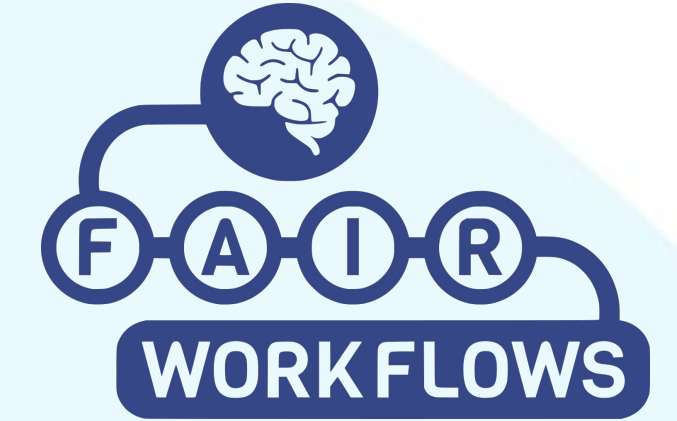




CONNECTING RESEARCH,  
IDENTIFYING KNOWLEDGE



# Towards metadata completeness

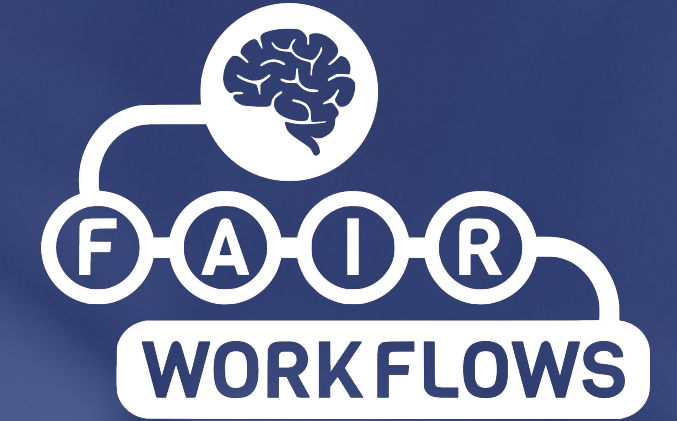
# FAIR Workflows

**Xiaoli Chen**

9. 11. 2022

DataCite Open Hour

# Agenda

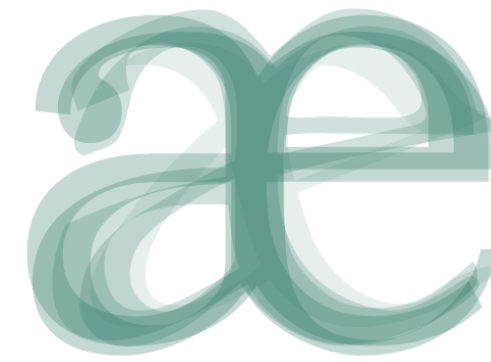


- Metadata considerations in the context of the project
- Capturing metadata
- Enriching metadata





TEMPLETON WORLD  
CHARITY FOUNDATION



Max-Planck-Institut  
für empirische Ästhetik

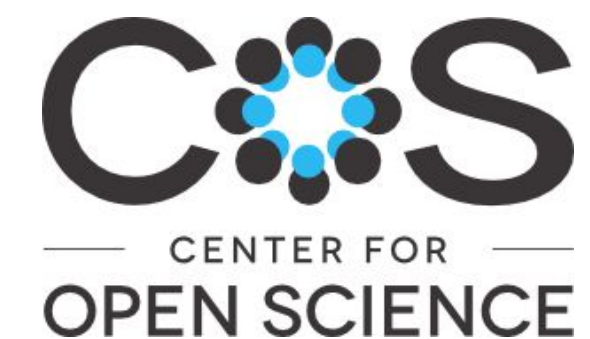


**DRYAD**

ChronosHub



CEDAR



# What makes a workflow FAIR



## FAIR Entities

- Uniquely identified resources associated to a project
  - Researcher (ORCID iD)
  - Research organization (ROR ID)
  - Funding agency (ROR ID) and grant (Grant-ID)

## FAIR Practices

- Sharing various types of interim outputs
  - Data Management Plan, Pre-registration, Protocol, Preprint, Code, Dataset, etc.

