# Andy Götz (ESRF)

## 7 December 2018

**A. Gotz – PaNOSC project – Laserlab workshop 7 December 2018.**

The European Synchrotron

The new **FAIR bible**

https://doi.org/10.2777/1524

The European Synchrotron

## 2.2 Definition of FAIR

https://doi.org/10.2777/1524

**The FAIR guiding principles:** https://doi.org/10.1038/sdata.2016.18

*To be Findable:*

F1. (meta)data are assigned a globally unique and persistent identifier

F2. data are described with rich metadata (defined by R1 below)

F3. metadata clearly and explicitly include the identifier of the data it describes

F4. (meta)data are registered or indexed in a searchable resource

*To be Accessible:*

A1. (meta)data are retrievable by their identifier using a standardized communications protocol

A1.1. the protocol is free, open and universally implementable

A1.2. the protocol allows for an authentication and authorization procedure, where necessary

A2. metadata are accessible, even when the data are no longer available

*To be Interoperable:*

I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation

I2. (meta)data uses vocabularies that follow FAIR principles

I3. (meta)data include qualified references to other (meta)data

*To be reusable:*

R1. (meta)data are richly described with a plurality of accurate and relevant attributes

R1.1. (meta)data are released with a clear and accessible data usage license

R1.2. (meta)data are associated with data provenance
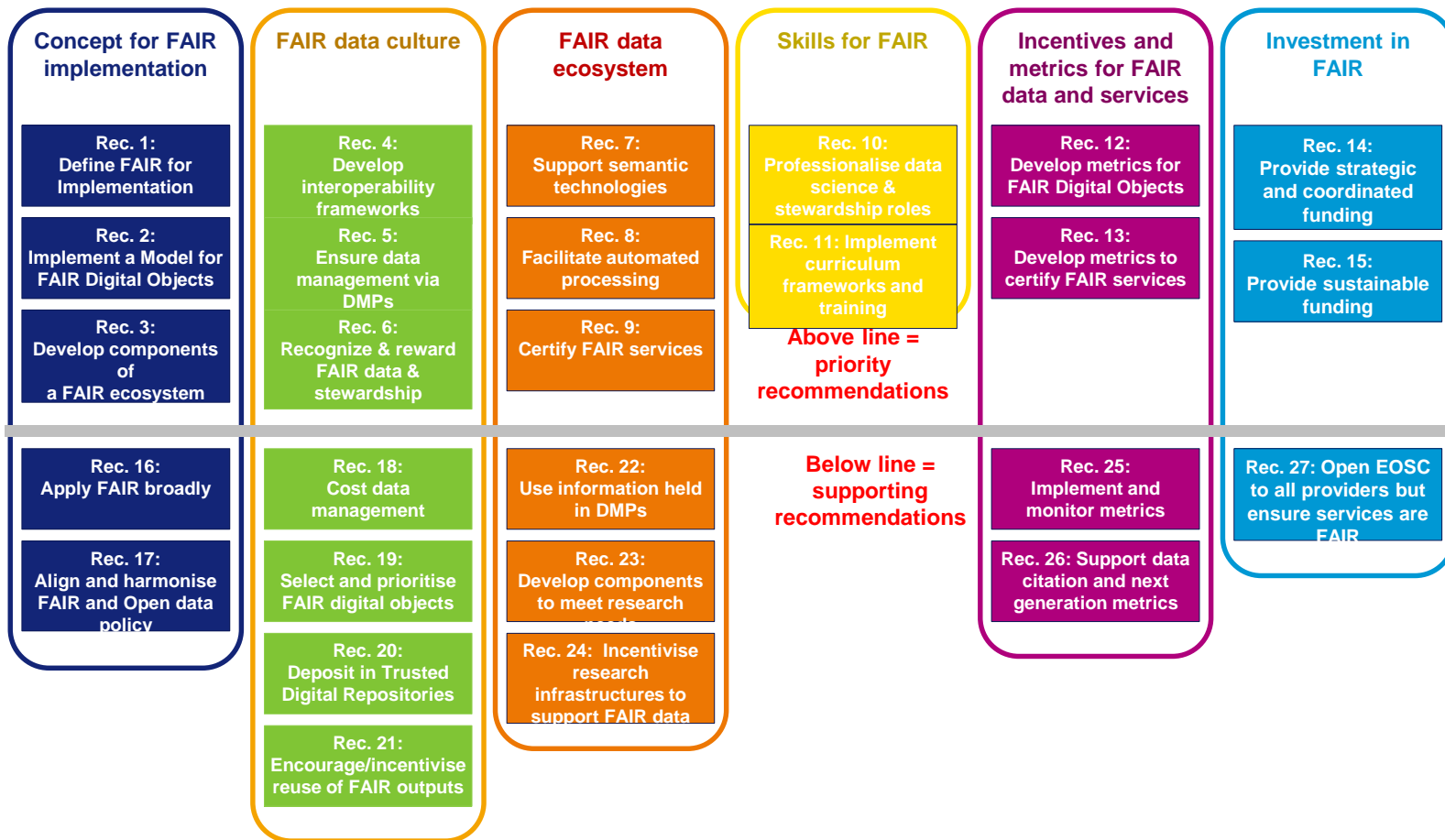
R1.3. (meta)data meet domain relevant community standards

