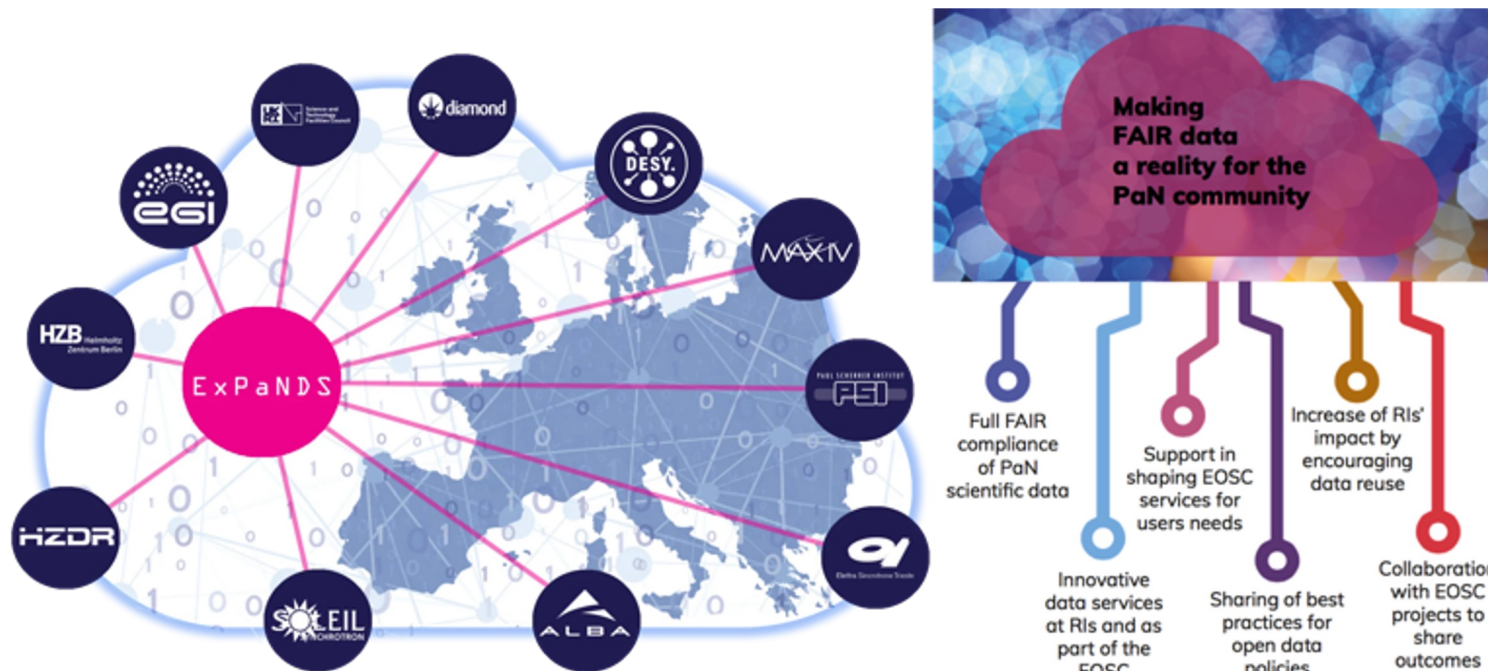
**HPC cluster** 252 core (504 thread) CPU and 2 TB di RAM for beamline online processing

**HCC NUP** (general) Storage + Virtualisation Cluster Sofa (scratch, online, bl acquisition and controls),

**Legacy Cluster** (administration)

PLANIMETRIA GENERALE

# FAIR Data Projects: ExPaNDS, PaNOSC

# Scientific Data Storage Architecture

| DATA CATALOGUE | Without Metadata | With Metadata | |
|---|---|---|---|
| STANDARD STATION | 1 (Data Uploader) | 2 | 3 |
| AUTOMATIZED STATION | | 1 | 2 |
| STORAGE SYSTEM | SCRATCH Can be organized Can be processed Can be transferred | ONLINE Can be organized Can be processed Can be transferred | OFFLINE Can be transferred |
| CAPACITY | 2PB (CEPH – Rep. 3) | 4PB (CEPH – Rep. 3) (48%) | 6 PB Up to 60PB (Replica 4) |

