



Metadata & Persistent identifiers



Espen Åberg
Data Steward
ELIXIR Norway/BioMedData

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?

“data about data”

“Data is content, and metadata is context”

“Metadata is a Love Note to the Future”

Metadata “all over the place”?



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



This is why you 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) <u>F2. data are described with rich metadata</u> 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. <u>meta(data) are richly described with a plurality of accurate and relevant attributes</u> <u>R1.1. (meta)data are released with a clear and accessible data usage license</u> <u>R1.2. (meta)data are associated with detailed provenance</u> <u>R1.3. (meta)data meet domain-relevant community standards</u>

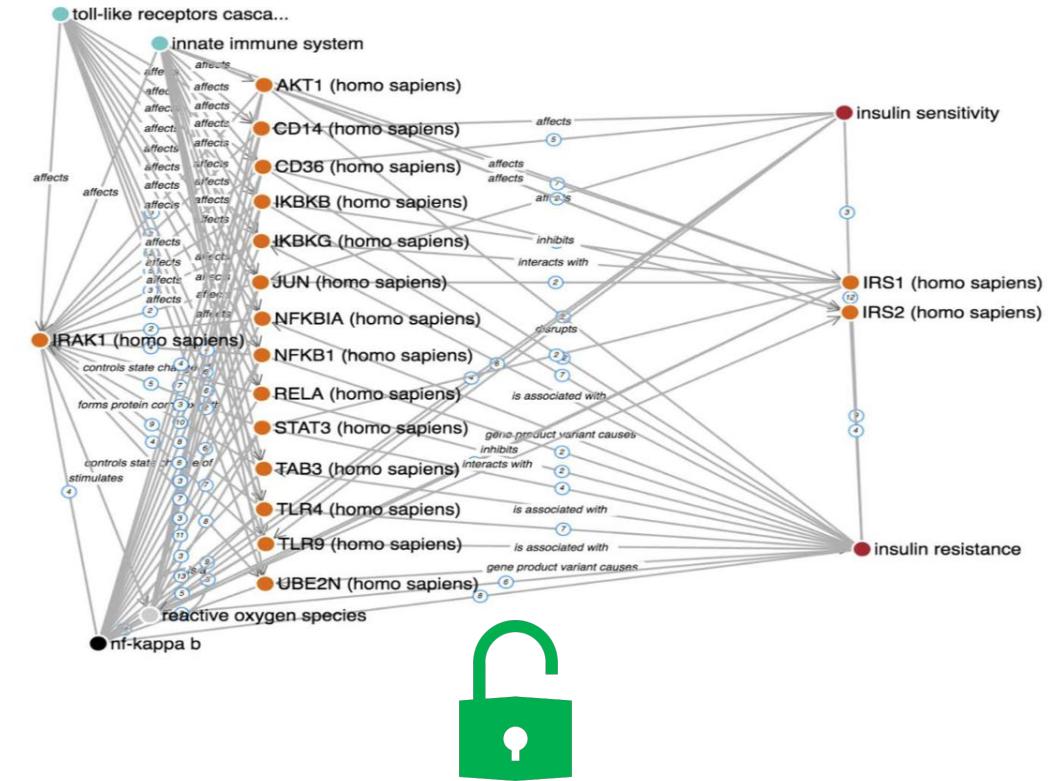
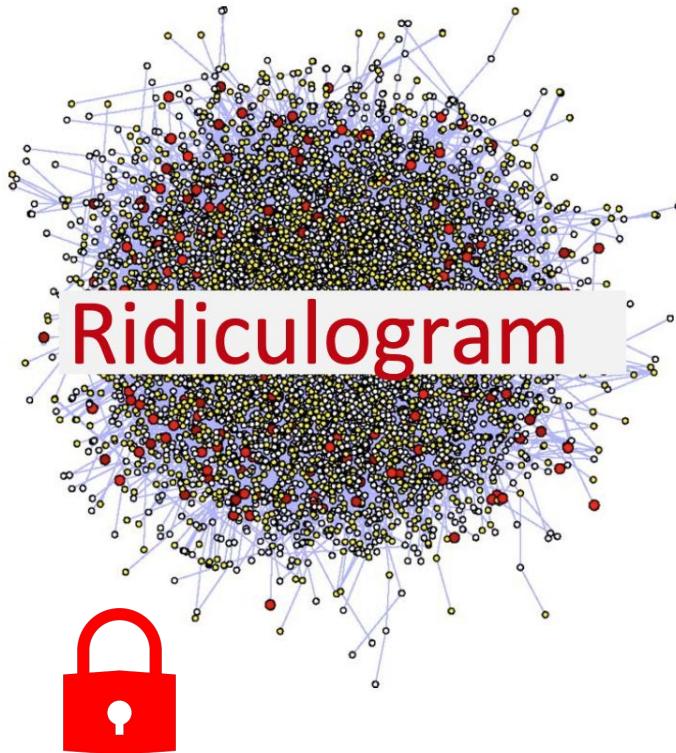
What is Metadata?

“Data”

“Metadata”

Outcome =

Helps to gain insight

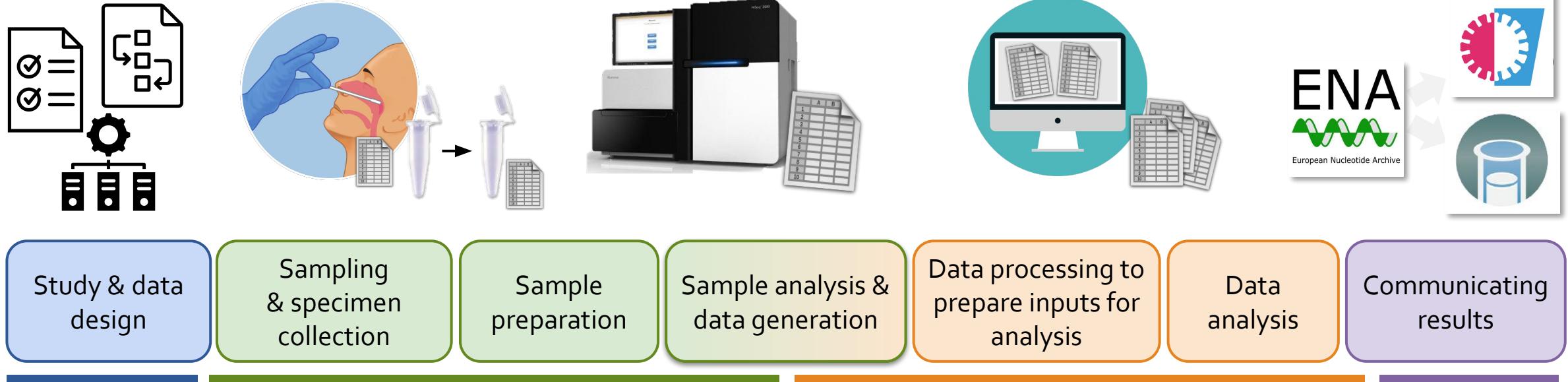


“If data is the new oil, metadata is the refinery”

— Adam Rauh

(Meta)data generation and collection

"Protocol" & "project plan" icons by Justin Blake, and "infrastructure" icon by Eko Purnomo, from thenounproject.com



Procedures

data protection,
ethics permit,
infrastructure,
standards,
protocols,
data dictionaries,
data access, ...

Biosamples and instruments

populations (statistical) and inclusion criteria,
physical processing steps,
working storage conditions,
long-term storage location,
sample quality assessment,
sample annotations,
reagents, instruments, kits, ...

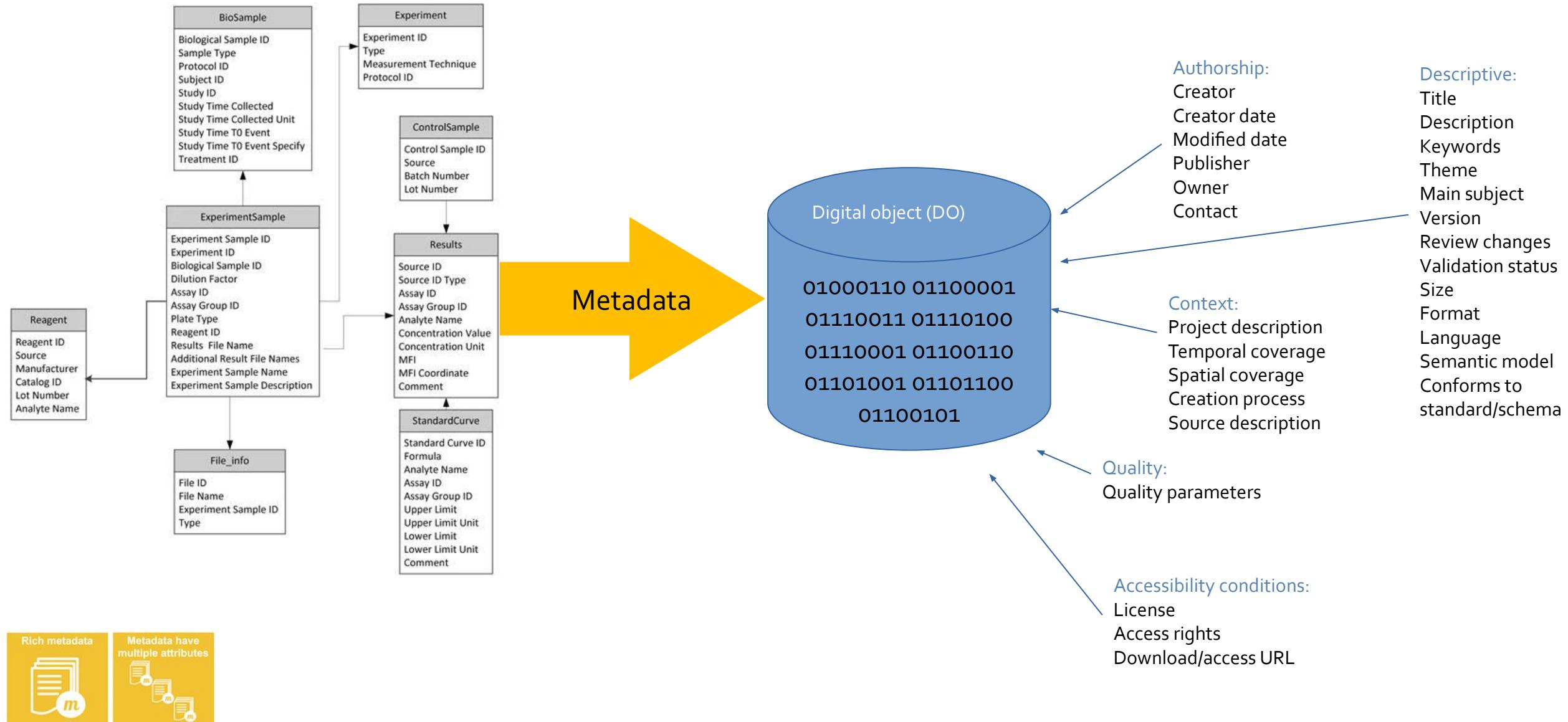
Data and computational workflows

digital processing steps,
working storage conditions,
long-term storage location,
data quality assessment,
sample/data annotations,
reference data,
analysis method...

