



ORCID DE Workshop

Organization identifiers

# Integration of organizations in the OpenAIRE Research Graph

Alessia Bardi

*National Research Council, Italy*



ORCID DE Workshop on organisation identifiers - 2 December 2020



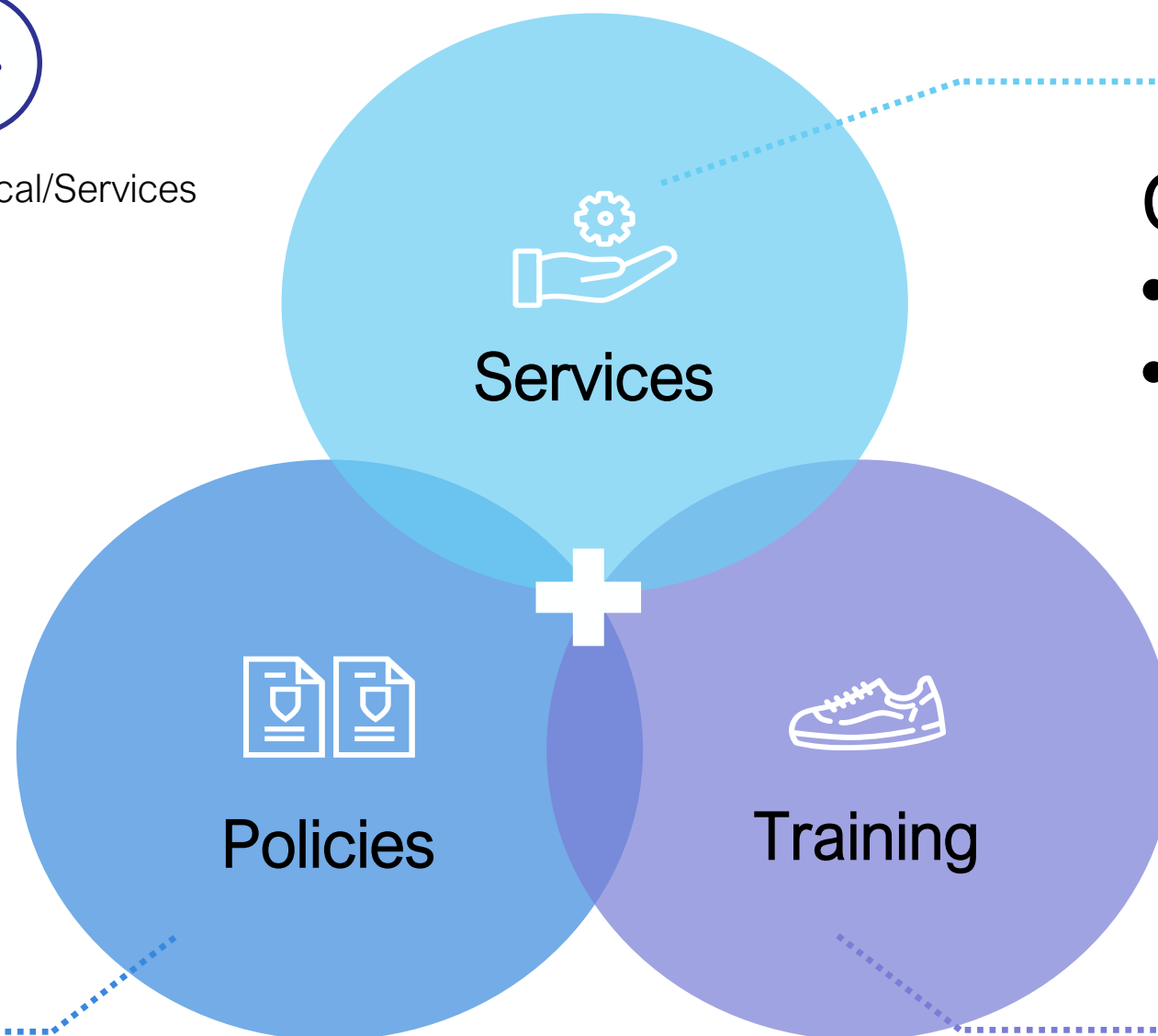
# OpenAIRE: three pillars of action

[www.openaire.eu](http://www.openaire.eu)



## Aligning

- *Standard*
- *Guidelines*
- *Practices*
- *Workflows*

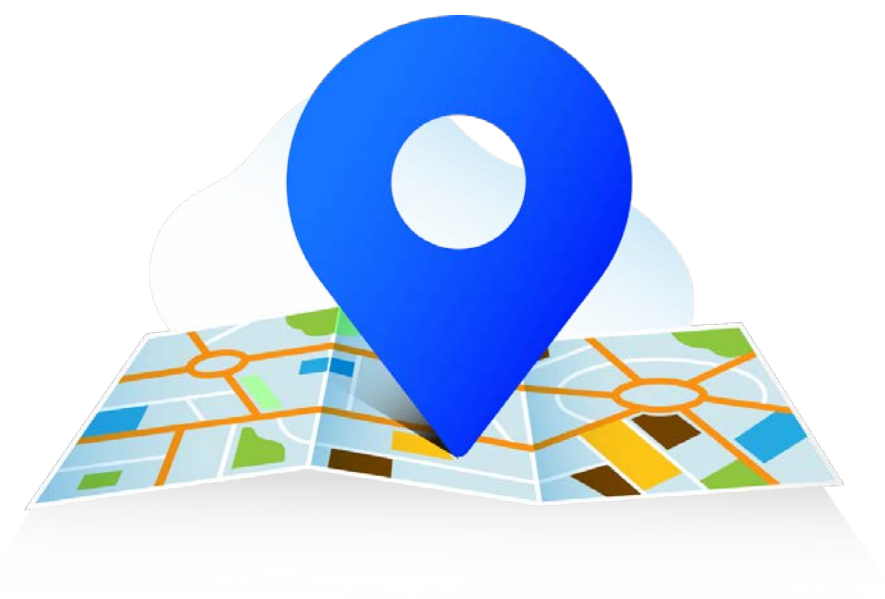


## Connecting

- *Establishing infrastructure*
- *Added value services*

## Empowering:

- *Open Science*
- *Open Access*
- *Policies*
- *Services*



## Tracking Open Science

Reproducibility and transparency require tracking of all outcomes of science and related "context"



## Monitoring Open Science

Monitoring quality, impact, and "open scienceness" of science should be a transparent, reproducible process for all, inclusive of research "context"



## Discovery Open Science

Discovery of reproducible science outcomes must find new ways, driven by "scientific intentions" that go beyond the "find articles related to a research topic".