Figure 2. The FAIR guiding principles

- **DOI**
- **Nexus**
- **E-logbook**
- **Metadata Catalogue**

- **http**
- **ResourceAsync**
- **AAI**
- **Metadata Catalogue**

- **Nexus**

- **CC BY**
- **ESRF DOI**

zon 2020 research and
The European Synchrotron

27 recommendations !

**Define** → **Implement** → **Embed and sustain**

**Concept for FAIR implementation**

- **Rec. 1:** Define FAIR for Implementation
- **Rec. 2:** Implement a Model for FAIR Digital Objects
- **Rec. 3:** Develop components of a FAIR ecosystem
- **Rec. 16:** Apply FAIR broadly
- **Rec. 17:** Align and harmonise FAIR and Open data policy

**FAIR data culture**

- **Rec. 4:** Develop interoperability frameworks
- **Rec. 5:** Ensure data management via DMPs
- **Rec. 6:** Recognize & reward FAIR data & stewardship
- **Rec. 18:** Cost data management
- **Rec. 19:** Select and prioritise FAIR digital objects
- **Rec. 20:** Deposit in Trusted Digital Repositories
- **Rec. 21:** Encourage/incentivise reuse of FAIR outputs

**FAIR data ecosystem**

- **Rec. 7:** Support semantic technologies
- **Rec. 8:** Facilitate automated processing
- **Rec. 9:** Certify FAIR services
- **Rec. 22:** Use information held in DMPs
- **Rec. 23:** Develop components to meet research needs
- **Rec. 24:** Incentivise research infrastructures to support FAIR data

**Skills for FAIR**

- **Rec. 10:** Professionalise data science & stewardship roles
- **Rec. 11:** Implement curriculum frameworks and training

**Above line = priority recommendations**

**Below line = supporting recommendations**

**Incentives and metrics for FAIR data and services**

- **Rec. 12:** Develop metrics for FAIR Digital Objects
- **Rec. 13:** Develop metrics to certify FAIR services
- **Rec. 25:** Implement and monitor metrics
- **Rec. 26:** Support data citation and next generation metrics

**Investment in FAIR**

- **Rec. 14:** Provide strategic and coordinated funding
- **Rec. 15:** Provide sustainable funding
- **Rec. 27:** Open EOSC to all providers but ensure services are FAIR

# FAIR – don't be put off …



People who say it cannot be done should not interrupt those who are doing it.

