# Training on Research Data Curation

Luigia Cristiano, HMC Hub Matter

Helmholtz Zentrum Berlin

# Agenda

Today

- Overview of FAIR data curation practices
- Start with data documentation
- Implementation of FAIR
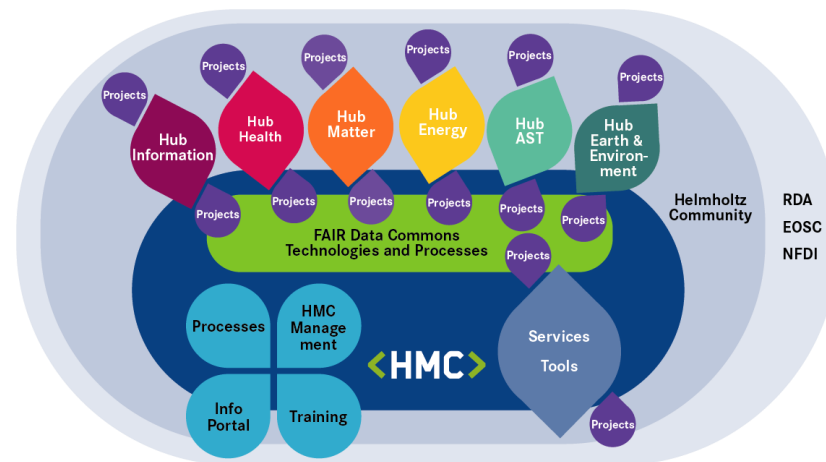
  - Metadata

  - PIDs

  - Repository and licences

Tomorrow

- HZB- HMC Hub Matter available tools for the data staging
- EMIL use case

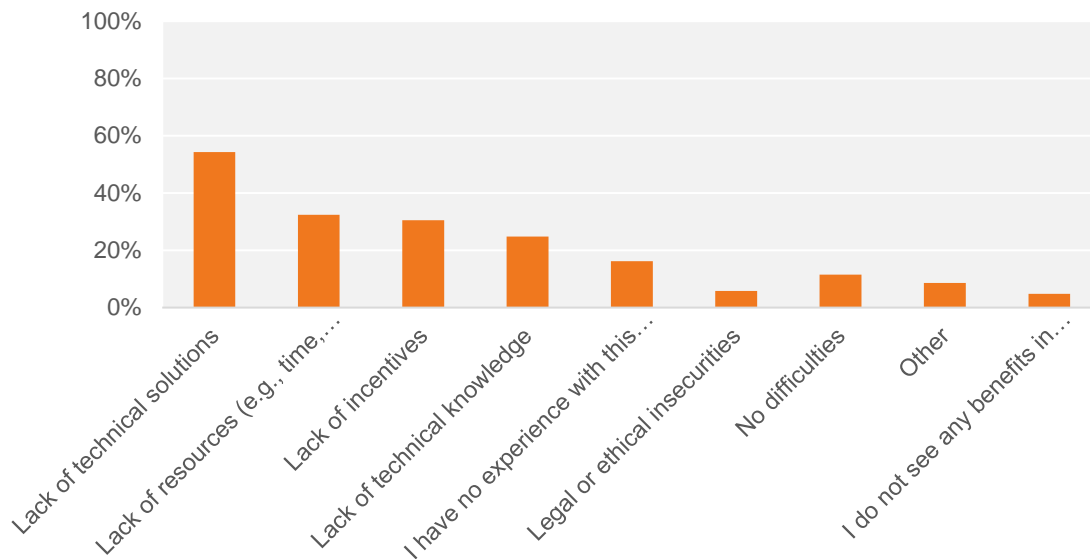**Familiarize end users with FAIR data curation**

- How to extract, format and standardize data description

- How to stage the data in institutional repository

- Discuss licences and data copyright

- Overall we want to show the benefits of FAIR data sharing and foster their use in the daily research activity

- Show the potential of the HMC-Hub Matter data staging tools

**We ease FAIR you use FAIR**



[1]