# OpenAIRE Open Science Monitoring



Research  
admins

Research impact

**MONITOR** **CONNECT**



Funders



Research  
communities

Open Access/Science  
trends

**MONITOR** **CONNECT**



Research  
Organization



Research  
Infrastructures



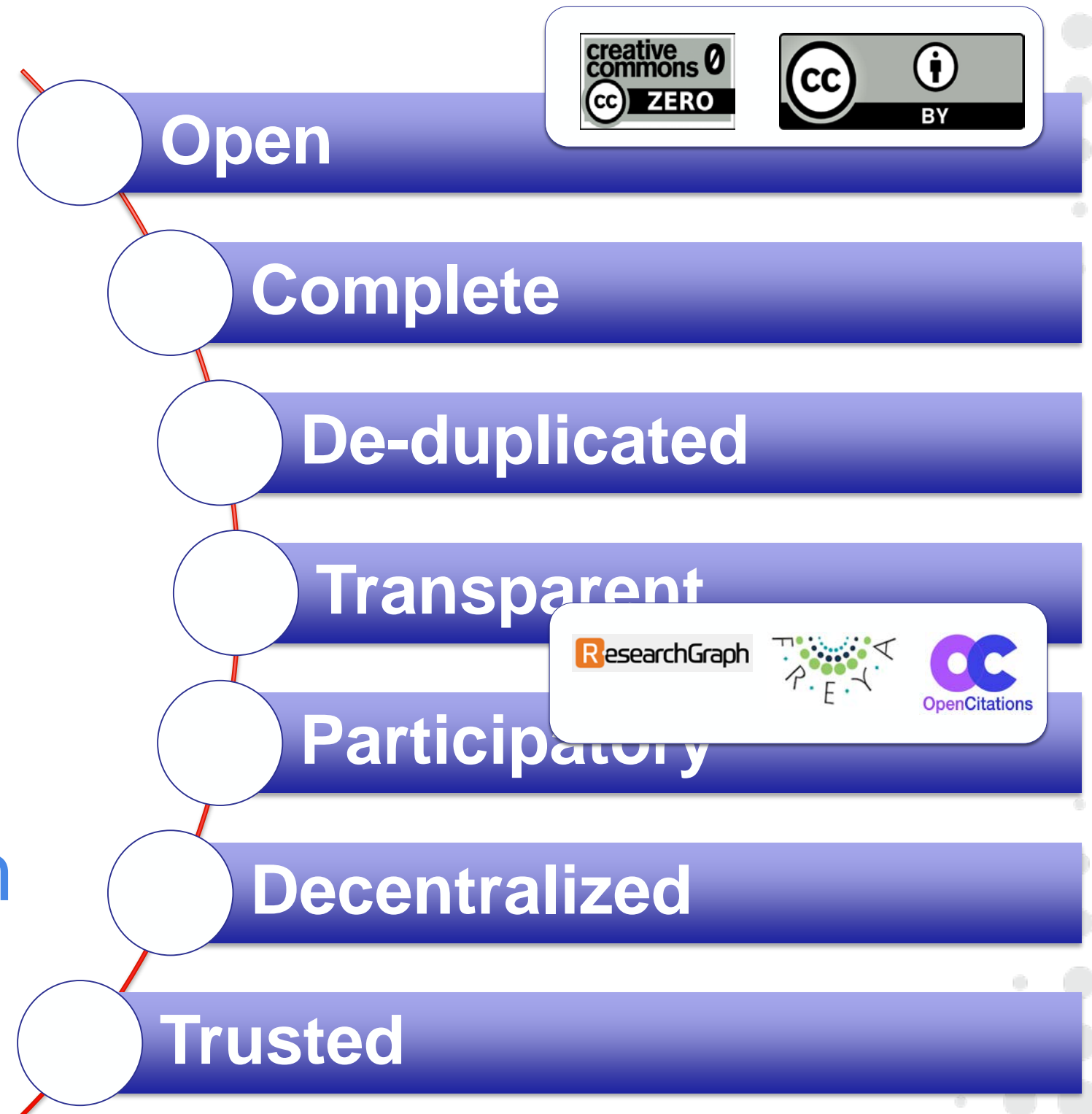
[www.openaire.eu](http://www.openaire.eu)

