# Metadata & Persistent identifiers

Espen Åberg
Data Steward
ELIXIR Norway/BioMedData

RDMkit

Reuse · Plan · Collect · Process · Analyse · Preserve · Share

Link to RDMkit: https://rdmkit.elixir-europe.org/

Has a useful purpose

Can be acted upon and processed by humans and machines

"Metadata is constructed, constructive, and actionable."

Definition from Karen Coyle, Digital Librarian and Author of Coyle's InFormation

"information about something"

# What is metadata?

Rich metadata

Metadata have multiple attributes

"data about data"

"Data is content, and metadata is context"

"Metadata is a Love Note to the Future"

# Why do I care?

Metadata **facilitates** organization, indexing, discovery, access, analysis, and use of data.

Metadata **presence and quality** (or the lack thereof) can significantly **help or hinder** time and money expenditures in research activities.

# Metadata helps make data FAIR

| | |
|---|---|
| Data should be **Findable** | F1. (meta)data are assigned a globally unique and persistent identifier (DOI) |
| | F2. data are described with rich metadata |
| | F3. metadata clearly and explicitly include the identifier of the data it describes |
| | F4. (meta)data are registered or indexed in a searchable resource |
| Data should 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 |
| Data should 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 |
| Data should 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 |

Medyckyj-Scott, David et al. (2016).

# Experimental design

"Data"

"Metadata"

Outcome = Treatment effect + Biological effect + Technical effects + Error

| Environment | Sex | Operator | Experimental |
|---|---|---|---|
| Compound | Age | Batch | Treatment |
| Infection | Weight | Plate | Sampling |
| Inhibitor | Litter | Cage | Measurement |
| siRNA | Genotype | Array | |
| sgRNA | Species | Flowcell | |
| Dose | Cell line | Instrument | |
| Time | | Day | |
| | | Order | |
| | | Source | |

# Helps to gain insight



"If data is the new oil, metadata is the refinery"

— Adam Rauh

# "Rich" Metadata

# Metadata templates/checklists



https://www.ebi.ac.uk/ena/browser/checklists

# Metadata Submission Workflow



Figure 2 from: Using Semantic Technologies to Enhance Metadata Submissions to Public Repositories in Biomedicine

# Make it visible

Digital object (DO)

With Rich Metadata

● TAG it with a PID

# PIDs helps make data FAIR

| | |
|---|---|
| Data should be **Findable** | F1. (meta)data are assigned a globally unique and persistent identifier (DOI) |
| | F2. data are described with rich metadata |
| | F3. metadata clearly and explicitly include the identifier of the data it describes |
| | F4. (meta)data are registered or indexed in a searchable resource |
| Data should 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 |
| Data should 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 |
| Data should 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 |

Medyckyj-Scott, David et al. (2016).

# Why don't I just use a link (URL)?



25. Supplemental data showing the predicted secondary structures of each construct (Fig. 3) and explaining the ligation activity of truncated ribozymes (Fig. 2B) are available at *Science* Online at www.sciencemag.org/feature/data/1050240.shl.

"Link rot"

404

Hmmm...

**This doesn't *look* like science.**

It seems you're in search of a page that doesn't exist, or may have moved. You can use the Back button in your browser to return to the page that brought you here, or **search for your missing page**.

# PID = PDI = GUID

PID = Persistent Identifier
PDI = Persistent Digital Identifier
GUID = Globally Unique Identifier

**Physical objects**: a dog, building, microscope, star, person etc

It doesn't "rot"

A persistent identifier (PID) is a long-lasting reference to a resource

Somebody commits to keep it alive

People AND computers can find it

**Digital Objects**: data, collections, metadata, software, publications, configurations, categories, workflows etc

globally unique string of characters

# A PID consists of two components:

**Visible** string of letters and/or numbers

1. A unique identifier

2. A service that locates the resource (or "**resolves**" it)

Behind the scene

# Persistent over time

.. by design



**11839/abc123**

ID is unique and always the same

**11839/abc123**

http://www.example.com/

http://www.moved.com/

URL may change over time

```
1110000
1000111
1
```

```
1110000
1000111
1
```

Supports access to resource as it moves from one location to another.

