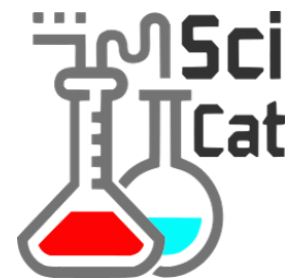Harvesting by EOSC repos, OAI-PMH, extended schemas

Gareth Murphy

European Spallation Source

# Harvesting vs authenticated search

- For EOSC/OpenAIRE/ B2FIND, we need to provide all our metadata to be "harvested". No login, anyone can access, truly open data

- For analysis/WP4 users, can query for their own (embargoed) metadata. Requires login, authentication, securing data and metadata



http://doi.org/10.17616/R31NJMKO
SciCat

http://doi.org/10.17616/R33H18
ILL Data Portal

## European Open Science Cloud (EOSC)

- Still a work in progress
- We want to provide our data and metadata to EOSC
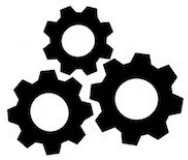- We don't know what EOSC will be like in its final form

# OpenAIRE, B2FIND

- Open access infrastructure for research in Europe
- We join them and they will connect to EOSC
-

# FAIR Digital Object

## DIGITAL OBJECT
### Data, code and other research outputs
*At its most basic level, data or code is a bitstream or binary sequence. For this to have meaning and to be FAIR, it needs to be represented in standard formats and be accompanied by Persistent Identifiers (PIDs), metadata and documentation. These layers of meaning enrich the object and enable reuse.*

## IDENTIFIERS
### Persistent and unique (PIDs)
*Digital Objects should be assigned a unique and persistent identifier such as a DOI or URN. This enables stable links to the object and supports citation and reuse to be tracked. Identifiers should also be applied to other related concepts such as the data authors (ORCIDs), projects (RAIDs), funders and associated research resources (RRIDs).*

## STANDARDS & CODE
### Open, documented formats
*Digital Objects should be represented in common and ideally open file formats. This enables others to reuse them as the format is in widespread use and software is available to read the files. Open and well-documented formats are easier to preserve. Data also need to be accompanied by the code use to process and analyse the data.*

## METADATA
### Contextual documentation
*In order for Digital Objects to be assessable and reusable, they should be accompanied by sufficient metadata and documentation. Basic metadata will enable data discovery, but much richer information and provenance is required to understand how, why, when and by whom the objects were created. To enable the broadest reuse, they should be accompanied by a plurality of relevant attributes and a clear and accessible usage license.*

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

http://doi.org/10.17616/R31NJMKO

SciCat

doi

iD ORCID
Connecting Research and Researchers

NeXus

HDF

Dublin Core Elements

| Rights | Contributor | Creator |
| Subject | Coverage | Title |
| Publisher | Identifier | Description |
| Type | Date | Source |
| Relation | Format | Language |

panosc
photon and neutron open science cloud

Wavelength  Chemical Formula
Start Date  Sample Name
Facility  Scientific Technique

