

Where's the data? A story of data discovery, cleaning, and equality

IASSIST 2023

Alicia Hofelich Mohr, University of Minnesota Joel Herndon, Duke University Cynthia Hudson Vitale, Association of Research Libraries

Realities of Academic Data Sharing (RADS) Study: Metadata Analysis

RADS Institutions:





- Cornell University
- Duke University
- University of Michigan
- University of Minnesota
- Virginia Tech
- Washington University in St. Louis











DATA CURATION NETWORK



RADS has been funded by NSF EAGER grant #2135874: Completing the Lifecycle: Developing Evidence Based Models of Research Data Sharing

RADS Research Questions



Where are funded researchers sharing their data and what is the quality of that metadata?



How are researchers making decisions about why and how to share research data?



What is the cost to the institution to implement federally mandated public access to research data policies?





ACT 1: Where do we look?

What we set out to do



Identify the location of published data between 2012–2022



Search DOI registries using APIs



Parse affiliation across our six organizations using RORs



Facet results by subject/keywords and funder ID

Done, easy, right?

The best laid plans...



Identify the location of published data between 2012-2022



APIs varied in accessibility & speed → needed to combine approaches



No widespread use of RORs → instead used text search for institution name



Limited & inconsistent use of funder fields and keywords → removed plans for this facet

The Search





Search institutions in creators.affiliation.name (n= 55,634)

publicationYear >= 2012
 (n= 51,053)

resourceTypeGeneral =
"Dataset" or "Software"
(n=31,946)

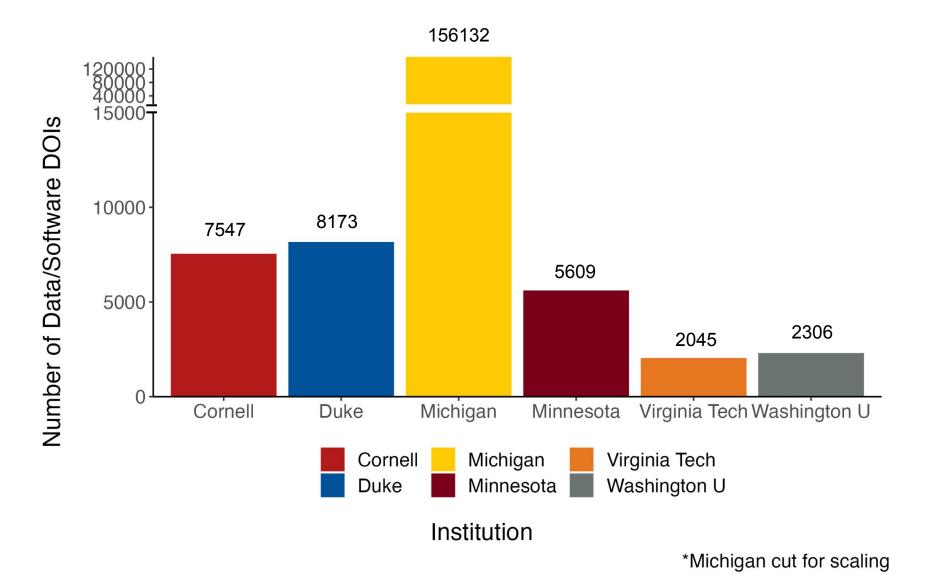
Author affiliation contained institution; type = "Dataset"; date-parts >= 2012 (n=152,376)

Remove non-relevant institutions (n=147,702)

n = 179,648 DOIs

ACT 2: What's here?

Finding all the data



Well, almost...

We knew researchers were publishing data in our institutional repositories...

	Michigan	Minnesota	Cornell	Virginia Tech	WashU	Duke
Institutional Data Repository Records	645	692	174	333	95	225

But we didn't find them

Institutional Repository 1 90 DOIs	34 111	16 0
------------------------------------	--------	------

Affiliation Woes

- None of us consistently entered "affiliation" metadata in DataCite/CrossRef.
- Solution: Search for our repositories



