# PERSISTENT IDENTIFIERS AND THEIR USE CASES

## IMPROVING FINDABILITY OF DATA

CHRISTINE STAIGER

COORDINATOR DATA STEWARDSHIP

MAY, 2019

Helis Academy

# Agenda

| | |
|---|---|
| 9:00 –9:15am | Arrival and Coffee |
| 9:15 – 10:45am | Persistent Identifiers |
| 10:45 – 11:15am | Coffee |
| 11:15 - 11:45am | Data sharing, publishing and archiving |
| 11:45 - 12:00pm | Intro exercise Dataverse |
| 12:00 - 12:30pm | Exercise Dataverse |
| 12:30 – 1:30pm | Lunch |
| 1:30 – 3:00pm | Data sharing, publishing and archiving |

# PIDs – Why?

- Managing increasing numbers of **data objects**
- **Sharing data from different sources** amongst researchers

- Data needs to be **(globally) identifiable and addressable** to ensure reuse of data
- Data citation
- Linking data from different sources

- Challenges
  - Object locations change over time
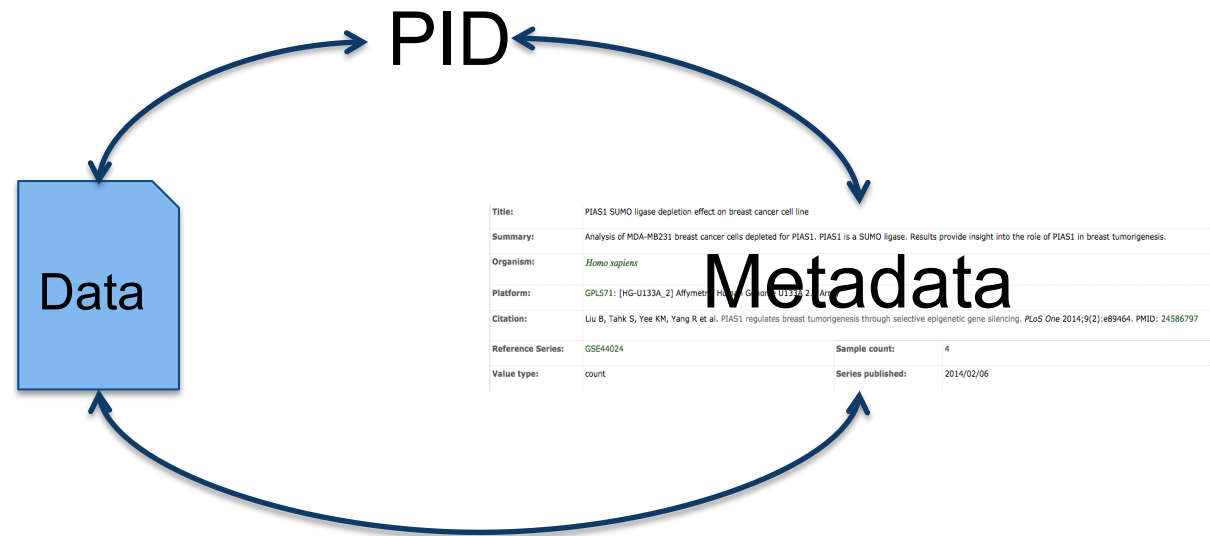  - Object migration between repositories

# What do we want from data?

- Findable – Easy to find by both humans and computer systems
  - → Expose Metadata

- Accessible – Stored for long term, accessed and/or downloaded with well-defined license and access

- Interoperable – Ready to be combined with other datasets by humans as well as computer systems;

- Reusable – Ready to be used for future research and to be processed further using computational methods.