# Persistent over time

.. by design



| Stable | 11839/abc123 | | 11839/abc123 |

http://www.example.com/ ──────────────► http://www.moved.com/

Update information
Redirection

**Responsibility** of the PID owner to keep it up-to-date when the resource changes

```
1110000
1000111
1
```

```
1110000
1000111
1
```

# Priciple:



TAG

PID

PID Service

URL

Electronic Resource

**Publication date:**
November 24, 2017

**DOI:** 10.5281/zenodo.1065991

**Keyword(s):**
FAIR, FAIRness, checklist, research data, Findable, Accessible, Interoperable, Reusable, PID, repository, DOI, metadata, licence, data sharing, research data management,

**Grants:**
European Commission:
- EUDAT2020 - EUDAT2020 (654065)

**License (for files):**
Creative Commons Attribution 4.0

https://www.clarin.eu/sites/default/files/pid-CLARIN-ShortGuide.pdf

01000110 01100001
01110011 01110100
01110001 01100110
01101001 01101100
01100101

# Example



**DOI**

10.25820/data.006172

**URL**

https://iro.uiowa.edu/esploro/outputs/dataset/9984240535802771

**Citation**

Van Benschoten, W. Z., & Shepherd, J. J. (2022). *Dataset for "Piecewise Interaction Picture Density Matrix Quantum Monte Carlo"* [Data set]. University of Iowa.
https://doi.org/10.25820/data.006172

**Landing Page**

DATASET | OPEN ACCESS

**Dataset for "Piecewise Interaction Picture Density Matrix Quantum Monte Carlo"**

William Z Van Benschoten and James J Shepherd

University of Iowa;
05/20/2022;
DOI: 10.25820/data.006172

View | Share | Export

Files and links

DATASET_VanBenschoten2022_... | 8.67 kB | Download | View
TXT | README | Description of the data and file overview, relationship bet

DATADICTIONARY_VanBenschot... | 7.23 kB | Download | View
CSV | Data Dictionary | Open Data Commons Attribution (ODC-By) V1.0,

DATASTRUCTURE_VanBenschoten... | 3.31 MB | Download | View
CSV | Describes file organization within the zipped folders of files. | Open Da

# Different systems

## Some Common Identifiers:

Digital Object Identifiers (doi:10.1186/2041-1480-3-9)
Handles (hdl:2381/12775)
URN (urn:isbn:0451450523)
Archival Resource Keys (ARK) (ark:/13030/tf5p30086k)
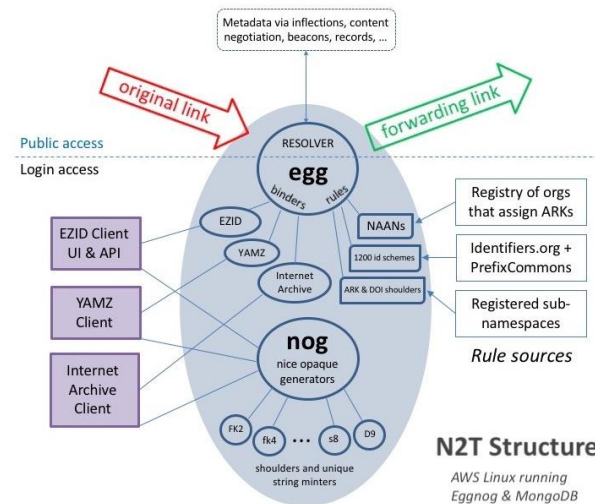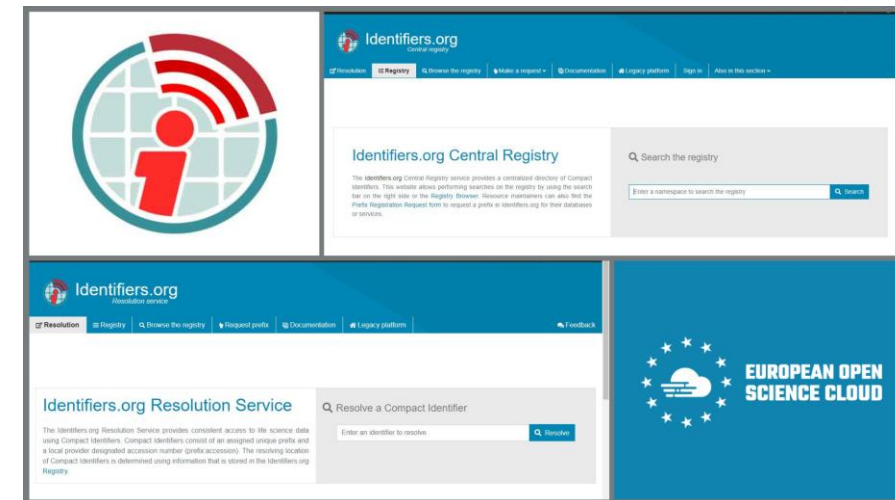Persistent Uniform Resource Locator (PURL)

