# Metadata Management in Scientific Research: an overview

**Rossella Aversa (KIT-SCC)**
iEntrance Advanced School, 22.02.2024

# Outline

| 1. Introduction | 2. Effective Metadata Management | 3. Metadata Management in practice | 4. Conclusions |
|---|---|---|---|
| Motivation and recap of the FAIR principles | Best practices and tools to implement them | Top 10 management steps in real projects | Summary and main takeaways |

# 1. Introduction

# Motivation



Powered by Bing Image Creator

February 22, 2024     Rossella Aversa

# Metadata

3. How manage metadata?

1. What is metadata?

Why is metadata needed?

Data describing other data

February 22, 2024    Rossella Aversa

# Metadata

1. What is metadata?

2. Why is metadata needed?

3. How to manage metadata?

To add context and meaning to data

February 22, 2024    Rossella Aversa

# Metadata

2. Why metadata needed?

3. How to manage metadata?

What is metadata?

We will see it throughout the next slides…

February 22, 2024    Rossella Aversa

# The FAIR Guiding Principles



https://www.go-fair.org/fair-principles/

# The FAIR Guiding Principles

| Findable | Accessible | Interoperable | Reusable |
|----------|------------|---------------|----------|

(Meta)data should be easy to find for both humans and computers

Globally unique persistent identifiers (PID)

https://www.go-fair.org/fair-principles/

# The FAIR Guiding Principles

| 🔍 Findable | 👆 **Accessible** | ⚙️ Interoperable | ♻️ Reusable |
|---|---|---|---|

It should be known how (meta)data can be accessed

📌 (Meta)data repositories, authorization & authentication

https://www.go-fair.org/fair-principles/

# The FAIR Guiding Principles

| Findable | Accessible | **Interoperable** | Reusable |
|----------|------------|-------------------|----------|

Data should be exchanged and interpreted by humans and computers

📌 Structured metadata (schemas, vocabularies)

https://www.go-fair.org/fair-principles/

# The FAIR Guiding Principles

| Findable | Accessible | Interoperable | **Reusable** |
|---|---|---|---|

It should be clear how data can be reused and/or replicated

Licences, rich (provenance) metadata

https://www.go-fair.org/fair-principles/

# 2. Effective Metadata Management

February 22, 2024    Rossella Aversa

# What to describe?



R1: Metadata should richly describe the data with a plurality of accurate and relevant attributes.

Sample

Sample Holder

Instrument

Research data

Images: courtesy of R. Thelen and M. Mail. Powered by Bing Image Creator

February 22, 2024    Rossella Aversa

# How to describe data?

**Minimal requirements:**

- The vocabulary and its terms have globally unique PIDs
- The vocabulary and its terms are documented
- The documentation is findable and accessible by users

12: Metadata use vocabularies that follow the FAIR principles.



https://skosmos.org

# How to represent metadata?

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

**Resource Description Framework:** metadata model to represent interconnected data.
https://www.w3.org/RDF/

**Simple Knowledge Organization System:** standard to represent knowledge organization systems using RDF
https://www.w3.org/2004/02/skos/

**Web Ontology Language:** computational logic-based language to represent complex knowledge.
https://www.w3.org/OWL/

# How to structure metadata?

R1.3: Metadata meet domain-relevant community standards or best practices.
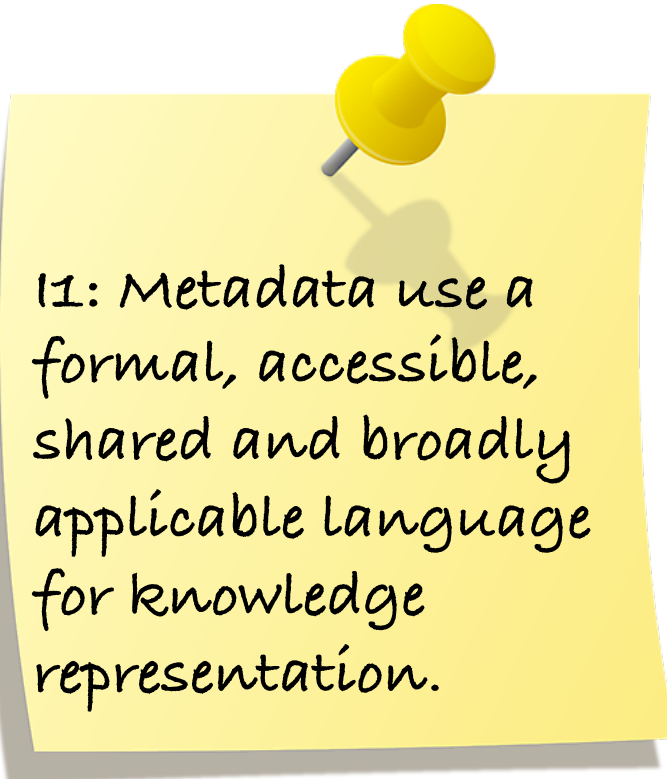
## General purpose

DublinCore — http://dublincore.org/schemas/

DataCite — http://schema.datacite.org

Schema.org — https://schema.org

## Neutron, x-ray, muon

NeXus — http://www.nexusformat.org

## Crystallography

CIF — https://www.iucr.org/resources/cif

February 22, 2024    Rossella Aversa

# How to represent structured metadata?

I1: Metadata use a formal, accessible, shared and broadly applicable language for knowledge representation.