→ Reference data and identify data
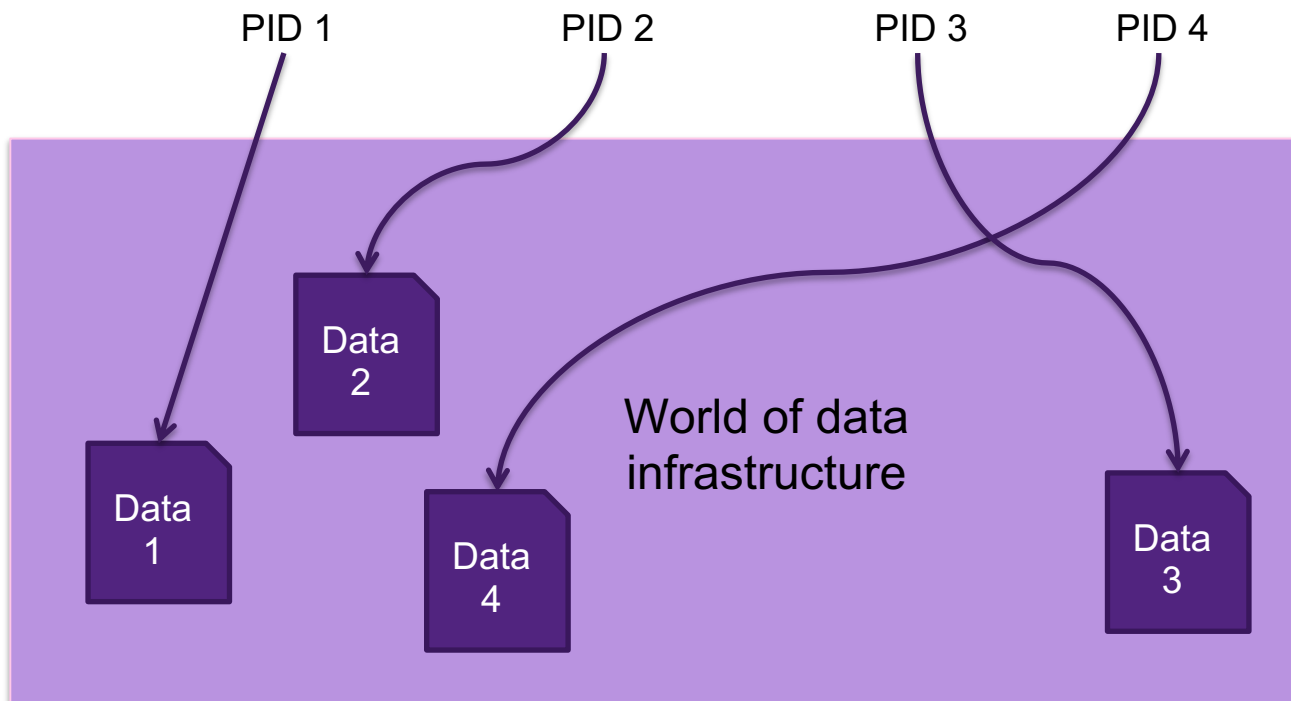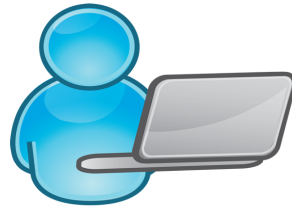→ Infrastructure should take care of some aspects

# Digital Object (DO)



- Persistent Identifier: reference and identify object, either metadata or data object

- Synchronise PID, Data and Metadata during creation, maintenance, update and deletion of a digital object!
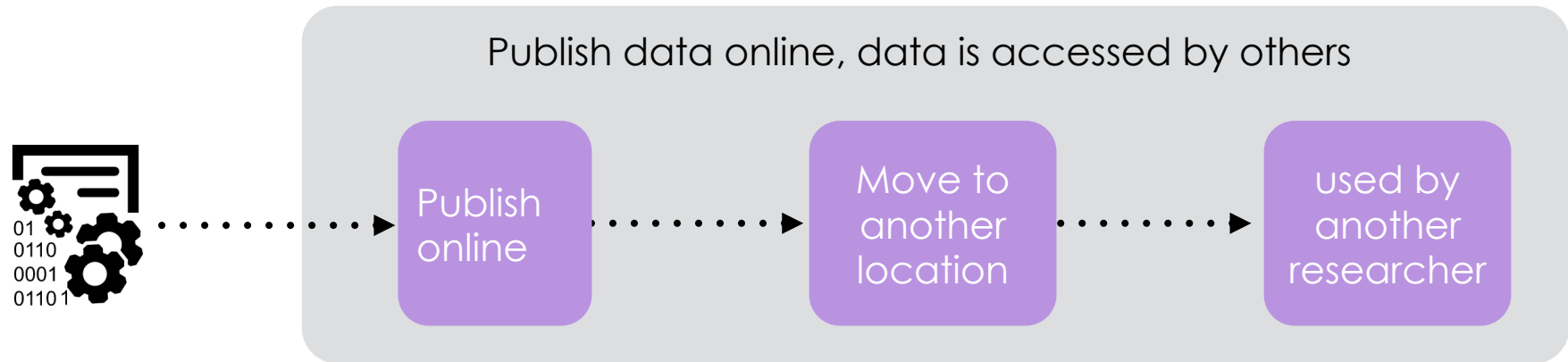
# Pro: PIDs are static