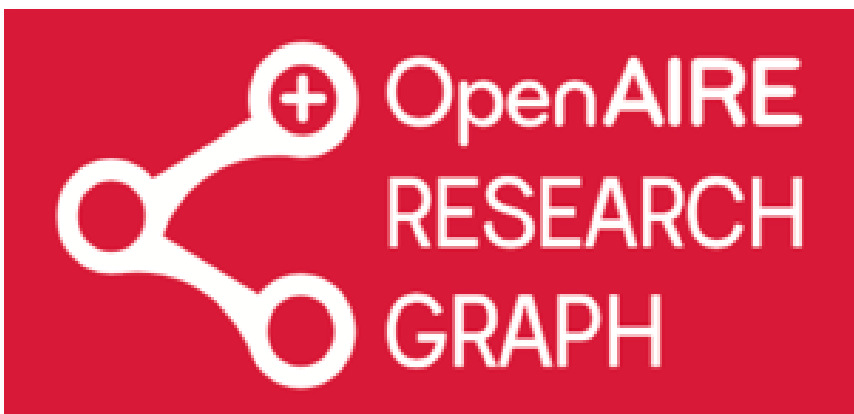
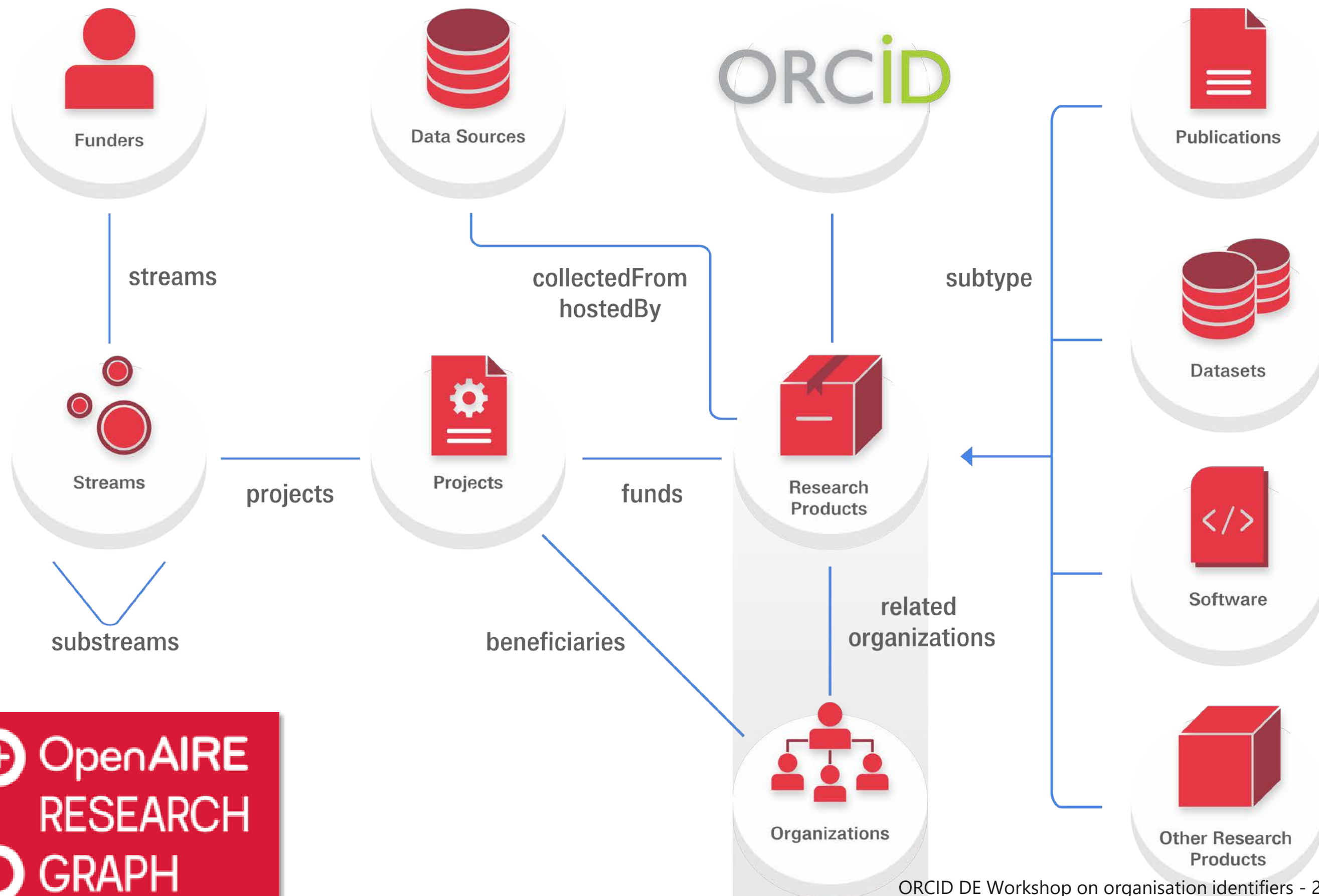
<https://monitor.openaire.eu>

