# Metadata & Persistent identifiers

Espen Åberg
Data Steward
ELIXIR Norway/BioMedData

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

"Metadata is constructed, constructive, and actionable."
Definition from Karen Coyle, Digital Librarian and Author of Coyle's InFormation

"data about data"

# What is metadata?



Rich metadata

Metadata have multiple attributes

"information about something"

"Data is content, and metadata is context"

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

— Adam Rauh

"Metadata is a Love Note to the Future"

"Metadata is constructed, constructive, and actionable."

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

"data about data"

# What is ... a?

"information about something"

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.

"Data is content, a...

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

— Adam Rauh

...etadata is a Love Note to the Future"

# Metadata helps make data FAIR

| | |
|---|---|
| **Data should be Findable** | F1. (meta)data are assigned a globally unique and persistent identifier (DOI)<br>F2. data are described with rich metadata<br>F3. metadata clearly and explicitly include the identifier of the data it describes<br>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<br>A1.1 the protocol is open, free, and universally implementable<br>A1.2 the protocol allows for an authentication and authorization procedure, where necessary<br>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.<br>I2. (meta)data use vocabularies that follow FAIR principles<br>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<br>R1.1. (meta)data are released with a clear and accessible data usage license<br>R1.2. (meta)data are associated with detailed provenance<br>R1.3. (meta)data meet domain-relevant community standards |

# Experimental design

**"Data"**                    **"Metadata"**

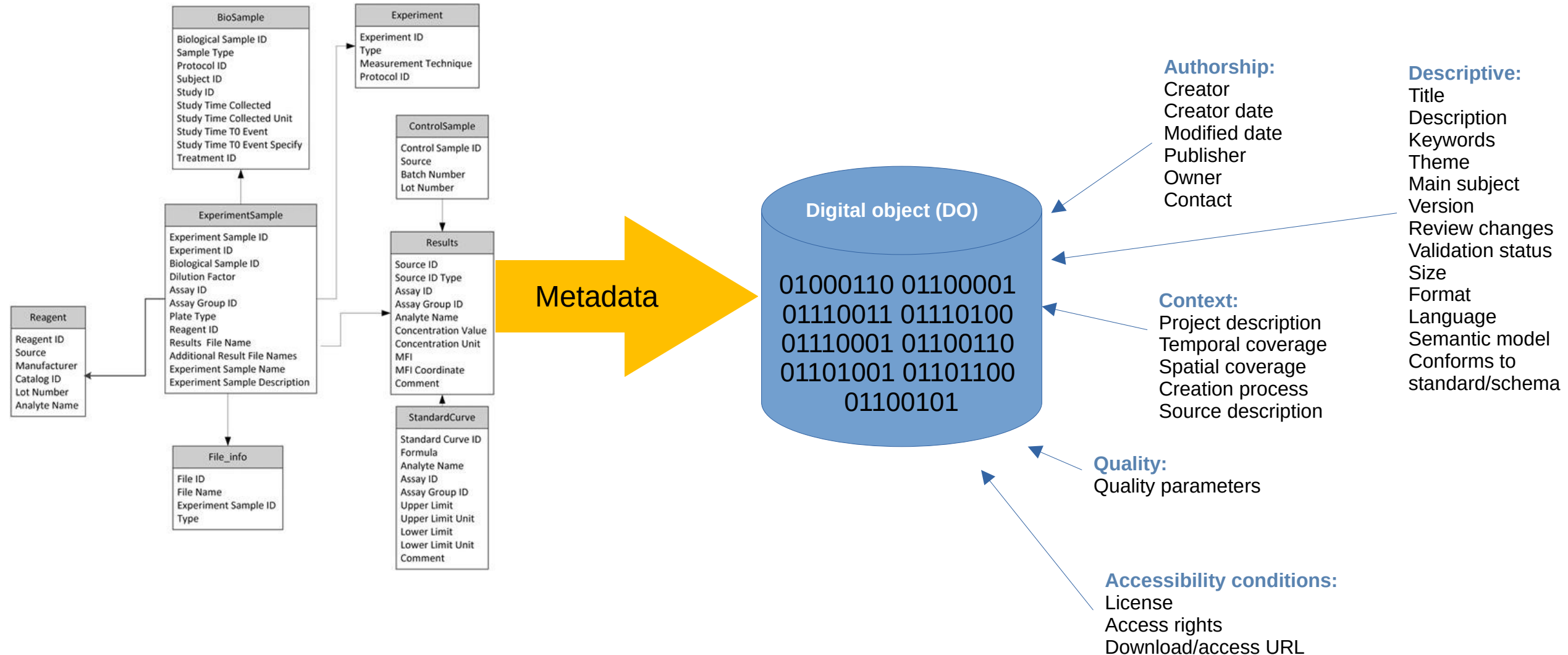**Outcome = Treatment effect + Biological effect + Technical effects + Error**

| | | | |
|---|---|---|---|
| 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 | |