according to emerging standards for trustworthiness and FAIR. The overall system and interactions between components and stakeholders are driven by metrics, incentives, investment and skills. In a European context, this FAIR ecosystem should be delivered primarily via the EOSC.
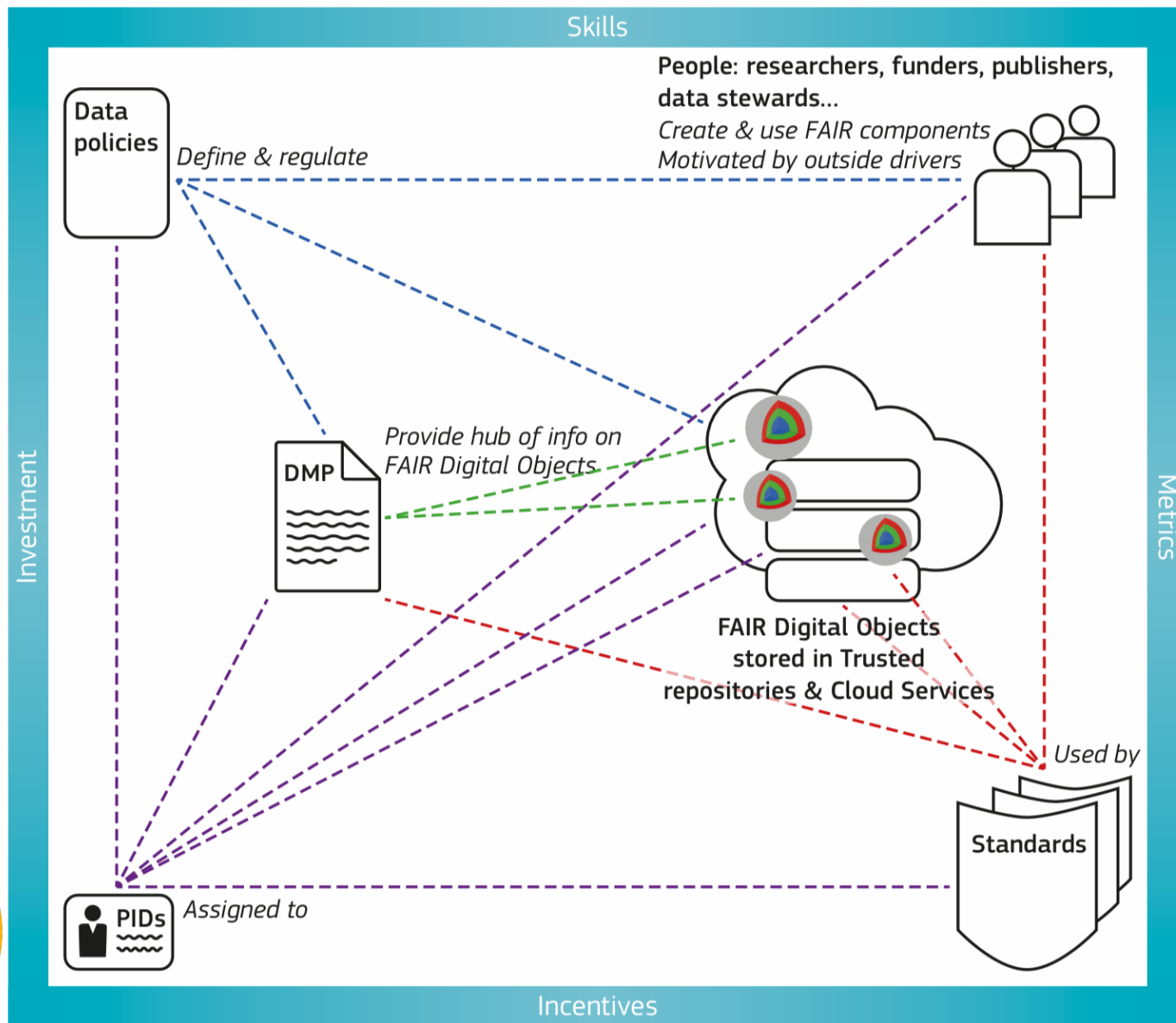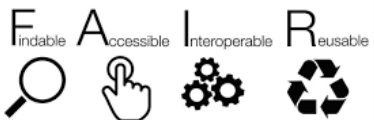


Figure 9. The interactions between components in the FAIR data ecosystem. Notes on this figure:

# Extending the schema

- Supported formats, (OAI) Dublin Core and PaN format

# OAI-PMH

- OAI-PMH (Open Archives Initiative - Protocol for Metadata Harvesting)
- 6 verbs
  - Identify
  - ListMetadataFormats
  - ListRecords
  - GetRecord
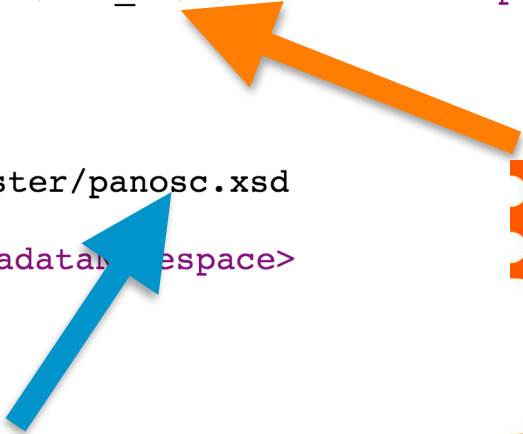  - ListSets
  - ListIdentifiers
- https://www.openarchives.org/pmh/