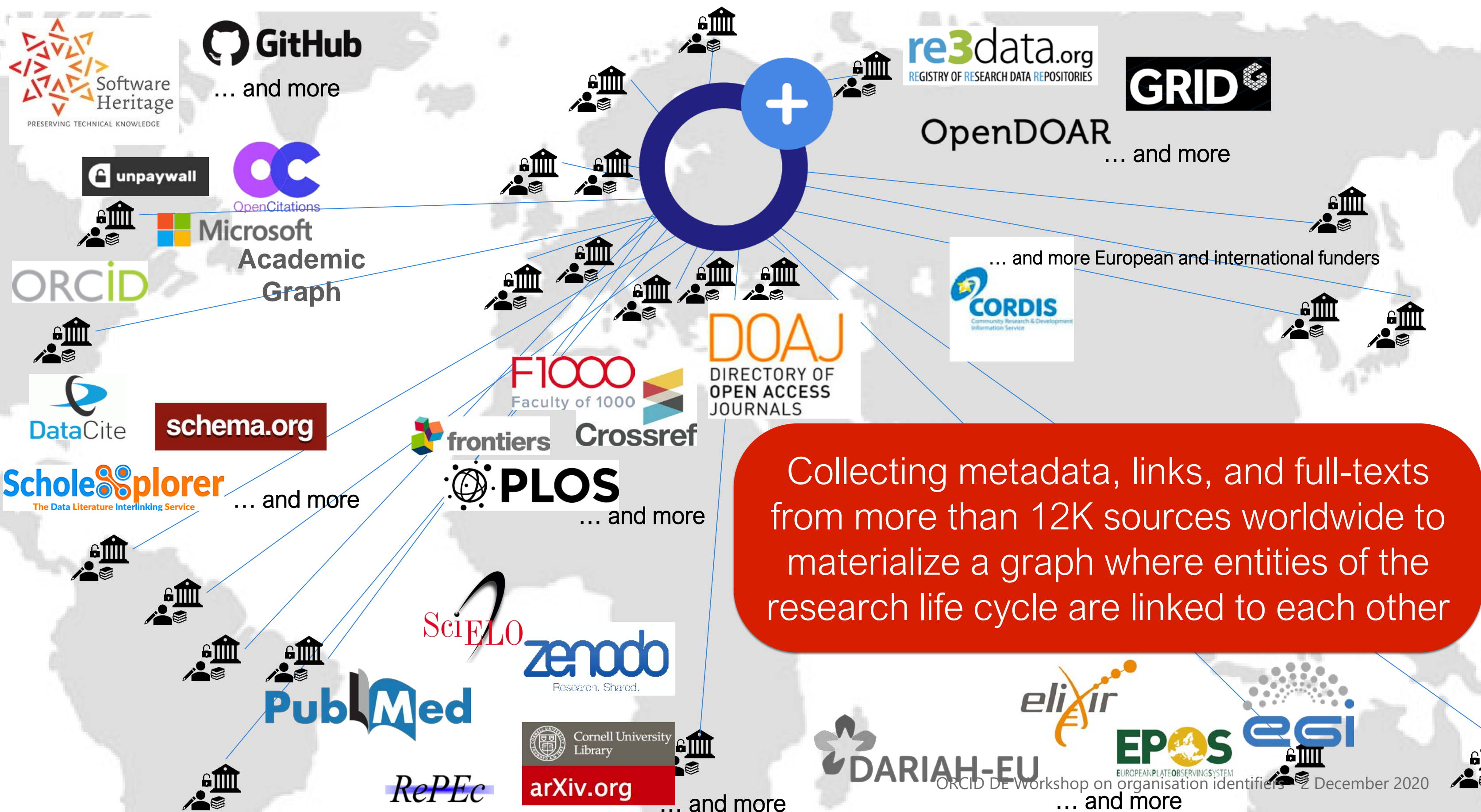
<https://connect.openaire.eu>



Providing an **open metadata** research graph of interlinked **scientific products**, with **Open Access** information, linked to **funding information** and **research communities**







Collecting metadata, links, and full-texts from more than 12K sources worldwide to materialize a graph where entities of the research life cycle are linked to each other

**Software Heritage**  
PRESERVING TECHNICAL KNOWLEDGE

**GitHub**  
... and more

**re3data.org**  
REGISTRY OF RESEARCH DATA REPOSITORIES

**GRID**

**OpenDOAR**  
... and more

**unpaywall**

**OpenCitations**

**Microsoft Academic Graph**

**ORCID**

... and more European and international funders  
**CORDIS**  
Community Research & Development Information Service

**DOAJ**  
DIRECTORY OF OPEN ACCESS JOURNALS

**F1000**  
Faculty of 1000

**DataCite**

**schema.org**

**frontiers** **Crossref**

**ScholarXplorer**  
The Data Literature Interlinking Service  
... and more

**PLOS**  
... and more

**SciELO** **zenodo**  
Research. Shared.

**PubMed**

**RePEc**

**arXiv.org**  
... and more

**DARIAH-EU**

**elixir**

**EPOS**  
EUROPEAN PLATE OBSERVING SYSTEM

**esi**

# PIDs in OpenAIRE

## Authors

ORCID



## Organizations

GRID

ROR

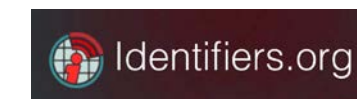


and other funders

## Products



arXiv.org



## Funders/Projects



FCT  
Fundação para a Ciência e a Tecnologia



and other funders

## Data sources

re3data.org  
REGISTRY OF RESEARCH DATA REPOSITORIES

OpenDOAR

+ journal lists from publishers

DOAJ  
DIRECTORY OF OPEN ACCESS JOURNALS



## Services (fortcoming...)



EUROPEAN OPEN SCIENCE CLOUD

...Clusters and RIs service catalogues





# OpenAIRE Research Graph: the supply chain

Data sources **12K+**  
 Products **400Mi+**  
 Harvested/mined links **490Mi+**  
 Full-texts **14.Mi+**

Data sources

Aggregation

Collection and harmonisation

Integration of additional scholarly communication sources

Metadata, relationships Relationships Metadata

Deduplication

Different records representing the same entity (results or organisation) are merged in one

Metadata records corresponding to equivalent objects are merged. Pre-print, post-print, published versions are considered equivalent for stats & monitoring purposes

Inference of new properties and links via full-text mining

Inference of new properties and links from meta-data

Enrichment Post-cleaning

Final cleaning step to harmonize values according to controlled vocabularies

Data sources **12K+**  
 Publications **102Mi**  
 Datasets **12Mi+**  
 Software **200K**  
 Others **7Mi**