Data:

- Files

- Folders

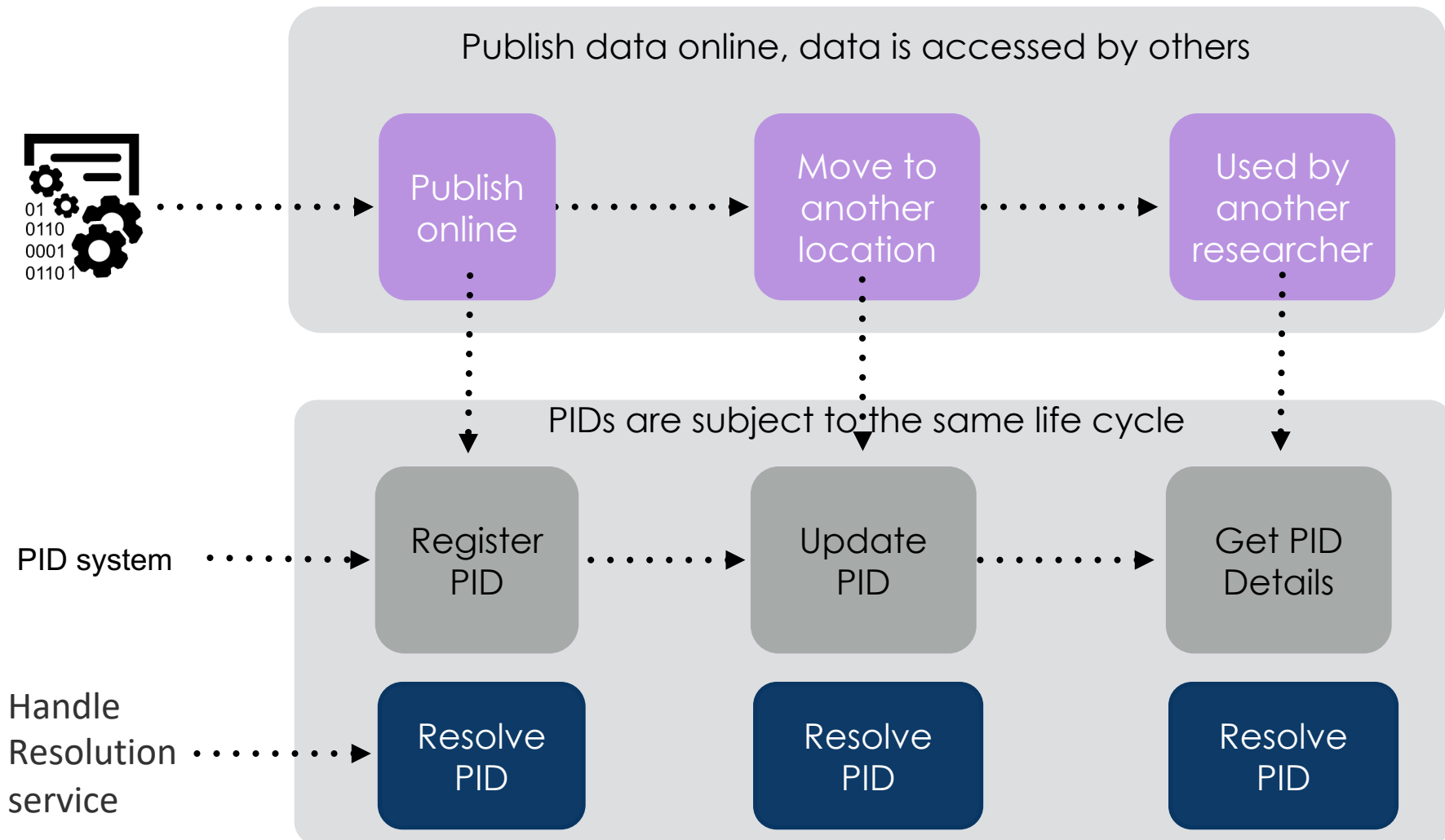- Webpages

- Sometimes even real world objects



PID 1    PID 2    PID 3    PID 4

Data 1

Data 2

Data 4

Data 3

World of data infrastructure

# Simple example of data sharing

Publish data online, data is accessed by others

Publish online → Move to another location → used by another researcher

- Published online: http://www.test.com/test.html

- Other users may cite, access, re-use this url

- Relocate the resource at http://www.example.com/

- Other users are not informed  -> 404

# Simple example of data sharing



Publish data online, data is accessed by others

| Publish online | Move to another location | Used by another researcher |

PIDs are subject to the same life cycle

PID system → Register PID → Update PID → Get PID Details

Handle Resolution service → Resolve PID · Resolve PID · Resolve PID

# Structure of a PID

11304/3265434c-4b34-11e4-81ac-dcbd1b51435e

Prefix:
- Denoting the owner of the PID
- One prefix → thousands of PIDs
- Unique in the world

Suffix:
- Specific for the thing that it identifies
- Prefix and Suffix together are unique in the world

Resolver:
- Maps PID to the target
- Web-browser compliant; HTTP redirect

http://hdl.handle.net/11304/3265434c-4b34-11e4-81ac-dcbd1b51435e

# PID Use Cases

## Use Case 1: Data publication

- PIDs point to landing page of the digital repository showing metadata

- "Real" data can be downloaded from this page with another link

- E.g. B2SHARE, FigShare, Zenodo, …

