

DATA CURATION NETWORK

Data Curation Fundamentals Workshop



This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

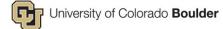
Ethical. Reusable. Better.

DATA CURATION NETWORK

datacurationnetwork.org

DATA CURATION NETWORK

UNIVERSITY OF MINNESOTA
Driven to Discover®



Ethical. Reusable. Better.

DATA CURATION NETWORK

Mission

Trusted, community-led
network of curators
advancing open research
by making data

Ethical. Reusable. Better.

datacurationnetwork.org

Code of Conduct

The Data Curation Workshop is dedicated to providing a harassment-free experience for everyone, regardless of gender, gender identity and expression, sexual orientation, disability, physical appearance, body size, race, age or religion. We do not tolerate harassment of participants in any form. Sexual language and imagery is not appropriate for any conference venue, including talks. Participants violating these rules may be sanctioned or expelled from the workshop at the discretion of the conference organizers.

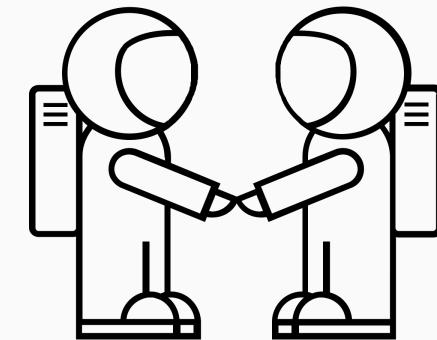
Our code of conduct and anti-harassment policy can be found at:

<https://datacurationnetwork.org/about/code-of-conduct/>

Code of Conduct (cont)

Examples of behavior that contribute to creating a positive environment include:

- Using welcoming and inclusive language
- Being respectful of differing viewpoints and experiences
- Gracefully accepting constructive criticism
- Focusing on what is best for the community
- Showing empathy towards other community members



Learning Outcomes

1. Increase understanding of data curation practices and tools in various disciplines, data types, and formats.
2. Share expertise and enhance curation capacity nationwide.
3. Meet like-minded colleagues who are interested in building and extending curation practices at their institutions.

The Value of Curation

Introduction to the workshop

Reasons for data sharing



Open data movement



Reproducibility



Funding agency mandates



Publisher policies

Research data have value beyond their original purpose.

But....

Data can be messy, incomprehensible, and lack context

Normal View Ready Sum = 0.296

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	0.296	-0.55	-0.243	-3.04	-2.885	0.071	1.847	-0.521	0.107	-3.528	-3.089	2.45	-1.46	2.404	-0.857	0.354	-0.923	1.487	-2.06
2	0.354	-0.487	-0.193	-2.974	-2.933	-0.029	1.951	-0.506	0.061	-3.643	-3.162	2.131	-1.438	2.482	-0.796	0.435	-1.075	1.448	-2.03
3	0.438	-0.449	-0.099	-2.979	-2.901	-0.09	2.051	-0.501	-0.005	-3.647	-3.34	1.713	-1.428	2.519	-0.681	0.548	-1.258	1.372	-2.10
4	0.519	-0.431	-0.018	-3.042	-2.831	-0.13	2.107	-0.452	-0.027	-3.616	-3.42	1.322	-1.434	2.535	-0.564	0.624	-1.438	1.351	-2.22
5	0.696	-0.37	0.023	-3.065	-2.832	-0.187	2.202	-0.415	-0.072	-3.648	-3.504	0.831	-1.403	2.529	-0.467	0.718	-1.651	1.315	-2.30
6	0.939	-0.332	0.083	-3.089	-2.815	-0.233	2.314	-0.342	-0.119	-3.603	-3.648	0.38	-1.333	2.551	-0.367	0.773	-1.856	1.243	-2.39
7	1.188	-0.295	0.171	-3.08	-2.753	-0.295	2.431	-0.197	-0.186	-3.598	-3.619	-0.07	-1.272	2.594	-0.263	0.867	-1.978	1.137	-2.48
8	1.503	-0.284	0.279	-3.129	-2.746	-0.363	2.52	-0.116	-0.298	-3.487	-3.5	-0.608	-1.175	2.591	-0.121	0.904	-2.033	1.019	-2.60
9	1.826	-0.288	0.36	-3.183	-2.743	-0.496	2.59	-0.012	-0.316	-3.318	-3.456	-0.989	-1.077	2.574	0.034	0.946	-2.054	0.915	-2.68
10	2.153	-0.289	0.369	-3.162	-2.632	-0.615	2.653	0.197	-0.345	-3.249	-3.388	-1.33	-0.985	2.529	0.16	0.989	-2.125	0.805	-2.75
11	2.59	-0.244	0.359	-3.205	-2.51	-0.761	2.761	0.412	-0.416	-3.204	-3.296	-1.58	-0.896	2.494	0.318	0.99	-2.223	0.734	-2.81
12	2.97	-0.196	0.319	-3.218	-2.463	-0.944	2.933	0.643	-0.421	-3.143	-3.089	-1.746	-0.836	2.486	0.426	1.002	-2.274	0.687	-3.00
13	3.269	-0.222	0.297	-3.148	-2.454	-1.045	3.051	0.904	-0.356	-2.983	-2.829	-1.813	-0.74	2.485	0.58	1.003	-2.314	0.659	-3.05
14	3.512	-0.266	0.274	-3.157	-2.429	-1.147	3.119	1.116	-0.286	-2.783	-2.595	-1.927	-0.659	2.505	0.75	0.954	-2.351	0.609	-3.18
15	3.684	-0.271	0.289	-3.214	-2.396	-1.255	3.052	1.222	-0.227	-2.627	-2.292	-2.081	-0.595	2.456	0.915	0.918	-2.353	0.549	-3.19
16	3.824	-0.275	0.233	-3.289	-2.4	-1.262	2.996	1.39	-0.16	-2.475	-2.019	-2.286	-0.449	2.324	1.062	0.856	-2.367	0.54	-3.25
17	3.889	-0.294	0.186	-3.295	-2.303	-1.306	2.961	1.545	-0.083	-2.293	-1.825	-2.461	-0.368	2.243	1.196	0.764	-2.4	0.463	-3.26
18	3.896	-0.295	0.158	-3.289	-2.266	-1.383	2.93	1.645	-0.095	-2.195	-1.606	-2.573	-0.234	2.222	1.344	0.701	-2.392	0.432	-3.29
19	3.838	-0.283	0.152	-3.286	-2.273	-1.352	2.876	1.615	-0.074	-2.086	-1.385	-2.672	-0.063	2.182	1.556	0.622	-2.289	0.405	-3.24
20	3.712	-0.338	0.139	-3.328	-2.23	-1.302	2.778	1.637	-0.007	-1.971	-1.328	-2.759	0.078	2.144	1.789	0.566	-2.074	0.358	-3.25
21	3.526	-0.363	0.125	-3.387	-2.198	-1.275	2.604	1.624	-0.019	-1.814	-1.35	-2.817	0.217	2.032	1.996	0.562	-1.771	0.293	-3.33
22	3.343	-0.393	0.112	-3.466	-2.107	-1.235	2.412	1.645	-0.008	-1.744	-1.385	-2.886	0.373	1.929	2.219	0.573	-1.502	0.225	-3.38
23	3.17	-0.397	0.04	-3.434	-1.975	-1.199	2.27	1.612	0.025	-1.692	-1.416	-2.868	0.528	1.793	2.473	0.514	-1.197	0.207	-3.46
24	3.011	-0.384	-0.04	-3.374	-1.829	-1.116	2.122	1.604	0.014	-1.607	-1.452	-2.679	0.624	1.66	2.757	0.527	-0.949	0.173	-3.46
25	2.758	-0.369	-0.067	-3.294	-1.669	-1.005	1.963	1.586	-0.099	-1.53	-1.546	-2.495	0.732	1.585	3.029	0.597	-0.721	0.163	-3.44
26	2.467	-0.313	-0.129	-3.247	-1.425	-0.874	1.744	1.469	-0.191	-1.427	-1.618	-2.295	0.804	1.482	3.248	0.663	-0.519	0.17	-3.39
27	2.172	-0.267	-0.23	-3.282	-1.103	-0.738	1.466	1.408	-0.267	-1.316	-1.674	-2.177	0.827	1.373	3.46	0.756	-0.4	0.178	-3.32
28	1.897	-0.224	-0.304	-3.319	-0.821	-0.656	1.237	1.361	-0.354	-1.228	-1.716	-2.019	0.877	1.216	3.548	0.872	-0.293	0.187	-3.34
29	1.643	-0.174	-0.372	-3.24	-0.6	-0.544	1.066	1.295	-0.409	-1.167	-1.694	-1.944	0.995	1.05	3.599	0.934	-0.215	0.234	-3.21
30	1.391	-0.133	-0.456	-3.215	-0.378	-0.465	0.91	1.237	-0.462	-1.151	-1.73	-1.808	1.111	0.924	3.666	1.023	-0.194	0.243	-3.18
31	1.152	-0.131	-0.533	-3.195	-0.257	-0.381	0.792	1.159	-0.466	-1.102	-1.711	-1.622	1.162	0.819	3.66	1.102	-0.133	0.244	-3.08
32	0.893	-0.159	-0.614	-3.167	-0.204	-0.363	0.65	1.048	-0.518	-1.1	-1.734	-1.394	1.251	0.679	3.67	1.171	-0.044	0.278	-3.01
33	0.635	-0.144	-0.643	-3.13	-0.092	-0.352	0.517	0.933	-0.579	-1.149	-1.763	-1.301	1.317	0.551	3.664	1.223	0.043	0.266	-2.93
34	0.415	-0.151	-0.681	-3.169	-0.075	-0.292	0.429	0.861	-0.619	-1.156	-1.785	-1.176	1.403	0.405	3.572	1.291	0.119	0.318	-2.86

