

(SII)  
Program  
Office.

## Gilles Landais (CDS/Vizier) DOI status in IVOA

2nd session after this morning

DOI metadata is managed by archives **outside of the IVOA**

Provides metadata that can be used to transform or create citation mechanisms.

Used to build network of resources of many resource types

Comparison of Registry and DOI metadata

- Lots of similarities. Resource links, etc.

DOI in Registry via VOResource 1.2 ("alt identifier")

Data Origins in VO result.

The state of the DOI Note by Schaaff et al. ... in 2018. Lots have changed.

- Interconnected network of resources.
- Preserved DOIs.

How to match editorial requirements of citations

Duplicates with different origins of metadata.

The image shows three presentation slides from IVOA meetings. The first slide, titled "DOI status in IVOA", discusses working groups (ICP, A. Schaaff et al., coord IVOA 2018), DOI structure, and examples of DOI implementation. The second slide, titled "DOI in the Open Science", lists challenges such as improving data & service citation, making the bridge between resources available in the registry and Open Sciences, and building an interconnected network of "visible" resources. The third slide, titled "Make the bridge with Open sciences", discusses the proliferation of catalogues, platforms, search engines, and the risk of non-integration for resources without DOIs.

## **Gus Muench (AAS): The utility of dataset DOIs in manuscript review and scientific publications**

VIA Mark Parsons:

An RDA recommendation on citing dynamic data: <http://dx.doi.org/10.15497/RDA00016>

Credit is complicated: <https://eos.org/opinions/credit-where-credit-is-due>

## **S.Peroni (Bologna university) DOI (and Beyond) for Publications and Other Citable Research Outcomes**

- Keeps track of Citation links. Keep in CC0 database.
- Data model: OpenCitations Model (based on SPAR)
- Two RDF/triple-store databases: Index and Meta
- Changed the perspective. The main entity is the citation itself.
  - Two metadata fields: the two objects.
  - By prioritizing the citation event then it doesn't matter what the citer/citee includes.
  - They are concerned with duplicates where an object has multiple names and have been cited twice in an article. Huh?
- Sources for the merged "COCI"
  - PubMed open citations
  - DOCI from Datacite
- They are concerned with
  - Duplicates. Citations to the "same" thing but look like distinct cites bc they use different identifiers. Huh?
  - Many things do not have DOIs or PIDs.
  - How to go beyond the DOI — assign their own PIDs via OpenCitations Meta.
    - Is there a threshold for when
- Provenance.
  - When ingested. Agents doing the ingest.
  - Enable trust from all those pieces.
- Journal articles are versioned! There are snapshots!
- An evolution of data (in time) is crucial to be able to reconstruct those steps.

Q/A

1. How do we take care of the deltas? (Data deltas and metadata deltas)

1. Just store all the snapshots.
2. Their citation entities are diff'd/snapshotted with their PID system.

## A.Accomazzi (ADS): DOI-Enabled Discovery and Credit: an ADS Perspective

3 attributes of DOIs

- Persistent
- Notion of registered metadata for the identifier
- Widely adopted and can be adopted easily.

**Context**

ADS is primarily a literature database, and as such it does not aim to be an index for all research data products.

However, one of ADS's goals is to make relevant data products discoverable from the literature, whenever feasible.

Some types of data which are of most interest to ADS:

- Datasets "close" to publications, either as DDF, supporting archival links, or citations, as they supplement the science presented therein; examples include 'Voilà!' catalogs, text-mined Zenodo links, archival data links, data citations
- Reference catalogs, collections, and services, which are highly used and (possibly) cited; examples include 2MASS, WISE, CGCG, etc.
- (to a somewhat lesser extent) Observations linked to proposals; examples include obs. proposals from JWST, CXC, XMM, etc.

Indexed versus Linked resources

Some typically "linked" resources are becoming "indexed" resources

- Software
- Some data products
- Some notebook products

**What's the difference**

**Indexed Dataset**

- ADS has a record corresponding to the dataset
- Dataset has higher level of discoverability (retrieved by e.g. ADS author search)
- Dataset has ADS metrics associated with it
- Data is accessible from paper via citation and data link

**Linked Dataset**

- ADS does not have a record corresponding to the dataset
- Papers associated with dataset typically part of a linked data collection (e.g. Chandra, IRSA, MAST)
- Only metrics available are via associated paper metrics
- Data is accessible from paper via data link

Then there are a collections of linked datasets, e.g., ycats, e.g, Chandra Obs

Ingestion policy is evolving.

**Ingestion Policy (still evolving)**

**What is/will be indexed in ADS**

- Curated, high-level datasets with good metadata (registered DOI) and clear authorship information
- Research data collections that have shown reuse value (initially via citations, i.e. `o:stats_cswot > 3`)

**What will not be indexed in (but possibly linked from) ADS**

- Every data product out there
- Every single version of a data product
- Data collections created for bundling purpose (e.g. MAST user generated DOIs)

**Food for thought**

- How does the data indexing & linking policy outlined here fit the needs of our community?
- How does it help you, as an data archivist / publisher / scientist?
- Is there a need for a disciplinary index of data products cited / mentioned in the literature beyond what is described here?

Q/A:

- Will ADS use related identifiers to append in other linking in their system?
  - They have an existing set of relations
  - And they have a mapping from DataCite to their system.
  - Maintain Citation graph and citations that "come" from data cite are not included right now.
- 

## **B.Cecconi (Obs Paris): Data Management and DOI implementation and lessons' learnt**

Point of view of data provider.

MASER service (see Radio WG later)

A rather varied typology of data types: data collections, meta collections, docs, catalogs, etc..

Do a hard test on Data Managements plans: both for release and for dev and for hardware costs.

Showing a two parallel structure for the organization of the DMP

- ObsParis can mint DOI with Datacite.
- One DOI per collection/dataset/document.
  - Landing Page content: title, citation, abstract, link to data, description, acknowledgments, references
  - Web-semantic annotation (schema.org)
- Current status: **manual process for**
  - creation / maintenance of DOIs (on Datacite portal)
  - creation / maintenance of Landing Page (SPDF)
  - creation / maintenance of annotations (JSON-LD)
 Only two persons authorised.

Fully manual process (2 ppl; "a lot of work")

Investigating [Recherche.data.gov](https://recherche.data.gov) for a French National Repository based on Dataverse.

## **M.Parsons (Nasa): the new NASDA DOI Registration guidelines and general guidelines process.**

Developing guidelines for creating citable PIDs for NASA data.

"Which identifier" and how.

They have some big team and survey across all 5 divisions.

Mapping a set of "local" ids to something "citable"

Full RFC & Short Version

Things important to them:

- How DataCite works and the archive members vs STI members
- Three scenarios: planned, provider request, user request.
- The responsible repository is responsible for DOI requests
- Repositories must meet metadata requirements (not STI? )
- Repos responsible for DOI maintenance.
- DOI Request: offload to STI (rare); use Consortium Membership (coherent focal point); Direct members keep being direct members.
  - But if the repo is responsible and registration is offloaded then how is the repo responsible?
  - Just a fiscal layer — who pays etc.
- Agency-wide of Registry DOIs
- Lessons Learned
  - Saying repositories are in the "middle" between data users and HQ
  - Friction is part of the game. A wheel turns bc of friction.