## FAIR Supporting Structures

- Tools and platforms that integrate PIDs and metadata workflow
  - DMP authoring tools, Metadata templates, Data repositories, Notebooks, Collaborative research platforms, etc.

## FAIR Outputs

- Assigning PIDs to outputs with rich metadata annotation
  - Essential descriptive and connection metadata
    - Connection between inputs and outputs
    - Relations between outputs
  - Domain specific metadata
    - Disciplinary ontological information
    - Experimental setup

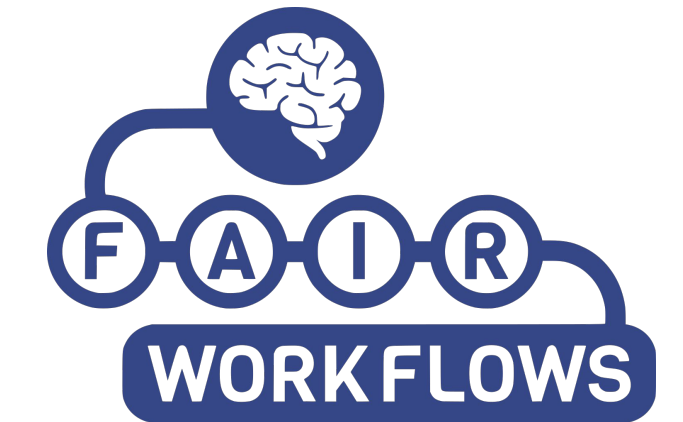
Identifier	Name
<a href="#">FsF-F1-01D</a>	<b>Data</b> is assigned a globally unique identifier.
<a href="#">FsF-F1-02D</a>	<b>Data</b> is assigned a persistent identifier.
<a href="#">FsF-F2-01M</a>	<b>Metadata</b> includes descriptive core elements (creator, title, data identifier, publisher, publication date, summary and keywords) to support data findability.
<a href="#">FsF-F3-01M</a>	<b>Metadata</b> includes the identifier of the data it describes.
<a href="#">FsF-F4-01M</a>	<b>Metadata</b> is offered in such a way that it can be retrieved by machines.
<a href="#">FsF-A1-01M</a>	<b>Metadata</b> contains access level and access conditions of the data.
<a href="#">FsF-A1-02M</a>	<b>Metadata</b> is accessible through a standardized communication protocol
<a href="#">FsF-A1-03D</a>	<b>Data</b> is accessible through a standardized communication protocol
<a href="#">FsF-A2-01M</a>	<b>Metadata</b> remains available, even if the data is no longer available.
<a href="#">FsF-I1-01M</a>	<b>Metadata</b> is represented using a formal knowledge representation language.
<a href="#">FsF-I1-02M</a>	<b>Metadata</b> uses semantic resources.
<a href="#">FsF-I3-01M</a>	<b>Metadata</b> includes links between the data and its related entities.
<a href="#">FsF-R1-01MD</a>	<b>Metadata</b> specifies the content of the data.
<a href="#">FsF-R1.1-01M</a>	<b>Metadata</b> includes license information under which data can be reused.
<a href="#">FsF-R1.2-01M</a>	<b>Metadata</b> includes provenance information about data creation or generation.
<a href="#">FsF-R1.3-01M</a>	<b>Metadata</b> follows a standard recommended by the target research community of the data.
<a href="#">FsF-R1.3-02D</a>	<b>Data</b> is available in a file format recommended by the target research community.