Digital file formats are constantly at risk



Most data never leave a laptop ⇒ benign neglect



RETRO INFORMATIQUE

Digital Equipment Corporation Rainbow 100
Date de l'édition : 2000
N° de ISBN : 2-85654-200-2 - ISSN : 0988-14-May-98
Mise en page : Gérard Pichot
Illustrations : Gérard Pichot
Photographies : Gérard Pichot
Conception graphique : Gérard Pichot
Révision : Gérard Pichot
Format : 24 x 32 cm
Papier : 100 g/m²
Imprimé par : Gérard Pichot
Imprimerie : Gérard Pichot
Relié : Gérard Pichot
Reliure : Gérard Pichot
Éditeur : Gérard Pichot
Distributeur : Gérard Pichot
Gérard Pichot

What is data curation?

“Data curation is the active and ongoing management of data through its lifecycle and interest and usefulness to scholarship, science, and education; curation activities enable data discovery and retrieval, maintain quality, add value, and provide for re-use over time.”

- University of Illinois’ Graduate School of Library and Information Science

What is data curation?

“Data curation is akin to work performed by an art or museum curator. Through the curation process, data are organized, described, cleaned, enhanced, and preserved for public use, much like the work done on paintings or rare books to make the works accessible to the public now and in the future.”

- ICPSR

Who assists in curation?

Data curation involves

- Automated curation actions taken by **repositories**
- Decisions and actions taken by the **curator**
- Engagement with the **researcher**

Definitions of Data Curation Activities used by the DCN

The value-add

- ★ Easier for fellow scholars and future collaborators to understand
- ★ More likely to be trusted
- ★ The research they represent are more likely to be reproducible
- ★ More likely to be properly cited
- ★ Represent potential cost-savings
- ★ Findable, accessible, interoperable, and reusable (FAIR)

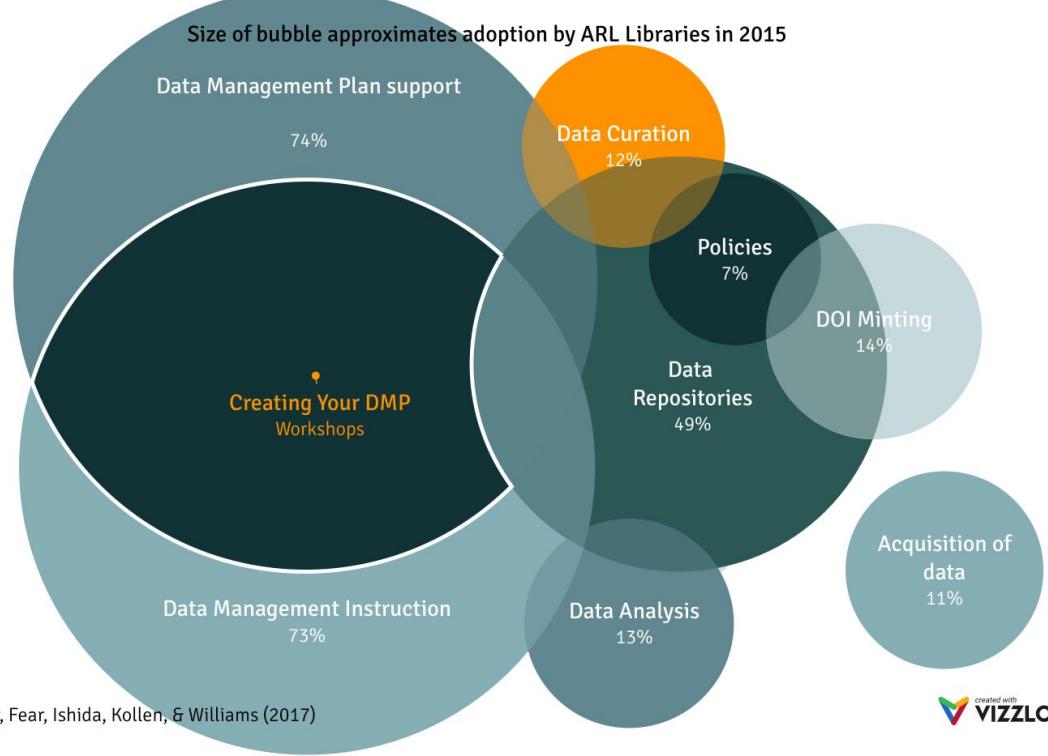
Remember... data curation is just **one piece** of the broader data services landscape

DATA CURATION NETWORK

Kouper, Fear, Ishida, Kollen, & Williams (2017)

Research Data Services

Size of bubble approximates adoption by ARL Libraries in 2015



Role of libraries in data curation

Libraries and academic-based data repositories are just one piece of the data repository landscape.

Baker, K., and Dueri, R. (2017). "Data and a diversity of repositories" in *Curating Research Data: A handbook of current practice* (L. R. Johnston, ed.). ACRL press.

TABLE 4.3
Examples of kinds of data repositories found in the United States.

Kind of Repository	Examples
Federally Funded Data Centers	NASA Distributed Active Archives (DAAC), NOAA National Centers for Environmental Information (NCEI), National Snow and Ice Data Center (NSIDC), USGS Earth Resources Observation Systems (EROS) Data Center (EDC)
Federally Funded Research and Development Centers (FFRDC)	National Center for Atmospheric Research (NCAR), Jet Propulsion Lab (JPL), Oak Ridge National Laboratory (ORNL)
National Libraries	National Library of Medicine (NLM), National Agricultural Library (NAL), Library of Congress (LOC)
State and Local Agencies	State geological surveys, County planning offices
Thematic Repository	Long Term Ecological Research Network Information System (LTER NIS), Andrews Forest LTER (AND), National Snow and Ice Data Center (NSIDC), Maria Rogers Oral History Program
Domain Repository	Global Biodiversity Information Facility (GBIF), Inter-university Consortium for Political and Social Research (ICPSR), DataOne, Interdisciplinary Earth Data Alliance (IEDA)
Institutional Repository	Purdue University Research Repository (PURR), Data Repository for the University of Minnesota (DRUM)
Replication Repository	Dryad Digital Repository, Pangaea Data Library
Software Repository	GitHub, SourceForge
Commercial Archives	DigitalGlobe, Aerial photography companies, Resource exploration companies, Figshare
Private Archives	Huntington Library, Getty Research Institute

Data Curation in Practice

Research performed by the Data Curation Network

DCN Research

- Identification of Data Curation Activities (2016)
- Understanding what data curation activities are important to researchers (2017)
- SPEC Kit for Data Curation Activities in ARL Institutions (2017)
- Value of Curation from the Repository and Depositor Perspective (2021-2022)

What is the Value of Curation?