# Why getting you on board

**Researchers difficulties in making data available in repositories**

- Allocate time to describe the data
- Finding trusted repository
- Lack of specific tools and curation expertise
- Insecurity in terms of misuse, licencing and authorship
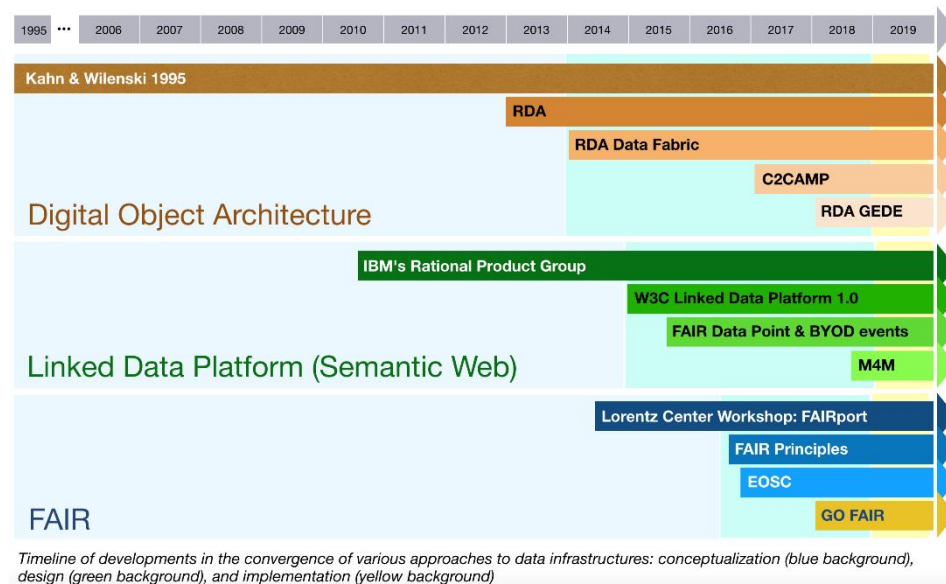- Data publication not as first class product



HMC-Hub matter focus is the development of use cases for the research data management and related training

HELMHOLTZ

# 01 **FAIR Principles and PaN data life cycle**

# FAIR: Turning data to knowledge.

## History of the FAIR initiative

- Lorentz workshop „jointly designing a Data FAIRport, 2014

- Mark D. Wilkinson et al., 'The FAIR Guiding Principles for Scientific Data Management and Stewardship,' Scientific Data 3 (March 15, 2016): 160018.)

- European task force /action plan document 2018

- Fair data maturity model, 2020



Timeline of developments in the convergence of various approaches to data infrastructures: conceptualization (blue background), design (green background), and implementation (yellow background)

[3]

HELMHOLTZ

# Why FAIR

- FAIR is not Open Access
- Different solutions and different paces of implementation are suggested for the different communities
- **Findable**

  Persistent ID

  Metadata

  Online Repository
- **Accessible**

  Community or generic repository
- **Interoperable**

  Community or generic standards

  Open file formats
- **Reusable**

  Rich, complete and clear documentation

[4] Licenced

Guidelines for promoting data visibility, reuse and ensure data quality

Wilkinson, M. D. et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci. Data 3:160018 doi: 10.1038/sdata.2016.18 (2016).

**Box 2 | The FAIR Guiding Principles**

**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 open, free, and universally implementable
A1.2 the protocol allows for an authentication and authorization procedure, where necessary
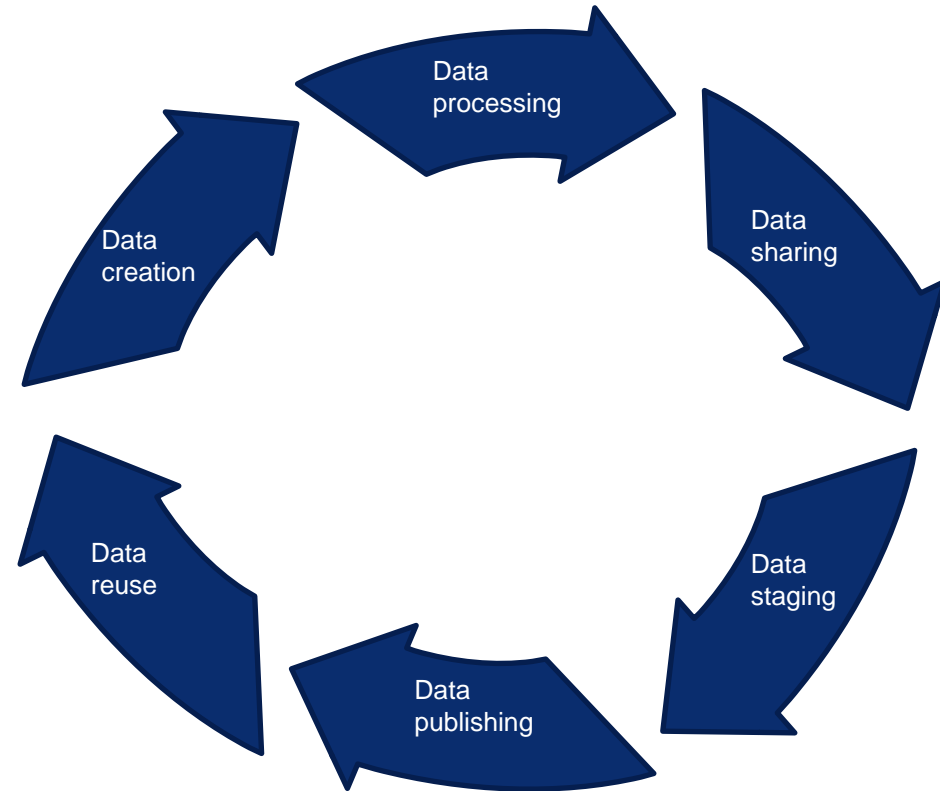A2. metadata are accessible, even when the data are no longer available