# What does complete mean?

- **Completeness is one aspect of quality**
  - Accuracy
  - Provenance/authority
  - Accessibility
  - Timeliness
  - ...
- **Completeness is only meaningful when measured against a standard**
  - Community recommendation
  - Use case requirements
  - Local context

Ad  
e

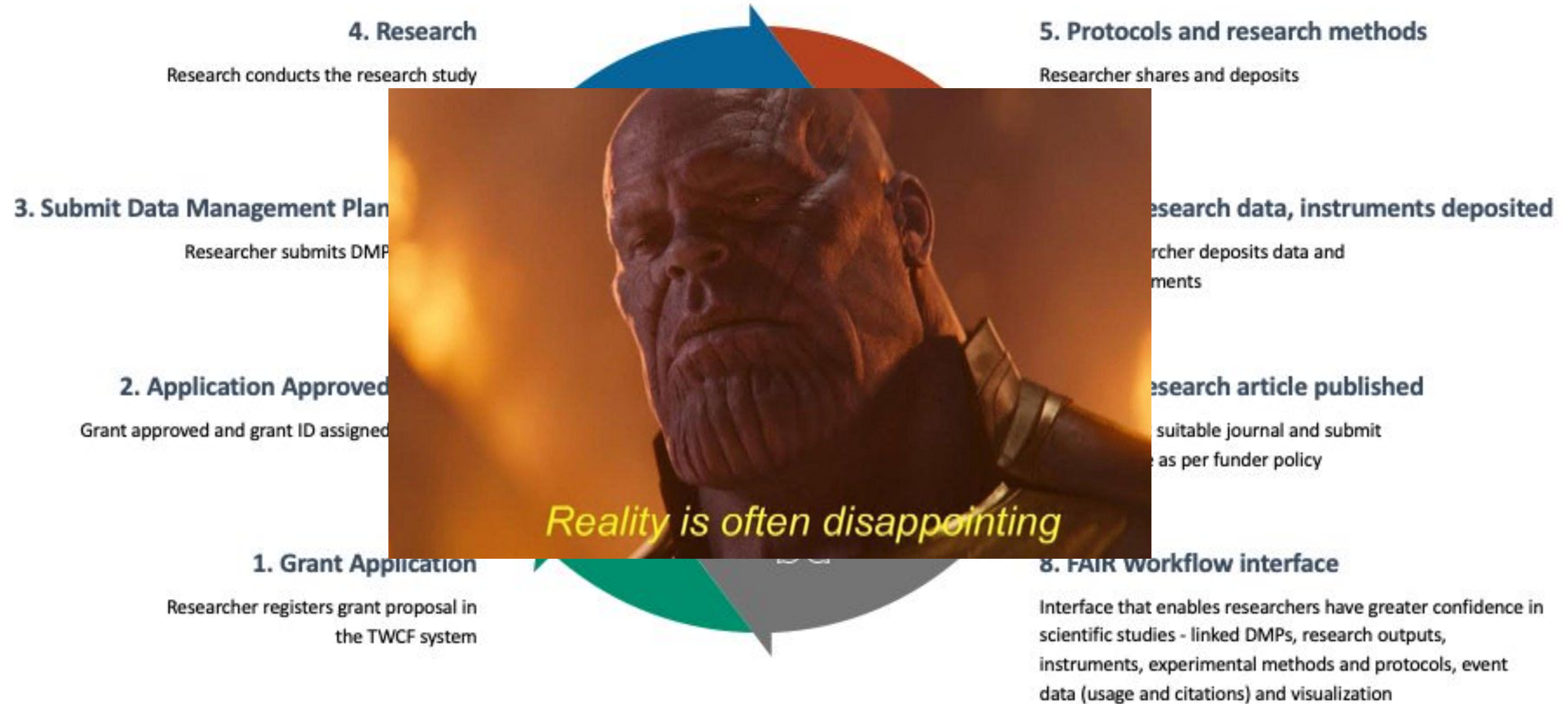


## **FAIR Data Maturity Model** Specification and Guidelines 2020

Proposed RDA Recommendation  
Produced by: **FAIR Data Maturity Model WG, 2019-2020**  
<https://www.rd-alliance.org/groups/fair-data-maturity-model-wg>

