



Data Policy at Large Research Infrastructure. Data Retention challenge.

WORKSHOP ON *WHEN SHOULD SMALL MOLECULE CRYSTALLOGRAPHERS PUBLISH RAW DIFFRACTION DATA?* 11-12 August 2021





Outline

- 1. Why FAIR and Open Access
- 2. FAIR data policy, DMP and Open Access
- 3. What FAIR brings/what data to publish/how to publish
- 4. ELI/PaNOSC use-cases
- 5. Conclusions





Why FAIR? Why OPEN ?

Science produces much more than Publications

Science produces Publications







"Ideas, theories, and their supporting intellectual frameworks should co a common good of all humanity, freely shared for our mutual enjoyment and

FAIR and Open Science Impact! "1



Scientific impact and cost effective:

The European Union commissioned a study of the cost to the research community from not having FAIR data¹.

The study estimates a cost impact to the entire European research community of over €10bn per year. This cost impact provides a market valuation of research data, that without FAIR, is not curated and has degraded quality.

More details: https://sciencebusiness.net/report/why-open-science-future-and-how-make-it-happen





No raw data, no science: another possible source of the reproducibility crisis¹

Computational physics – only 12 results reproduced out of over 300?²

Is there a Reproducibility Crisis? What are the causes? How can it be addressed?

Nature533, 452–454 (2016)Cite this article5320Accesses1225Citations3871AltmetricMetrics

Reproducibility crisis???



HTTPS://EN.WIKIPEDIA.ORG/WIKI/REPLICATION_CRISIS HTTPS://PHYS.ORG/NEWS/2017-03-SCIENCE-CRISIS.HTML What are the factors (<u>https://www.nature.com/articles/d42473-019-00004-y</u>): (ex comes from biomedical but conclusions are accurate!)

- A lack of access to methodological details, *Taw data*, and research materials
- Use of **misidentified**, cross-contaminated, or over-passaged cell lines and microorganisms.
- Inability to manage complex datasets
- Poor research practices and experimental design
- A competitive culture that rewards novel findings and undervalues negative results (But are there really NEGATIVE Results or is it just a different perspective?)

It all leads to (RAW) Data!

1. https://pubmed.ncbi.nlm.nih.gov/32079532

2. https://royalsocietypublishing.org/doi/10.1098/rsta.2020.0069





FAIR principles, present in the Data Policy is becoming THE community standard!



PaN Partners are also approaching the Open Access challenge!

ISIS neutron and Muon facility	no guarantee - long term archive	3 years	https://tinyurl.com/f3z hnpw3
ESRF	5 years minimum, 10 years expected	3 years	https://tinyurl.com/3rp e9vk6
ILL	5 years minimum, 10 years expected	5 years	https://tinyurl.com/2af uk755
Sol eil	5-10 years	3 years	https://tinyurl.com/48v b9f73
DESY	Not specifically defined	Not Specifically defined	https://tinyurl.com/hrr 4nzpb
FRMII	10 years	Not Specifically defined	https://tinyurl.com/tdk n67y9
HZDR	10 years	5 years	https://tinyurl.com/4br vdtuv
HZB	10 years	5 years	https://tinyurl.com/n6 2tnv62

PaNOSC - of the surveyed facilities with a specific policy for scientific data, 15 facilities have specific terms for open access to data.

There's still some work to be done:

- The existence of a given facility's data policy does not necessarily guarantee that the policy terms are fully implemented.
- This requires considerable resources for both staff and in capital investment, which are tensioned against other aspects of a facility's operation.

What are the challenges:

- What's the cost, how you make the
- Know your DATA becomes the most critical aspect



eli

FAIR Data Challenges

1. FAIRification of existing data and upgrading systems for facilities already operating



Fig. 1 The FAIRification workflow for health data. The blue boxes come from the GO FAIR (an initiative implementing FAIR data principles) process, while the green boxes are newly introduced steps to meet the specific challenges of health data. FAIR, findability, accessibility, interoperability, and reusability.

Costs challenges:

- There are no clear cost drivers PaNOSC preliminary analysis = different IT strategies => different cost impact
- Adds extra IT&C requirements
- Data Acquisition/Data Ingestion pipelines will need upgrades
- CS upgrades are surely needed.



FAIR Data Challenges

2. Some facilities are using the timing advantage, aiming to **be FAIR by Design.**

FAIR is easier for new RIs:

- Using FAIR as standard in designing and implementing Scientific Data Management Tools
- Increases the load on CS and DaQ groups.

Challenges:

- Computing models are not yet well-known
- New users, new physics, new requirements
- Raw data is challenging. What is raw data?

How it is for ELI?

We aim to start with FAIR data for FAIR experiments!

- Data policy framework¹ – PaNOSC

- PaNOSC guidelines on Best Practices implementing a research data poliy²

13(1) 'Data' refers to all information collected by USERS and the staff while performing scientific experiments under the ACCESS FOR USERS Policy and performing operations of the ELI FACILITIES.

13(2) **Open Access to FAIR data sets and metadata stored in Open Access repositories shall be favored** for data collected as a result of the use of the ELI FACILITIES and, to the extent possible in case of software and computer programs created by the ELI ERIC and the ELI FACILITIES; open source principles shall be considered.