— George Bernard Shaw —

AZ QUOTES

The European Synchrotron
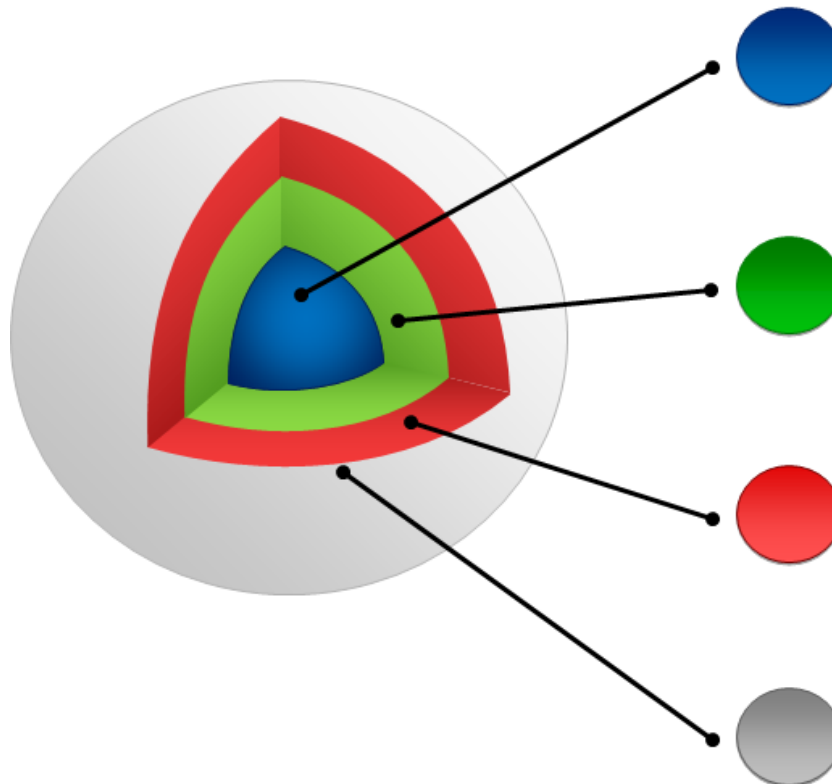
# FAIR – in simple words

**GOAL:** Making your data FAIR means managing your data in a professional manner so that it can be found, accessed, understood and re-used in the future (10 years or more) by scientists who were not involved in producing the data.

**ADVANTAGES:** By making data FAIR (1) the quality of the metadata are improved, (2) the data can be referenced by journals and data search machines, (3) the data are better organized, (4) data integrity is ensured, (5) fraud is more difficult, (6) services can be built on top of data repositories, (7) re-use of data is enabled, (8) new algorithms can be developed more easily, (9) the origin of the data are correctly acknowledged, (10) metrics on data can be tracked.

**DISADVANTAGES: implementing FAIR requires human and infrastructure resources.**

The European Synchrotron

# FAIR – Data Objects

**DATA:** are not simply data in files any more but are now data objects with rich metadata, following standards, linked to software codes and identified by persistent identifiers



**DATA**
**The core bits**
*At its most basic level, data is a bitstream or binary sequence. For data to have meaning and to be FAIR, it needs to be represented in standard formats and be accompanied by Persistent Identifiers (PIDs), metadata and code. These layers of meaning enrich the data and enable reuse.*

**IDENTIFIERS**
**Persistent and unique (PIDs)**
*Data should be assigned a unique and persistent identifier such as a DOI or URN. This enables stable links to the object and supports citation and reuse to be tracked. Identifiers should also be applied to other related concepts such as the data authors (ORCIDs), projects (RAIDs), funders and associated research resources (RRIDs).*

**STANDARDS & CODE**
**Open, documented formats**
*Data should be represented in common and ideally open file formats. This enables others to reuse the data as the format is in widespread use and software is available to read the files. Open and well-documented formats are easier to preserve. Data also need to be accompanied by the code use to process and analyse the data.*