CONNECT EXPLORE DEVELOP

indexing stats analysis

MONITOR

# APIs integration with third-party services

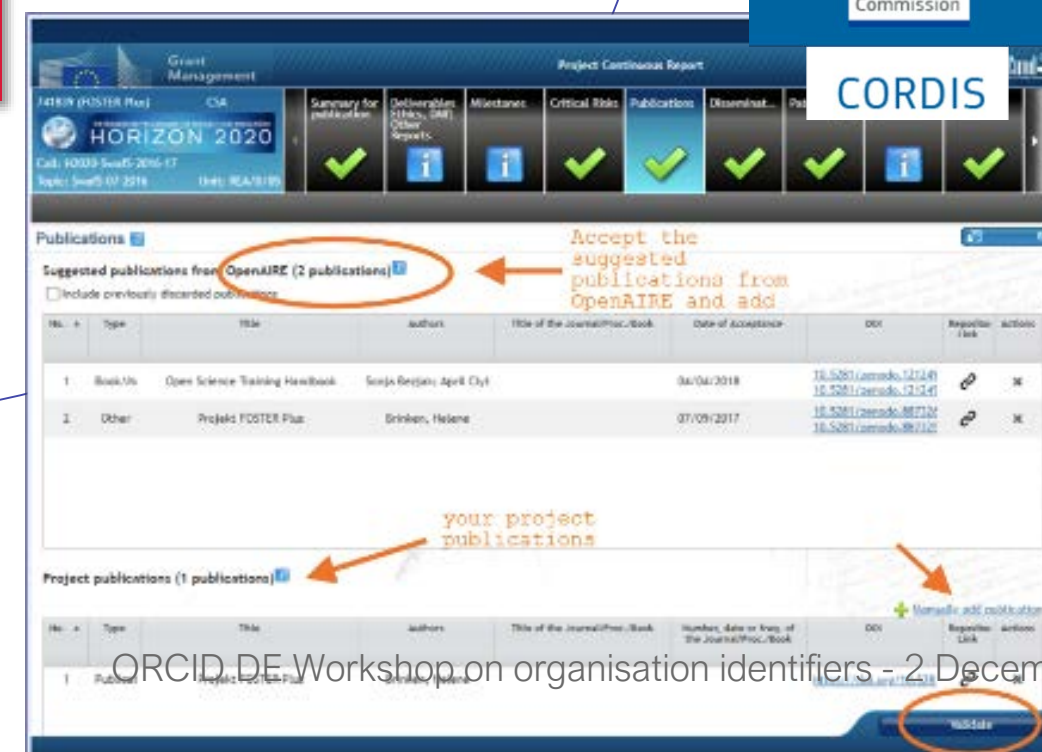
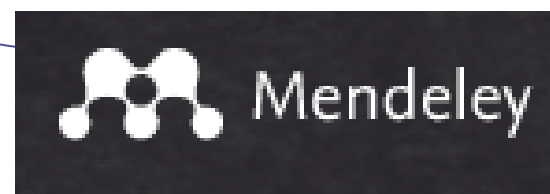
<https://develop.openaire.eu>



The Graph as EOSC Resource Catalogue



Research infrastructures



# OpenAIRE Research Graph: Dumps

- OpenAIRE Research Graph
- OpenAIRE COVID-19
- OpenAIRE communities
- DOIBoost
- Scholexplorer

zenodo Search Upload Communities admin@openaire.eu

## OpenAIRE Research Graph

Recent uploads

Search OpenAIRE Research Graph

November 11, 2020 (4.0) Dataset Open Access View

**OpenAIRE ScholeXplorer Service: Scholix JSON Dump**  
La Bruzzo, Sandro; Manghi, Paolo;  
This dataset contains the GZ-compressed dump of the Scholix links (schema Version 3) exposed by the OpenAIRE ScholeXplorer service. The dataset consists of 445+Mi bi-directional links (i.e. 890+Mi directed links) between literature-dataset and dataset-dataset involving 17+ Mi literature objects  
Uploaded on November 12, 2020  
3 more version(s) exist for this record

November 9, 2020 (v1) Poster Open Access View

**OpenAIRE Usage Counts. The analytics service of OpenAIRE Research Graph**  
Dimitris Pierrakos; Andreas Czerniak; Jochen Schirrwagen;  
Usage metrics for all types of scholarly output are one of the measures to assess Open Access impact and are a value added service of Open Access repositories. OpenAIRE has a successful record in providing usage metrics for for a large number of repositories from around the world. OpenAIRE  
Uploaded on November 11, 2020

November 3, 2020 (1.0.0) Other Open Access

OpenAIRE Research Graph: Json schem

Community

New upload

OpenAIRE RESEARCH GRAPH

OpenAIRE Research Graph  
This Zenodo community serves as container of the dumps of the OpenAIRE Research Graph.  
The OpenAIRE Research Graph is one of the largest open scholarly record collections worldwide, key in fostering Open Science and establishing its practices in the daily research activities. Conceived as a public and transparent good, populated out of data sources trusted by scientists, the Graph aims