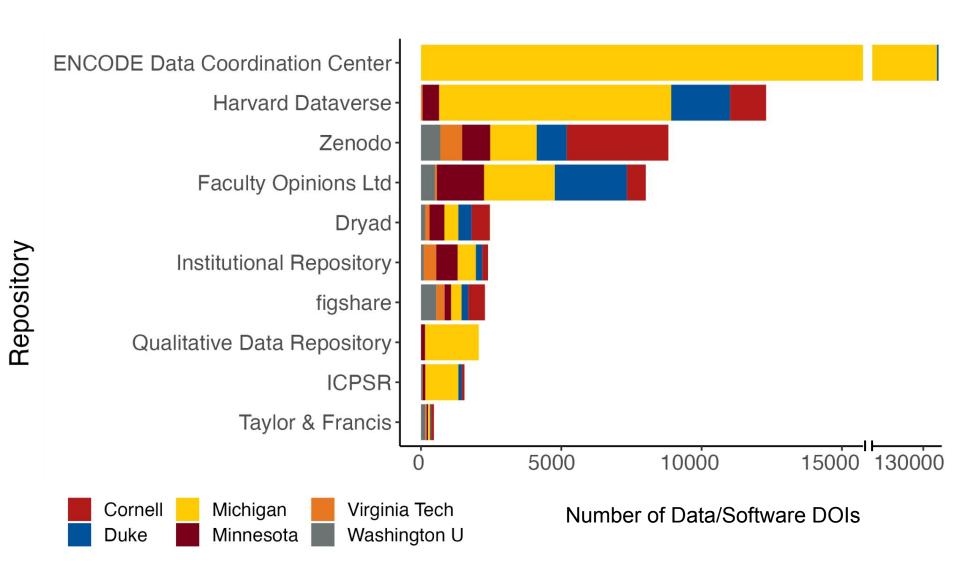
Search five institutional repositories in publisher; publicationYear >= 2012; generalResourceType = dataset or software (n = 1,939)



Search Duke's member prefix for year >=2012; type = data (n = 225)

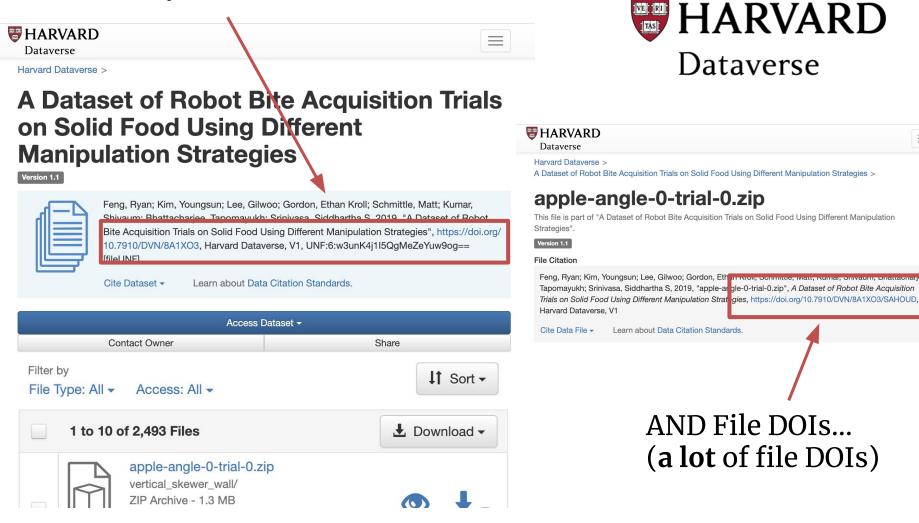
Additional n = 2,164 DOIs

Now, where's the data (Top 10)?

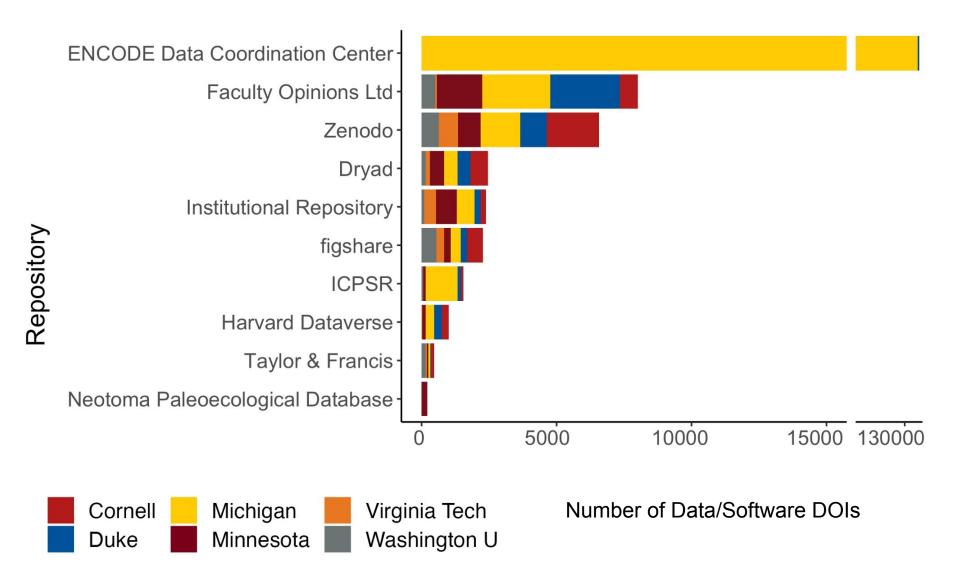


A closer look...