**METADATA**
**Contextual documentation**
*In order for data to be assessable and reusable, it should be accompanied by sufficient metadata and documentation. Basic metadata will enable data discovery, but much richer information and provenance is required to understand how, why, when and by whom the data were created. To enable the broadest reuse, data should be accompanied by a 'plurality of relevant attributes' and a clear and accessible data usage license.*
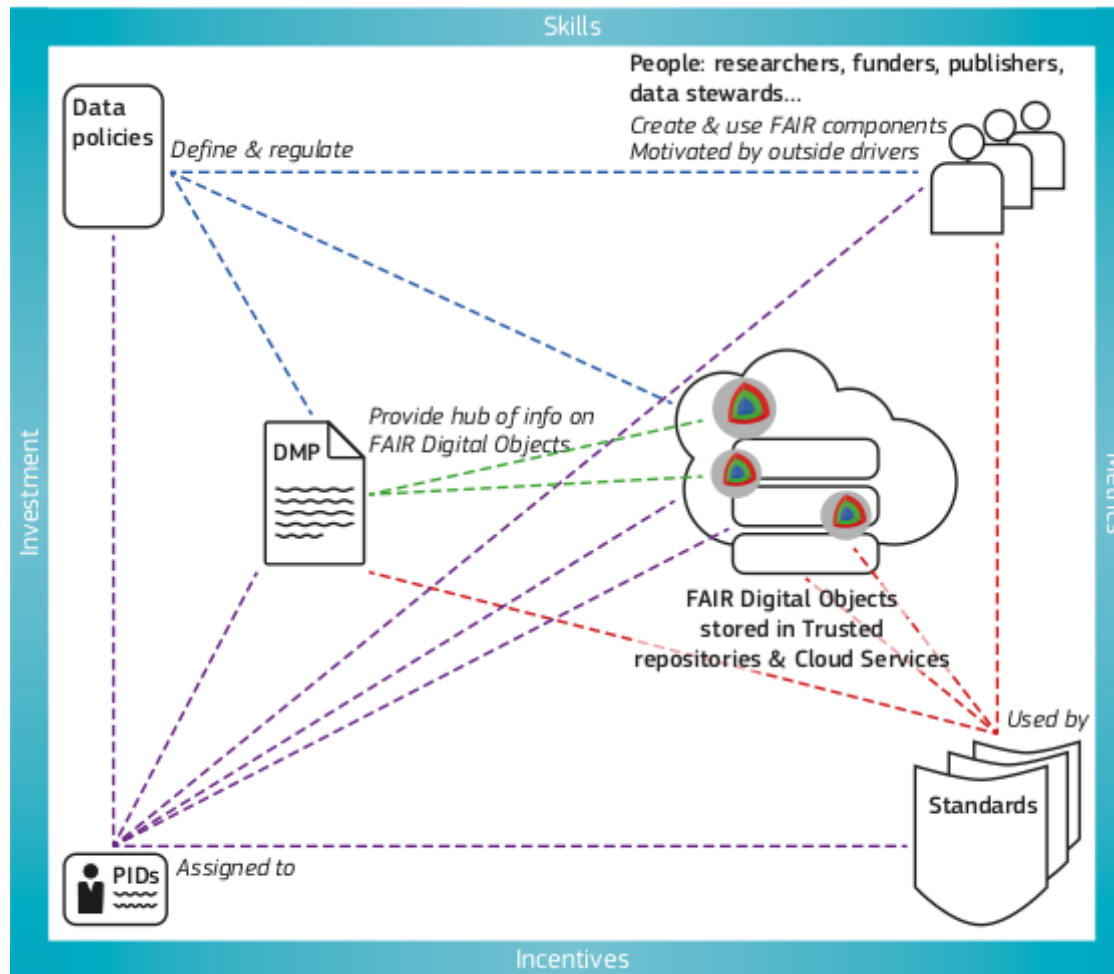
0 research and
bean Synchrotron

Figure 9. The interactions between components in the FAIR data ecosystem. Notes on this figure:

The European Synchrotron

# ESRF's long and winding road to FAIR



**2011 -** Developed Data Policy framework

**2015 -** Data Policy endorsed

**2014 -** Hired Data manager

**2014 -** Built a prototype

**2015 -** Tested on beamline

**2018 -** Implement DOIs

**2015-2018 -** Rollout on ½ beamlines

**2017 -** Hired 2nd Data manager

**2018 -** Developed e-logbook

**2019 -** Train scientists

**2021 –** Rollout on all beamlines

**2021-future** – Develop data services

The European Synchrotron

DOI > 10.15151/ESRF-DC-142893590

Data collection

Dataset  Open access

## STRUCTURAL EVIDENCE FOR A ROLE OF THE MULTI-FUNCTIONAL HUMAN GLYCOPROTEIN AFAMIN IN WNT TRANSPORT

Andreas Naschberger ; Matthew W. Bowler ; Bernhard Rupp.