XML (eXtensible Markup Language)

JSON (JavaScript Object Notation)

JSON Schema    https://json-schema.org

JSON-LD (JSON for Linked Data)

February 22, 2024    Rossella Aversa

# How to publish (meta)data?

F1: (Meta)data are assigned globally unique and persistent identifiers.

## Cite this article

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

### Details

**DOI**
DOI 10.5281/zenodo.7778338

**Resource type**
Dataset

**Publisher**
Zenodo

**Languages**
English

ORCID
Connecting research and researchers

https://orcid.org/
**0000-0003-2534-0063**
Preview public record

# How to publish metadata?

A2: Metadata should be accessible even when the data is no longer available.

Not Found

The requested ... found on this ser...

Apache Server

Something Went Wrong

:/
Oops, this URL is invalid...

Error 404 – the requested URL

Bad Request - Invalid URL

HTTP Error 400. The request URL is invalid.

DATA    METADATA

Dataedo /cartoon

Piotr @Dataedo

Metadata Repositories!

# How to find data from metadata?

F3: Metadata clearly and explicitly include the identifier of the data they describe.

Details

**DOI**
DOI 10.5281/zenodo.7778338

**Resource type**
Dataset

**Publisher**
Zenodo

**Languages**
English

February 22, 2024     Rossella Aversa

# How to find data from metadata?



R1.2: Metadata are associated with detailed provenance.

Data acquired from measurement

Measurement performed on sample

Sample placed on holder

February 22, 2024    Rossella Aversa

# How to reuse the data?

R1.1: (Meta)data are released with a clear and accessible data usage licence.

# Should FAIR data be open?

A1.2: The protocol allows for an authentication and authorization procedure where necessary.

**Open data:** "can be freely used, modified, and shared by anyone for any purpose"
https://opendefinition.org

**FAIR data:** "as open as possible, as closed as necessary"

# FAIR or open?

| | |
|---|---|
| My data is copyright protected | **FAIR** |
| My dataset can be used only by a specific group of scientists | **FAIR** |
| An image is shared on a public website | **Open** |
| A dataset is published on Zenodo with an open licence | **FAIR**  **Open** |
| A data file is on my Dropbox | **None** |

February 22, 2024    Rossella Aversa

# 3. Metadata Management in practice

February 22, 2024    Rossella Aversa

# The projects



Nanoscience Foundries and Fine Analysis – Europe Pilot (NEP)
https://nffa.eu

Access to nanoscience research infrastructure Synthesis, growth of nanostructures, fine analysis, theory and simulation

February 22, 2024    Rossella Aversa

# The projects



Joint-Lab Model and Data driven
Materials Characterization (MDMC)
of the Helmholtz Association
https://jl-mdmc-helmholtz.de

Platform for multiscale and
multidimensional characterization,
analytics and simulation methods

February 22, 2024    Rossella Aversa

# The projects



National Research Data Infrastructure for Materials Science and Engineering (NFDI-MatWerk)
https://nfdi-matwerk.de

Infrastructure for the digital representation of materials and their relevant process

# The projects



**Common aims:**

- Implement (meta)data management practices following the FAIR principles
- Develop tools and infrastructure solutions guided by community requirements
- Agree on common descriptions
- Collaborate on interoperable results

# 1. Definition of terms



## Glossary of Terms

- High-level description of experimental and computational workflows
- Framed in the management infrastructure of the projects
- Allows to track the provenance information
- Adopts existing terms

MDMC-NEP Glossary of Terms. DOI: 10.5281/zenodo.10663833

**Research User**

Person, usually member of a **Project**, who conducts any part of the **Study**, in order to collect and/or analyse **Research Data** or is interested in reusing **Research Data** by a third party (e.g., **Reference Data**) with the final aim to extract insights that support the answer to some specific research question (i.e., **Conclusions**). **Research Users** may be assigned with a role (data curator, instrument scientist, team leader, team member).