# "Rich" Metadata

# Metadata templates/checklists



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

# Metadata Submission Workflow
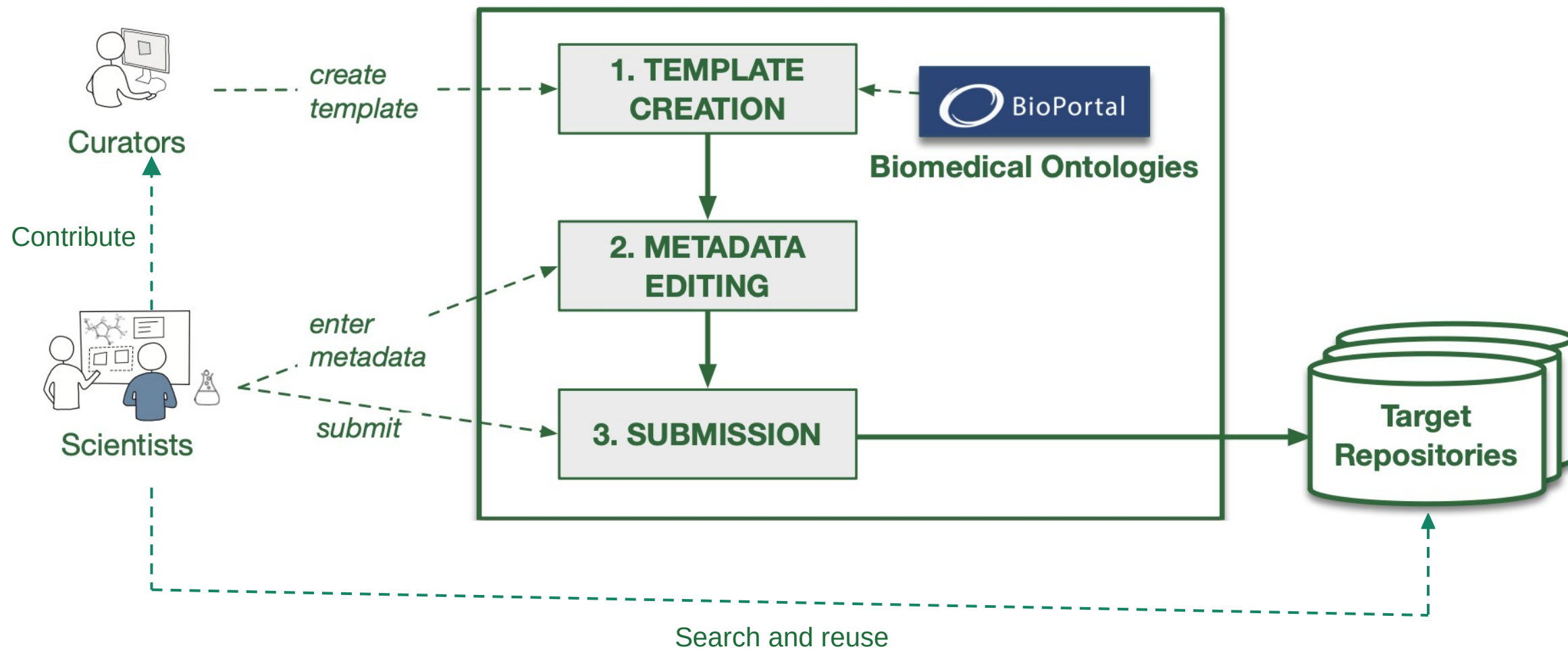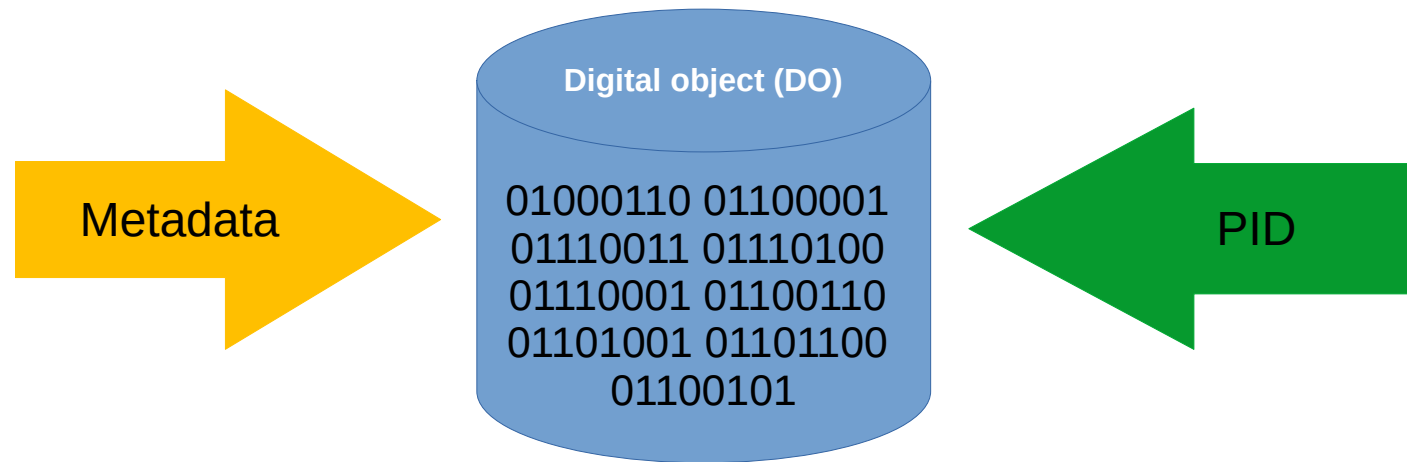


Figure 2 from: Using Semantic Technologies to Enhance Metadata Submissions to Public Repositories in Biomedicine

**Metadata** → **Digital object (DO)**
01000110 01100001
01110011 01110100
01110001 01100110
01101001 01101100
01100101
← **PID**

# PIDs helps make data FAIR

| | |
|---|---|
| **Data should be Findable** | F1. (meta)data are <u>assigned a globally unique and persistent identifier</u> (DOI) |
| | F2. data are described with rich metadata |
| | F3. metadata <u>clearly and explicitly include the identifier of the data</u> it describes |
| | F4. (meta)data are registered or indexed in a searchable resource |
| **Data should be Accessible** | A1. (meta)data are <u>retrievable by their identifier</u> 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 |

# Why not just use a URL?

domain may change

resource may be relocated

URL may change