Perspective of repositories *doing* curation



Perspective of depositors *receiving* curation



Ethical. Reusable. Better.

DATA CURATION NETWORK

JOHNS HOPKINS
UNIVERSITY

datacurationnetwork.org

Methodology

- US-based data repositories
- Non-probabilistic sampling
 - Recruitment: multiple listservs
 - Risk of bias and data skewness towards higher level of curation
 - Limited generalizability
- Open for 3 weeks (Jan. 2021)



Demographics



34 Directors
52 Staff
4 Depositors
5 Users

95 respondents



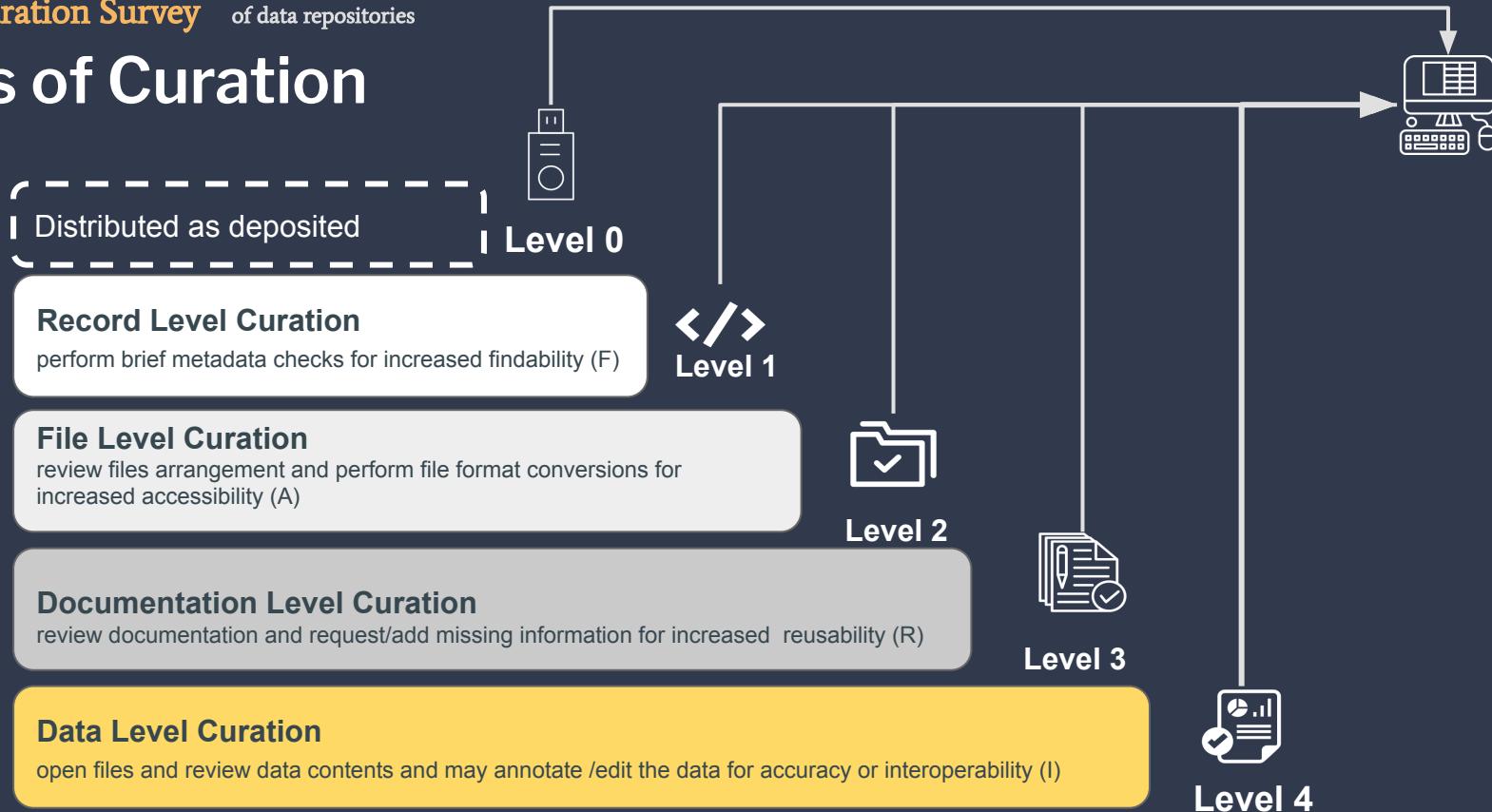
31 Disciplinary
25 Institutional
3 Generalist

59 repositories



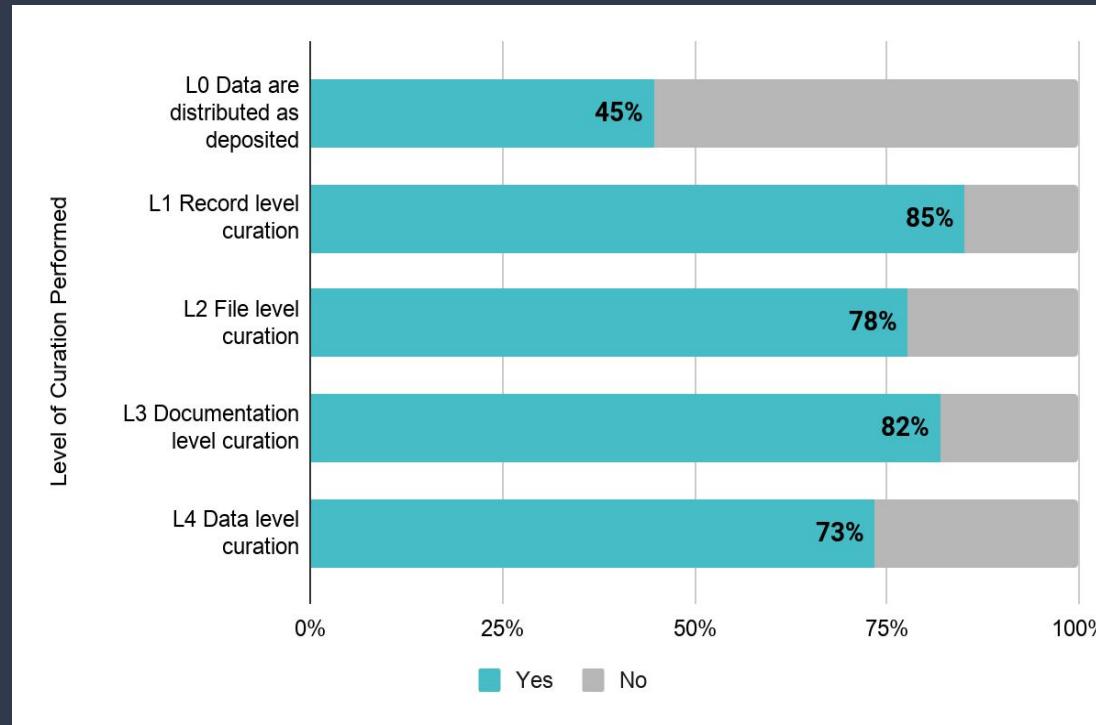
US Repositories
from 23 states
- 10 CoreTrustSeal
- 11 DCN