## Resolver Services

N2T (Name-to-Thing)
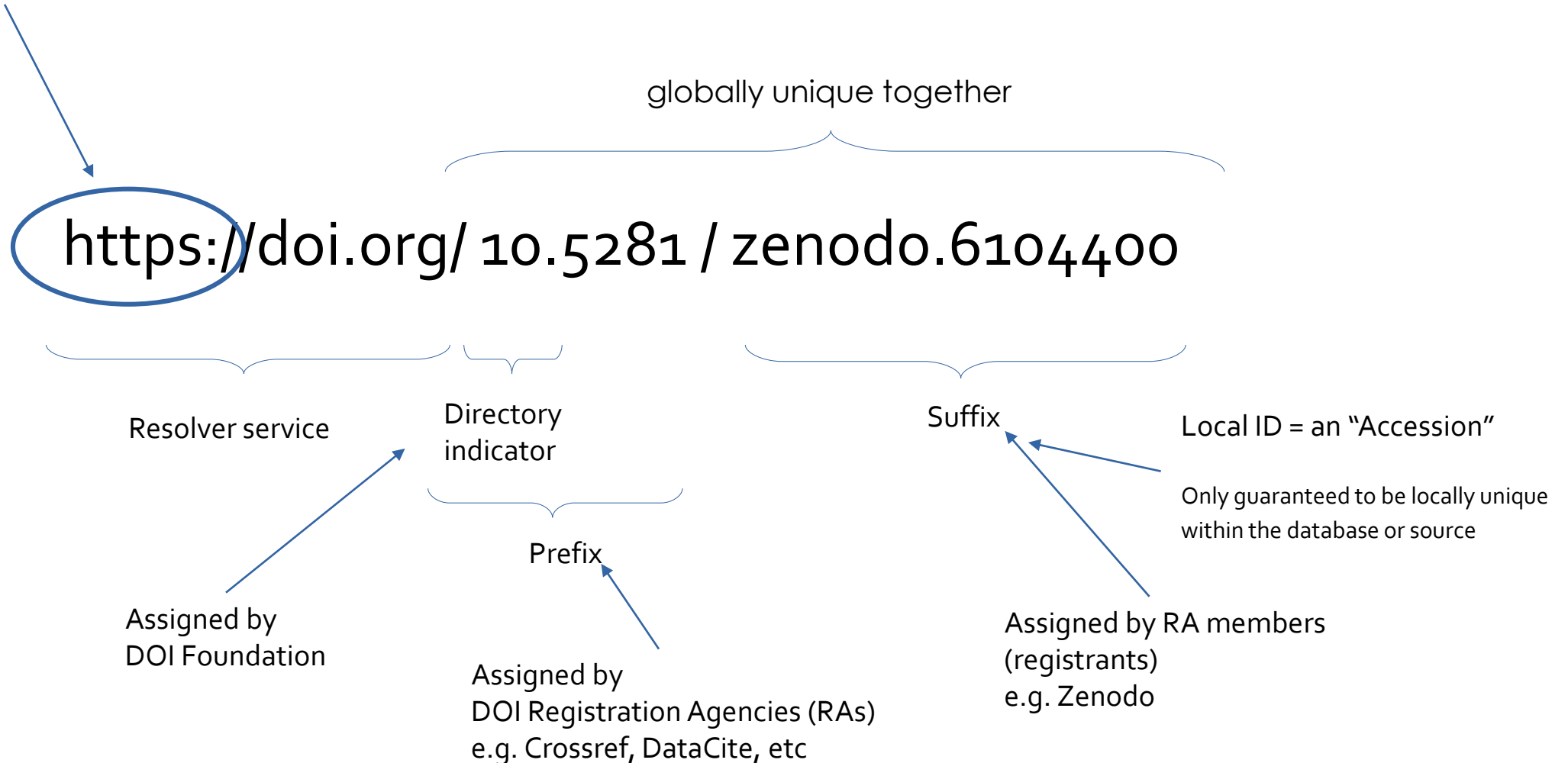Identifiers.org



https://arks.org/about/n2t-global-resolver/

https://eosc-portal.eu/news-and-events/news/identifiers-ensuring-robust-and-reliable-access-life-sciences-data