A PaN institute, e.g. ESS, Soleil

This XML file does not appear to have any style information associated with it. The document tree is shown below.

```xml
<OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-
instance" xmlns:dc="http://purl.org/dc/elements/1.1/"
xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/ http://www.openarchives.org/OAI/2.0/OAI-
PMH.xsd">
  <responseDate>2019-08-29T10:12:03.626Z</responseDate>
  <request verb="ListMetadataFormats">http://scicat.esss.se/scicat/oai</request>
  <ListMetadataFormats>
    <metadataFormat>
      <metadataPrefix>oai_dc</metadataPrefix>
      <schema>http://www.openarchives.org/OAI/2.0/oai_dc.xsd</schema>
      <metadataNamespace>http://www.openarchives.org/OAI/2.0/oai_dc/</metadataNamespace>
    </metadataFormat>
    <metadataFormat>
      <metadataPrefix>panosc</metadataPrefix>
      <schema>
        https://github.com/panosc-eu/fair-data-api/blob/master/panosc.xsd
      </schema>
      <metadataNamespace>http://scicat.esss.se/panosc</metadataNamespace>
    </metadataFormat>
  </ListMetadataFormats>
</OAI-PMH>
```

**Dublin Core**

Dublin Core Elements

| Rights | Contributor | Creator |
| Subject | Coverage | Title |
| Publisher | Identifier | Description |
| Type | Date | Source |
| Relation | Format | Language |

panosc
photon and neutron
open science cloud

Wavelength   Chemical Formula
Start Date   Sample Name
Facility   Scientific Technique

```xml
▼<OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
 xmlns:dc="http://purl.org/dc/elements/1.1/" xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
 http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd">
    <responseDate>2019-08-29T10:12:03.626Z</responseDate>
    <request verb="ListRecords" metadataPrefix="panosc">http://scicat.esss.se/scicat/oai</request>
  ▼<ListRecords>
    ▼<record>
      ▼<header>
          <identifier>10.17199/BRIGHTNESS/NMX0001</identifier>
          <datestamp>updatedAt</datestamp>
        </header>
      ▼<metadata>
        ▼<panosc:panosctype xmlns:panosc="http://www.openarchives.org/OAI/2.0/oai_dc/"
          xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
          xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/oai_dc/ https://raw.githubusercontent.com/panosc-eu/fair-
          data-api/master/panosc.xsd">
            <panosc:id/>
            <panosc:name>Sample Data from NMX</panosc:name>
          ▼<panosc:description>
              https://github.com/ess-dmsc/ess_file_formats/wiki/NMX
            </panosc:description>
            <panosc:owner>Dorothea Pfeiffer</panosc:owner>
            <panosc:contactEmail/>
            <panosc:orcidOfOwner/>
            <panosc:license/>
            <panosc:embargoEndDate/>
            <panosc:startDate/>
            <panosc:path/>
            <panosc:technique/>
            <panosc:sampleName/>
            <panosc:chemicalFormula/>
            <panosc:size>12496253739</panosc:size>
            <panosc:wavelength/>
          </panosc:panosctype>
        </metadata>
```

panosc

# OAI-PMH pros and cons

**Pros:**

✅ Well supported by e.g. OpenAire

✅ Lots of versions available

✅ Little implementation work required

**Cons:**

❌ OAI-PMH doesn't scale

❌ In order to change some entries, harvesters have to harvest everything again

❌ Change not supported

# ResourceSync

- Upgrade/rewrite of OAI-PMH
- Developed to address missing parts of PMH
- Supports Change List, Change Dump, Versioning,
- Sitemap technology so XML files broken into 50 MB chunks
- http://www.openarchives.org/rs/toc
- Works for any search engine