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.



Science — Contents ▾ News ▾ Careers ▾ Journals ▾

Read our COVID-19 research and news.

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**.

"Link rot"

farm3.staticflickr.com

# PID?

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

globally unique string of characters

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

# A PID consists of two components:

1. a unique identifier


2. a service that locates the resource over time even when it's location changes

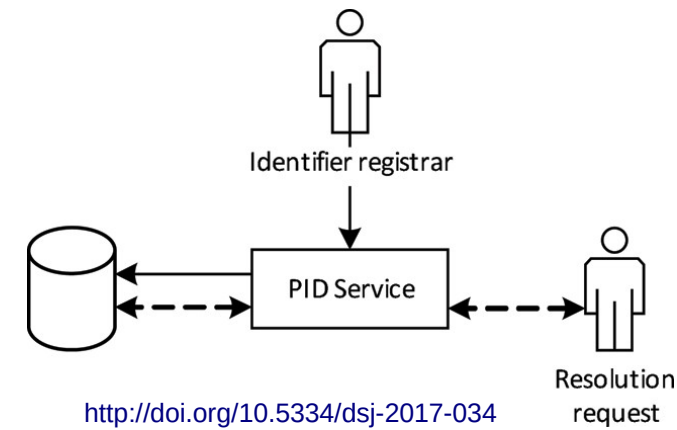# Persistent over time

.. by design

today ... ... ... 2030

**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

today ... ... 2030

| Stable | 11839/abc123 | 11839/abc123 |

http://www.example.com/

Update information
Redirection

http://www.moved.com/

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

```
1110000
1000111
1
```

```
1110000
1000111
1
```

# Examples for digital objects

Digital Object Identifiers **doi**

Handles            Handle.Net®

Archival Resource Keys (ARK)