# How do I recognize a PID?

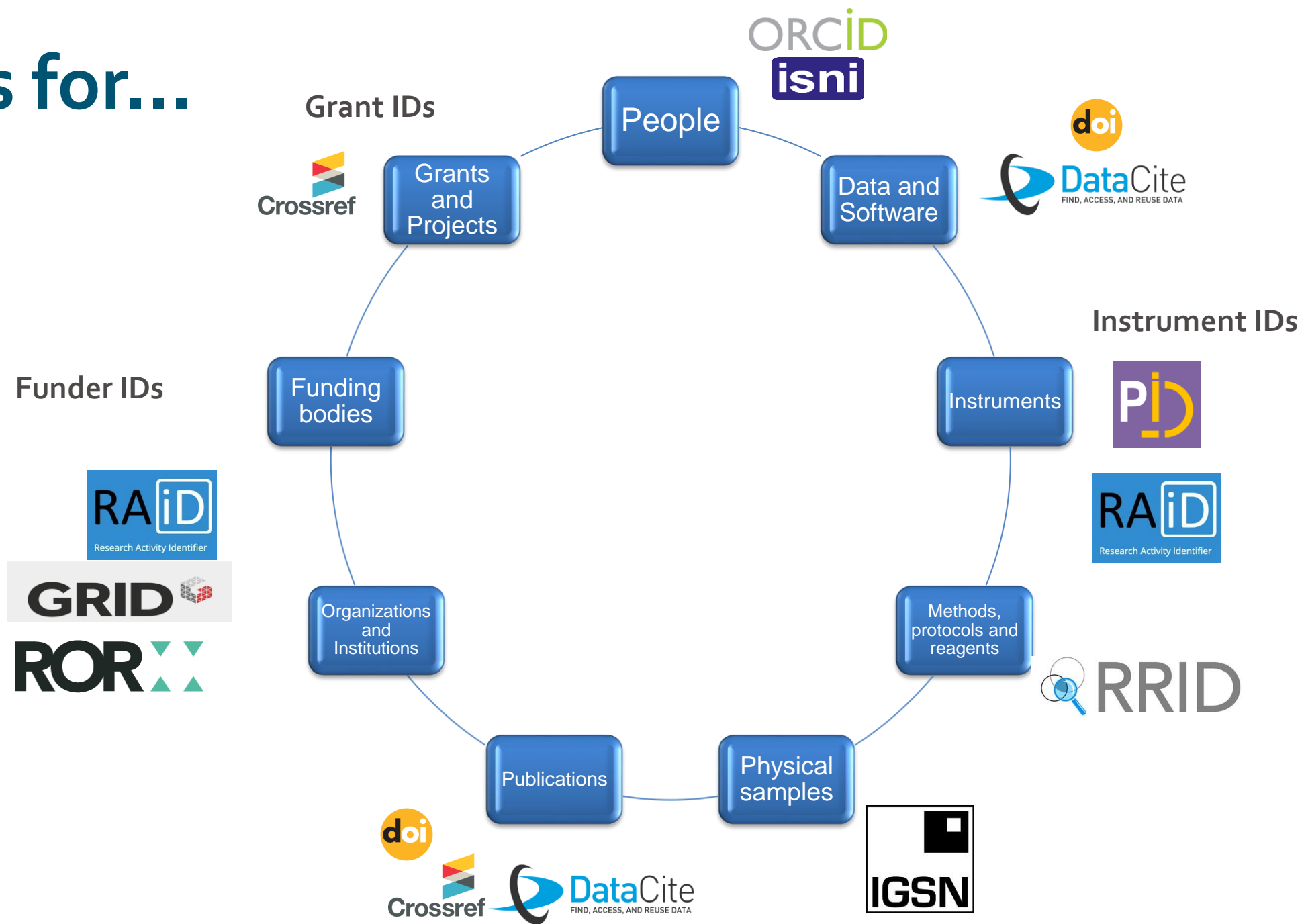