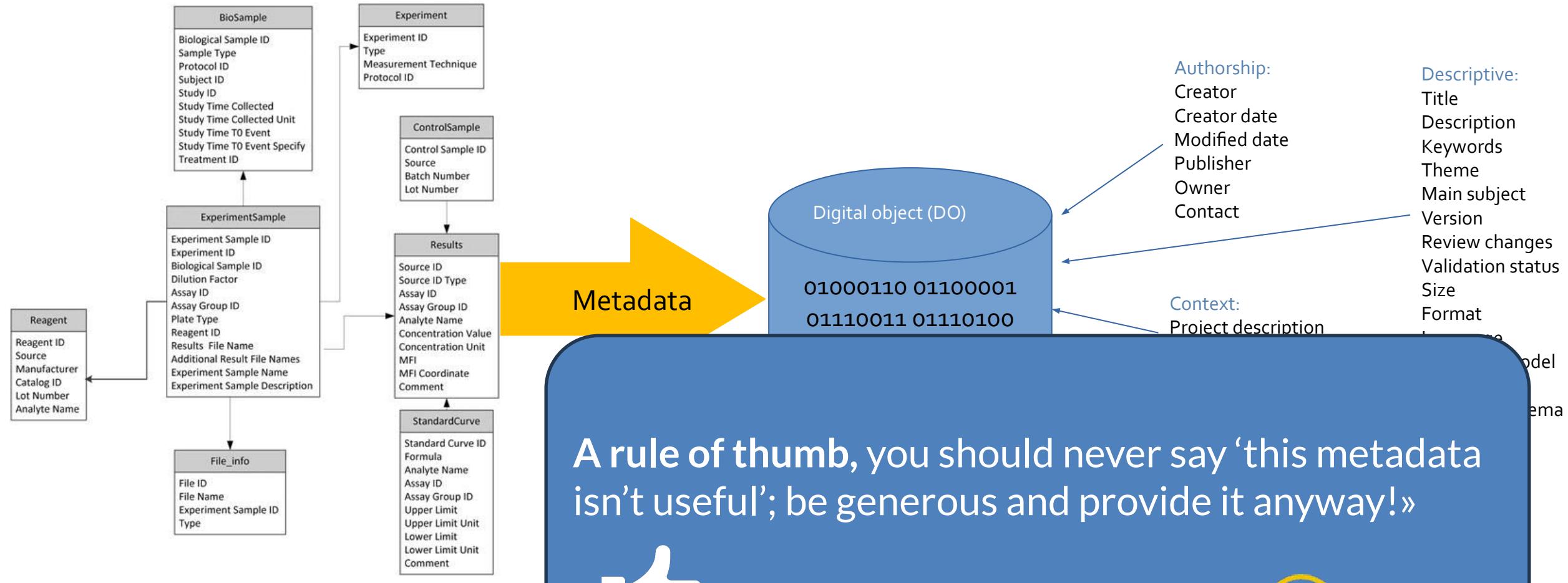
Outputs

publications,
data,
tools,
workflows,
reports,
dashboards, ...

"Rich" Metadata



"Rich" Metadata



A rule of thumb, you should never say 'this metadata isn't useful'; be generous and provide it anyway!»



Download/access URL



Metadata collection (Simple solution)

Not optimal but can be fitted to a set of rules (schema) and allowed values (controlled vocabulary)

Metadata templates/checklists



European Nucleotide Archive

Home | Submit | Search | Rulespace | About | Support

Enter text search terms Search 

Examples: histone, BN000065

Enter accession View 

Examples: Taxon:9606, BN000065, PRJEB402

Sample Checklists

There is a minimum amount of information required during ENA sample registration and all samples must conform to a defined checklist of expected metadata values. The most suitable checklist for sample registration depends on the type of the sample.

Filter checklists... 

Accession	Name	Description
ERC000012	GSC MIxS air	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000013	GSC MIxS host associated	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000014	GSC MIxS human associated	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000015	GSC MIxS human gut	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000016	GSC MIxS human oral	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000017	GSC MIxS human skin	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...
ERC000018	GSC MIxS human vaginal	Genomic Standards Consortium package extension for reporting of measurements and observations obtain...

These sample checklists have been developed to meet the needs of different research communities. Different communities have different requirements on the minimum metadata expected to describe biological samples.

Checklist: ERC000031

GSC MIxS built environment

Genomic Standards Consortium package extension for reporting of measurements and observations obtained from the environment where the sample was obtained. By choosing the environmental package, a selection of fields can be made from a relevant subsets of the GSC terms.

Checklist Fields

Filter fields... 

Filter by type: 

Field Name	Field Format	(Field Restriction)	Requirement	(Units)
relative air humidity	restricted text	regular expression	mandatory	%
absolute air humidity	restricted text	regular expression	mandatory	kg
surface humidity	restricted text	regular expression	optional	%
air temperature	restricted text	regular expression	mandatory	°C
surface temperature	restricted text	regular expression	optional	°C
surface moisture	restricted text	regular expression	optional	options
surface moisture pH	restricted text	regular expression	optional	
dew point	restricted text	regular expression	optional	°C

Metadata Submission Workflow

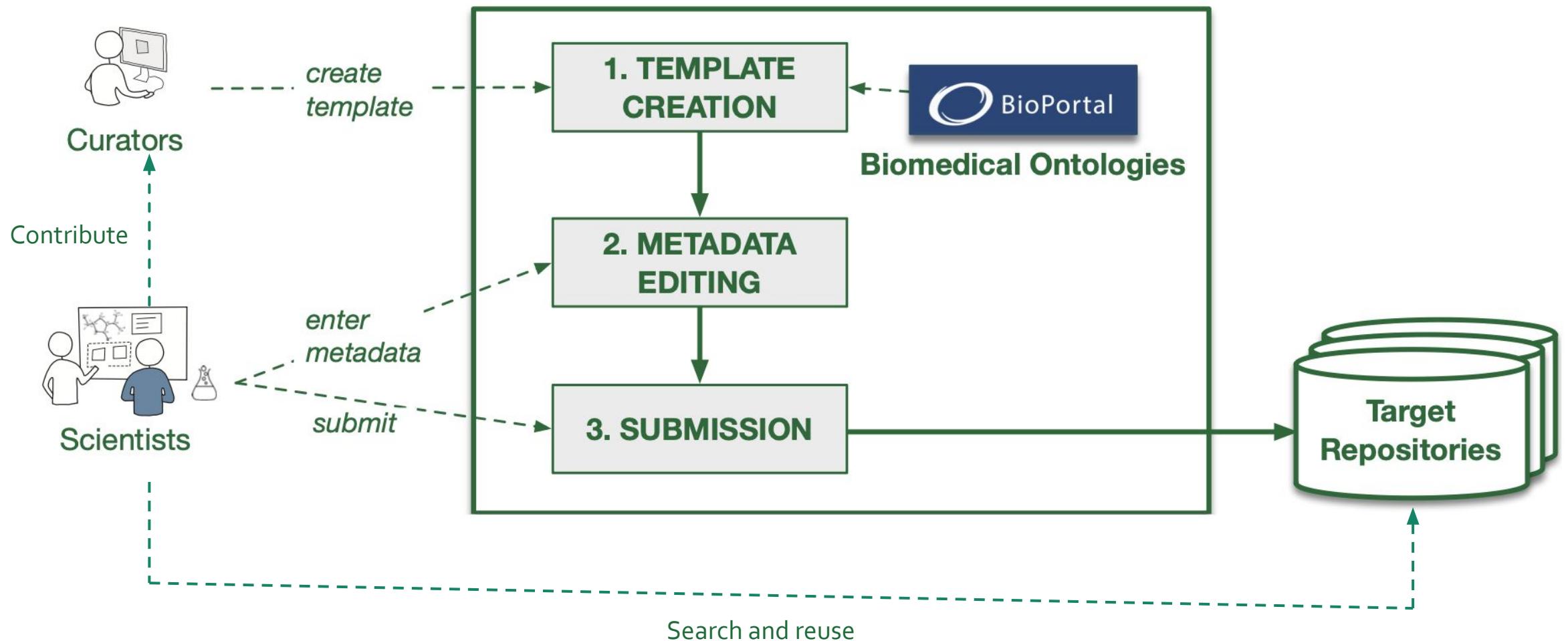
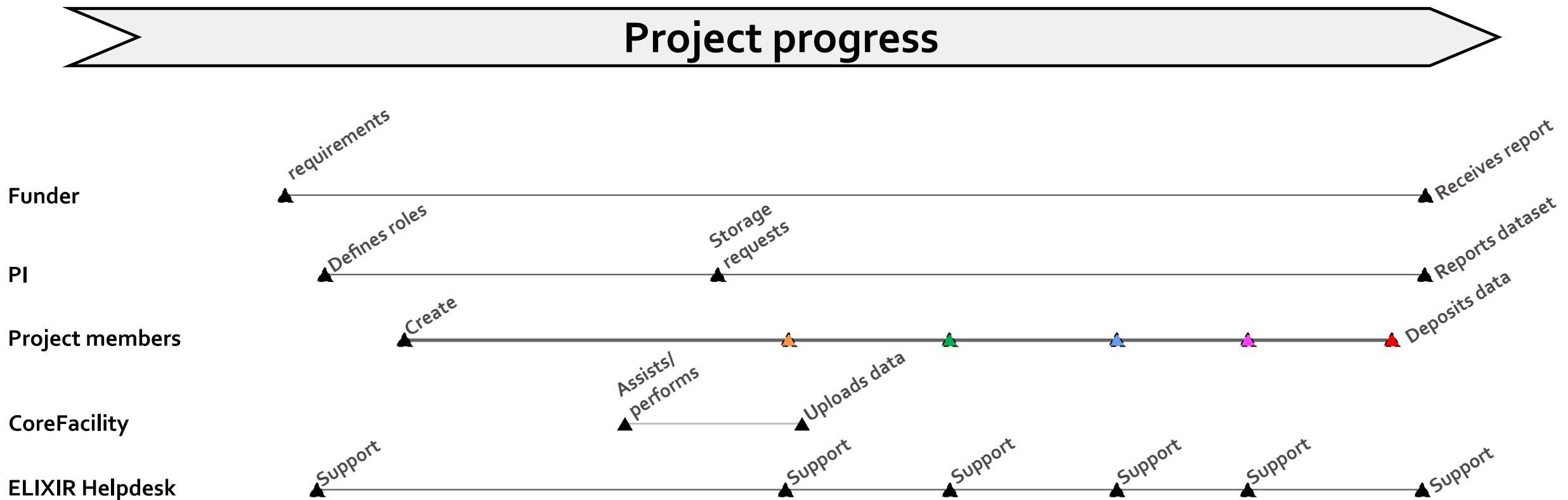
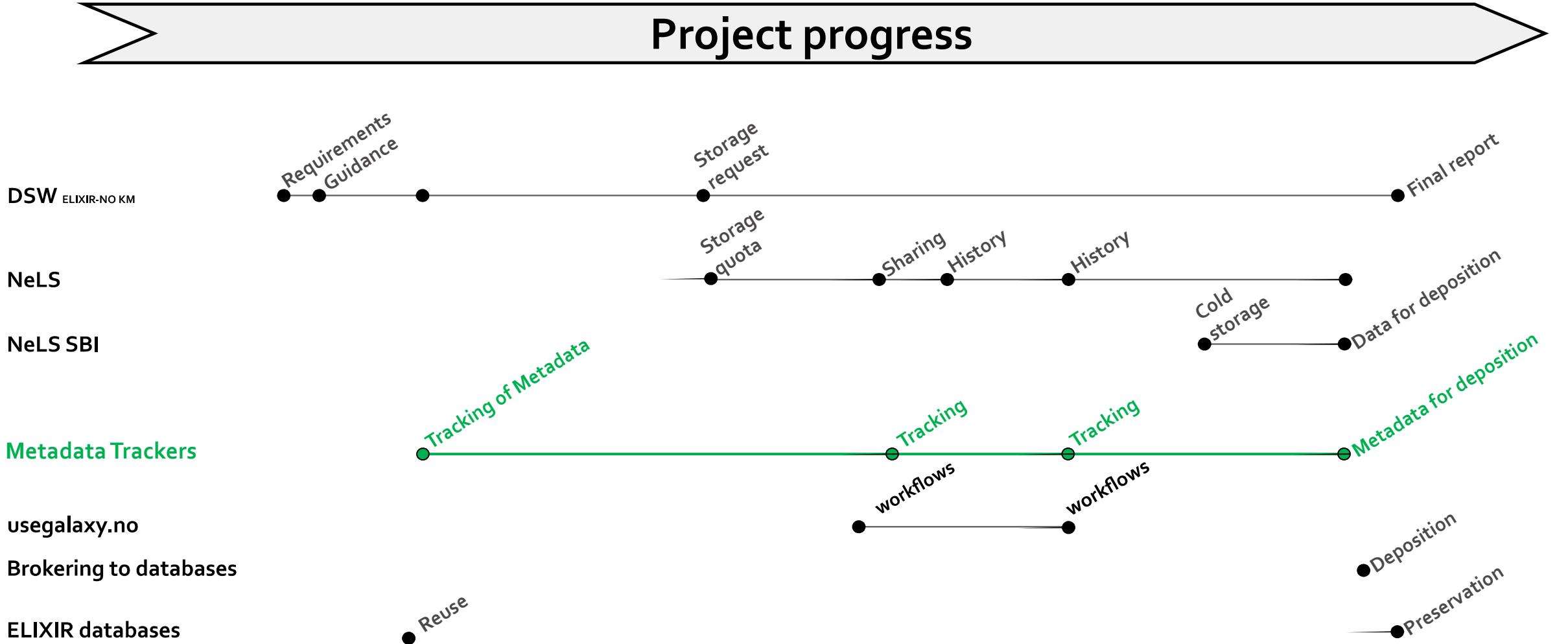