MDMC-NEP Glossary of Terms. DOI: 10.5281/zenodo.10663833

# 2. Terms in a Vocabulary Service

| PREFERRED TERM | **User Role** 🗗 |
|---|---|
| NARROWER CONCEPTS | Data Curator |
| | Instrument Scientist |
| | Team Leader |
| | Team Member |
| URI | http://matwerk.datamanager.kit.edu:8001/DemoTerms-1/en/page/userrole 🗗 |
| DOWNLOAD THIS CONCEPT: | RDF/XML TURTLE JSON-LD |

**Reseach User**

User Name

Rossella Aversa

User Role

✓ Data Curator
Instrument Scientist
Team Leader
Team Member

## EVOKS Vocabulary Service: Collaborative online vocabulary editor
- RDF model + SKOS model
- Persistent identifier to each term
- Can be resolved in interfaces, websites, automatic processes…
- Centrally maintained
- Public read-only Skosmos instance

*Work in Progress*

# 3. Metadata schemas

- Describe inputs/outputs of processes
- JSON schema
- Adopt existing solutions
- Avoid proliferation of schemas



https://xkcd.com/927/

MDMC-NEP Glossary of Terms. DOI: 10.5281/zenodo.10663833

# 3. Metadata schemas



Minimal

System

Sample description

Input → Fabrication → Precursor → Sample Preparation → Sample ← Sample Component → Measurement → Raw Data

Ongoing collaboration with CNR-IOM and FBK (see talk by L. Ferrario)

Work in Progress

Adopted from the Materials Data Vocabulary DOI: 10.5334/dsj-2021-018

SEM
TEM
MRI
STM
SEM/FIB Tomography
Nano CT/micro CT

# 4. Metadata schemas and documents

**Metadata Schema:** outline of the overall structure of the metadata (elements, value types, rules, …)

**Metadata Document:** structured information about a data resource

```
"instrumentID": {
    "type": "string"
},
"instrumentManufacturer": {
    "type": "object",
    "properties": {
        "manufacturerName": {
            "type": "string"
        },
        "modelName": {
            "type": "string"
        },
```

```
ent": {
trumentID": "425590",
trumentManufacturer": {
"manufacturerName": "Bruker B
"modelName": "Biospec 152/11"
```

MRI schema, DOI: 10.5445/IR/1000159552

# From data to metadata



### Raw Data

### Metadata Schema

```
"instrumentID": {
    "type": "string"
},
"instrumentManufacturer": {
    "type": "object",
    "properties": {
        "manufacturerName": {
            "type": "string"
        },
        "modelName": {
            "type": "string"
        },
```

### Metadata Document

```
ent": {
trumentID": "425590",
trumentManufacturer": {
"manufacturerName": "Bruker B
"modelName": "Biospec 152/11"
```

Image from Magnetic Resonance Imaging Copper Sulfate Dataset. DOI: 10.5281/zenodo.6107720

# 5. Mapping service



- Online service
- Input: data file(s)
- Extract unstructured metadata
- Map them to the metadata schema
- Output: structured metadata



**Mapping Service UI**
Extract metadata and map it to json

The Mapping Service is a tool designed to extract metadata from different kinds of data produced by instruments, and map this metadata to published metadata schemas. Show More

**Choose a suitable mapping from available options**

| SEM to TXT | MRI to JSON | SEM to JSON |
|---|---|---|
| ...on based tool extracts ...from machine generated ...microscopy images in the ...at and generates a TXT ...ining a summary of the ...metadata. Last edited: 19.12.2023 | Takes a single .dcm or zipped directory of .dcm files and maps to the MRI schema returning a JSON metadata document. LU: 01.02.2024 | This plugin is able to handle variety of SEM images and processes them using the Hype library. A resulting metadat document in JSON format is th created. LU: 06.02.2024 |
| Select | Select | Select |

Drag & Drop your files or Browse

https://matwerk.datamanager.kit.edu/mapping-service-ui.html

# 6. Metadata Editor



- Local service connected to the metadata repository
- Load schema from registered ones
- Load existing metadata documents
- Manually edit metadata documents
- Download metadata documents
- Register metadata documents
- Create the provenance file



https://metadata-editor.gitlab.io/documentation/

# 7. ELN and LIMS



- ▪ Electronic Lab Notebooks
- ▪ Lab Information Management Systems
- ▪ Metadata schemas as templates
- ▪ Ongoing collaborations:
  - ▪ KIT
  - ▪ CNR-IOM
  - ▪ FBK (see talk L. Ferrario)

# 8. MetaStore



- Metadata repository
- Register/find metadata schemas
- Register/find metadata documents
- Validate metadata documents
- Versioning
- Access control management
- User authentication

https://metarepo.nffa.eu/

https://github.com/kit-data-manager/metastore2

# 9. Link metadata to data



February 22, 2024     Rossella Aversa

# 9. Link metadata to data

February 22, 2024   Rossella Aversa

# 10. Find data from metadata

# 10. Find data from metadata

- Use the content of metadata documents to search for relevant data
- Private vs Public resources
- What is the data about? Is it useful for my needs?
- Full-text search
- (basic, customizable) faceted search

*Work in Progress*

# 4. Conclusions

February 22, 2024    Rossella Aversa

February 22, 2024    Rossella Aversa

**Contacts:** rossella.aversa@kit.edu

**Acknowledgements to:**
G. Abdildina, N. Blumenröhr, F. Ernst, V. Haltmann, M. Inkmann, T. Jejkal, A. Kirar, E. Vitali, NEP JA6, JL-MDMC Metadata WG.