- PID
http://hdl.handle.net/11304/3265434c-4b34-11e4-81ac-dcbd1b51435e

- resolves to landing page

https://b2share.eudat.eu/records/feafb12e810c489b9e878949c6c35345

## Use Case 2: Enabling compute workflows

Molecular profiling dataflow in TraIT

# Use Case 2: Enabling compute workflows



Figure 8 Common ontology structure consists of [...] and data. CED stands for the undivided data, on [...] multiple CEDs in the other; therefore this flexibil[...]

Figure 10 The structure of RDF graph : yellow colours stand for the conceptual terms; grey colours stand for the physical location terms; the blue stand for the conceptual predicates and the green stand for the physical predicates; the dotted terms stand for the non-core terms for the structure to be compatible with different stages of realizations

Zhang, Abeln, Bijlard, Staiger: https://dx.doi.org/10.12688/f1000research.12168.1

# Use Case 3: Labelling code

- Execute program hidden behind a PID

- A way to refer to workflows → reproducibility

- Example: Identification and resolving:

```
In [16]: prefix = "841"

In [17]: suffix = "/5f6fb451-5841-11e4-9665-14109fe83170"

In [18]: ec.getValueFromHandle(prefix, "URL", suffix)
Out[18]: '/Users/christines/PIDs/helloWorld.py'

In [19]: pid = subprocess.Popen([sys.executable, ec.getValueFromHandle(prefix, "URL", suffix)])

In [20]: Hello World!
```

# PID systems

## Hands-on: Resolution

- Resolve the PIDs

- What happens if you resolve a PID
  with a foreign resolver?