# ResourceSync

```xml
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9"
        xmlns:rs="http://www.openarchives.org/rs/terms/">
<url>
    <loc>http://example.com/res1</loc>
    <lastmod>2013-01-02T09:07:00Z</lastmod>
    <rs:md change="updated"
            hash="md5:1584abdf8ebdc9802ac0c6a7402c03b6"
            length="8876"
            type="text/html"/>
  </url>
  <url>
  …
  </url>
</urlset>
```

# Google Dataset Search

- They will scrape your metadata if it is in JSON-LD or RDF-a format

- https://datasetsearch.research.google.com/

# Roadmap

1. Deploy OAI-PMH at 6 PaNOSC institutes
2. Register at re3data
3. Provide data to B2FIND/ OpenAire
4. Upgrade to ResourceSync



European XFEL Data Portal



SciCat



ILL Data Portal

# Roadmap

| | Task | Target | CERIC | ELI | ESRF | ESS | ILL | XFEL |
|---|---|---|---|---|---|---|---|---|
| 1 | Deploy OAI-PMH | March 2020 | ✅ | | | ✅ | ✅ | |
| 2 | Re3data registration | March 2020 | ✅ | | ✅ | ✅ | ✅ | ✅ |
| 3 | OpenAIRE registration | March 2020 | ✅ | | | ✅ | ✅ | |
| 4 | Upgrade to ResourceSync | 2021 | | | | | | |

panosc

provide.openaire.eu

EXPLORE    PROVIDE    CONNECT    MONITOR    DEVELOP

DASHBOARD        SOURCES        COMPATIBILITY        CONTENT        METRICS

OpenAIRE | PROVIDE

GARETH CHARLES MURPHY

## Register

Register data sources in the OpenAIRE infrastructure

## Validate

Validate data sources against OpenAIRE guidelines

MY DATASOURCES AT A GLANCE

**Retrieving your datasources...**

## Notifications

View notifications to enrich the metadata and the content

## Metrics

View aggregated, cleaned usage statistics for repository access

nosc

**DASHBOARD**    **SOURCES**    **COMPATIBILITY**    **CONTENT**    **METRICS**

OpenAIRE | PROVIDE

**GARETH CHARLES MURPHY**

| | | | | |
|---|---|---|---|---|
| Property Identifier & identifierType (M) | The type of the Identifier.<br>View guideline | 5 | 0/456 | ✕<br>View Errors |
| Property Identifier (M) | The Identifier is a unique string that identifies a resource.<br>View guideline | 5 | 0/456 | ✕<br>View Errors |
| Property Language (R) | The primary language of the resource.<br>View guideline | 3 | 0/456 | ⚠<br>View Warnings |
| Property PublicationYear (M) | The year when the data was or will be made publicly available.<br>View guideline | 5 | 0/456 | ✕<br>View Errors |

The name of the entity that holds, archives, publishes, prints, distributes, releases, issues,

nosc

provide.openaire.eu

DASHBOARD    SOURCES    COMPATIBILITY    CONTENT    METRICS

OpenAIRE | PROVIDE

# Previous validations

**GARETH CHARLES MURPHY**

**Filter by job type:**    --none selected--

**Filter validation jobs:**    All jobs (9)    successful (2)    failed (7)    ongoing (0)    **Jobs per page:**    10

‹ Previous    page 1 of 1    Next ›

| REPOSITORY | VALIDATION TYPE | STATUS | SCORE | STARTED | GUIDELINES | ACTIONS | |
|---|---|---|---|---|---|---|---|
| https://scicat.esss.se/scicat/oai | OAI Content | finished | 100 | 2020-01-14 | For Data | View Results | ● |
| | OAI Usage | finished | 100 | 10:36:03 | Archives (2.0) | › | |
| | | | | | | Resubmit Job | |
| | | | | | | ↻ | |

21

nosc

EXPLORE    PROVIDE    CONNECT    MONITOR    DEVELOP

DASHBOARD    SOURCES    COMPATIBILITY    CONTENT    METRICS

OpenAIRE | PROVIDE

# Manage your datasources

GARETH CHARLES MURPHY

SciCat

Congratulations! Your repository was successfully registered in OpenAIRE. You can download this logo to use in your site.

Validated

Download

OpenAIRE

**Dashboards**

**Support**

**Updates**

Explore

NOADs

News

nosc

22