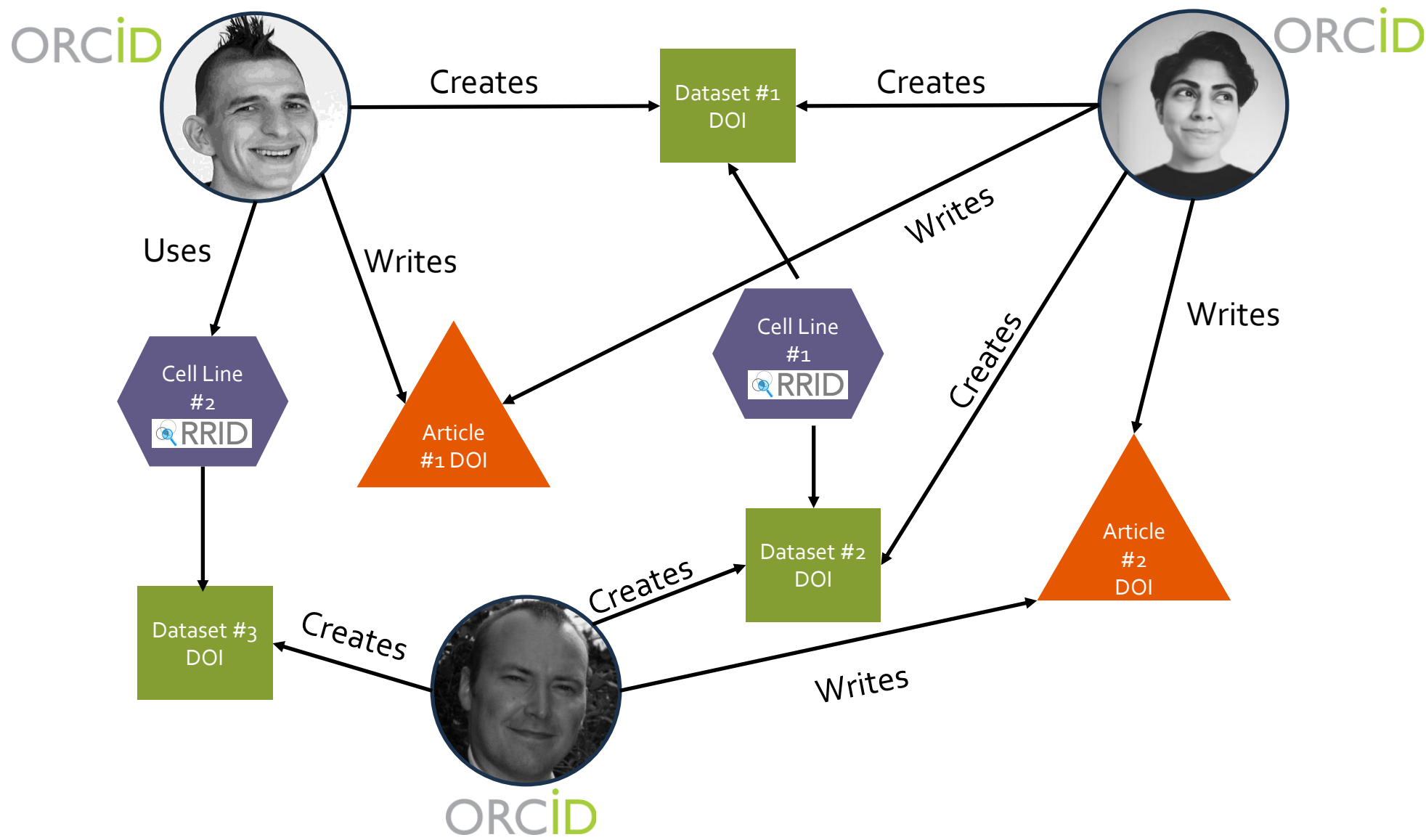DOI: 10.5281/zenodo.6104400