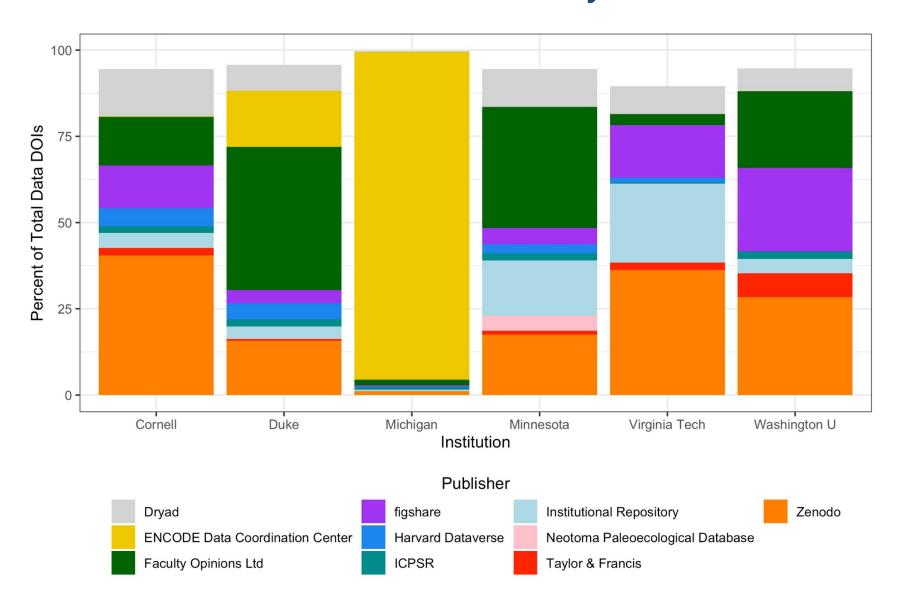
Study DOI



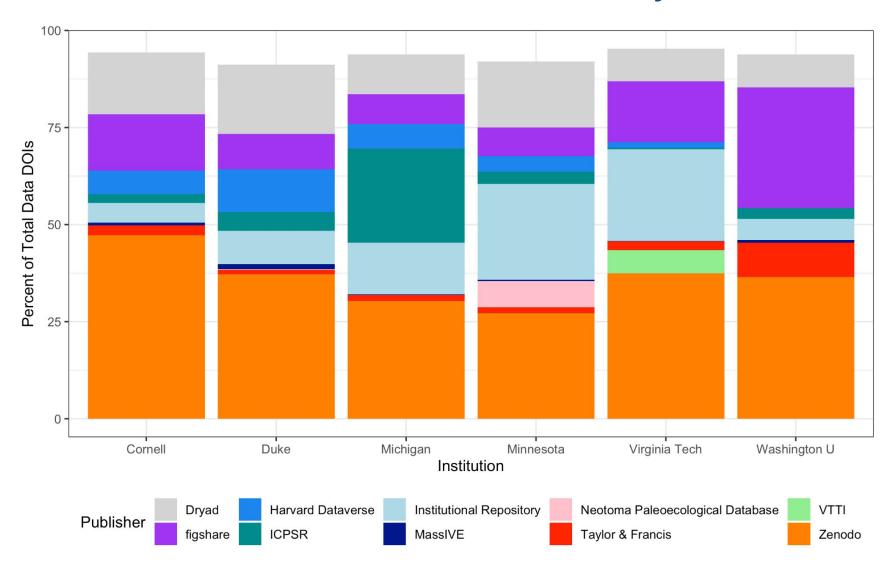
Ok... now where's the data (Top 10)?



Distribution of Publisher by Institution



Without ENCODE and Faculty LTD



ACT 3: What's missing?