**To be Interoperable:**
I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
I2. (meta)data use 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 detailed provenance
R1.3. (meta)data meet domain-relevant community standards

**HELMHOLTZ**

# Research Data Management (RDM)

- Living Protocols and actions to optimize the research data usage

- The implementation of RDM:

   technical, organization and legal aspects must be accounted

- Consultancy is offered from data managers to researchers on:

 - Optimization of data staging in distributed locations

 - Development of tools and workflows for data curation

 - Guideline to implement data management

 - Scale down the complexity

 - Identify in advance the bottlenecks

 - solutions for the estimated the data volume and data types

   - data access solutions

[5]

Data processing

Data sharing

Data staging

Data publishing

Data reuse

Data creation

**HELMHOLTZ**

**Why a good data management in research is important**

- Encourage high quality data
- Foster/enable the data discovery and reuse and support efficient data usage with machine readability and data linking
- Enhance visibility of the research results and increase funding and cooperation opportunities
- Access to data and processing codes is requested for paper submission
- EU projects are incentivating FAIR by funding the efforts
- Ensure trasparency and reduce costs

| Data Statement | Funding Agency | | | |
|---|---|---|---|---|
| | H2020 | ERC | DFG | BMBF |
| Open Access Policy | as open as possible, as closed as necessary | as open as possible, as closed as necessary | as open as possible, as closed as necessary | as open as possible, as closed as necessary |
| DMP requested | within first six months of project | as part of the proposal | as part of the proposal | as part of the proposal |
| Template available? | yes | yes | yes | yes |
| Which data should be made available? | all data and metadata | data collections and metadata | reusable raw and structured (meta)data | |
| When should data be made available? | as soon as possible, embargo possible | as soon as possible, embargo possible | as soon as possible | as soon as possible, at the latest 9 months after completion of project |

https://dataservices.gfz-potsdam.de/portal/drr.html

[6]

# How to FAIR

**Implementation Challenges for Researchers**

- How to plan

- How to create suitable structures to host data, ensure the access, the privacy and licensing of the data and the versioning

- How to allocate time and resources for it

**Practices for data curation:**

**Data description and staging**

- Write a clear documentation for the data interpretation and reuse

- Use community standards for data and metadata when available

- Use repositories for staging the data

- Automatize the data processing and curation

- Assign PIDs to datasets, software and cite them in your publication

**HELMHOLTZ**

# How to start

## DO FAIR

Data management can determine the impact of your research data

Reduce the complexity by using machine readable metadata and **standard** formats,

Ensure the accessibility to your data and implement practice to data reuse (**licences**)

Use a Trusted **data repository** and **author** your data assigning **PID**s

**HELMHOLTZ**

02 **Helmholtz Metadata Collaboration**

- **The FAIR principles will be guiding** the development of data management services

The mission of the Helmholtz Metadata Collaboration (HMC) is to facilitate the discovery, access, machine readability, and reuse of research data of the Helmholtz Association.

HMC promotes a coordination with the scientific community for the services  development in order to establish widely accepted practices in the handling of research data.