Figure 2 from: [Using Semantic Technologies to Enhance Metadata Submissions to Public Repositories in Biomedicine](#)

Different Persona



Infrastructure tools



Metadata tracking platforms

Domain specific:

COPO for plant sciences



MOLGENIS for biobanking



...

MOLGENIS

Adaptable (configuration requires domain knowledge):

Proprietary ELNs/LIMS - often poor support for ontologies



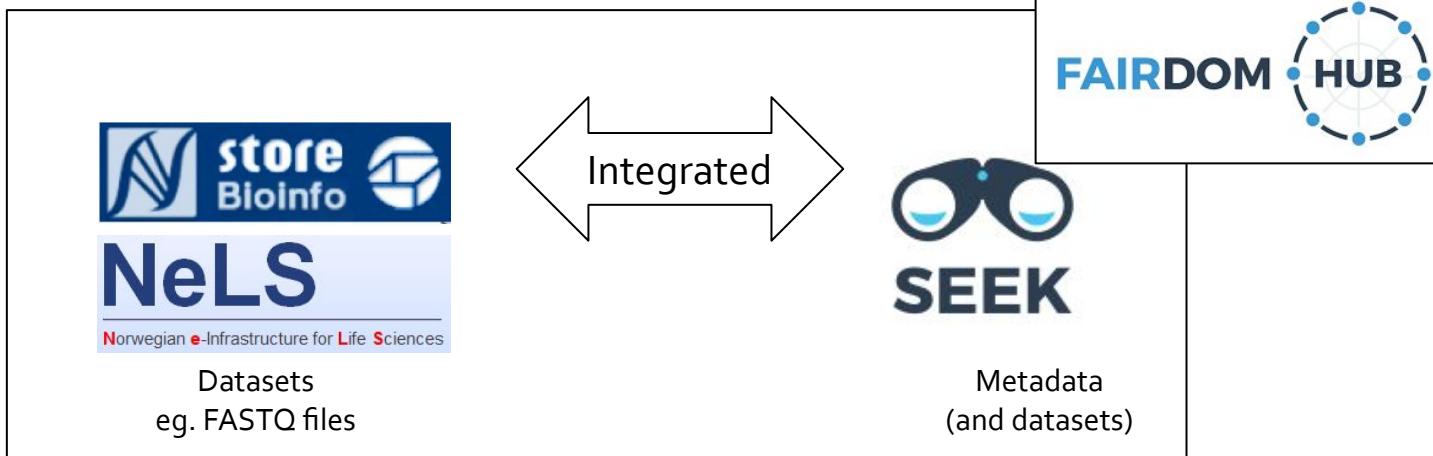
openBIS - open source ELN/LIMS

SEEK

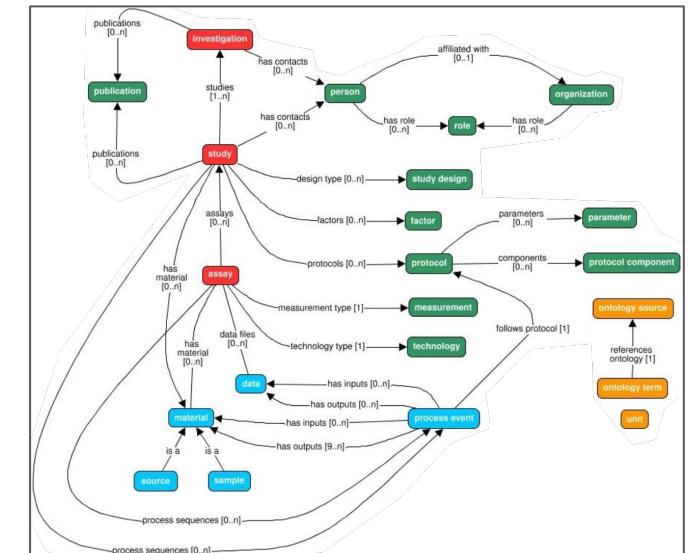
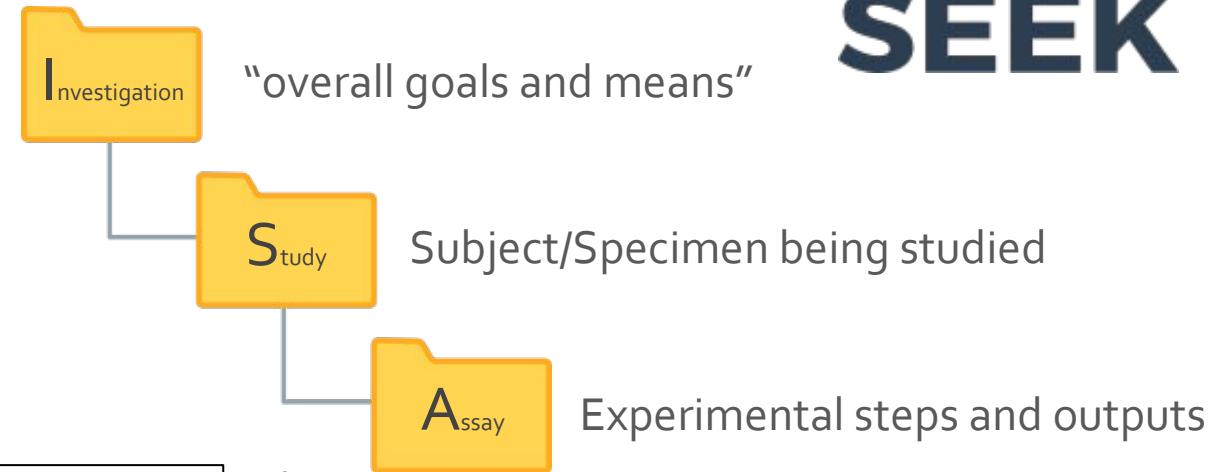


FAIRDOM SEEK

The SEEK platform is a web-based tool for organising and storing data, for sharing, exploring and annotating data