# FAIR Workflows Project

Potential workflow



*This project was made possible through the support of a grant from Templeton World Charity Foundation, Inc. The opinions expressed in this publication are those of the author(s) and do not necessarily reflect the views of Templeton World Charity Foundation, Inc.*



# Capturing Metadata

**Put in place mechanism and practices for metadata submission and capture.**



## Researchers

- Plan for output management routine
- Build protocols for not only individuals but also teams, foster lab-wide FAIR culture
- Preserve and share research outputs
- Select tools and platforms with intention.

## Repositories

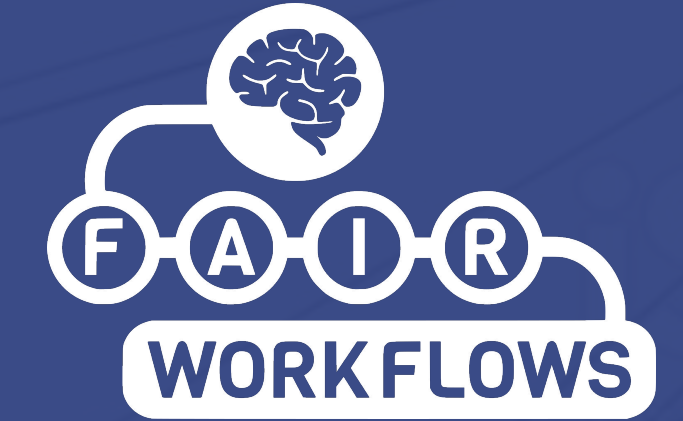
- Adopt PIDs
- Use standardized semantic resources
- Build workflows that capture rich metadata and continuously enrich metadata over time
- Foster an open environment and interoperability by providing API endpoints

## Tools and platforms

- Adopt PIDs when possible
- Integrate PID enabled features

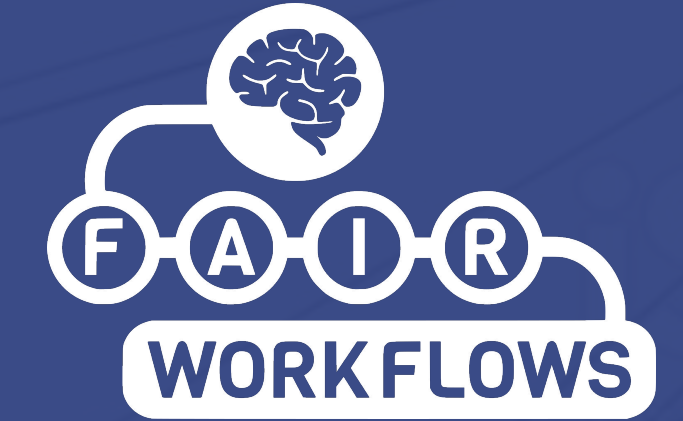


# System Integrations





# Enriching Metadata



**Improve the quality of metadata entered / submitted to the metadata commons.**

## **Semantic resources**

- Ontologies
- Classification
- Controlled list
- Standardized format
- ...

## **Comprehensive Crosswalk**

- Maximize number of submittable metadata fields when registering resource
- Introduce new metadata fields