Implementation of standards for the data description to favour the data interoperability and discovery are also goals of our initiative

[7]

# FAIR implementation Profile

## Researcher and data provider

- Curate data description
- Use standard formats
- Follow the FAIR guidelines and contribute to the services optimization
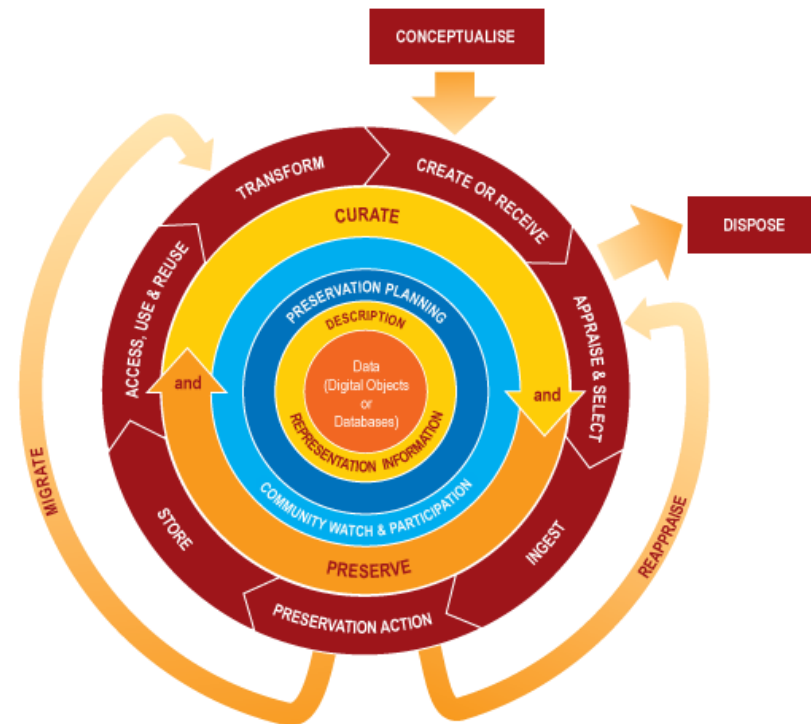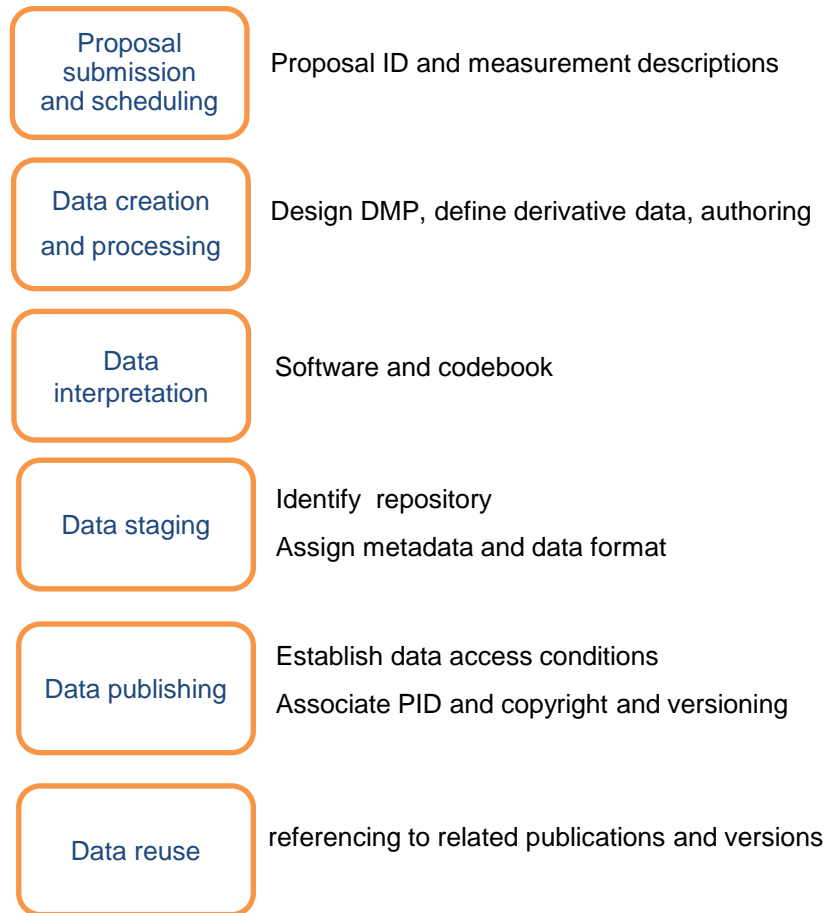- Cite datasets using DOI

**HMC and reserchers community**

**work together for a GO BUILD**

## Institutional data manager

- Raise awareness of community standards
- Support researchers with tailored services
- Ensure standard protocols to data access and the compliance to data policy
- Associate an ID to the data
- Enable the data visibility
- Promote good practices for research data management
- Scale down the complexity

**Heike Görzig head of the FDM group at HZB**

[8]