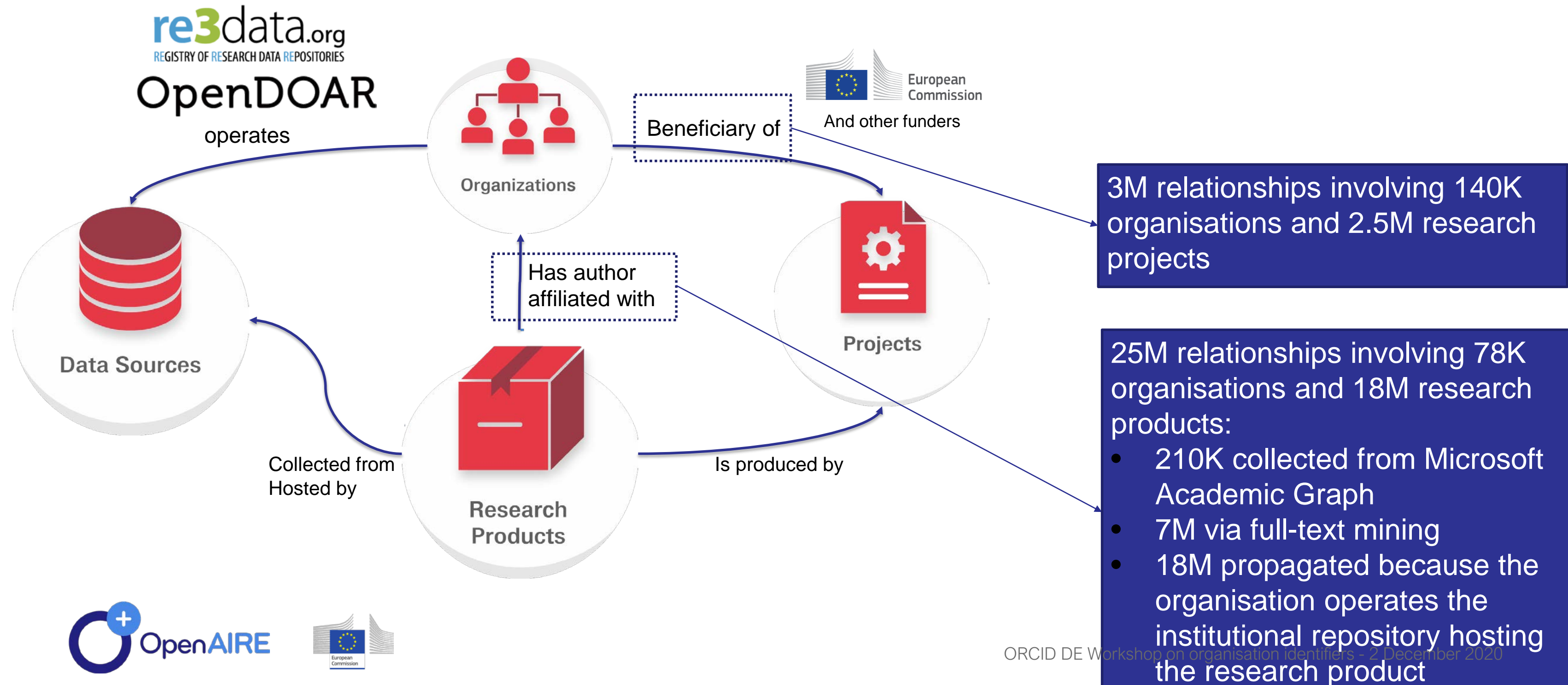
**November 2020**  
**11000+ views 3000+ downloads**

<https://zenodo.org/communities/openaire-research-graph>

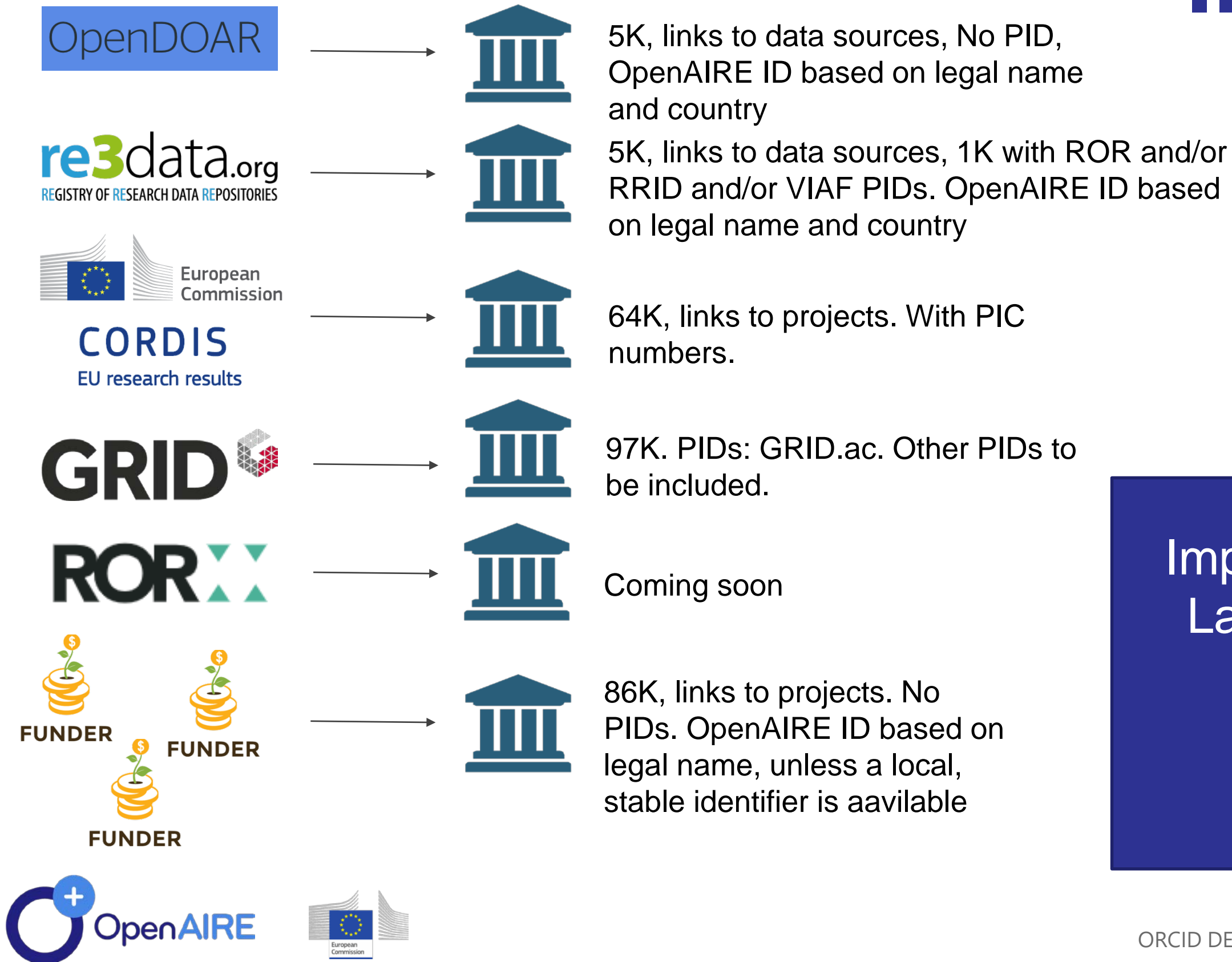


# Organisations in OpenAIRE

# Organizations in OpenAIRE



# Import phase



Imported 330K organisations  
Lack of common identifiers  
PIDs not referenced  
=  
Duplication