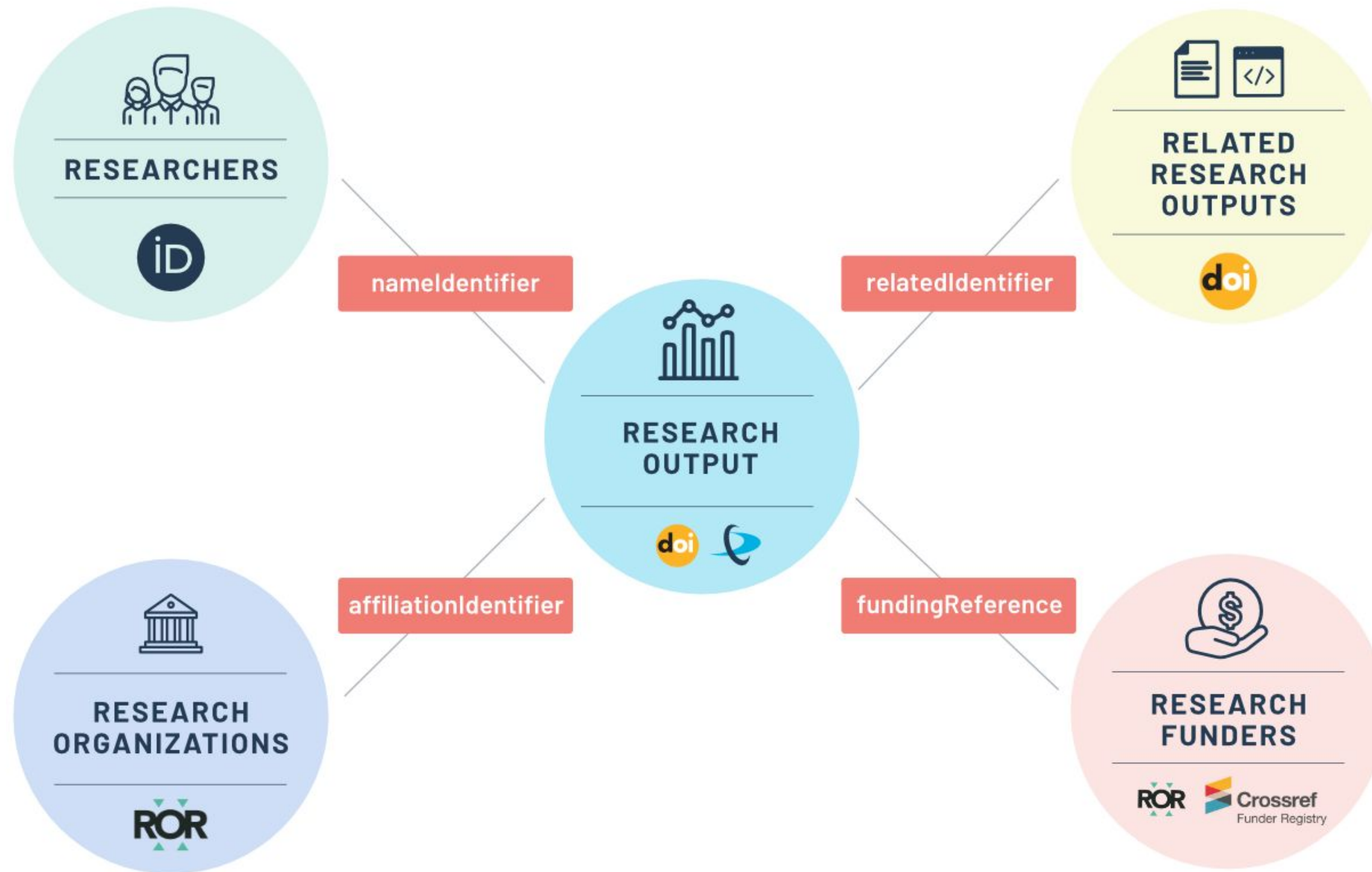
## **Connection metadata**

- Formulate connection information into standardized metadata
- Share connection metadata