**HELMHOLTZ**

# PaN data life cycle

**Proposal submission and scheduling** — Proposal ID and measurement descriptions

**Data creation and processing** — Design DMP, define derivative data, authoring

**Data interpretation** — Software and codebook

**Data staging** — Identify repository
Assign metadata and data format

**Data publishing** — Establish data access conditions
Associate PID and copyright and versioning

**Data reuse** — referencing to related publications and versions

[9]

**HELMHOLTZ**

# Data documentation

## Study level documentation

- Project description and scientific background
- Data acquisition settings
- Data acquisition workflow
- Samples preparation and tracking
- History of the tests and implemented changes
- Define the derivative products
- Assign quality assessment procedure
- Funding program

## Data level documentation

- Content of the data and how can be reused
- Authors, acquisition time, usage licence and data location
- Data creation procedure, indication of raw or derived data
- Instruments characteristics
- Quality assessment
- Link to related datasets

[10]

# Data description, tabular : Mandatory- Optional-Ancillary

- Readme for the datafile
- Name/institution/email/ORCID
- Contributors and authors, contact person
- Acquisition time and parameters
- Geographic location/ sample, source, instrument info
- Language information
- Keywords to describe data topic
- Information about funding program and agencies
- Versioning
- Link to related data collection
- Information about data type, variables definition, unit of the measured quantities
- Explanation on acquisition settings and proc. soft.
- Uncertainty associated, gaps
- Short abstract

**Tools**

- Elab book
- Parameter files
- Activity reports
- Files header
- Data processing workflow and logfiles
- Codebooks (DDI)
- Data acquisition software
- Data management software
- Instruments log files

**Automated generation of data documentation along with data staging workflows (from measurements to publication)**

[11]

# Data formats and data structures

## Formats

- Prefer open
- Community standards
- Look at the repository recommendation
- Machine readable (metadata)

## Structures

- Keep the versioning in the filenaming
- File versioning system
- Codes for format conversion
- Associate readme, codebook, data dictionary and data list

[12]

# Data documentation workflow

- **Attach-embed-link** to the data all the information relative to the data production and data analysis.

- **Log and parameter** files help the tracking of the analysis settings and facilitate the reproducibility of the data processing

- Keep **updated** the documentation with an history of the data processing

- Use of scripting to **automatize** the update and associate quality parameters to the data analysis

- Use templates or define one. Use structured doc

- Optimize the data exploration and the data processing

- Which structure is optimal for your data set depends on the analysis tools, the data type and the data processing you need to perfom

- The access protocol and also the access terms to the data (raw or processed) plays also a key role on deciding which database structure is the most suitable for your research work

- **SQL versus noSQL databases**

  (training module on this topic is planned, register to our mailing list if interested)

- **Svn vs Git repository**

  (training module on this topic is planned, register to our mailing list if interested)

# Data organization

**General guidelines for the data organization:**

Raw data---optimized for the access and back-up

**Projectname_Date_time_instrument_analysis**

Processed data:

- Clone of the raw data structure

- Folder naming associated to the analysis method

- Versioning of the file, parameter and log files saved in the folder

- Readme and vocabulary for abbreviations

- Keep the syntax machine readable

**Output_folder**

**V?_XSF_..._..._...**

**V1.log**

**Parameters.txt**

**Instrument_analysis_v1.txt**

**Define a file name convention facilitates the data discovery**

**Complete data description should be available in each folder**

**Scripts**

# How prepare data to FAIRness

- Identifier to the data, repository, metadata, licences
- Standard or not proprietary formats
- Keep data documented and versioned
- Project or community repository
- Cross references to related data

# Data about data

**How to ensure your data are understable and usable by others**

- Machine and human redable data documentation: structured information

- Metadata organized in keys and values are descriptors to data discovery, interpretation and analysis

- Metadata keep info on data provenance, copyrights and access restrictions

- Metadata can be generic or discipline specific

**Metadata are remaining even when the data are gone**

**Metadata: different levels of informations**

| |
|---|
| **General discovery: discipline, authors, institute** |
| **Cross domain keys** |
| **Domain specific keys** |

- Minimum set of data should allow the discovery and citation- Mandatory

[13]

**HELMHOLTZ**

# Metadata categories

**Administrative**    Authors ID, dataset id, time frame, embargo, versioning

**Descriptive**    Methodological and geog. Info, data quality metrics

**Structural**    Discipline dependent

**Ancillary informations**    Additional material, or derived metrics

[14]

**HELMHOLTZ**

# Metadata Structure

Metadata elements (categories/keys) and associated values Defined by ontologies

Hierarchical structure supported by a taxonomy

Dictionary to extract the components (xml,json)