# How to identify the duplicates?

## A matter of data quality and completeness

OpenDOAR



re3data.org  
REGISTRY OF RESEARCH DATA REPOSITORIES



European Commission



CORDIS  
EU research results

GRID



FUNDER FUNDER FUNDER



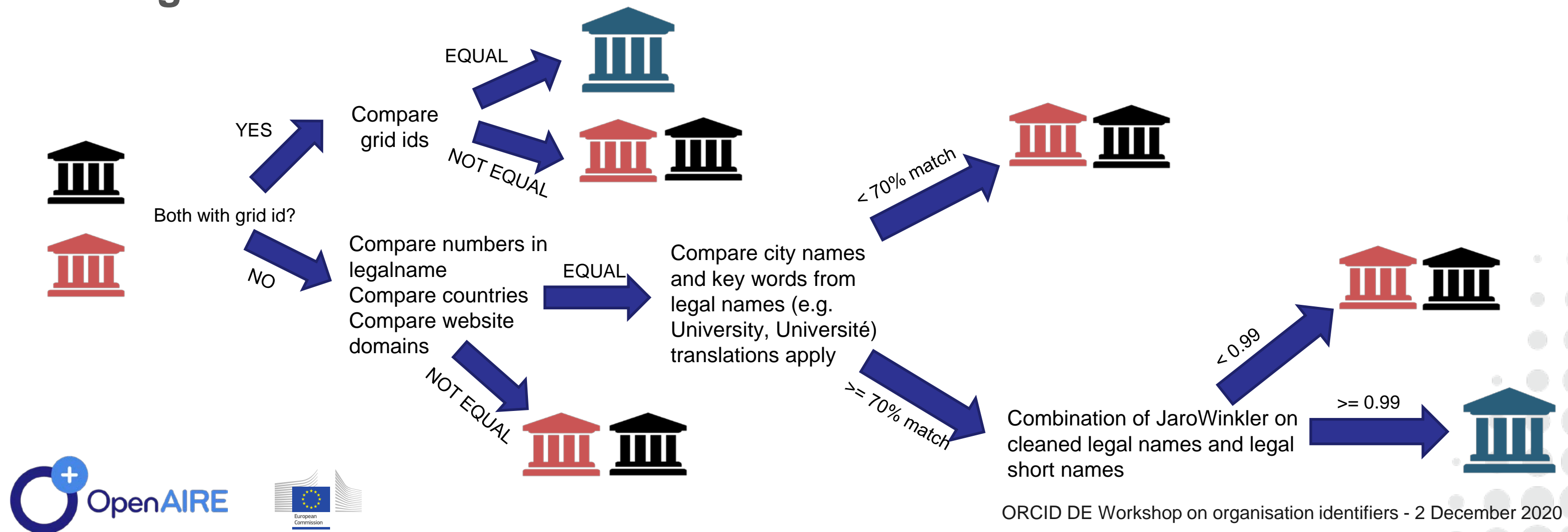
OpenAIRE



organization field completeness		
field	missing	% missing
legalname	254	0.1%
legalshortname	156,923	43.2%
alternativenames	319,825	88.0%
country	82,002	22.6%
websiteurl	142,047	39.1%
pid	184,536	50.8%

# The deduplication algorithm in a nutshell

- Cluster organisations by legalname and by website URL (4 clustering functions)
- For each cluster we perform pairwise comparisons with early exit strategies.