# Anatomy of a DOI

Means that it is actionable: you can paste in a web browser address bar and be taken to the identified source.

globally unique together

https://doi.org/ 10.5281 / zenodo.6104400

Resolver service

Directory indicator

Suffix

Local ID = an "Accession"

Only guaranteed to be locally unique within the database or source

Assigned by DOI Foundation

Prefix

Assigned by DOI Registration Agencies (RAs) e.g. Crossref, DataCite, etc

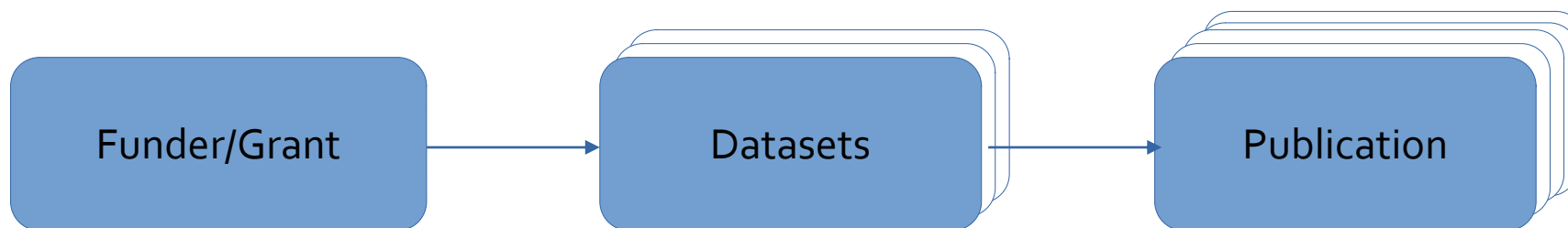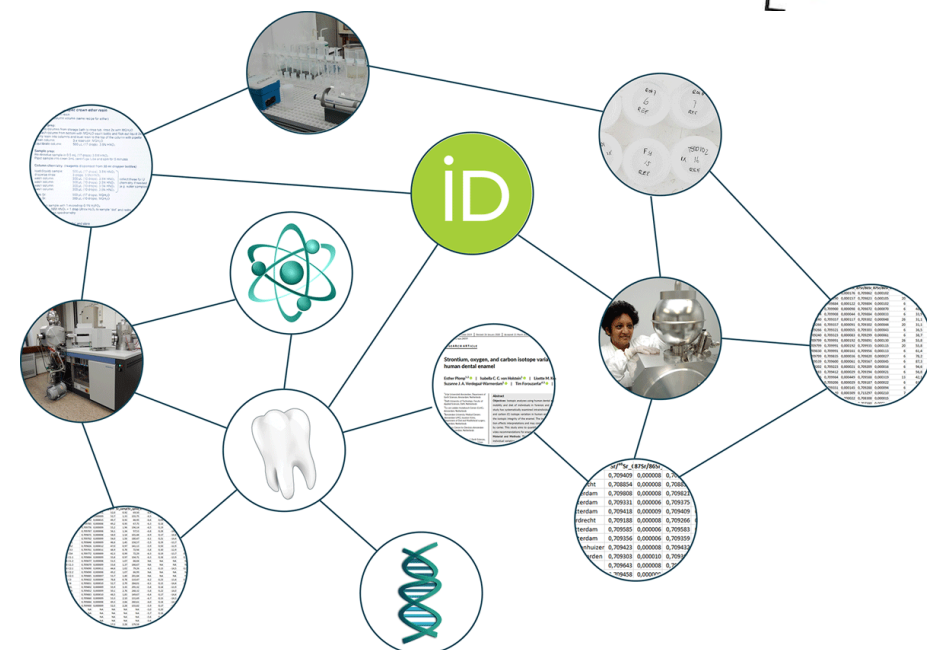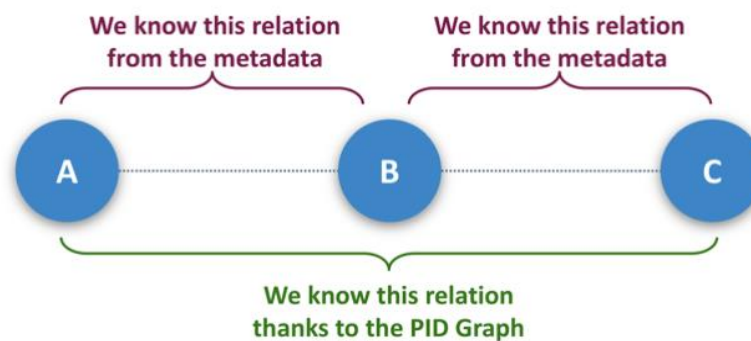Assigned by RA members (registrants) e.g. Zenodo

# PIDs for...

# PIDs connect different entities in research

# PID graphs



"I want to see all datasets funded by RCN cited by this article"

We know this relation from the metadata

We know this relation from the metadata

We know this relation thanks to the PID Graph

Funder/Grant → Datasets → Publication