

DATA CURATION NETWORK

Update for the Canadian Data Curators Forum

• • •

Hamilton, ON
October 17, 2019

Presented by



Lisa R. Johnston
Research Data Management Curation Lead
Principal Investigator, DCN
University of Minnesota



Cynthia Hudson Vitale
Head, Research Informatics & Publishing
Principal Investigator, DCN Education
Pennsylvania State University

Why data curation?

Most data is less than good

(Initially)

We generate a lot of data in our research. Many societies, funders, and publishers encourage data sharing, but...

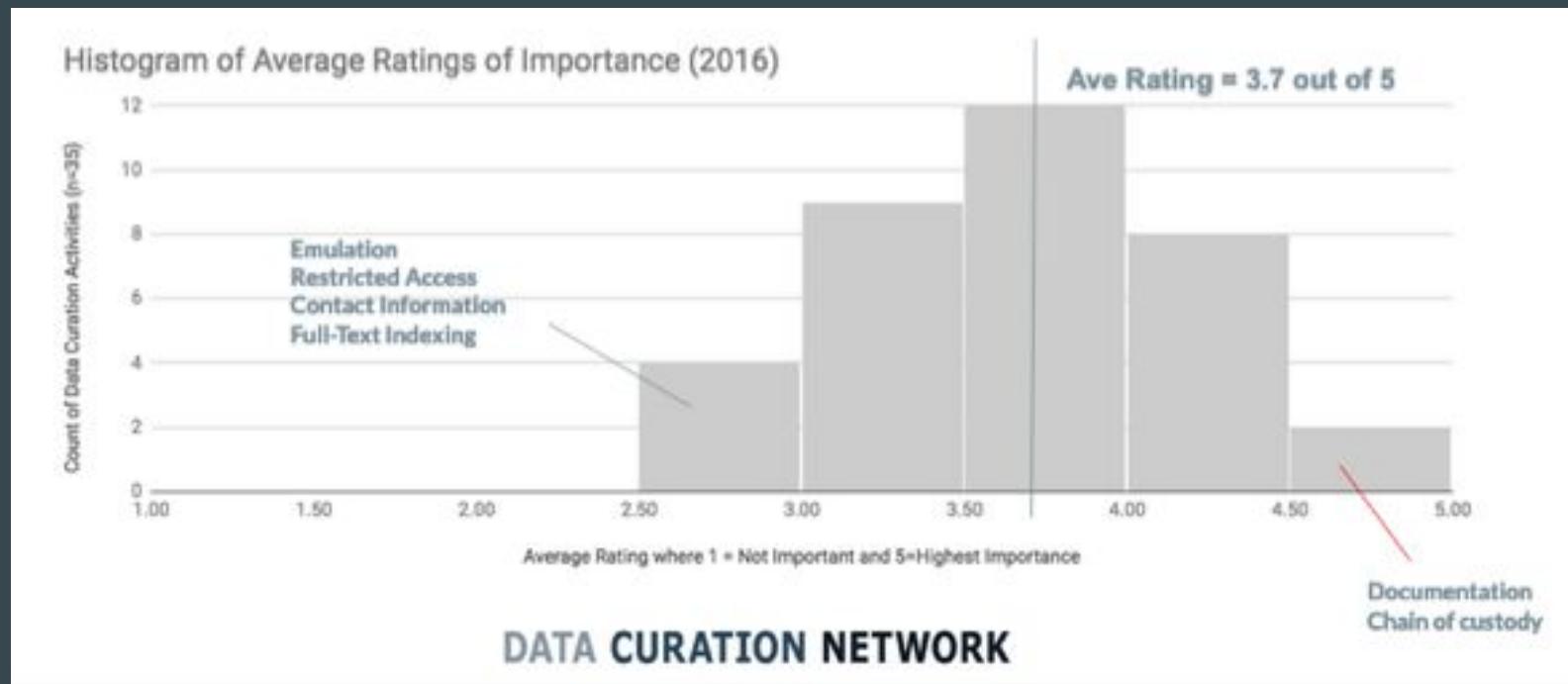
- Data are messy (lack context!)
 - Digital file formats are constantly at risk
 - Most data never leave their author's laptop ⇒ benign neglect
-

47%

Data sets with no documentation*

*2017 study of 175 data sets across 6 academic repositories in report: “Data Curation Network: A Cross-Institutional Staffing Model for Curating Research Data” (2017), <http://hdl.handle.net/11299/188654>.

Focus Group Results (n=91 across 6 institutions)



Focus Group Results (n=91 across 6 institutions)



Most Important Activities* (4 out of 5)

- (Create) Documentation (4.6)
- Secure Storage (4.4)
- Quality Assurance (4.3)
- Persistent Identifier (4.3)
- Software Registry (4.1)
- Data Visualization (4.0)
- File Audit (4.0)
- (Create) Metadata (4.0)
- Versioning (3.9)
- Contextualization (3.9)
- Code Review (3.9)
- File Format Transformations (3.9)

Focus Group Results (n=91 across 6 institutions)



Most Important Activities* (4 out of 5)

- (Create) Documentation (4.6)
- Secure Storage (4.4)
- Quality Assurance (4.3)
- Persistent Identifier (4.3)
- Software Registry (4.1)
- Data Visualization (4.0)
- File Audit (4.0)
- (Create) Metadata (4.0)
- Versioning (3.9)
- Contextualization (3.9)
- Code Review (3.9)
- File Format Transformations (3.9)

Does this happen for your data?