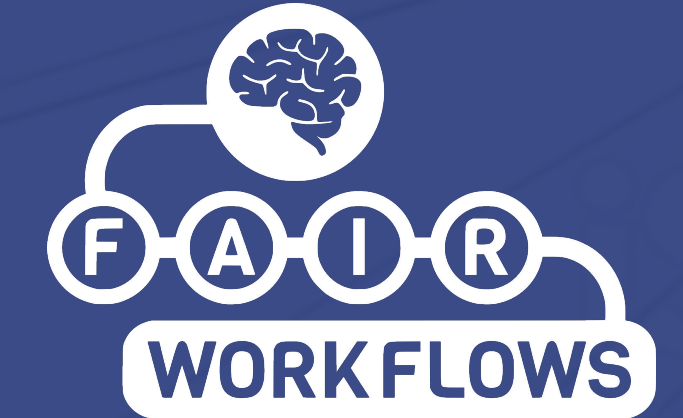
## DataCite Connection Metadata

Connect DataCite DOIs to every part of research ecosystem





# Domain specific Metadata Template



If we want to have FAIR data, we need good metadata. Good metadata need:

- **Reporting guidelines**—like MIAME—to provide a uniform structure
- **Ontologies** to provide controlled terms
- **Technology** to make it easy to author good metadata in the first place

A **metadata template** can ensure compliance with all investigator-controlled FAIR principles, including:

- Making metadata “rich”
- Using metadata vocabularies that follow the FAIR principles
- Meeting domain-relevant community metadata standards

A screenshot of a web-based metadata form titled 'BioSample Human'. It shows a tree-like structure of fields. The main section 'BioSample Human' contains fields for Sample Name (056), Organism (Homo sapiens), Tissue (skin of body), Sex (Male), Isolate (N/A), Age (74), and Biomaterial Provider (Life Technologies). Below this are three expandable sections: 'Attribute (1)' with Name 'disease' and Value 'dermatitis'; 'Attribute (2)' with Name 'description' and Value 'Cell line was cultured until the 5th passage'; and 'Attribute (3)' with Name 'treatment' and Value '350mg brodalumab'.