Levels of Curation



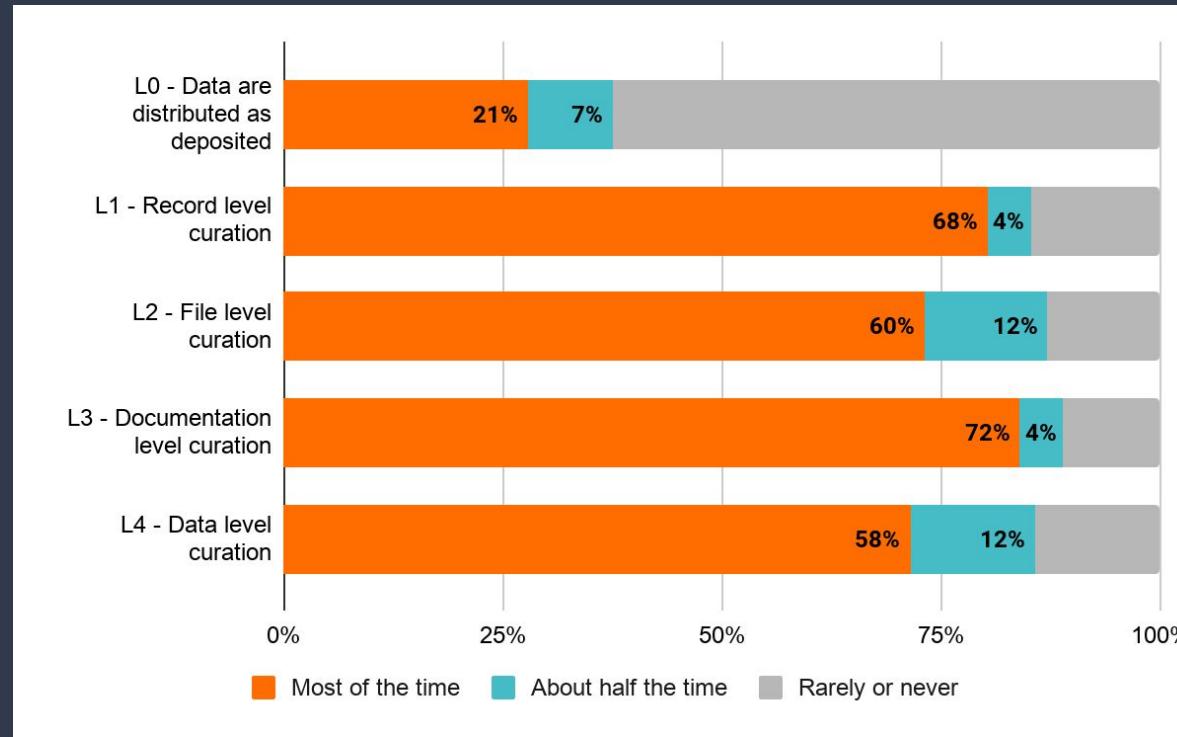
Results - Levels of Curation Performed

L1 and L3 are the most performed curation levels



*Multiple responses allowed

Results - Levels of Curation (Frequency)



*multiple responses allowed.

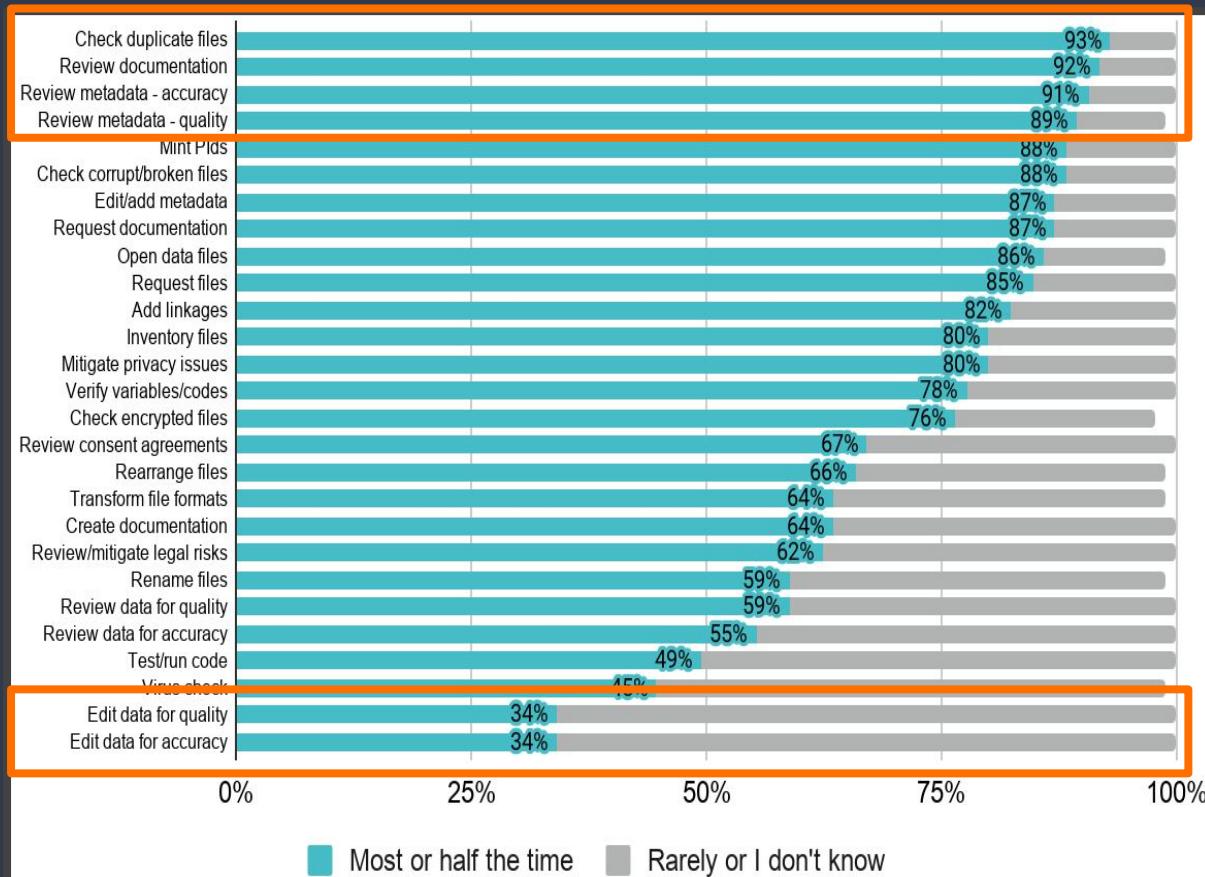
Results - Curation Actions Performed

~90%

- Check for duplicate files
- Review documentation and metadata for accuracy and quality

~35%

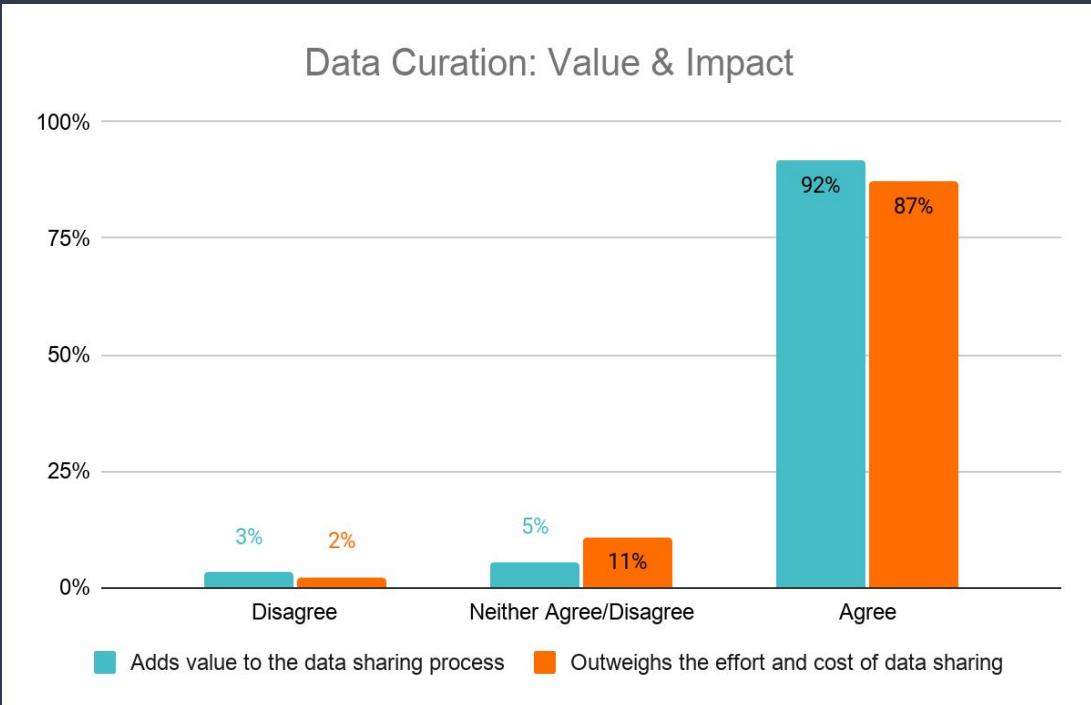
- Edit data for quality and accuracy



Results - Perceived Value & Impact of Curation

Top ranked impacts of data curation (out of 14)

- The ability for others to...
- #1 **find** the data
 - #2 **understand** the data
 - #3 **use** the data
 - #4 **access** the data
 - #5 **preserve** the data



What is the Value of Curation?