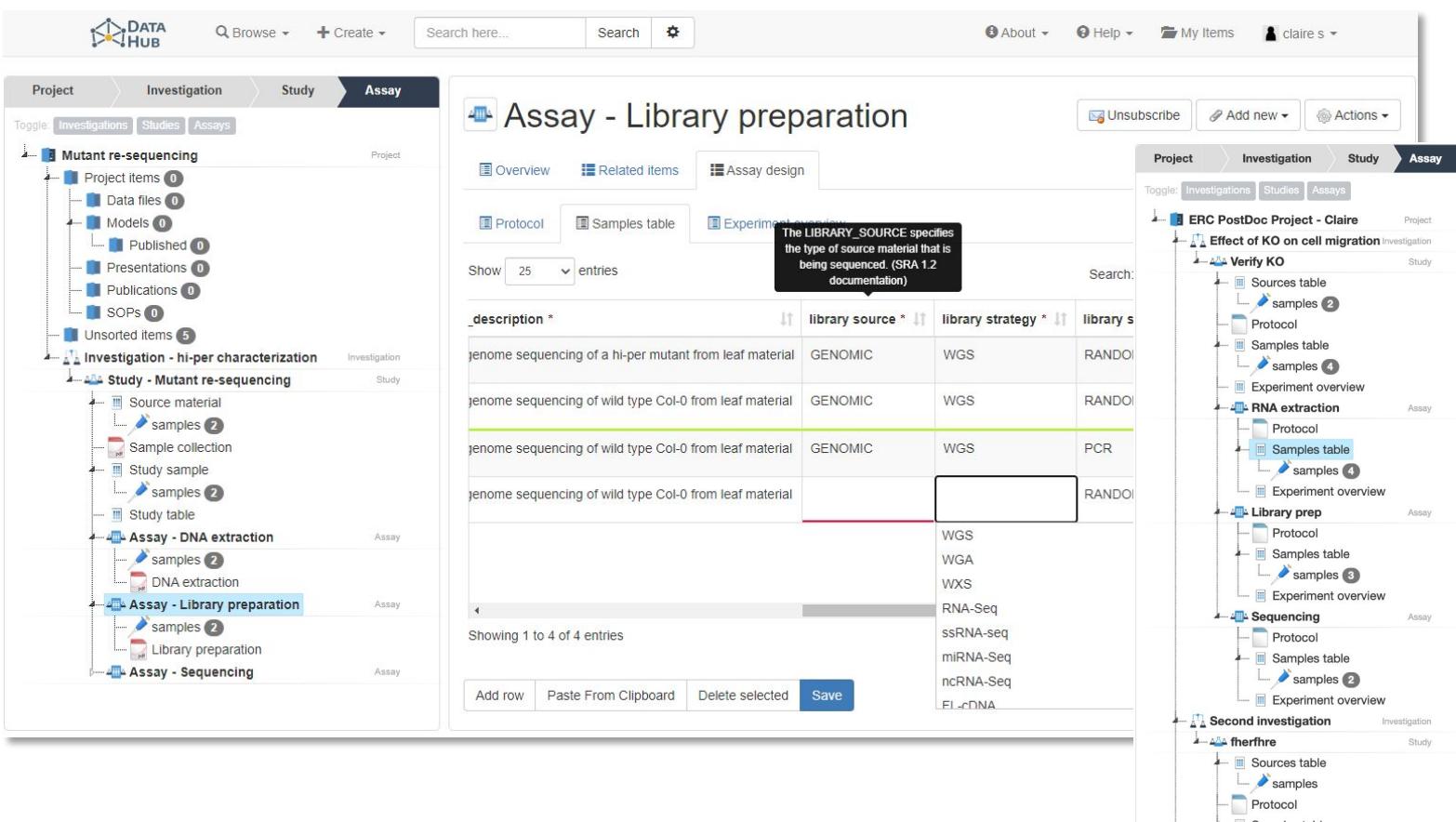


ISA model



FAIRDOM SEEK in DataHub

Dynamic table for samples metadata



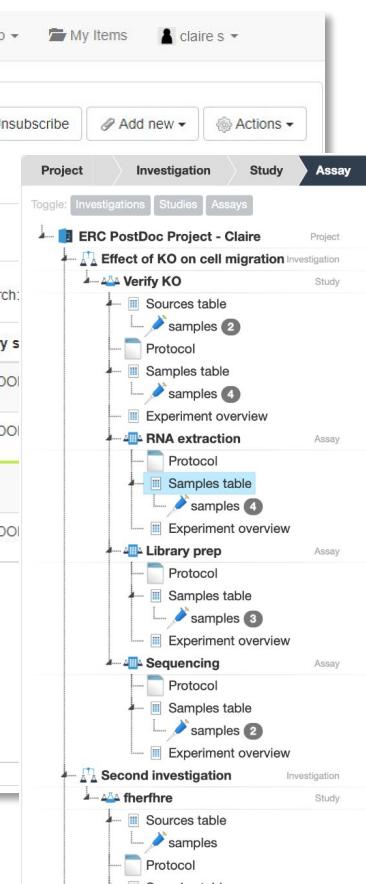
Assay - Library preparation

The LIBRARY_SOURCE specifies the type of source material that is being sequenced. (SRA 1.2 documentation)

description *	library source *	library strategy *	library s
genome sequencing of a hi-per mutant from leaf material	GENOMIC	WGS	RANDO
genome sequencing of wild type Col-0 from leaf material	GENOMIC	WGS	RANDO
genome sequencing of wild type Col-0 from leaf material	GENOMIC	WGS	PCR
genome sequencing of wild type Col-0 from leaf material			RANDO
			WGS
			WGA
			WXS
			RNA-Seq
			ssRNA-seq
			miRNA-Seq
			ncRNA-Seq
			EI -cDNA

Show 25 entries

Add row Paste From Clipboard Delete selected Save



RNA extraction

FAIRDOM-SEEK id	uuid	Input (Sample Name) *	nucleic acid extraction *	Extract Nan
101	12897400-2221-013c-43e0-7a163e608de1	x WT1	kit	WT1_RNA
102	72fa5100-2223-013c-43e0-7a163e608de1	x WT2	kit	WT2_RNA
103	730558f0-2223-013c-43e0-7a163e608de1	x KO1	kit	KO1_RNA
104	730fdc10-2223-013c-43e0-7a163e608de1	x KO2	kit	KO2_RNA

Show 1 to 4 of 4 entries

Add row Paste From Clipboard Delete selected Batch sharing permissions Batch download to Excel Save

Upload excel spreadsheet:
Choose File no file selected Upload

Batch sharing permissions for samples from the dynamic table

FAIRDOM SEEK in DataHub

Download to Excel

samples 4

Experiment overview

RNA extraction

- Protocol
- Samples table
 - samples 4
- Experiment overview

Assay

Library prep

- Protocol
- Samples table
 - samples 3
- Experiment overview

Assay

Sequencing

- Protocol
- Samples table
 - samples 2
- Experiment overview

Assay

Second investigation

fherfhe

- Sources table
- Samples
- Protocol
- Samples table

Investigation

Study

FAIRDOM-SEEK id	uuid	Input (Sample Name) *	nucleic acid extraction *	Extract Name
<input checked="" type="checkbox"/> 101	12897400-2221-013c-43e0-7a163e608de1	WT1	kit	WT1_RNA
<input checked="" type="checkbox"/> 102	72fa5100-2223-013c-43e0-7a163e608de1	WT2	kit	WT2_RNA
<input type="checkbox"/> 103	730558f0-2223-013c-43e0-7a163e608de1	KO1	kit	KO1_RNA
<input type="checkbox"/> 104	730fdc10-2223-013c-43e0-7a163e608de1	KO2	kit	KO2_RNA

Showing 1 to 4 of 4 entries

Previous 1 Next

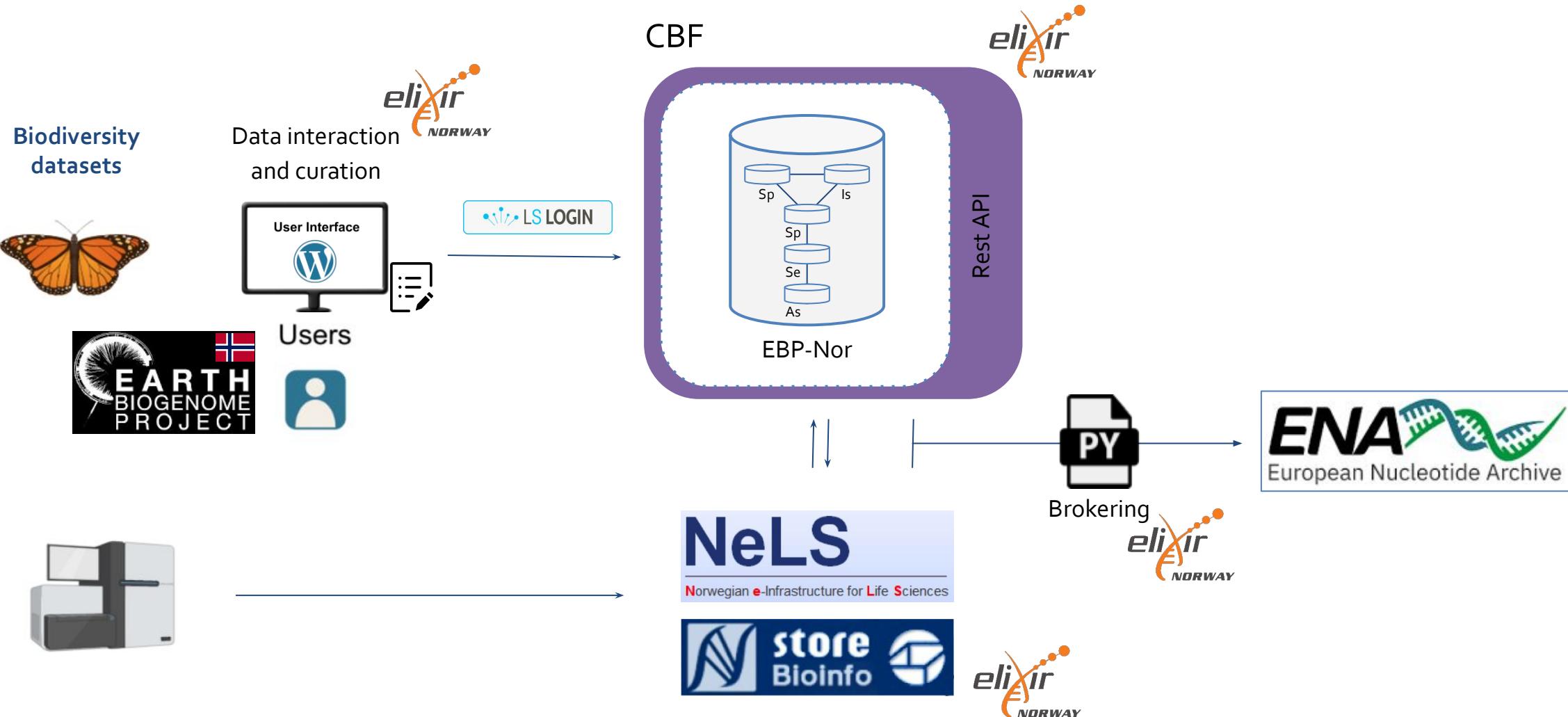
[Add row](#) [Paste From Clipboard](#) [Delete selected](#) [Batch sharing permissions](#) [Batch download to Excel](#) [Save](#)

Upload excel spreadsheet:
 no file selected [Upload](#)