Thesaurus for the definition of categories

A rough illustration of the semantic gradient

Stronger semantics

**Ontologies**

**Taxonomies**

**Controlled vocabularies**

**Data models**

**Thesauri**

**Glossaries**

Weaker semantics

Modified from McCreary D (2006)
Patterns of Semantic Integration. CC 2.5

The metadata can be distributed attached to the data as 1. „header" or as
2. separated file.
Both options offer advantages and disadvantages

[15]

**HELMHOLTZ**

# Metadata templates

**HELMHOLTZ METADATA COLLABORATION**

## Generic

Dublin Core, DataCite, Schema.org,

EUDAT Core

[The Simple Dublin Core Metadata](#)

1. Title
2. Creator
3. Subject
4. Description
5. Publisher
6. Contributor
7. Date
8. Type
9. Format
10. Identifier
11. Source
12. Language
13. Relation
14. Coverage
15. Rights

## Formats

Xml, csv, json, html

```
-->
-<xs:schema targetNamespace="http://datacite.org/schema/kernel-4" elementFormDefault="qualified" xml:lang="EN">
  <xs:import namespace="http://www.w3.org/XML/1998/namespace" schemaLocation="include/xml.xsd"/>
  <xs:include schemaLocation="include/datacite-titleType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-contributorType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-dateType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-resourceType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-relationType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-relatedIdentifierType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-funderIdentifierType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-descriptionType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-nameType-v4.xsd"/>
  <xs:include schemaLocation="include/datacite-numberType-v4.xsd"/>
  -<xs:element name="resource">
    -<xs:annotation>
      -<xs:documentation>
        Root element of a single record. This wrapper element is for XML implementation only and is not defined in the Da
        within this schema.
      </xs:documentation>
      <xs:documentation>No content in this wrapper element.</xs:documentation>
    </xs:annotation>
    -<xs:complexType>
      -<xs:all>
        <!--REQUIRED FIELDS-->
        -<xs:element name="identifier">
          -<xs:annotation>
            -<xs:documentation>
```
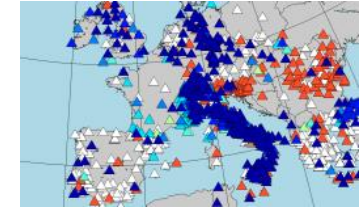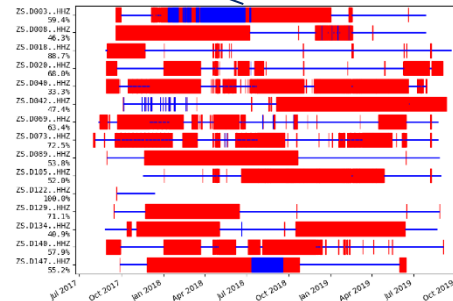
**Discipline dependent** NXDL

https://github.com/nexusformat/definitions/blob/main/nxdl.xsd

[16]

[https://schema.datacite.org/meta/kernel-4.1/example/datacite-example-full-v4.1.xml](#)
[https://www.fdsn.org/xml/station/fdsn-station-1.1.xsd](#)

https://schema.datacite.org/meta/kernel-4.4/metadata.xsd

**HELMHOLTZ**

# How can be used the metadata information



**Data analysis and quality control**

**Data geographic distribution**

**Data exploration**

**Data completness and availability**

```
<FDSNStationXML schemaVersion="1.1">
  <Source>SeisComP</Source>
  <Sender>GFZ</Sender>
  <Created>2021-11-20T01:53:36.667865</Created>
  <Network code="ZS" startDate="2017-01-01T00:00:00" endDate="2019-12-31T00:00:00" restrictedStatus="closed">
    <Description>
      AlpArray Seismic Network (SwathD) temporary component
    </Description>
    <Identifier type="DOI">10.14470/MF7562601148</Identifier>
    <Comment id="0">
      <Value>Grant GIPP201717</Value>
    </Comment>
    <Station code="D159" startDate="2018-10-23T16:28:00" endDate="2019-12-31T00:00:00" restrictedStatus="closed">
      <Latitude>46.966342</Latitude>
      <Longitude>14.152825</Longitude>
      <Elevation>1035</Elevation>
      <Site>
        <Name>Metnitz, Austria</Name>
      </Site>
      <CreationDate>2018-10-23T16:28:00</CreationDate>
      <Channel code="HHE" startDate="2018-10-23T16:28:00" endDate="2019-12-31T00:00:00" restrictedStatus="closed" locationCode="00">
        <Latitude>46.966342</Latitude>
        <Longitude>14.152825</Longitude>
        <Elevation>1035</Elevation>
        <Depth>0</Depth>
        <Azimuth>90</Azimuth>
```
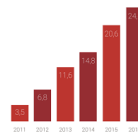
[17]

# Formats converter

https://github.com/G-Node/gogs/blob/master/conf/datacite/datacite.yml

German Neuroinformatics Node, guidelines for DOI request



**Several routines offer format conversions:**

Yaml file offer the possibility to add text blocks

It has a structure in which embedding free text environment.