Perspective of repositories *doing* curation



Ethical. Reusable. Better.

DATA CURATION NETWORK

Perspective of depositors *receiving* curation

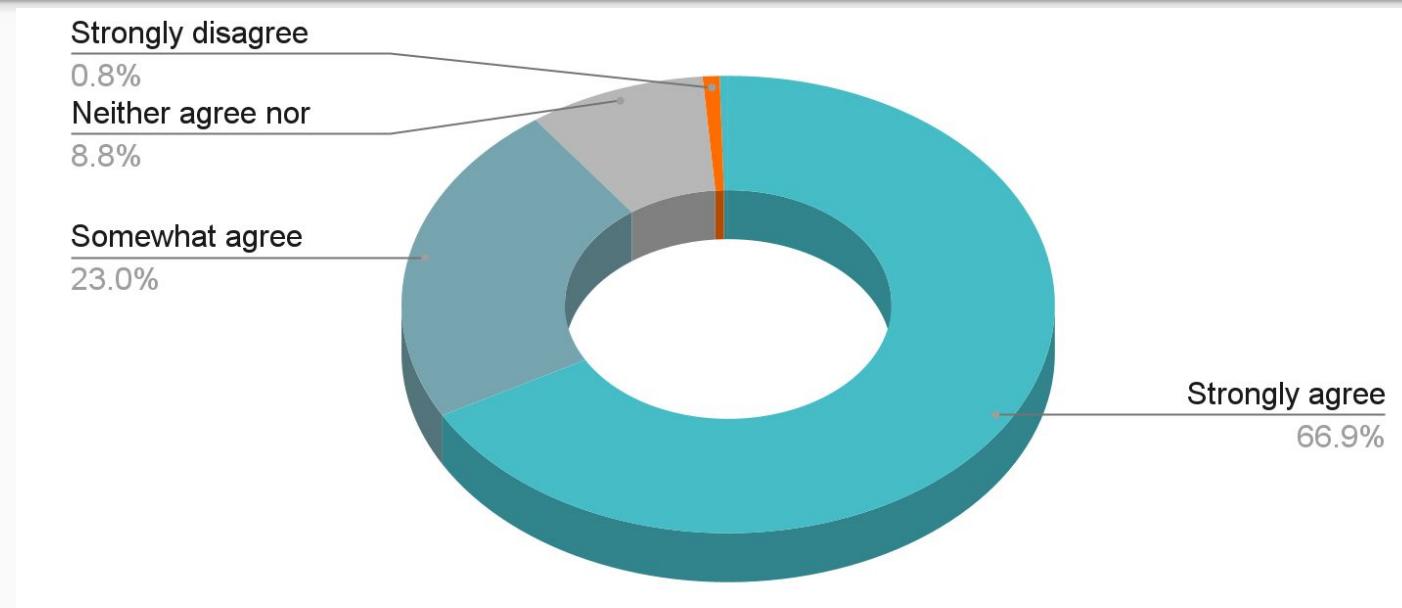


datacurationnetwork.org

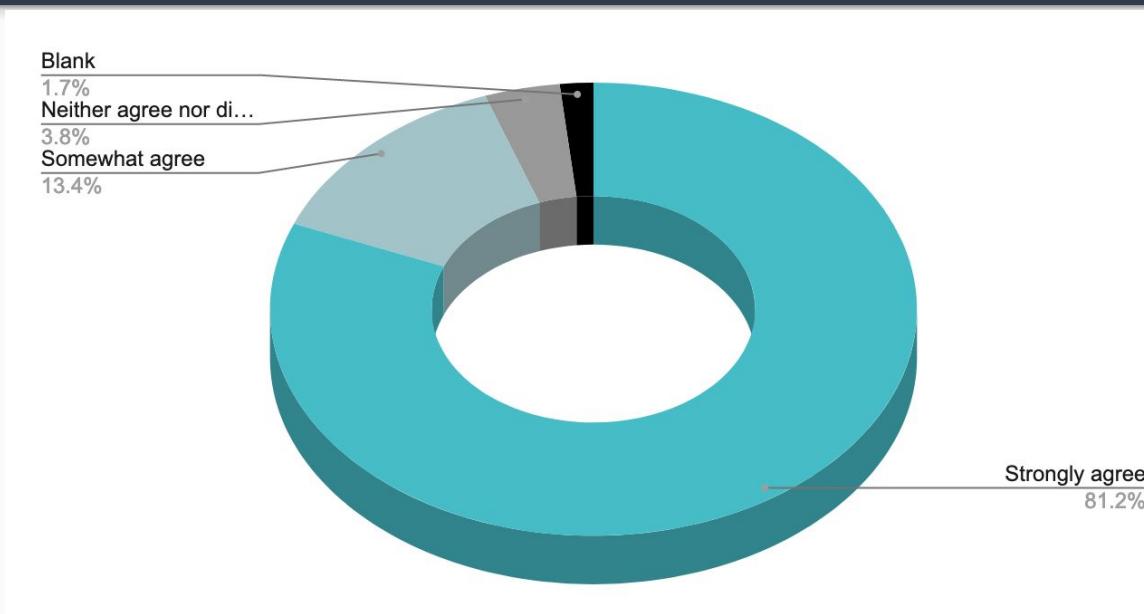
Survey Response rate = 40% (n=227)

	Date	Distribution Count	Response Count	Response Rate
Minnesota	4-26-21	197	82	42%
Cornell	4-26-21	34	18	53%
Michigan	5-11-21	130	63	48%
Duke	5-05-21	54	19	35%
Illinois	5-18-21	121	45	37%
Johns Hopkins	6-17-21	32	12	38%
Total		568	227	40%

Due to the curation process I felt more confident sharing my data. (N = 227)



Data curation by this repository adds value to the data sharing process. (N = 227)



What is the most "value-add" curation action taken by this repository? (N = 172)

caught mistakes - increased dataset quality!

Helping me prepare the data/metadata for long term access. The curator's expertise in this area was very helpful as I was not as familiar with what metadata was required to ensure easy, long term utility/access of my dataset.

The curation process makes the data more accessible and readable for other scientists. This amplifies the impact of our research.