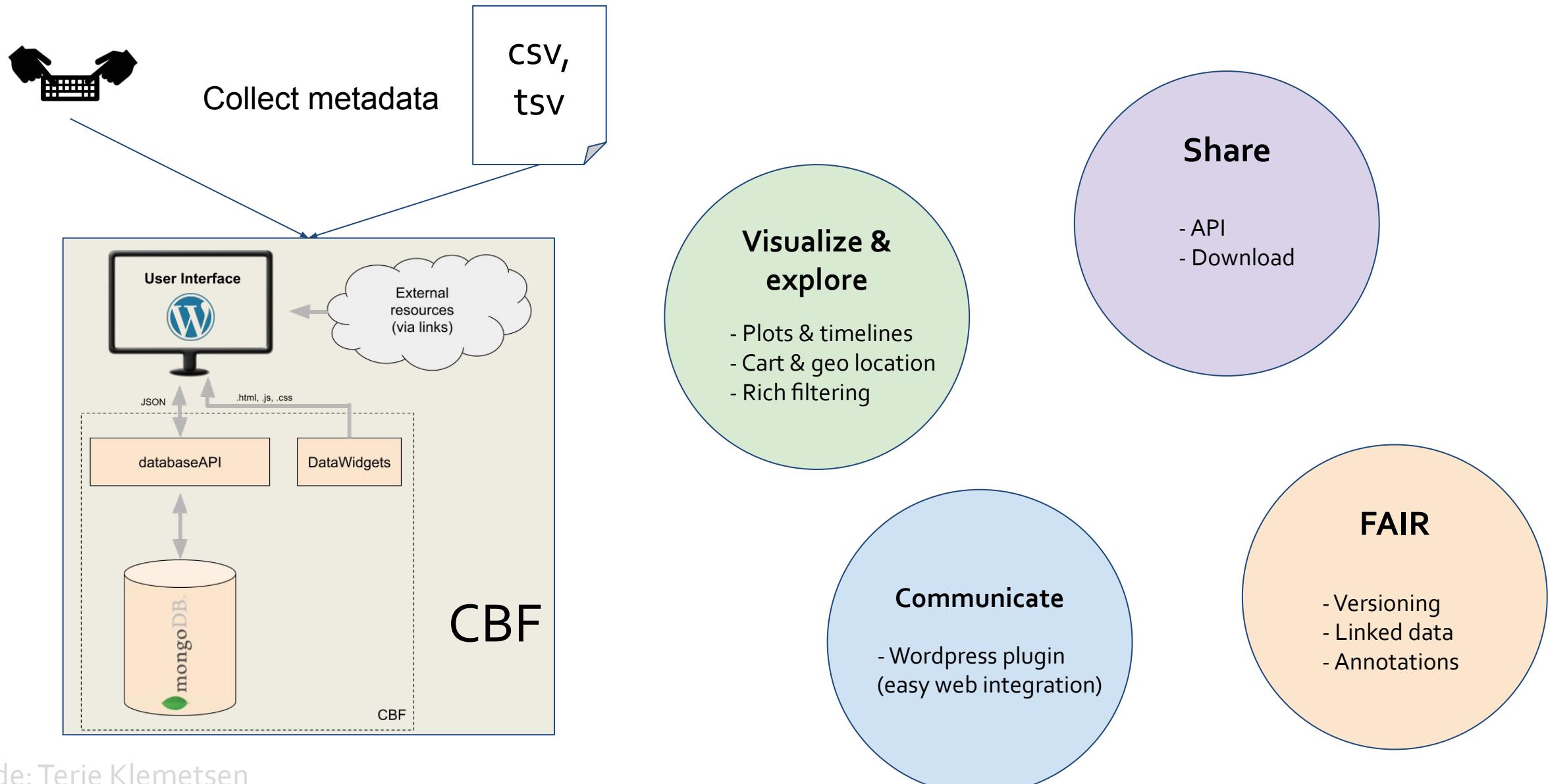
A	B	C	D	E	F	G
id	uuid	Input (Sample Name) *	nucleic	Extract Na	Material typ	*
101	12897400-2221-013c-43e0-7a163e608de1	WT1	kit	WT1_RNA	RNA	
102	72fa5100-2223-013c-43e0-7a163e608de1	WT2	kit	WT2_RNA	RNA	
<i>Locked cells</i>						

Material type
 Choose a valid option. This value is REQUIRED!

EBP-Nor use case - Brokering



Contextual Biodata Framework (CBF)





CBF

Species

1.1
Showing records 1 to 25 out of 131



Species name	Unique ID	Common name (nor.)	Common name (engl.)	Family	Isolate	Specimen	Sequencing	Assembly
Bathyraja spinicauda								
Amblyraja hyperborea								
Arctogadus glacialis								
Anelasma squamicola								
Sebastes norvegicus								
Zootoca vivipara								
Strongylocentrotus								
Sphagnum arcticum								
Vulpes lagopus								
Tolypella normaniana								
Strobilurus esculentus								
Sphagnum annulatum								
Sphagnum jensenii								
Sphagnum balticum								
Sphagnum compactum								
Usnea longissima								
Triturus cristatus								
Umbellopsis vinacea								
Thraustochytrium sp.								
Synocium turgens								
Trichopterus loricatum								
Schistidium bryophilum								

Marine Metagenomics Portal



MarDB

MarDB

1.6
Showing records 1 to 25 out of 25322



SARS-CoV-2 Database



SARS-CoV-2 Browser

2.14
Showing records 1 to 25 out of 3026244

Text search...



Log in

Page size 25

ID	Accession	Isolate name	WHO label	Pango lineage	Nextstrain clade	Anatomic part	Collection date	Isolation country	Geographic location
SFB_COVID19_OV883982	OV883982.1	- Hcov-19/switzerland/vd-chuvg-en10348/2022 - Hcov-19/switzerland/vd-chuvg-en10348/2022	-	BA.1.1	-	-	2022-01-31	Switzerland	Canton of Vaud (VD)
SFB_COVID19_OV883981	OV883981.1	- Hcov-19/switzerland/so-uhb-43327720/2022 - Hcov-19/switzerland/so-uhb-43327720/2022	-	BA.1.1	-	-	2022-02-04	Switzerland	Canton of Solothurn (SO)
SFB_COVID19_OV883980	OV883980.1	- Hcov-19/switzerland/bs-uhb-43322790/2022 - Hcov-19/switzerland/bs-uhb-43322790/2022	-	BA.1.1	-	-	2022-02-02	Switzerland	Canton of Basel-Stadt (BS)
SFB_COVID19_OV883979	OV883979.1	- Hcov-19/switzerland/bs-uhb-	-	BA.1	-	-	2022-02-05	Switzerland	Canton of Basel-Stadt

Functionalities - Editing an entry

EBP-Nor

Howdy, Terje Klemetsen

EBP-Nor DB

Showing records 1 to 10 out of 127

Scientific name	Common name (Nor)	Common name (Eng)	Family	NCBI taxon ID	Sampling status	Sequencing status	Assembly status	Annotation status
Aurantiochytrium sp.	Encellet marin thraustochytrid	Unicellular marine thraustochytrid	Thraustochytriaeae	1689870	Not collected	-	-	-
Colura calyptrifolia	Levermose	Epiphytic liverwort	Lejeuneaceae	1187000	Not collected	-	-	-
Cortinarius alpinus	Fjellsøpp	Ectomycorrhizal fungus	Cortinariaceae	858854	Not collected	-	-	-
Cortinarius caperatus	Rimsøpp	Gypsy mushroom	Cortinariaceae	75324	Not collected	-	-	-
Cortinarius cordatae	Ladegårdssøpp	Fungus	Cortinariaceae	-	Not collected	-	-	-
Cortinarius osloensis	Oslosøpp	Fungus	Cortinariaceae	418151	Not collected	-	-	-
Cunninghamella blakesleeana	Oljemugg	Fungus	Cunninghamellaceae	155726	Not collected	-	-	-
Cupophylus canescens	Tinnvokssøpp	Agaric fungus	Hygrophoraceae	889777	Not collected	-	-	-
Dasya adela	Rød alge	Red algae	Dasyaceae	-	Not collected	-	-	-
Epirrita autumnata	Fjellbjørkemåler	Autumnal moth	Geometridae	201501	Not collected	-	-	-

Showing records 1 to 10 out of 127

Page size 10

Edit

elixir NORWAY

You are watching a draft version of this record.

Cortinarius osloensis 1.0

Overview Sample Sampling event Preservation Sequencing Assembly GoT output Record information

General

① Scientific name	Cortinarius osloensis
① Common name (Eng)	Fungus
① Common name (Nor)	Oslosøpp
① Sampling status	Not collected
① Sequencing status	-
① Assembly status	-
① Annotation status	-

Cortinarius osloensis 1.0

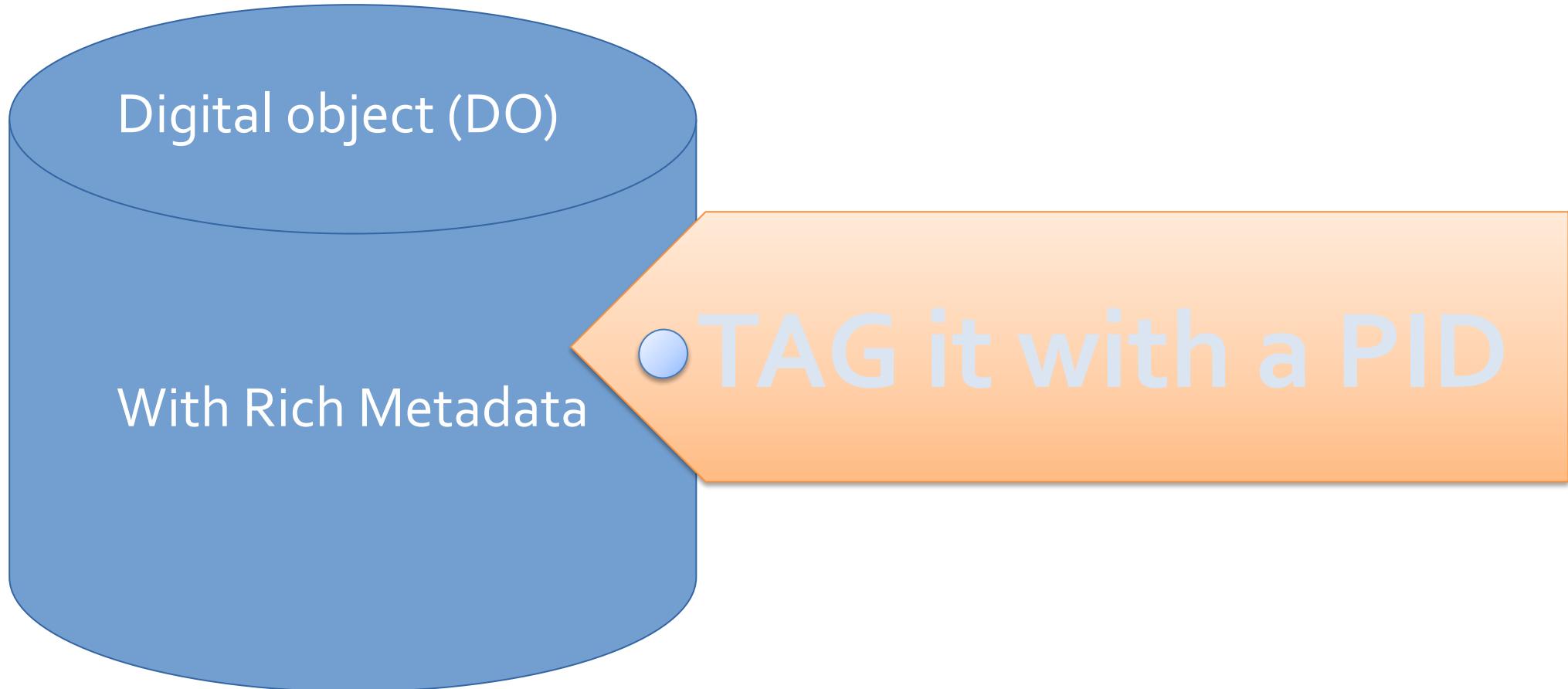
Overview Sample Sampling event Preservation Sequencing Assembly GoT output Record information

General

① Scientific name	Cortinarius osloensis
① Common name (Eng)	Fungus
① Common name (Nor)	Oslosøpp
① Sampling status	Not collected
① Sequencing status	✓ - INSDC open INSDC submitted
① Assembly status	In progress
① Annotation status	Published Sample acquired Uploaded to NeLS

Cancel Save

Make it visible

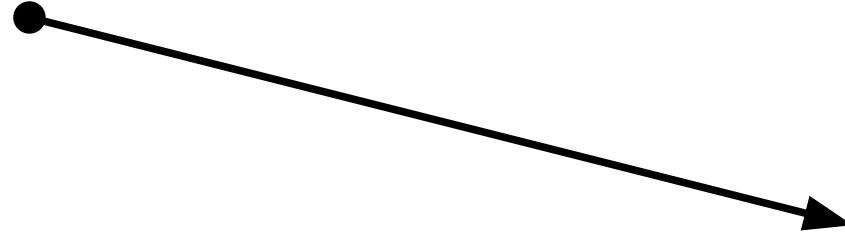


PIDs helps make data FAIR

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

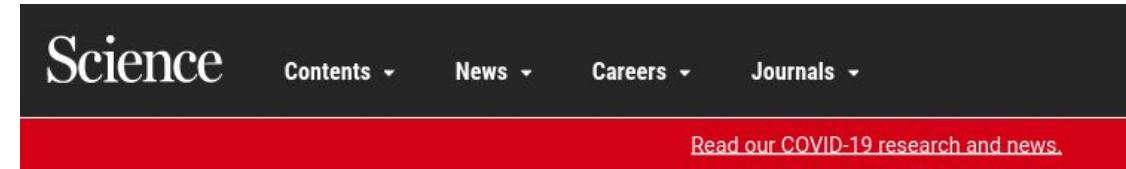
Why don't I just use a link (URL)?