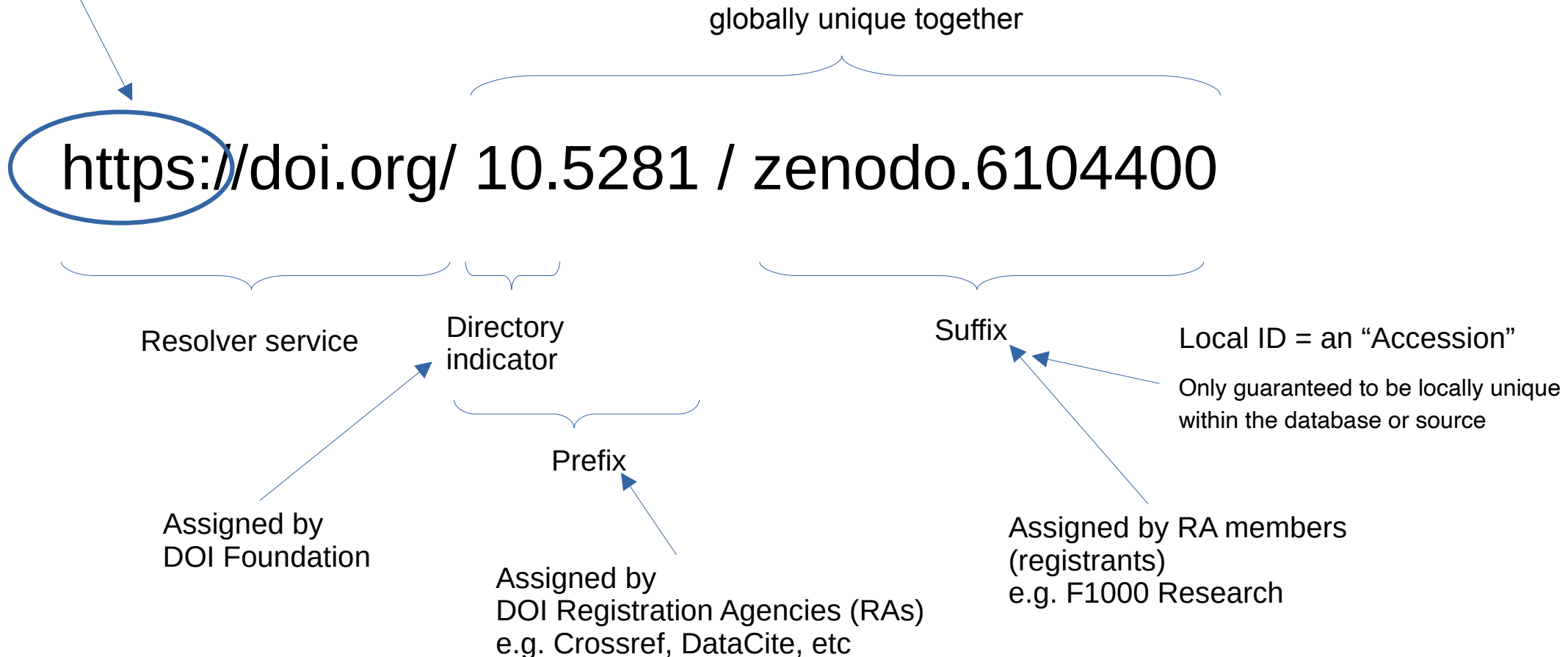Persistent Uniform Resource Locator (URL)