✓ When online datasets can have an associated DOI and you can search the datasets in the DataCite, PaNOSC and EOSC portals

# In summary what happens now …

✓ **Before proposal submission**

- The new scientific data policy is available on the web

✓ **During proposal submission**

- The principal investigator has to accept the new scientific data policy

✓ **Before beamtime**

- A chat group is be created 6 months before the first experiment

✓ **During beamtime**

- chat and all local and remote support services are available to whole experimental team

✓ **After beamtime**

- 2 weeks after the end data is copied offline to the tape library,
- BEST (the day after), achievements (after 30 days) check DOI text, generate **DOI**, generate **DMP**, chat 1 year after the last experiment closes is closed

✓ **Data become open access:** 3y[+1y[+1y]] or on request by the principal investigator

✓ **When a data access request arrives**: the principal investigator is informed

✓ **Moreover** … scratch is pruned when space occupation is above 75%, on-line when above 75%, when deadline expires (10 years …) and when above 75% off-line space will be pruned while metadata is kept forever … best effort.

# FAIR: scientist feedback



✓ During a BLEX-IT meeting and while doing beamtime together we discovered that users need something:

- Similar to **Zenodo**
- Allows for **public data access** but **not anonymous**
- Hosted in **Elettra**
- Compatible with major **scientific Journals** and EU projects
- In accordance with the Elettra Scientific **Data Policy**
- **Open** and FAIR
- Provides the data owner with useful **information and control**

You upload or select your data and get two links, one to the VUO where you can access the data (http://dx.doi.org/10.34965/i10166) and the other one to a public "search" portal.

# KPIs to monitor progress and achievements

| Beamline | ta | Percentage of Open Access Datasets | Data Stored on Online Storage | Percentage of Online Storage used | Percentage of Scratch Storage used | Offline Data Archiving | Number of Archived Datasets | Percentage of Archived Data |
|---|---|---|---|---|---|---|---|---|
| ALOISA | | | | | | | | |
| APE-HE | Yes | 0% | Yes | 1% | 1% | Yes | 1 | 33.33% |
| APE-LE | Yes | 0% | Yes | 1% | 1% | Yes | 1 | 16.67% |
| BACH | | | | | | | | |
| BAD_ELPH | Yes | 0% | Yes | 18% | 17% | Yes | 408 | 99.51% |
| BEAR | Yes | 0% | Yes | 3% | 3% | Yes | 10 | 90.91% |
| CIRCULARPOLARIZATION | Yes | 0% | Yes | 56% | 31% | Yes | 21 | 58.33% |
| DIPROI | Yes | 0% | Yes | 82% | 9% | Yes | 46058 | 85.21% |
| EIS-TIMER | Yes | 0.06% | Yes | 74% | 22% | Yes | 8404 | 77.21% |
| EIS-TIMEX | Yes | 0% | Yes | 74% | 35% | Yes | 11298 | 73.12% |
| ESCAMICROSCOPY | Yes | 1.52% | Yes | 19% | 17% | Yes | 253 | 96.2% |
| GASPHASE | Yes | 0.44% | Yes | 94% | 74% | Yes | 62 | 27.43% |
| IUVS | Yes | 0% | Yes | 1% | 1% | Yes | 6 | 100% |
| LDM | Yes | 0% | Yes | 95% | 9% | Yes | 67495 | 77.01% |
| MATERIALS SCIENCE | Yes | 0% | Yes | 0% | 0% | Yes | 418 | 100% |
| MCX | Yes | 0% | Yes | 18% | 30% | Yes | 174 | 91.58% |
| MagneDyn | Yes | 1.21% | Yes | 82% | 61% | Yes | 8508 | 69.36% |
| NANOSPECTROSCOPY | Yes | 0% | Yes | 31% | 57% | Yes | 2244 | 99.34% |
| NanoESCA | | | | | | | | |
| SAXS | Yes | 0% | Yes | 69% | 73% | Yes | 22906 | 83.24% |
| SISSI-BOFF | Yes | 0% | Yes | 29% | 85% | Yes | 39 | 97.5% |
| SISSI-Chem - Life Sci | Yes | 0.67% | Yes | 29% | 59% | Yes | 146 | 98.65% |
| SPECTROMICROSCOPY | Yes | 0% | Yes | 2% | 3% | Yes | 105 | 99.06% |
| SUPERESCA | Yes | 0% | Yes | 1% | 3% | Yes | 262 | 96.68% |
| SYRMEP | Yes | 0.14% | Yes | 85% | 77% | Yes | 5890 | 97.5% |