Python functions/packages e.g.

```
80 lines (65 sloc) | 2.27 KB

1   # Metadata for DOI registration according to DataCite Metadata Schema 4.1.
2   # For detailed schema description see https://doi.org/10.5438/0014
3
4   ## Required fields
5
6   # The main researchers involved. Include digital identifier (e.g., ORCID)
7   # if possible, including the prefix to indicate its type.
8   authors:
9     -
10      firstname: "GivenName1"
11      lastname: "FamilyName1"
12      affiliation: "Affiliation1"
13      id: "ORCID:0000-0001-2345-6789"
14    -
15      firstname: "GivenName2"
16      lastname: "FamilyName2"
17      affiliation: "Affiliation2"
18      id: "ResearcherID:X-1234-5678"
19    -
20      firstname: "GivenName3"
21      lastname: "FamilyName3"
22
23   # A title to describe the published resource.
24   title: "Example Title"
25
26   # Additional information about the resource, e.g., a brief abstract.
```

[18]

**How to link data and metadata/data description**

- 1. Keeping the metadata embedded in the data file
- 2. Having them stored in separated files

1. Raw data : no update on the acquisition settings
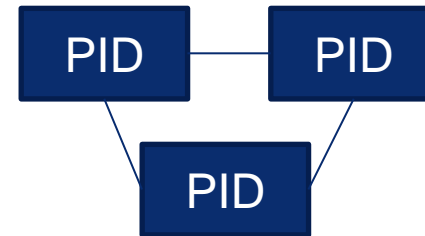1. Derivative data to preserve the link of the info

# xml editing tools

**Making research workflow FAIR:**
02 Resource Identification **PIDs**

- **Persisten identifiers** are **long lasting unique** digital reference to an object, contributor and organization.

- They act as pointers/entry points between the reference to the object and its actual location

-  They are character strings globally unique

Labeling  a dataset i.e. with a PID support the data findability: independent on the data physical location

It is managed and updated with new location in the registry

 and resolvable

```
  PID  ——  PID

        PID
```

PID attribution  (ARK, arXiv, DOI, ePIC, Handle, URN) to people and digital object.

Object for Instruments, software, workflows and labs procedures and samples

[19]

# ORCID and DOI

- **ROR** **research organization** Registry. It stores metadata about the organization

- **DOI** is built on the top of the Handle system and maintains a registry of metadata related to the object, It has local Registration agencies (DataCite for **datasets**, CrossRef for **publications** (chapters, proceedings, EIDR for audio objects). It is widely used in the publication registry and has specific metadata schemata, business model and assignment practices

- Associate the ID number to the URL hosting the data


- **ORCID (Open Researcher and Contributor ID)** registry for reseachers. Provides an ID to associate publications, grants and employments to individuals and desambiguate. It acts as record of professional activities. Journals requires that at least the corresponding authors has an ORCID

  https://orcid.org/

[20]

**HELMHOLTZ**

# Making research workflow FAIR:
## 03 **Repositories** and **Licenses**

- **Offer permanent archives and access to the data, provide persistent Identifiers (FAIR)**

- **Look for an institutional repository, project repository**

- **Generalist/ community/institutional repository:** zenodo, figshare, eudat, dryad home/ CXIDB, EMPIAR, EMDataResource, NOMAD (https://nomad-lab.eu/services/repo-arch)

- **FAIRsharing** registry of repositories https://fairsharing.org/ structured overview

- RatSWD https://www.konsortswd.de/datenzentren/alle-datenzentren/

  Collection of german data centers

- Re3data.org international registry of research data repository

https://www.re3data.org/search?query=&subjects%5B%5D=30701%20Experimental%20Condensed%20Matter%20Physics

- Measure of trustworthy of repository given (e.g. certificate, PIDs, metadata)

# License

- **Terms of use of your data**
- **Data proper attribution**
- **Data and databases for the licence**