FsF-I1-02M

**Metadata uses semantic resources**

FsF-R1-01MD

**Metadata specifies the content of the data**

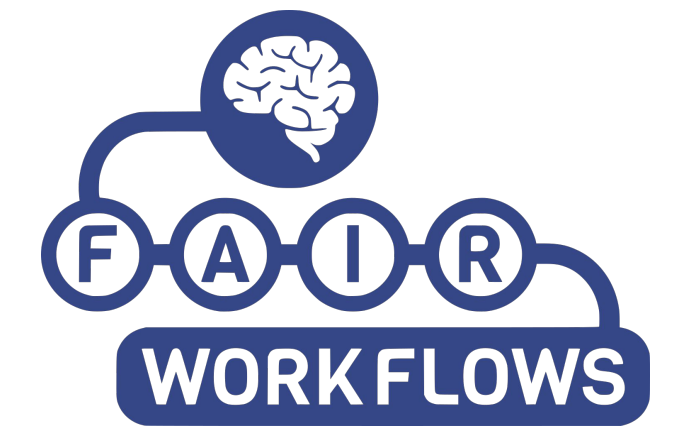
FsF-R1.3-01M

**Metadata follows a standard recommended by the target research community of the data**

FsF-R1.3-02D

**Data is available in a file format recommended by the target research community**



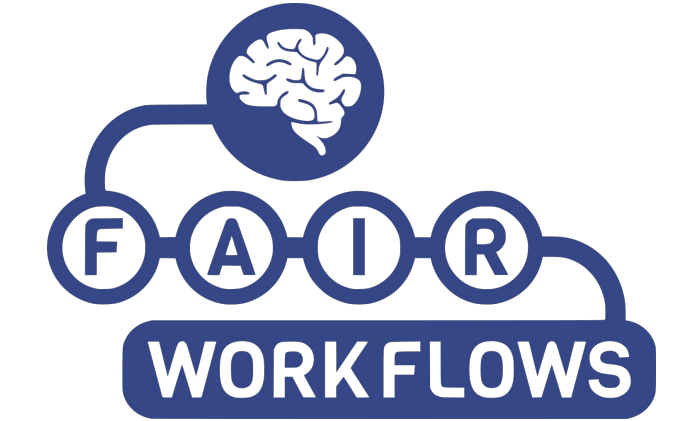




**Thanks!**



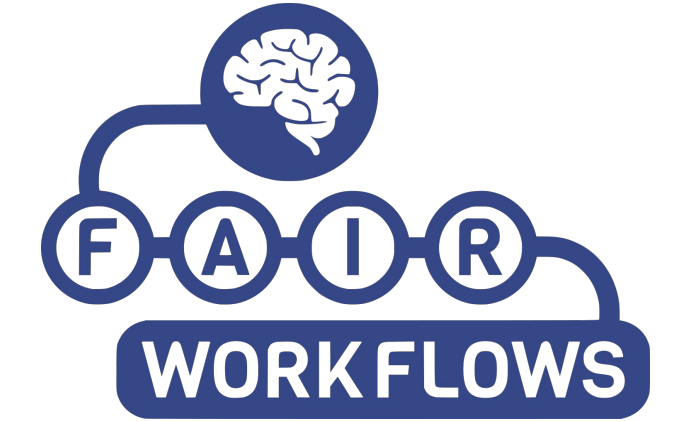
# Project advisory committee



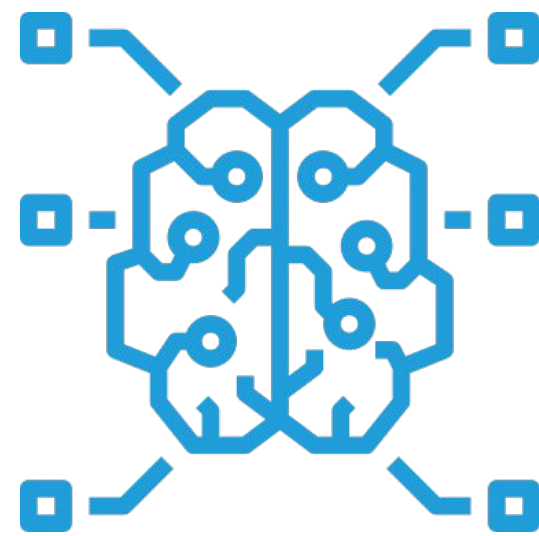
Internal advisors	<b>Adrian Burton</b>	Australian Research Data Commons
	<b>Jennifer Kemp</b>	Crossref
	<b>Nici Pfeiffer</b>	Center for Open Science
External advisors	<b>Russell Poldrack</b>	OpenNeuro
	<b>Franco Pestilli</b>	University of Texas
	<b>Rachael Kotarski</b>	British Library
	<b>Helena Ledmyr</b>	International Neuroinformatics Coordinating Facility
	<b>Dylan Roskams-Edris</b>	Tanenbaum Open Science Institute
	<b>Jean-Baptiste Poline</b>	Brain Imaging Centre Neuroinformatics, McGill University
	<b>Bryan Lawrence Caron</b>	NeuroHub (Lead PI), NeuroDataScience Director, INCF chair of scientific council