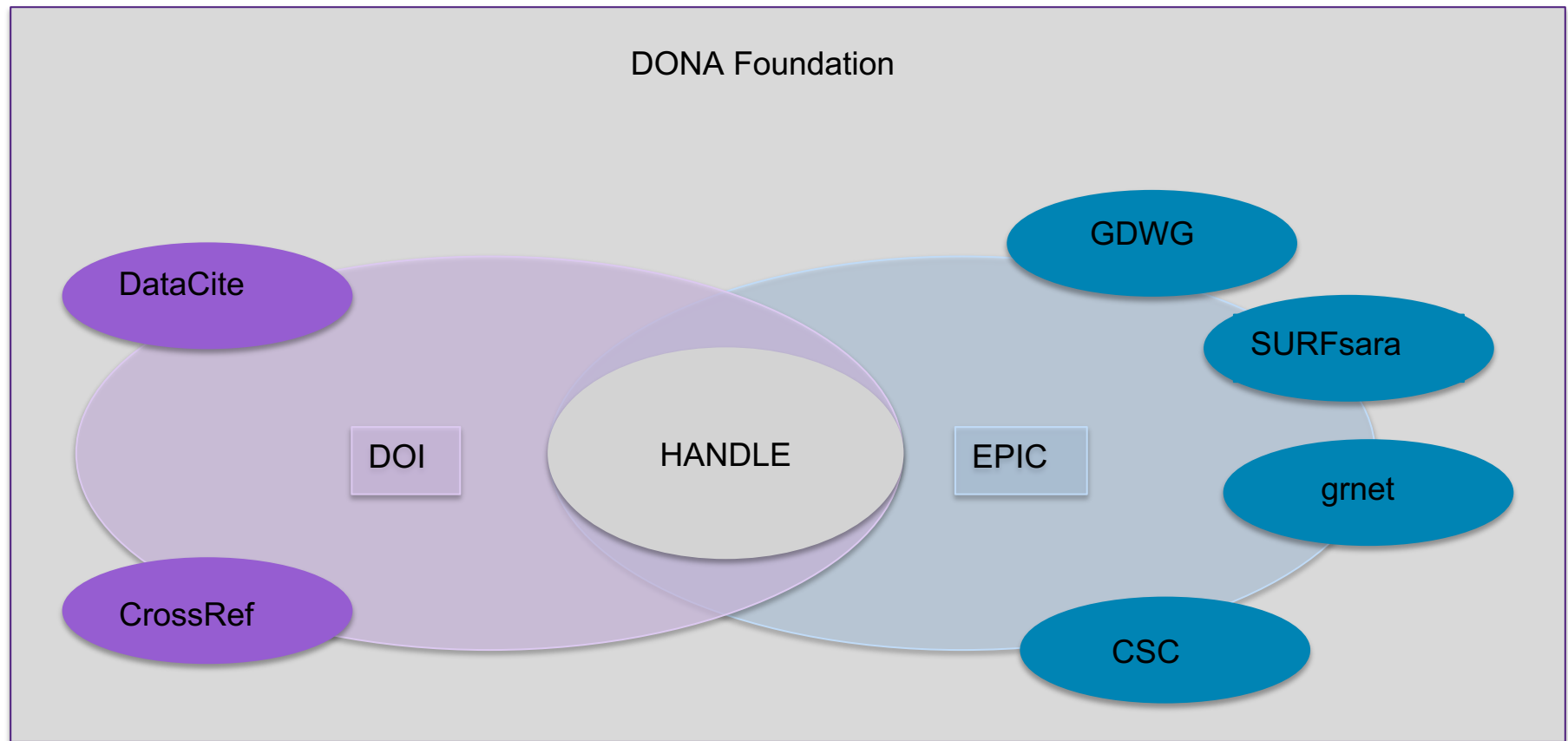# http://hdl.handle.net/21.T12995/PID-training

Exercise: Warming up!

# PID systems and issuing authorities

# PID systems and issuing authorities

- URN:NBN
    - Policies: PID is persistent and the data it is dereferenced to
    - Wants to be independent from transfer protocols
    - Currently all identifiers are compliant with http

- DOI
    - Policies: PID is persistent, data not
    - Based on the handle system
    - Datacite, Crossref are prefix issuing authorities
    - Requires extra metadata, stored in another database

- Both:
    - PIDs point to a landing page, not the file itself
        - →Taylored towards data citation
    - User needs to provide a **minimum set of metadata** (Dublin Core)

## PID systems and issuing authorities

- ePIC (European PID consortium)
    - Policies: PID is persistent, data is not
    - PIDs can point to anything
    - Based on the handle system
    - Taylored towards data identification and resolving

- DONA foundation (www.dona.net)
    - Maintains global handle registry
    - Partners:
        - CNRI (developer of the handle system)
        - GDWG (main partner in ePIC)
        - International DOI foundation (IDF)

## The Handle system – For whom?

- Metadata: You can create your own keyword-value pairs and store them with the PID

- PIDs allow to make a **distinction between data users and data managers**

- **Data users get a PID** and can directly access the data, or the metadata stored with the PID
- Pipelines can programmatically access the metadata and start specific applications

- **Requires some serious thoughts** about data organisation and developing the **code to put data policies into practice**, including code maintenance

→ For **bigger research groups or consortia** working in a distributed data environment

→ For **repositories** who are in need of a host for their PIDs

# Demo: Step-by-Step minting PIDs

- Register data with a Handle

- GET the details of a Handle

- Modify a Handle record

- Link two files on PID level

- Reverse look-up