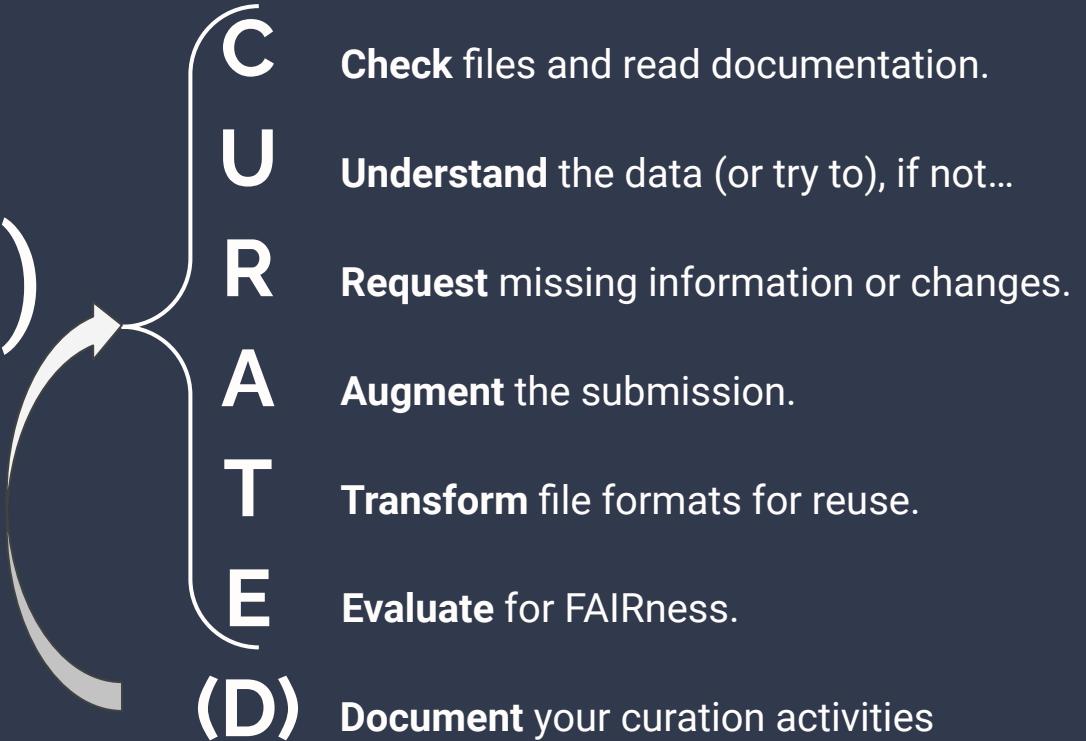
What do researchers value?

“The data curation was helpful in that it forced me to consider the **viewpoint of the end user** of the data. The curators nudged me to add pieces that would help third parties navigate the data (e.g. readme file) months or years in the future. **In retrospect** this is a good thing; at the time I was focused on getting my thesis chapter out and published, and didn’t pay attention to this aspect as much as I needed to.”

CURATE(D) Model



CURATE(D) Steps



C ➡ U ➡ R ➡ A ➡ T ➡ E ➡ (D)



Check files

CHECK step - action

- Review and inventory the content of the data files (e.g., open the files)
- Verify all metadata provided by the author and review the available documentation
- Ensure content of dataset is in scope
- Look for obvious ethical/sharing red flags

CHECK step - curator checklist

- Begin Curator Log to track curation decisions
- Open the related article and supporting information if available
- Inventory the dataset
 - Identify file formats
 - Review file organization, hierarchy, and naming convention(s)
 - Extract zip files when possible)
- Create working copy of files for formal inventory and testing
- Examine code for obvious errors/missing components, etc.
- Check that metadata quality is rich, accurate, and complete to institutional requirements.
- Check documentation
 - Complete
 - Needs work
 - If missing, document for the “Request” step
- Check whether human subject data. If so,
 - Request consent form / participation agreement if not present

CHECK step - curator checklist

- Check the accessibility of all files
 - Ensure there are robust descriptions in plain text of data files and any images.
- Check whether any visualization(s) of data are easily accessible
 - Review alt-text and visualization descriptions. Ensure these describe, but do not interpret, associated visualizations.
- Check data visualizations follow accessible color contrast guidelines
- Recommend graphical representation

- Recommend web-accessible surrogate

CHECK step – key ethical considerations

Review participant agreement and data use agreements; examine potential impacts of sharing this data. Consider:

- Individuals and communities represented
- Representativeness of diverse human populations
- Protection or endangerment status of species
- Geographic locations (e.g., specificity, contested boundaries, historical and current political situations)
- Animal research ethics and approval

Is it possible that the data deposit may impact a specific group?

Does this data deposit follow compliance and institutional policy?

CURATED Training:

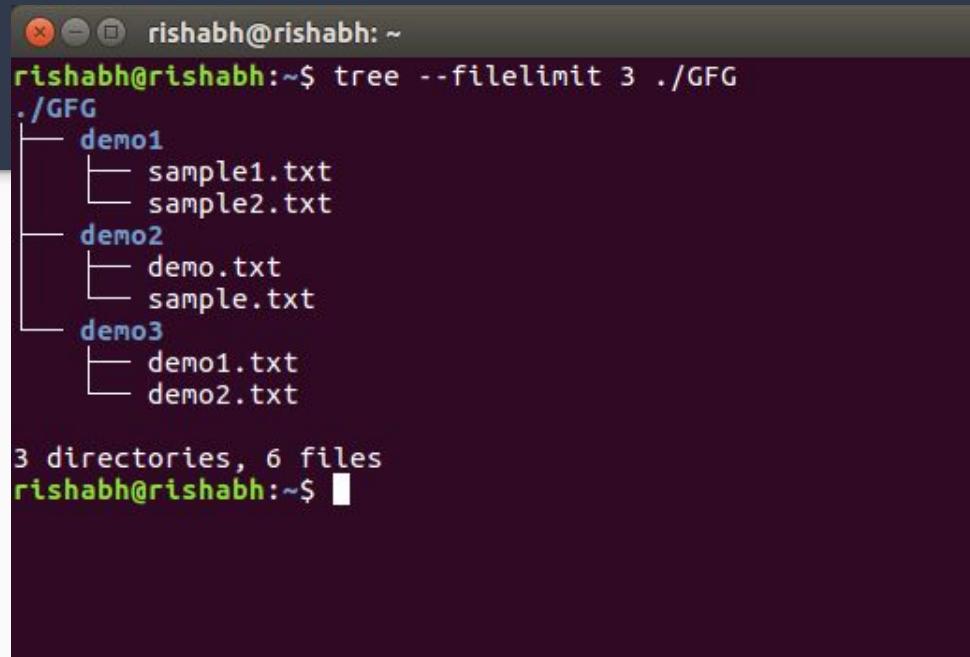
<https://datacurationnetwork.github.io/CURATED/>

Software Roundup

1. Appraisal Tools

a. Tree and Print Window

Tree /f /a > FileToPrint



The screenshot shows a terminal window with a dark background and light-colored text. The title bar reads "rishabh@rishabh: ~". The command entered is "tree --filelimit 3 ./GFG". The output displays a file tree structure:

```
rishabh@rishabh:~/GFG
./GFG
├── demo1
│   ├── sample1.txt
│   └── sample2.txt
└── demo2
    ├── demo.txt
    └── sample.txt
└── demo3
    ├── demo1.txt
    └── demo2.txt

3 directories, 6 files
```

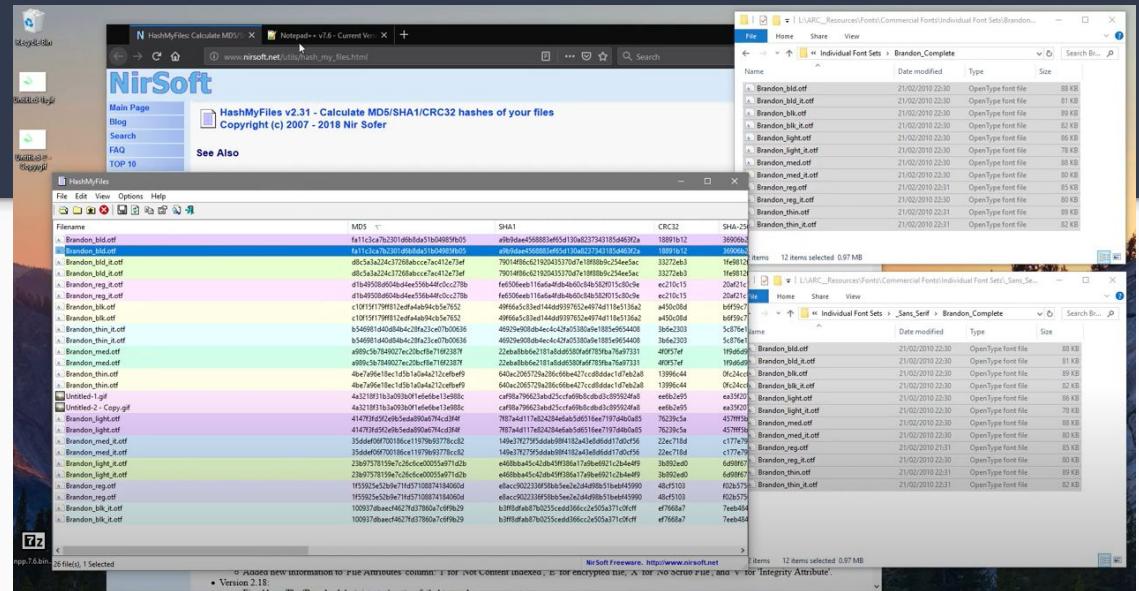
At the bottom, the prompt "rishabh@rishabh:~\$ █" is visible.