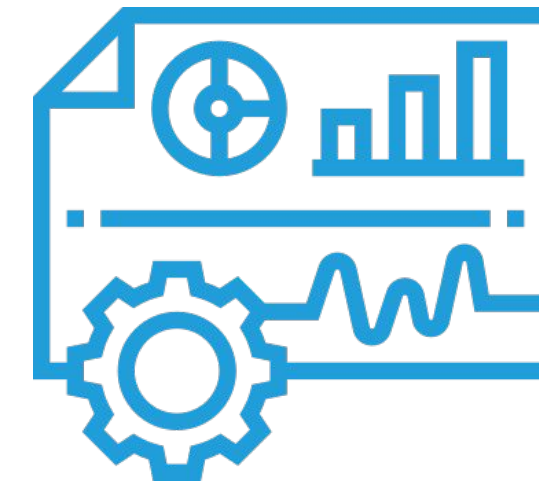
# Project work packages



**WP1**  
**Workflow  
development**



**WP2**  
**Application in  
research**



**WP3**  
**PID graph &  
dashboard**



**WP4**  
**Adoption &  
dissemination**

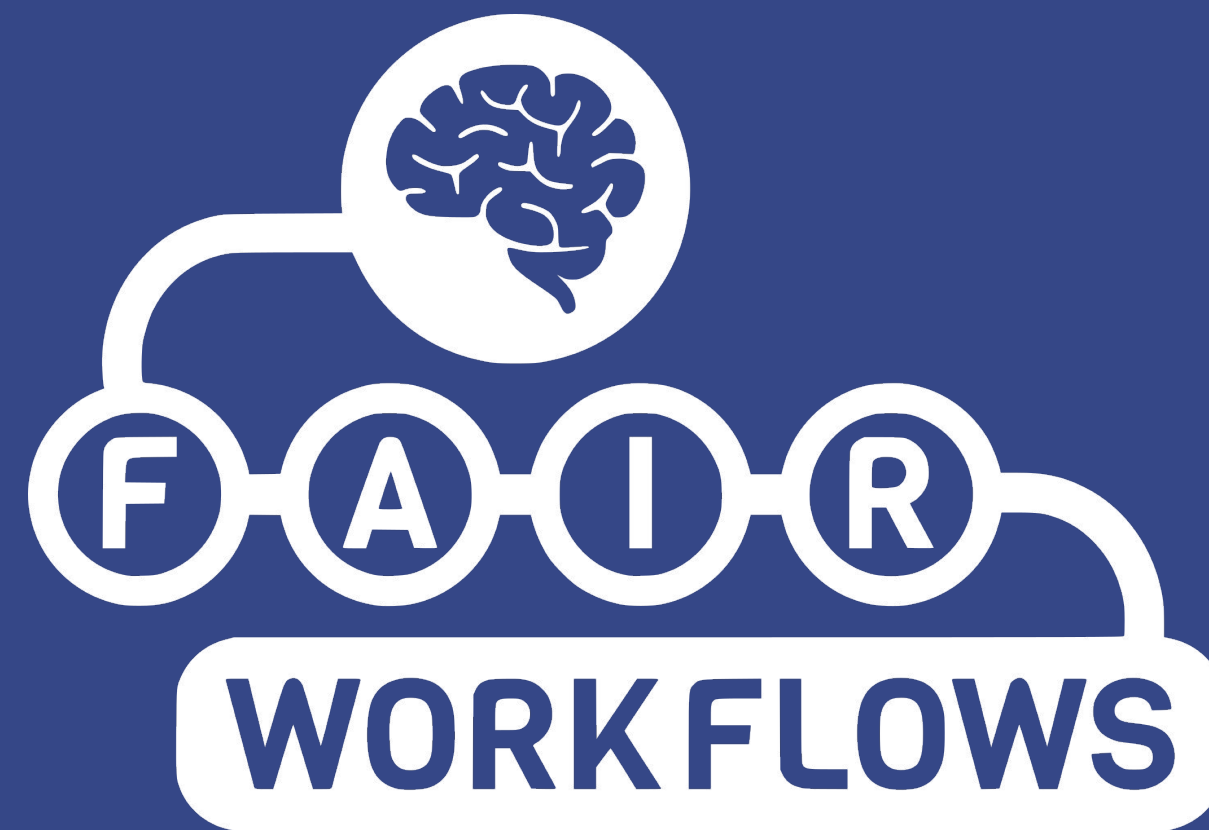


# Immediate Next steps



-





**Keep in touch!**



[Email](#)



[@DataCite](#)



[slack](#)



[webpage](#)



**DataCite**

**Chapter cover 02**



***“In order to make it more dynamic, we highlight pages as well as we highlight text.”***