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”

[farm3.staticflickr.com](#)



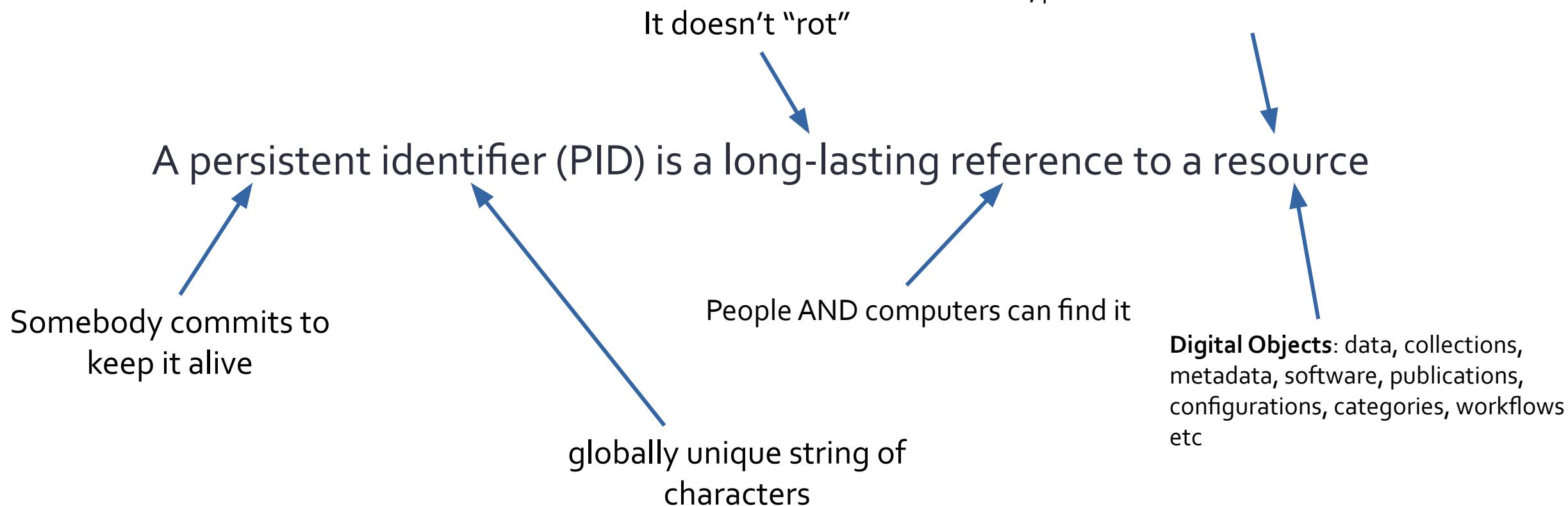
Hmmm...

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



A PID consists of two components:

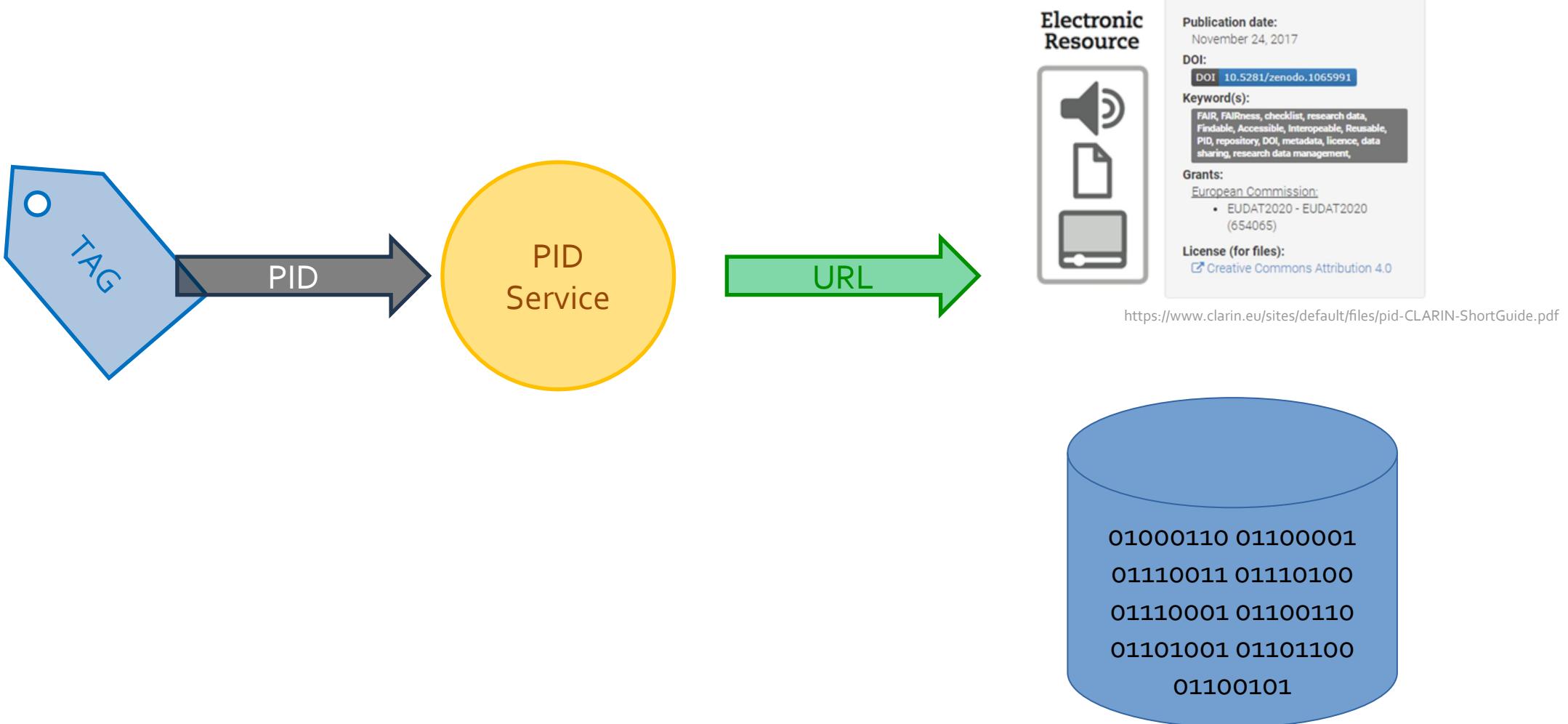
Visible string of letters and/or numbers

1. A unique identifier

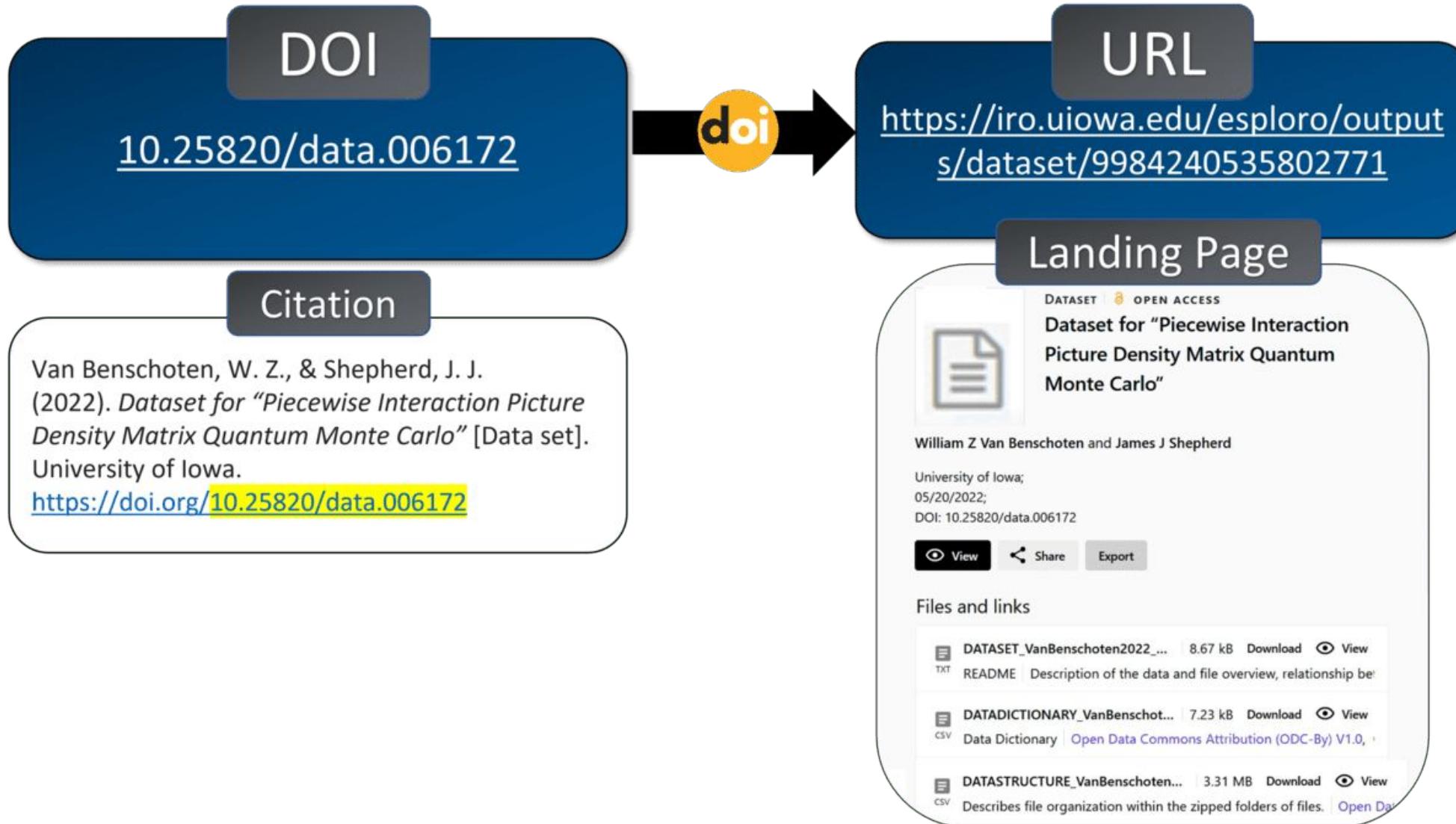
2. A service that locates the resource (or “resolves” it)