**DOI**
DOI  10.15151/ESRF-DC-142893590
**Licence (for files)**
Creative Commons Attribution 4.0

### Abstract

Afamin, a human plasma glycoprotein and putative transporter of hydrophobic molecules, has been shown to act as extracellular chaperone for poorly soluble, acylated Wnt proteins, forming a stable, soluble complex with functioning Wnt proteins. The 2.1-Å crystal structure of glycosylated human afamin reveals an almost exclusively hydrophobic binding cleft capable of harboring large hydrophobic moieties. Lipid analysis confirms the presence of lipids, and density in the primary binding pocket of afamin was modeled as palmitoleic acid, presenting the native O-acylation on serine 209 in human Wnt3a. The modeled complex between the experimental afamin structure and a Wnt3a homology model based on the XWnt8-Fz8-CRD fragment complex crystal structure is compelling, with favorable interactions comparable with the crystal structure complex. Afamin readily accommodates the conserved palmitoylated serine 209 of Wnt3a, providing a structural basis how afamin solubilizes hydrophobic and poorly soluble Wnt proteins.

| Proposals | Beamlines | Publication year |
|-----------|-----------|------------------|
| OPID-1 | ID30A1 | 2018 |

---

### Experimental report

There is currently no experimental report.

---

### Experimental data

The data can be accessed by clicking on the link below

Access data

# FAIR – Persistent Identifiers

**PID :** A persistent identifier is required to be able to refer to a data object in a permanent way i.e. independent of changing urls

## We did it this way

1. Chose datacite as PID provider (**datacite.org**)

2. Setup a contract with Datacite (a PID cost between 3 – 20 cents / year)

3. Chose and implemented a long term archiving solution (tape library for 90 PBs costs roughly 100 000 euros / yr)

4. Defined and collected metadata for experimental techniques

5. Chose and deployed a metadata catalogue (**icatproject.org**)

6. Setup a workflow to archive data with the correct metadata

7. Setup a web services for creating the landing page for the PIDs (**doi.esrf.fr**)

The European Synchrotron

# FAIR – Metadata

# What is a DMP?

A short plan that outlines:

- what data will be created and how

- how it will be managed (storage, back-up, access...)

- plans for data sharing and preservation

The European Synchrotron

# FAIR – Data Repositories

**Data Repository**: where your data is stored and are FAIR

**Public repositories** : zenodo, figshare, …

**Community repositories**: CXIDB, EMPIAR …

**Institute repositories**: ESRF, ILL, XFEL, …

**Private repositories**: are not a repository e.g. Dropbox

The European Synchrotron

# FAIR – Resources

**To implement FAIR Data Management you need at least one or more of the following human resources:**

1. **Scientist –** who produces data and does science

2. **Data scientist –** defines metadata for techniques

3. **Data manager –** IT specialist who implements data policy

4. **Data archiver –** DevOps specialist who implements data infrastructure

5. **Management –** who understands the need for Data Policy

The European Synchrotron

# FAIR – Open Data

**Some scientists misinterpret FAIR data policy to be a way to "steal their data" :**

1. **Open Data –** is data which the Scientist has made open

2. **Embargo Data–** data which is under restricted access

3. **Publicly funded –** data which was obtained for free at a publicly funded site

4. **Proprietary Data –** data the user has paid for and is private

5. **Data Services –** only make sense if there is a Data Policy and data will be open at some time

The European Synchrotron

**Implementing FAIR data is a big challenge but is worth it because it improves the quality of the data and enables Open Science**

**There are many resources out there – use them!**

**Hire a data manager !**

The European Synchrotron