Assumptions

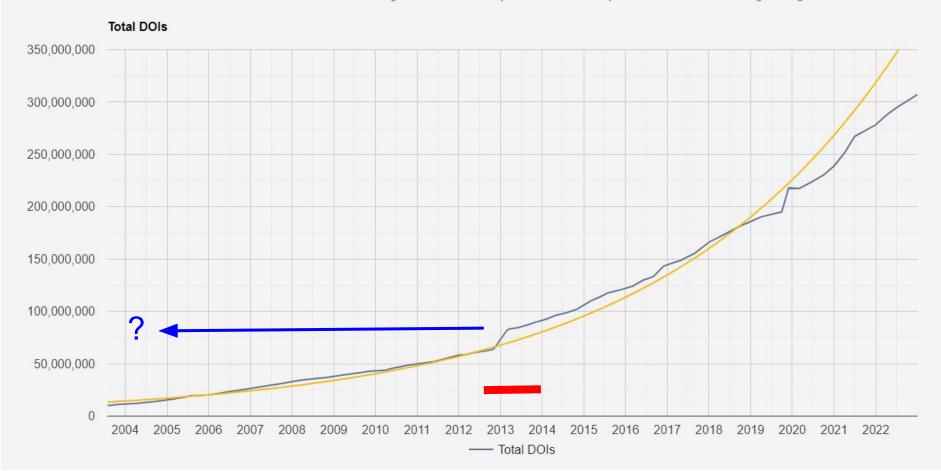
DOI as THE data/software identifier - but is this true?



DOI Growth Over Time

TOTAL DOIS

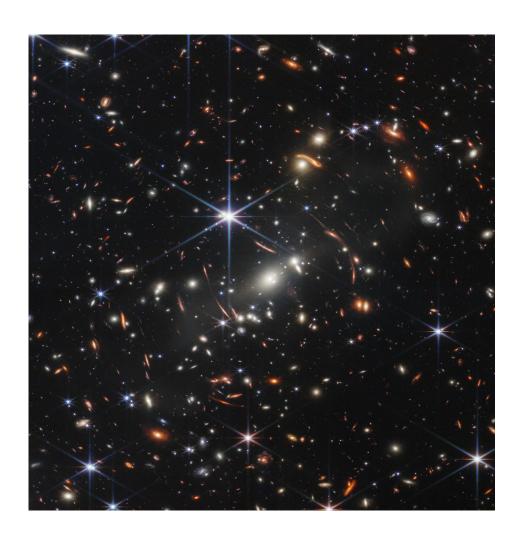
The cumulative number of DOI registered in the DOI System over time; the yellow line shows the moving average



https://www.doi.org/the-identifier/what-is-a-doi/

but... the data sharing space is vast...

- Handles, ARKs, other local identifiers for data
- Accession numbers (medical fields)
- No registered identifiers at all (including some linked data)



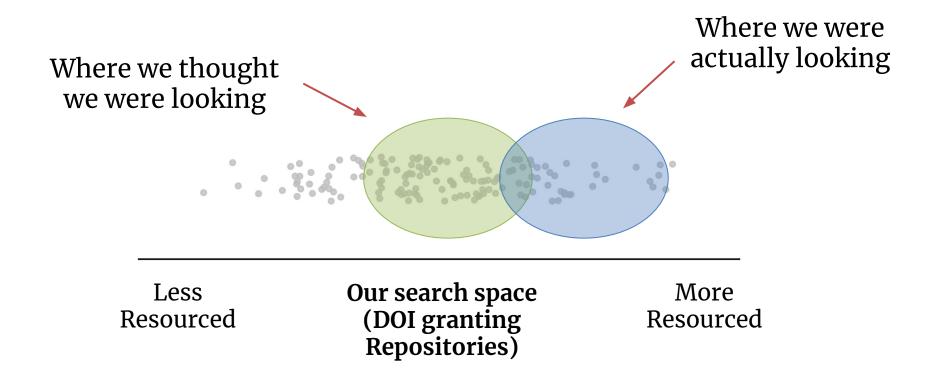
Source: NASA, bit.ly/43cf5mc

... and (even) DOI's present assessment challenges:

- Multiple registration agencies
- Metadata completeness (required vs "suggested")
- Granularity ("Itemness")/Versioning
- Software (sometimes coded as "data")

But there was bias within our search of DOIs

We are likely "losing" information from smaller, less-resourced repositories.



But there was bias within our search of DOIs

Taking a step back, we really only see a VERY small portion of the data sharing picture

 Metadata findability is highly biased towards well resourced repositories What we found





Data outside repositories

Data in repositories



The struggle is real!

Metadata decisions are ... complicated ...

- Capacity and resource availability
- Balancing description and discovery
- Multiple (evolving) standards
- Correction is time consuming, costly

Realities of Academic Data Sharing (RADS) Study: Metadata Analysis

RADS Institutions:





- Cornell University
- Duke University
- University of Michigan
- University of Minnesota
- Virginia Tech
- Washington University in St. Louis















Thank You!

Contact Us

Find our work on github: ajhmohr/rads metadata