Software Roundup

1. Appraisal Tools

- a. Tree command
- b. HashMyFiles

Download HashMyFiles:
https://www.nirsoft.net/utils/hash_my_files.html



Screenshot from HashMyFiles tutorial on YouTube

Software Roundup

2. Speciality Tools (Processing and Review)

- a. ArcGIS, QGIS, GeoNetwork
- b. RStudio
- c. Matlab
- d. AliView (genomics data viewer)
- e. Omero (microscopic images)
- f. ChemDraw
- g. MZmine - mass spec data
- h. Jena - 3d molecular structure/Crystallographic files

Discussion
**What tools do you use
in curation?**

C ➡ U ➡ R ➡ A ➡ T ➡ E ➡ (D)



Understand the data

Understand step - action

- Check for quality assurance and usability issues
- Try to detect and extract any “hidden documentation”
- Determine if the documentation of the data is sufficient for a user with similar qualifications to the author’s
- Assess whether the data requires a deeper dive due to the sensitivity of the data

Understand step - curator checklist

- ❑ Examine files, organization, and documentation more thoroughly.
Are there changes that could enhance the dataset?
 - ❑ Are there missing data?
 - ❑ Could a user with similar qualifications to the author's understand and reuse these data and reproduce the results?
 - ❑ Are the data, documentation and/or metadata presented in a way that aids in interpretation? (e.g., readme)
- ❑ Record all questions and concerns in a curation log.

Understand step - curator checklist

Tabular data sub-checklist (e.g., Microsoft Excel)

- Organization of data well-structured
 - Not rectangular
 - Split tables into separate tabs
- Headers/codes clearly defined
 - Define headers
 - Clarify codes used _____
 - Clarify use of “blanks”
 - Clarify units of measurement
- Quality control clearly defined
 - Unclear quality control
 - Update/add Methodology

[Full checklist](#)

UNDERSTAND step – key ethical considerations

If working with human data, is this research done with and not on communities and populations involved? (You may wish to review data sources, researchers, and their connections to the communities and subjects they are serving to facilitate further conversation with researcher(s).)

- Are there authoritative group representatives who should be contacted in the next (request) step?

Are there labels or other descriptive indicators that could be applied to better represent or protect an identified group of people impacted by this dataset? (Example: TK labels)

CURATED Training:

<https://datacurationnetwork.github.io/CURATED/>

Documentation can be for...

- assessing data for reuse
- assistance in actual reuse
- explaining and/or maintaining consistency of data
- training new people
- efficiency in archiving



```
<?xml version="1.0"?>  
  
<metadata  
    xmlns="http://example.org/myapp/"  
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"  
    xsi:schemaLocation="http://example.org/myapp/ http://example.org/myapp/schema.xsd"  
    xmlns:dc="http://purl.org/dc/elements/1.1/">  
  
    <dc:title>  
        UKOLN  
    </dc:title>  
    <dc:description>  
        UKOLN is a national focus of expertise in digital information  
        management. It provides policy, research and awareness services  
        to the UK library, information and cultural heritage communities.  
        UKOLN is based at the University of Bath.  
    </dc:description>  
    <dc:publisher>  
        UKOLN, University of Bath  
    </dc:publisher>  
    <dc:identifier>  
        http://www.ukoln.ac.uk/  
    </dc:identifier>  
  
</metadata>  
  
Note that the http://example.org/myapp/schema.xsd XML schema does not exist - this is a fictitious example.
```

Informal ReadMe

Formal Schema

Lower-Barrier
Fast
Easy

Lower-Potential
Irregular
Incomplete

Higher-Potential
Standardized
Machine actionable

Higher-Barrier
Slow
Skilled

Readme files as documentation

Recommended content

- Data and file overview
- Sharing and access information
- Methodological information
- Data-specific information
- General information



Cornell Un
Library

Some issues with CHECK

- Can't open the files/code
- Too many files to open
- Can open the files, but don't know what I'm looking for
- No (or not enough) documentation

Some issues with UNDERSTAND

- Not enough experience to evaluate these data
- Documentation is over my head
- Still can't tell what's typical in this field
- Required software is too hard to find/use/justify \$\$\$

C ➞ U ➞ R ➞ A ➞ T ➞ E ➞ (D)
↓

Request missing information

The “Request” Step

What information do you need to address the gaps that you have identified?

How will you get this information?



Goals for the “Request” Step



- Goal #1 - Get a response
- Goal #2 - Establish a rapport
- Goal #3 - Get the information that you need

Photo by [rawpixel](#) on [Unsplash](#)

Determine what information is needed



Image from Taskworld Blog:

<https://medium.com/taskworld-blog/build-on-top-of-taskworld-with-public-api-fca3d9e110d2>

- Be as specific as possible
- Prioritize your needs
 - Essential
 - Important
 - Supplemental
- Do you need information or permission to take action?

Advice: Make it Easy



Image from Staples website:
https://www.staples.com/Staples-Easy-Button/product_606396

- Limit to 4 asks
- Be specific, but try to keep it short
- Keep it as simple as possible (Ideally, response could be “yes or no”)
- Provide resources where useful
- If possible, do the work yourself and ask for approval

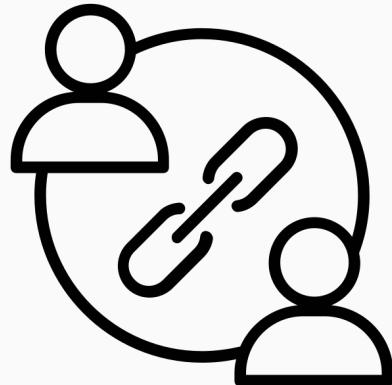
Request step - action

- Triage and prioritize issues. Highlight those with the highest data reuse implications
- Convey a sense of urgency, as it becomes more difficult to get responses from researchers as time passes
- Collaborate with the researcher(s) to make necessary changes
- Communicate any changes you, the curator, will make on their behalf
- Pause and consider how best to frame and communicate requests
- This should be the start of a conversation

REQUEST step - key ethical considerations

- Consider asking researchers if their participants will be notified that their data are being shared
- If you feel uncomfortable about sharing the data in its current state and/or it does not meet your institution's requirements, reserve the right not to publish
- Consider asking researcher(s) if there are limitations to how data could/should be used. Include any limitations in the documentation

Rapport: It may take more than one email



Created by Vectors Market
from Noun Project

- May be able to ask for additional information once rapport is established
- Offer to meet in person if needed
- Put the depositor first
 - What are their needs in depositing?
 - A little empathy goes a long way

REQUEST step - curator checklist

- ❑ Ask about additional data contributors, beyond publication authors. Consider using the Contributor Roles Taxonomy to communicate this: <https://casrai.org/credit/>
- ❑ Summarize conversations / outreach efforts in Curator Log

Long version at: [Request Step](#)

Request information

“I'm impressed with how well you seem to understand the data and how thorough your review was!”

- Researcher response to Data Repository for the U of Minnesota (DRUM) curator, Feb 2017

C $\Rightarrow\!\!>$ U $\Rightarrow\!\!>$ R $\Rightarrow\!\!>$ A $\Rightarrow\!\!>$ T $\Rightarrow\!\!>$ E $\Rightarrow\!\!>$ (D)



Augment the submission

Augment step - action



- Apply the information that you have received from the depositor
- Add metadata as appropriate to enhance discovery and access
- Improve documentation to make data more understandable, interoperable and reusable

Image: pxfuel

<https://www.pxfuel.com/en/free-photo-qfeat>

Ethical. Reusable. Better.