- <u>https://www.panosc.eu/data/panosc-data-policy-framework/</u>
- 2. https://www.panosc.eu/news/guidelines-for-implementing-a-research-data-policy-released/

Research Data Management Plans During The Project Life Cycle, describes Embargo and impacts the Data Retention!

Why should you create a DMP?

Check you have the necessary equipment and support in place
Decreasing the risk of duplication, data loss or data security breach
Ensuring that the data are complete, accurate and reliable
Saving time and energy (e.g., when searching the data, writing up papers, etc.)
Some funders require a DMP (e.g., Horizon 2020)

Most importantly, the DMP helps you prepare and understand the Data Challenges!

•Do you have enough storage, or will you need to include charges for additional services?

•Who will be responsible for backup and recovery?

•What are the risks to data security and how will these be managed?

•How will you ensure that collaborators can access your data securely?

The Data Policy shows what's available, first user requirements are collected in a DMP!

eli

What FAIR brings, what to publish? Science Article "Narrative" In diffraction applications, models such as molecular structures are then derived from the processed data. (For HEP = users' data @CERN, sometimes referred to as data analysis results) Derived Data Processed (also sometimes termed "reduced") to correct for artefacts and convert to Processed Data scientifically meaningful units "Raw" experimental data are read directly from sensors or detector, with no or little conversion. Raw Experimental Data (Ground Truth) Fig1 – Data Pyramid

Union of Crystallography described the importance of quality "the essential component of openness is that the data supporting any scientific assertion should be

• complete (i.e. all data collected for a particular purpose should be available for subsequent re-use); and

• **precise** (the meaning of each datum is fully defined, processing parameters are fully specified and quantified, statistical uncertainties evaluated and declared)."

Fig1 Zenodo. https://doi.org/10.5281/zenodo.5155882 - Gotz, Andy, Helliwell, John, Richter, Tobias, & Taylor, Jonathan. (2021). The vital role of primary experimental data for ensuring trust in (Photon & Neutron) science.

Sharing is more than Data! Software tools! The six principles of open science are

https://open-science-training-handbook.gitbook.io/book/open-sciencebasics/open-research-software-and-open-source

Users develop their own data analysis tools:

- Some are available we need to collect (DMP?)
- Some not open could be translated/improved
- Again, it all depends on the scientists!

- Open methodology
- Open source
- Open data
- Open access
- Open peer review
- Open educational resources

By Robbielan Morrison - Own work, CC BY 4.0,

https://commons.wikimedia.org/w/index.php?curid=100144897

Open Science – driven by the users!

"The goal is to turn data into information, and information into insight." – Carly Fiorina (former CEO HP)

Carly Strasser, who oversees the foundation's Data-Driven Discovery Initiative. "Open science, data sharing, software sharing is the future of science," she says. "It's only going to get more difficult to engage in science, without the ingcopenticles/nj7584-117a

Figure 40: Open data policies by country

United Kingdom United States Canada

> Ireland Sweden

Scientists are stakeholders but, at the s time, they are the **KEY CONTRIBUTOR!**

Figure 7: Have you done any of the folle

Figure 42: Journals by code sharing modality

eli

Number of open data policies, by type of manuate and

COUNTRY Source: Sherpa-Juliet – Reference date: October 21st, 2019

https://ec.europa.eu/info/sites/default/files/research_and_innovation/knowledge_publications_tools_and_data/documents/ec_rtd_open_science_monitor_fina l-report.pdf

What FAIR and Open Data requires? Champions and use-cases!

This use case is to link raw data in the ESRF data repository to the PDBe entries as shown in the

example here.

PDB 6gv0 coloured by chain and viewed from the front

Experimental raw data

Links to raw experimental data available for this entry are listed below

Raw experimental data related to PDB entry 6gv0:

Data DOI: 10.5281/zenodo.4456817

Dataset type: diffraction image data

Generalisation of the use case

All PaNOSC and ExPaNDS partners could add the adapter for linking PDB entries to raw data. Thereby making PDB entries FAIRer.

Resources

•ESRF data repository + PDBe adapter

<u>pDB 6gv0 structure summary</u>

Description of needs

Modify PaNOSC data repositories to add support for linking PDBe entries to raw data.

Use case action flow

1.ESRF data repository – add javascript adapter for PDBe

2.PDBe repository – test adapter

3. Scientists – add DOI to raw data to PDBe entry

Impacts from the implementation

Make Protein Data Bank entries FAIR by providing access to raw data. This will enable results to be verified and structures to be refined with new software. This ICUR CommDat committee has been a strong advocate of making raw data accessible for crystallography [1]

Summary and conclusions.

PaNOSC

- 1. Data Policy updates ads an interesting change the extension to processed data for publications!
- 2. Fairification requires time and planning(new challenges/new requirements/new cost drivers). Easier for "younger" facilities.

What is needed to approach FAIR/Open:

- RAW Data, Curated-processed Data, Standard Data formats,...
- It all starts with Data and Use Cases!

Science produces much more than Publications