| | Name | Descriptic | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| [Select] | Data Management Plan - (TWINMIC) | Data Mana | | | | | | | | |
| [Select] | Electronic Logbook - (TWINMIC) | Electronic l | | | | | | | | |
| [Select] | Experiment Chat - (TWINMIC) | Experiment | | | | | | | | |
| [Select] | General Metadata - (TWINMIC) | General Me | | | | | | | | |
| [Select] | Experiment Metadata - (TWINMIC) | Experiment | | | | | | | | |
| [Select] | Raw Data in HDF5 - (TWINMIC) | Raw Data i | | | | | | | | |
| [Select] | Number of Datasets - (TWINMIC) | Total numb | | | | | | | | |
| [Select] | DOI Minting - (TWINMIC) | DOI Mintin | | | | | | | | |
| [Select] | Number of minted DOI - (TWINMIC) | Number of | | | | | | | | |
| [Select] | Open Access to Data - (TWINMIC) | Open Acces | | | | | | | | |
| [Select] | Percentage of Open Access Datasets - (TWINMIC) | Percentage | | | | | | | | |
| [Select] | Data Stored on Online Storage - (TWINMIC) | Data stored | | | | | | | | |
| [Select] | Percentage of Online Storage used - (TWINMIC) | Occupied s | | | | | | | | |
| [Select] | Percentage of Scratch Storage used - (TWINMIC) | Occupied space on scratch storage | Restricted | Alessandra GIANONCELLI | % | 30 | Automatic | 42% | 24/01/2024 |
| [Select] | Offline Data Archiving - (TWINMIC) | Offline Data Archiving | Restricted | Alessandra GIANONCELLI | Y/N | 365 | Automatic | Yes | 25/12/2023 |
| [Select] | Number of Archived Datasets - (TWINMIC) | Total number of archived datasets | Restricted | Alessandra GIANONCELLI | | 30 | Automatic | 2348 | 24/01/2024 |
| [Select] | Percentage of Archived Data - (TWINMIC) | Percentage of archived data | Restricted | Alessandra GIANONCELLI | % | 30 | Automatic | 98.78% | 24/01/2024 |
| [Select] | Bytes per hour - (TWINMIC) | Bytes produced per hour by the beamline instrumentation | Restricted | Alessandra GIANONCELLI | | 30 | Automatic | 240861121 | 24/01/2024 |
| [Select] | Number of Workstations - (TWINMIC) | Total number of Workstations | Restricted | Marco DE SIMONE | | 365 | Manual | 8 | 08/11/2022 |
| [Select] | Number of remotely accessible Workstations - (TWINMIC) | Total number of remotely accessible Workstations | Restricted | Marco DE SIMONE | | 365 | Manual | 6 | 08/11/2022 |
| [Select] | Percentage of remotely accessible Workstations - (TWINMIC) | Percentage of remotely accessible Workstations | Restricted | Marco DE SIMONE | % | 365 | Automatic | 75% | 25/12/2023 |
| [Select] | Available Bandwidth - (TWINMIC) | Bandwidth available for the beamline | Restricted | Marco DE SIMONE | Gbit/s | 365 | Manual | 10 Gbit/s | 08/11/2022 |
| [Select] | Long Term Storage Duration - (TWINMIC) | Duration in years of the Long Term Storage for the beamline | Restricted | Marco DE SIMONE | Years | 365 | Manual | 10 Years | 08/11/2022 |
| [Select] | Number of CPUs - (TWINMIC) | Total number of CPUs available for the beamline | Restricted | Marco DE SIMONE | | 365 | Manual | 252 | 11/11/2022 |

Thanks!
Questions?

www.elettra.eu