DATA CURATION NETWORK

datacurationnetwork.org

Areas to consider



- Discoverability
 - Full text indexing
 - File Ordering
 - Data Set Description
 - File Description
 - Keywords
- Access
 - File / Folder Compression

Photo by [Louis Reed](#) on [Unsplash](#)

Ethical. Reusable. Better.

DATA CURATION NETWORK

datacurationnetwork.org

Areas to consider



- Linkages
 - Articles & other outputs
 - Other sources or inputs
 - Related datasets
- Persistent Identifiers
 - DOIs, ORCIDs, RORs

Photo by Chris Leipelt on [Unsplash](#)

Common metadata elements needed for datasets

- Author(s)
- PID (DOIs)
- Title
- Keywords
- Abstract / Description
- Language
- Citation
- Link(s) to related publications
- Licensing information
- Funder information

Brief recap of “C” & “U”

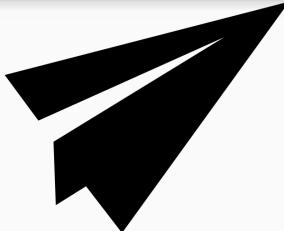


Image Source: <https://www.linkedin.com/learning/the-science-of-sales>

These are the steps taken to learn about the data set

- What is it you have?
- What should (or could) it be?
- What are the gaps between what you have and what you want?

Brief recap of “R” & “A”



Created by Briyan Design
from Noun Project



Created by NeMaria
from Noun Project

These are the steps taken get more information and improve the metadata

- What are the priority issues?
- What are 4 things you need to know?
- How can you enhance the metadata using what you've learned?
- What will make it more findable and useable?

C ➡ U ➡ R ➡ A ➡ T ➡ E ➡ (D)



Transform file formats

Example File Format Transformations

Native Software or Format	Suggested Formats or Transformations ➤	Transformation Tools and Notes
CZI (microscope images)	TIFF, JPG, FITS	Use “export”; Omero, Bioformats; WikiData tracks software and file formats for preservation
Microsoft Excel / XLS, XLS	CSV, TSV	Use “save as”; Use Excel Archival Tool to preserve formulas
Chemdraw / CDX	CDXML, MOL, JPG, 001, OPJ, TRI	Retain original. Some conversions will result in loss of information.
PDF	PDF-A	Use “save as”
MP4, MOV, WMV	Uncompressed AVI or MOV + captions	No information is gained going from a “lower” resolution image to a “higher” one, but long-term access may be improved. Use YouTube, Vimeo, Kaltura or other tools for captioning.
.SHP (geocoded xls)	CSV + extracted metadata	Retain .SHP. Use FME Tool or ArcGIS
Webpages	WARC, TIFF	Link to Internet Archive. Provide screen shots.

Transform step - curator checklist

- Preferred file formats in use
- Recommend conversion
 - from _____
 - to _____
- Retain original formats
- Software needed is readily available
- Unclear version of software
- Unclear software used
- Visualization of data easily accessible
- Recommend graphical representation _____
- Recommend web-accessible surrogate _____

C $\Rightarrow\!\!>$ U $\Rightarrow\!\!>$ R $\Rightarrow\!\!>$ A $\Rightarrow\!\!>$ T $\Rightarrow\!\!>$ E $\Rightarrow\!\!>$ (D)



Evaluate for...

What do we mean by Evaluate?



C ➞ U ➞ R ➞ A ➞ T ➞ E ➞ (D)



Evaluate for curation success

What do we mean by Evaluate?

Did the curation partnership result in FAIR data?

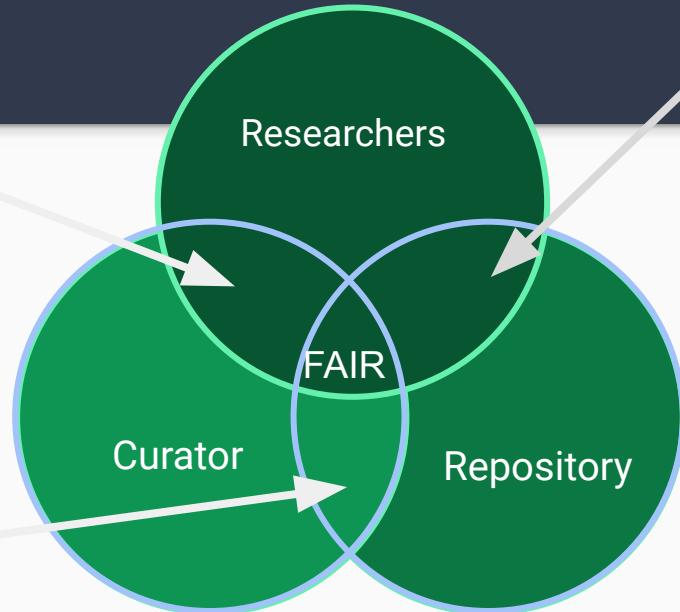
- Communication
- Expertise
- Buy-in?

Does the platform work for FAIR data?

- Platform Features
- Technology
- Standards

Will this repository make my data FAIR?

- Services
- Policies
- Transparency



C ➡ U ➡ R ➡ A ➡ T ➡ E ➡ (D)



Evaluate for curation FAIRness

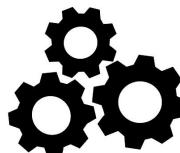
F
indable



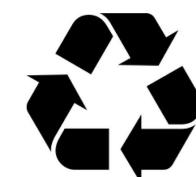
A
ccessible



I
nteroperable



R
eusable



Findable

*To be **findable (F)** or discoverable, data and metadata should be richly described to enable attribute-based search*

- (meta)data are assigned a globally unique and eternally persistent identifier
- data are described with rich metadata
- (meta)data are registered or indexed in a searchable resource
- metadata specify the data identifier

Accessible

*To be broadly **accessible (A)**, data and metadata should be retrievable in a variety of formats that are sensible to humans and machines using persistent identifiers*

- (meta)data are retrievable by their identifier using a standardized communications protocol
- the protocol is open, free, and universally implementable
- the protocol allows for an authentication and authorization procedure, where necessary
- metadata are accessible, even when the data are no longer available

Interoperable

To be **interoperable (I)**, the description of metadata elements should follow community guidelines that use an open, well defined vocabulary.

- (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation
- (meta)data use vocabularies that follow FAIR principles
- (meta)data include qualified references to other (meta)data

Reusable

*To be **reusable (R)**, the description of essential, recommended, and optional metadata elements should be machine processable and verifiable, use should be easy and data should be citable to sustain data sharing and recognize the value of data.*

- meta(data) have a plurality of accurate and relevant attributes
- (meta)data are released with a clear and accessible data usage license
- (meta)data are associated with their provenance
- (meta)data meet domain-relevant community standards

Evaluate step - curator checklist for FAIR

Findable -

- Metadata exceeds author/ title/ date,
- Unique PID (DOI, Handle, PURL, etc.).
- Discoverable via web search engines.

Accessible -

- Retrievable via a standard protocol (e.g., HTTP).
- Free, open (e.g., download link).

Interoperable -

- Metadata formatted in a standard schema (e.g., Dublin Core).
- Metadata provided in machine-readable format (OAI feed).

Reusable -

- Data include sufficient metadata about the data characteristics to reuse
- Contact info displayed if the direct assistance of the author needed.
- Clear indicators of who created, owns, and stewards the data.
- Data are released with clear data usage terms (e.g., a CC License).

C ⇒ U ⇒ R ⇒ A ⇒ T ⇒ E ⇒ (D)



Document your curation activities

Document your curation activities

Using whatever practices / tools in place at your institution, be sure you've documented and recorded everything you've done to curate this dataset



- **What to document**
 - Accessioning and deposit records
 - Repository dataset cataloging metadata
 - Provenance/change logs
 - Service workflow
 - Preservation packaging metadata



Stay in touch!

<https://datacurationnetwork.org/>

<https://datacurationnetwork.org/monthly-newsletter/>