Behind the scene

Principle:

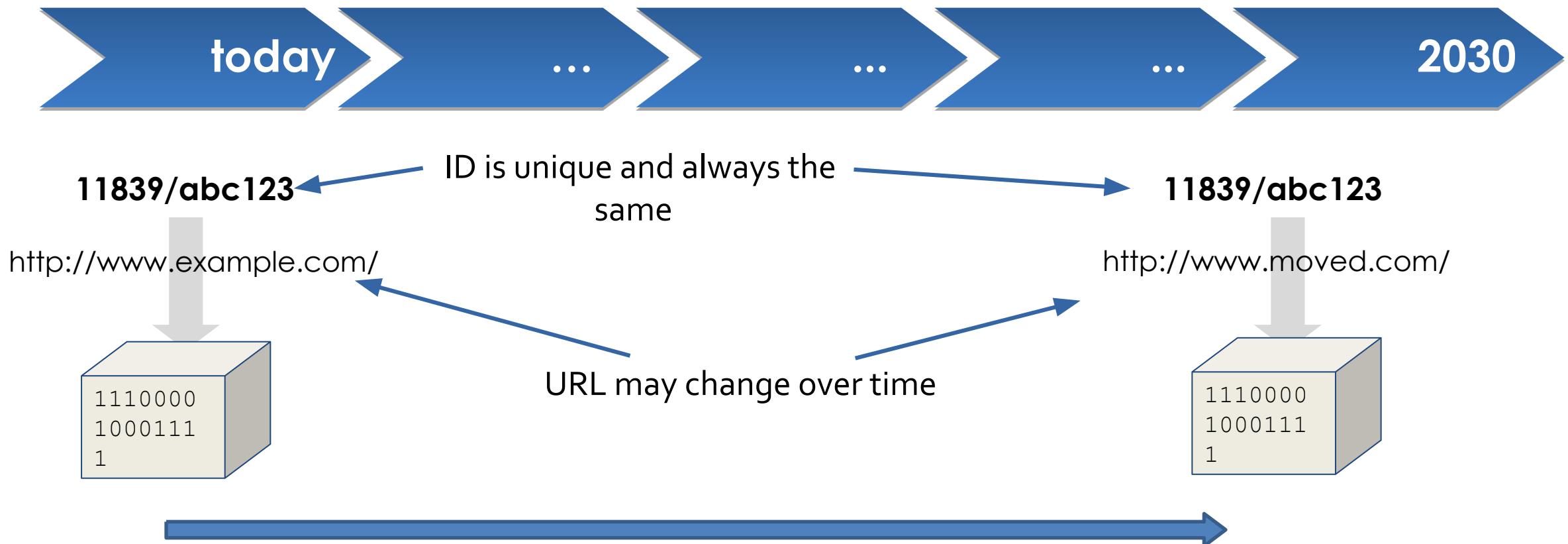


Example



Persistent over time

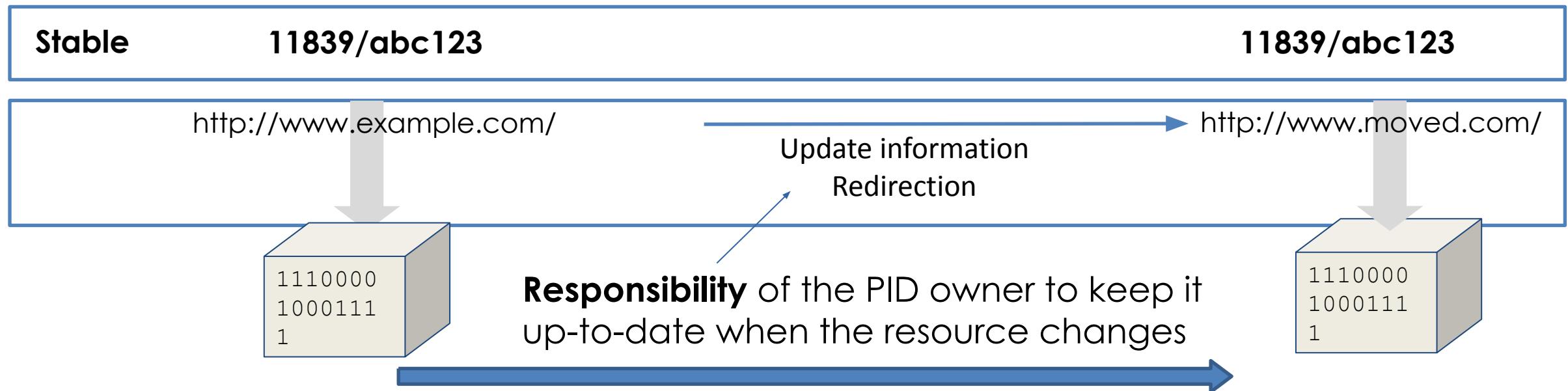
.. by design



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

Persistent over time

.. by design



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)

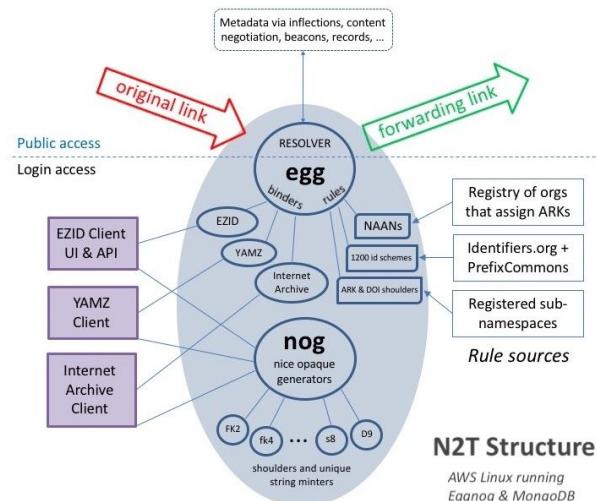


Handle.Net®

Resolver Services

N2T (Name-to-Thing)

Identifiers.org



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

The image shows the Identifiers.org Resolution Service interface. It features a large red "i" icon with a signal pattern. The top navigation bar includes links for "Resolution", "Registry", "Browse the registry", "Make a request", "Documentation", "Legacy platforms", "Sign in", and "Also in this section". The main search bar at the top right contains the placeholder "Enter a namespace to search the registry" with a "Search" button. Below the search bar, there are sections for "Identifiers.org Central Registry" and "Identifiers.org Resolution Service". The bottom right corner features the "EUROPEAN OPEN SCIENCE CLOUD" logo.

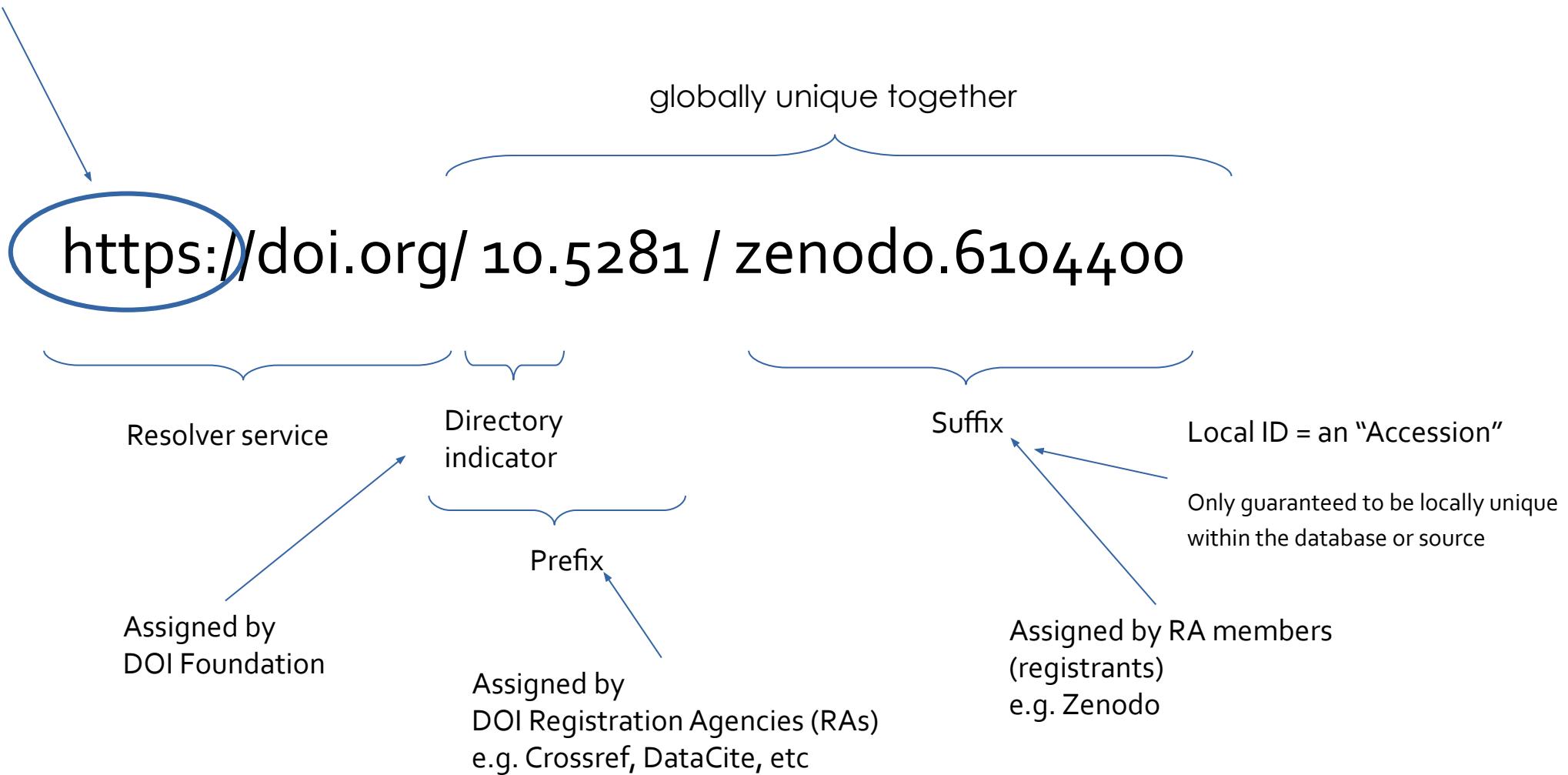
<https://eosc-portal.eu/news-and-events/news/identifiers-ensuring-robust-and-reliable-access-life-sciences-data>

How do I recognize a PID?