# Drawbacks



OpenDOAR



re3data.org  
REGISTRY OF RESEARCH DATA REPOSITORIES



CORDIS  
EU research results

GRID



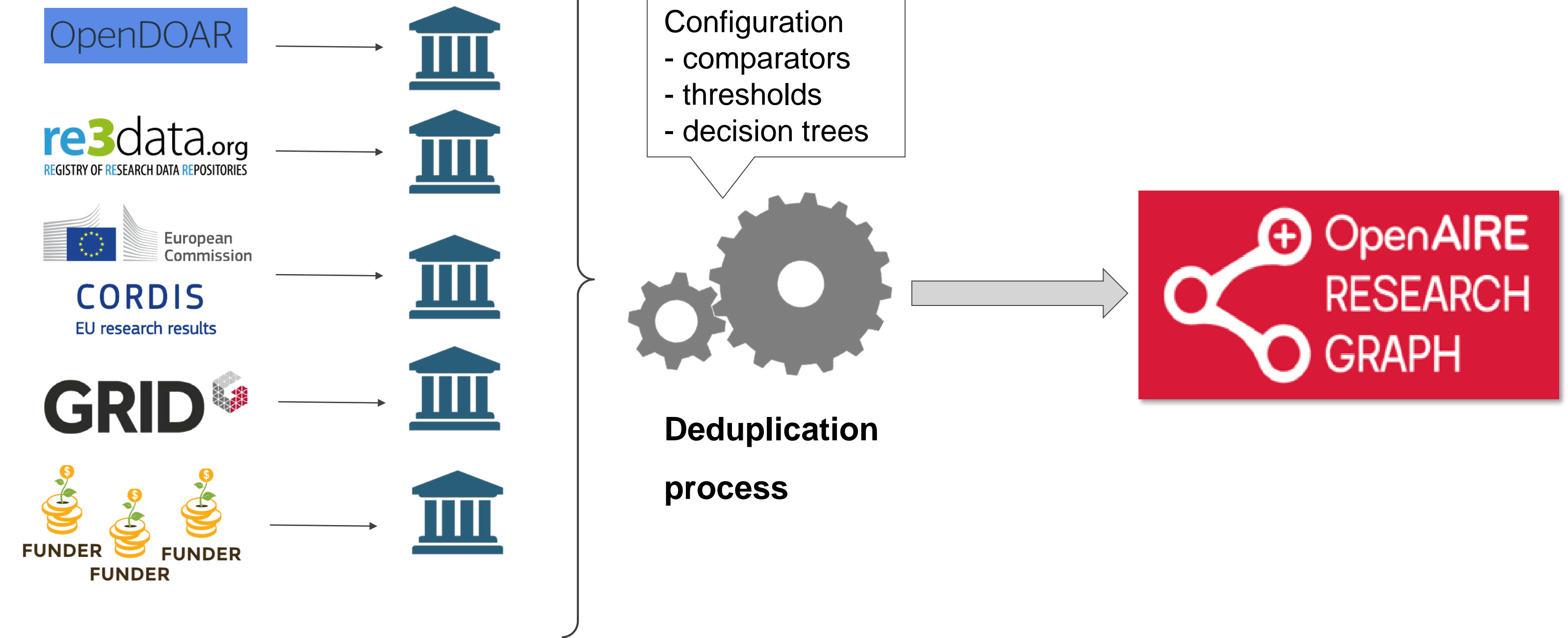
## False positives

- Mess up statistics (e.g. for institutional and country monitors)
- Mess up search functionality (e.g. finding only the US Concordia University and not the Canadian one)

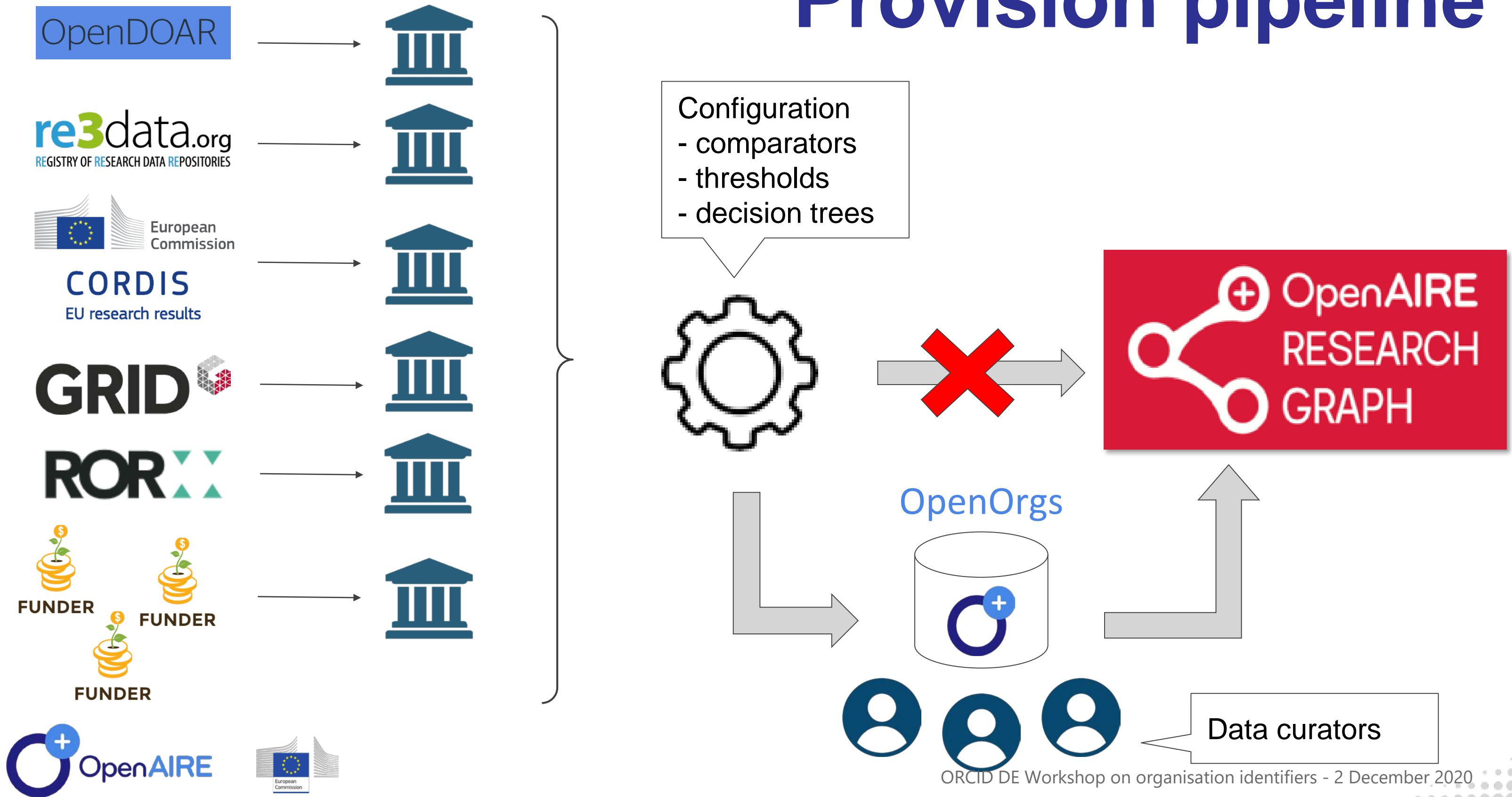
## False negatives

- One instance is linked to the papers, one to the project...but it is the same organization!

# Provision pipeline



# Provision pipeline



# Next steps

- **Operational OpenOrgs for OpenAIRE National Open Access Desks**
- **Import organisation mentions from repository records (CRIS systems, repository compliant to OpenAIRE guidelines 4)**
- **Notify missing PIDs to OpenAIRE sources (via the OpenAIRE Broker)**
- **Provide the sub-graph of entities with PIDs and bridge among different PID systems**

# Thank you!

**Alessia Bardi (CNR-ISTI)**

[Alessia.bardi@isti.cnr.it](mailto:Alessia.bardi@isti.cnr.it)

[graph.openaire.eu](http://graph.openaire.eu)