**Open Data Commons group**

  3 standard licenses and additional community norms

PDDL: Public domain dedication

ODC-By: free user under the provision of referencing to the source

ODC-Odbl: any use of the database must refer to the source, any derivative product is to be distributed under same terms

https://creativecommons.org/licenses/

**Creative Commons**

CC-BY-SA attribution share alike and same licencing of related work

CC-BY- Attribution only

CC0 : Public domain dedication

PDM:  Public Domain mark – free of know  restrictions

https://eu-datenschutz.org/

[23]

- **HZB data policy**

https://www.helmholtz-berlin.de/pubbin/vademecumdatei?did=326

https://www.helmholtz-berlin.de/pubbin/vademecumdatei?did=131

- Raw data :

  The raw data is under embargo for 5 years and hosted for a least 10 ys at HZB

  The embargo can be extended under request and can be shortened according to project needs

   The raw data are licences with CC0

  Results data are not affected

paN-data Europe common policy framework (february 2011) ---→ extended to FAIR by Pansoc

- Open access to raw data and metadata

- Curation of raw data supported by the facility

- Data catalogue to make data accessible

[24] - Embargo on the raw data

# Data publication guidelines

- Check the repository guidelines

- Check the publication service guidelines

- Check previous slides on data curation and metadata for publication

Data publication as supplement of peer reviews paper

Data publication on istitutional or generic repository

**+**

DOI to data , mentioned in the publication

Findable Accessible Interoperable Reusable

**HELMHOLTZ**

# Software publication guidelines

- Prepare a README file, usage description and dep.

- Provide usage examples and test integration

- Prepare a MAKE file or a setup.py file to prepare the repository clone or the installation with all dependencies included

- Associate a usage licence

- Enable the installation on OSs and give the compatibilities with the dependencies versions

- Obtain a DOI for the software e.g. DataCite and a snapshot of the gitlab

- Versioning the software

- Offer help support and possibility to start the bug fixes

- Use when possible the Gitlab /GIThub with institutional maintenance to guarantee the long term access to the software repository



https://github.com/

They both offer continuous integration and delivery and support the licences of your software



[26]

# Software publication guidelines



- No widely accepted FAIR principles as for research data (status 2019)

- Software citation principles force11.org (https://doi.org/1025490/a97f-egyk

- DataCite for software citation FORCE11.org principles mapped

**MIT Licence, Apache 2.0**

A short and simple permissive license requiring preservation of copyright and license notices.

https://github.com/G-Node/gogs/blob/master/LICENSE

Indication of modification respect the original version

**GNU General Public License v3.0** for free software www.gnu.org



https://github.com/

They both offer continuous integration and delivery and support the licences of your software



[27]

- Reconstruct the whole data life cycle

[28]

**Thanks.**

**HMC Hub Matter**

luigia.cristiano@helmholtz-berlin.de
heike.goerzig@helmholtz-berlin.de
gerrit.guenther@helmholtz-berlin.de
rolf.krahl@helmholtz-berlin.de
markus.kubin@helmholtz-berlin.de
oonagh.mannix@helmholtz-berlin.de

Helmholtz-metadata.de/en/pages-helpdesk
Hmc-matter@helmholtz-berlin.de

www.helmholtz-metadata.de

- [7] https://www.rd-alliance.org/system/files/PID-report_v6.1_2017-12-13_final.pdf; PIDs Beginners guide https://zenodo.org/record/4574566#.YXKYiedCRoR

- [9] https://zenodo.org/record/4312825#.YXFqSOdCRzp

- [1] Data management workshop –Abigail Mc Birnie

- [2] Stuart, David; Baynes, Grace; Hrynaszkiewicz, Iain; Allin, Katie; Penny, Dan; Lucraft, Mithu; et al. (2018): Whitepaper: Practical challenges for researchers in data sharing. figshare. Journal contribution. https://doi.org/10.6084/m9.figshare.5975011.v1

- Kubin et al., Report on community survey, HMC 26.10.21

- 10. S.Jones The Future of FAIR. N8 library management workshop. www.geant.org

- www.fairsfair.eu

- 28. A.Gonzalez-Beltran, Large-scale facilities experimental lifecycle and FAIRness, FAIR workshop 1-2 October 2020

- 23 creativecommons.org; opendatacommons.org

- 24 A.Mc Birnie, 2021. ExPaNDS Develops a Data Policy Framework for National Photon & Neutron Research Infrastructures. 10.5281/zenodo.5040078

- [25]  D. Salvat, 2010. Draft recommendation for FAIR Photon and Neutron Data Management. 10.5281/zenodo.4312825