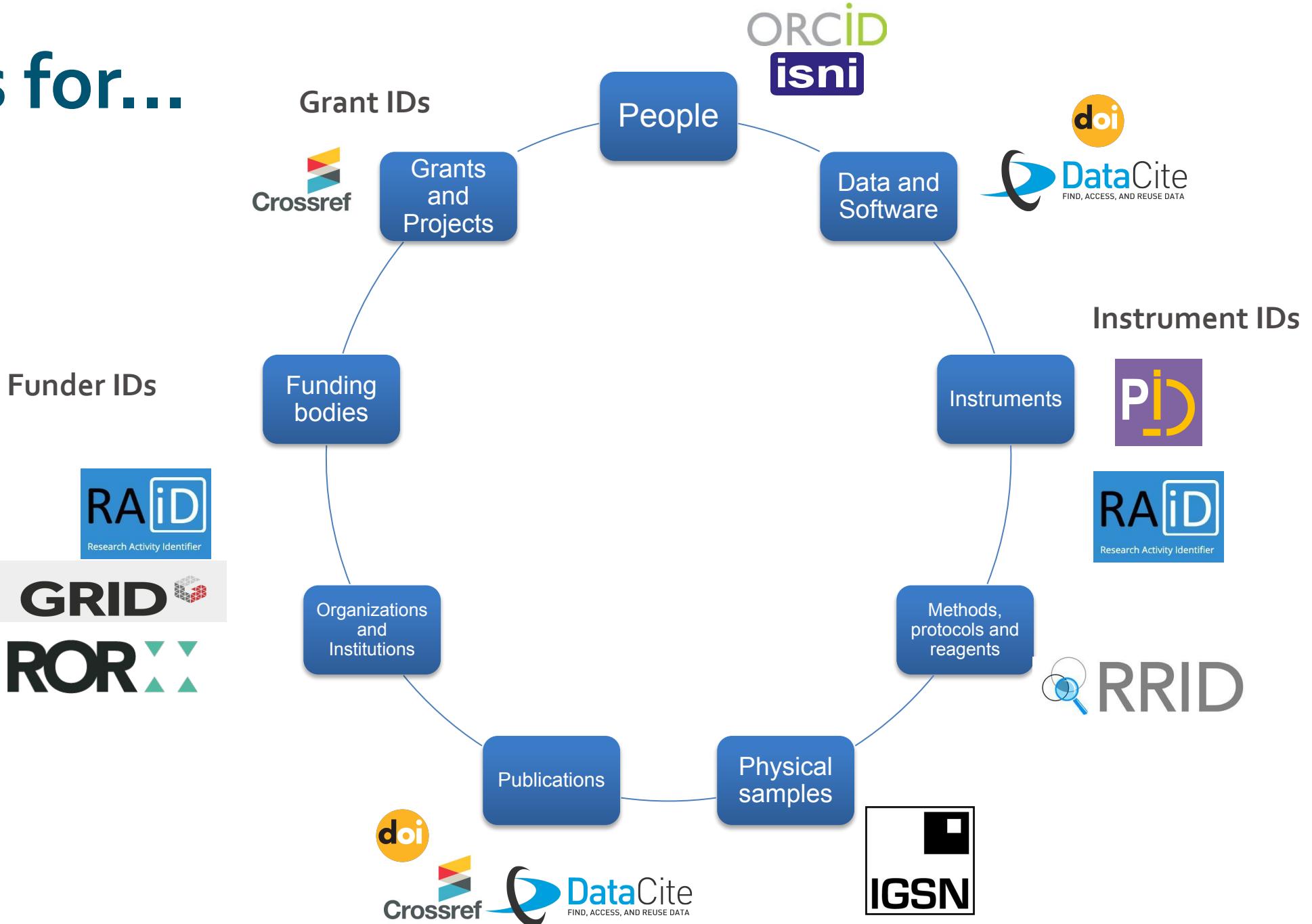
[DOI: 10.5281/zenodo.6104400](https://doi.org/10.5281/zenodo.6104400) 

Anatomy of a DOI

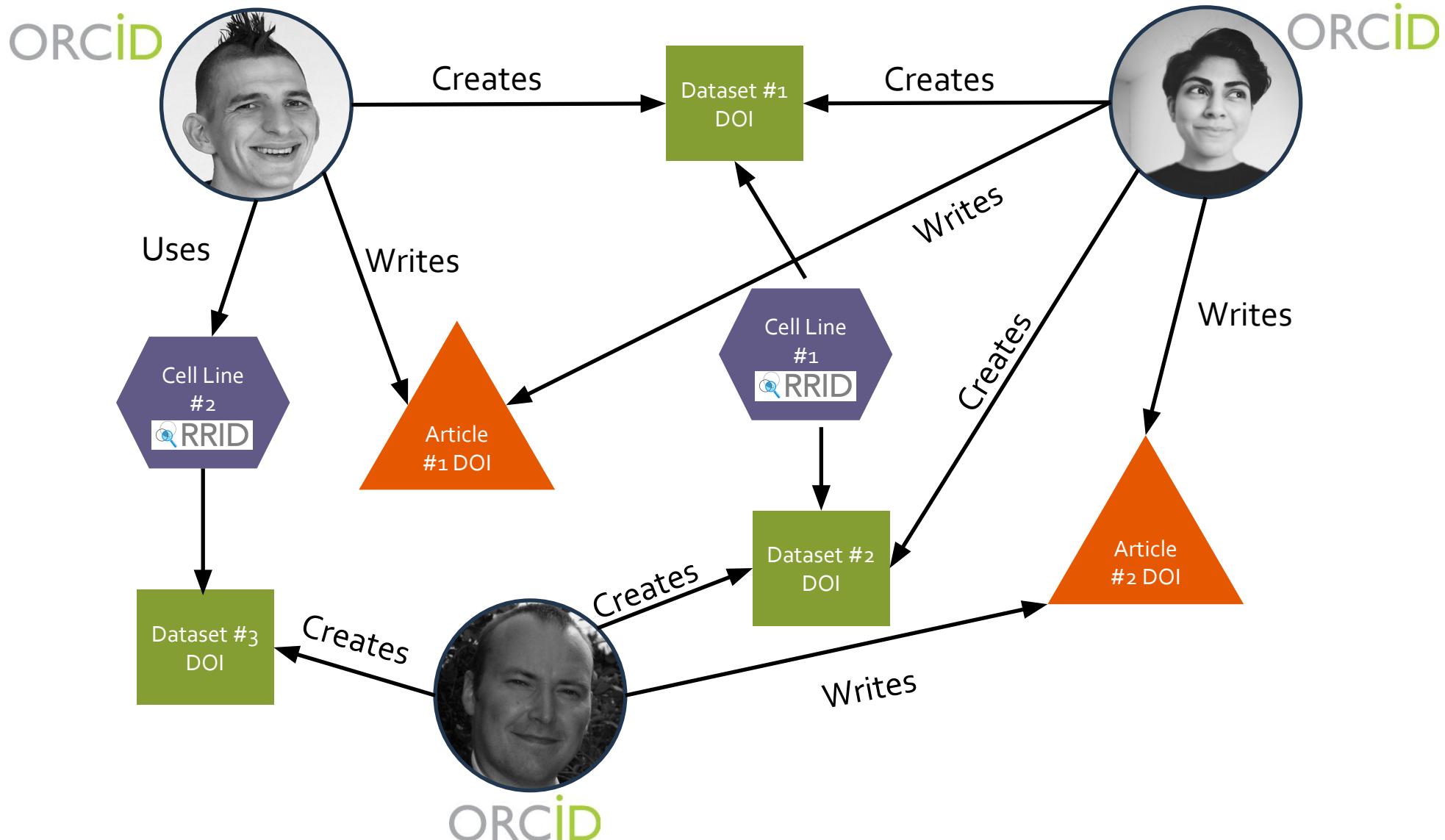
Means that it is actionable: you can paste in a web browser address bar and be taken to the identified source.



PIDs for...



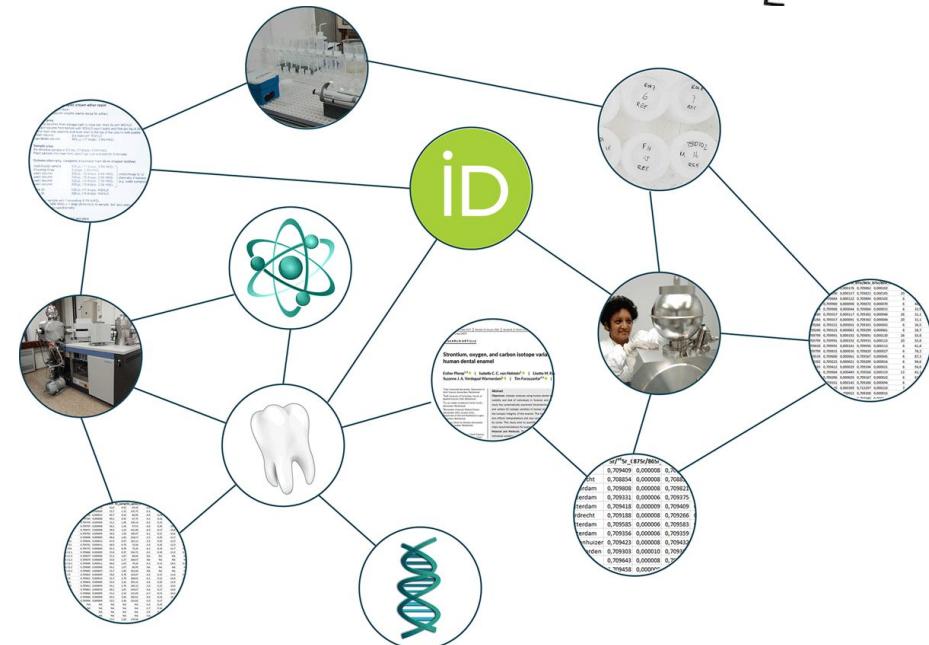
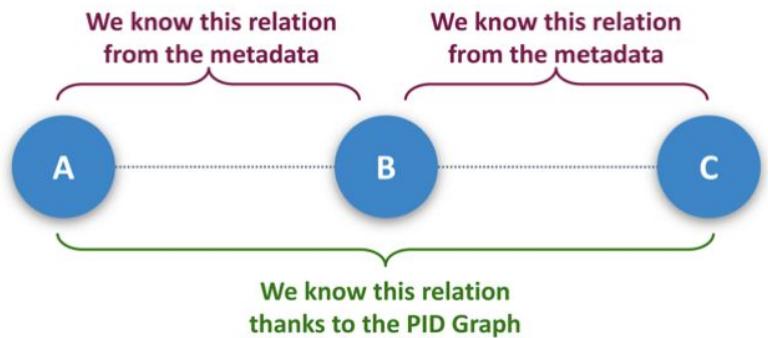
PIDs connect different entities in research





PID graphs

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





Except where otherwise noted, this work is licensed under:

<https://creativecommons.org/licenses/by/4.0/>