Identifiers.org

# How to recognize a PID
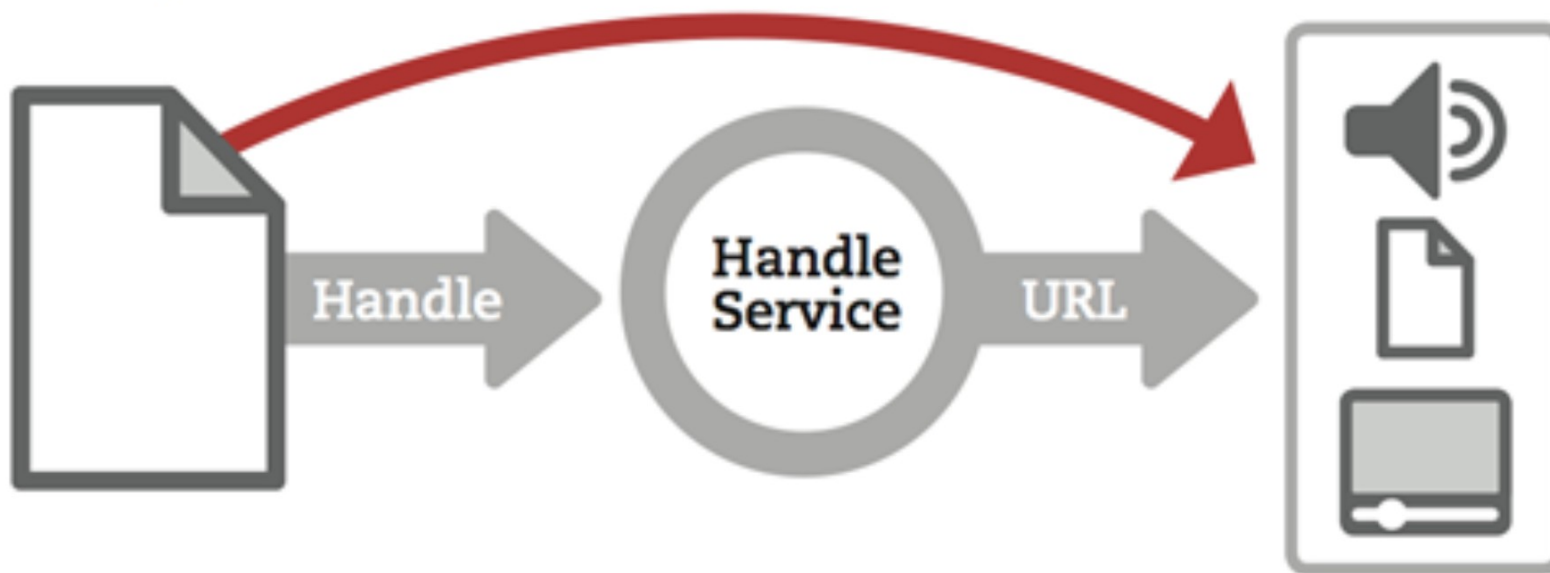
DOI: 10.5281/zenodo.6104400

# 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. F1000 Research

**Metadata Description** → **URL** → **Electronic Resource**

Handle → Handle Service → URL

Publication date:
November 24, 2017

DOI:
DOI 10.5281/zenodo.1065991

Keyword(s):
FAIR, FAIRness, checklist, research data, Findable, Accessible, Interopeable, 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

# PIDs for...

**People**    ORCiD isni

**Funding bodies**    Crossref
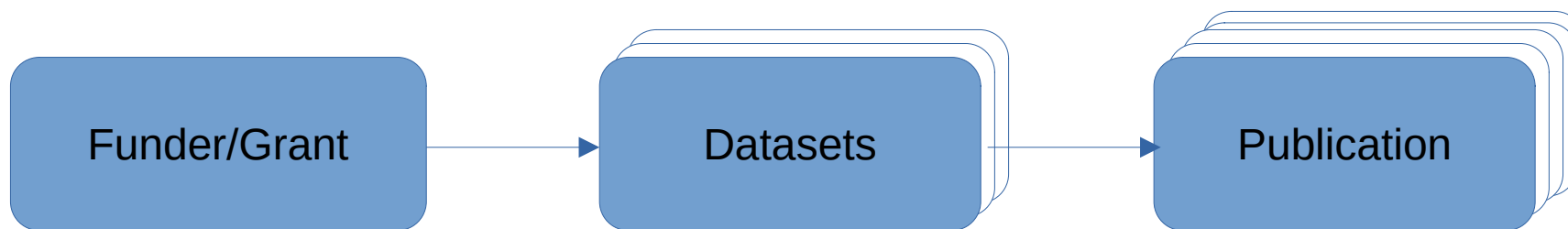
**Institutions**    GRID ROR

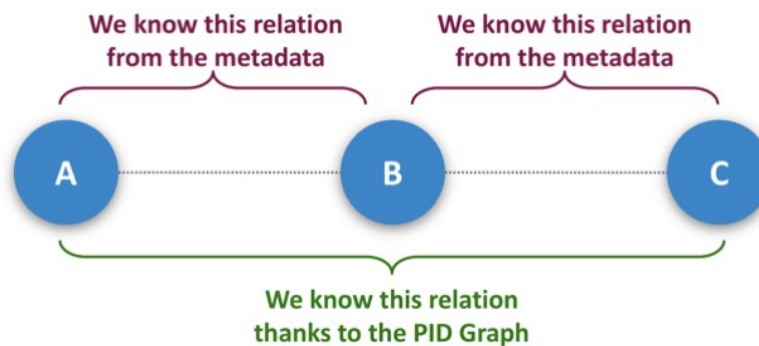**Instruments (soon)**    PiD

# PIDs assembled into 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

A — B — C

We know this relation thanks to the PID Graph

Funder/Grant → Datasets → Publication

Crossref

DataCite
FIND, ACCESS, AND REUSE DATA

Crossref