- Persistent Identifier (37% happens)
- Software Registry (41% happens)
- File Audit (16% happens)
- Contextualization (38% happens)
- Code Review (38% happens)

Focus Group Results (n=91 across 6 institutions)



Most Important Activities* (4 out of 5)

- (Create) Documentation (4.6)
- Secure Storage (4.4)
- Quality Assurance (4.3)
- Persistent Identifier (4.3)
- Software Registry (4.1)
- Data Visualization (4.0)
- File Audit (4.0)
- (Create) Metadata (4.0)
- Versioning (3.9)
- Contextualization (3.9)
- Code Review (3.9)
- File Format Transformations (3.9)



Is so, are you satisfied with the result?

- Documentation (26% satisfied),
- Secure storage (38% satisfied),
- Quality Assurance (14% satisfied),
- Data Visualization (12.5% satisfied),
- Metadata (29% satisfied)
- Versioning (13% Satisfied)
- File Format Transformations (29% satisfied)

Data Curation

The encompassing work and actions taken in order to provide meaningful and enduring access to data.

- ✓ Finding and adding missing files and documentation
- ✓ Screening for privacy disclosure risk
- ✓ Detecting and fixing code and other quality assurance issues
- ✓ Transforming file formats for long term access
- ✓ Arranging and describing files
- ✓ Reviewing and augmenting metadata

Impact of Data Curation

Researchers

Scholars and researchers trust professionally-curated data that are findable accessible, interoperable, and reusable (FAIR).

Curators

The actions taken by a data curator result in ethical, reusable, and better datasets for research and education.



Public

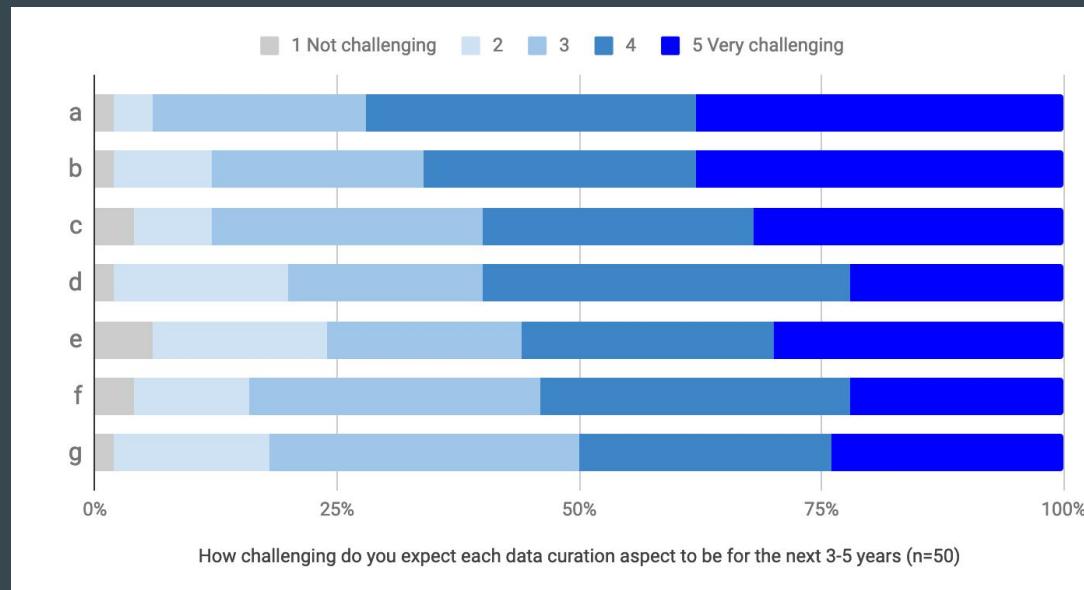
Everyone benefits from accessible, transparent, and reproducible research.

Repositories

Data repositories provide technical access and preservation services to publish high-quality data sets.

Barriers to Well-Curated Data

- A. Expertise in domain data
- B. Scaling with increased demand
- C. Training & retooling existing staff
- D. Outreach/marketing
- E. Recruiting & retaining staff
- F. Keeping up with technology changes
- G. Keeping up with data sharing
requirement changes



*2017 study of 80 ARL Institutions in: Hudson-Vitale, C, Imker, H, Johnston, LR, Carlson, J, Kozlowski, W, Olendorf, R and Stewart, C. **SPEC Kit #354: Data Curation**. (2017). Association of Research Libraries (ARL). May 2017. <https://doi.org/10.29242/spec.354>.

Data Curation Network

Mission

Trusted community-led network
that enables researchers to openly
share data in ways that are

Ethical. Reusable. Better.

DATA CURATION NETWORK



Alfred P. Sloan
FOUNDATION

DATA CURATION NETWORK

Community Led

10 Partners

28 data curators

43 domains

26 speciality file format types



Alfred P. Sloan
FOUNDATION



University of Minnesota
Lead: Lisa Johnston, PI
Liza Coburn, Project Coordinator
Curators: Katie Wilson, Alicia Hofelich Mohr, Shanda Hunt, Melinda Kernik, Wanda Marsolek, Alexis Logsdon
Admin: Janice Jaguszewski



Penn State University
Lead: Cynthia Hudson-Vitale
Curators: Xuying Xin, Seth Erickson



University of Illinois
Lead: Hoa Luong
Curator: Ashley Hetrick
Admin: Heidi Imker



Duke University
Lead: Joel Herndon
Curators: Jen Darragh, Sophia Lafferty-Hess
Admin: Timothy M. McGahey



Washington University in St. Louis
Lead: Jennifer Moore
Curator: Dorris Scott

University of Michigan
Lead: Jake Carlson
Curators: Susan Borda, Rachel Woodbrook

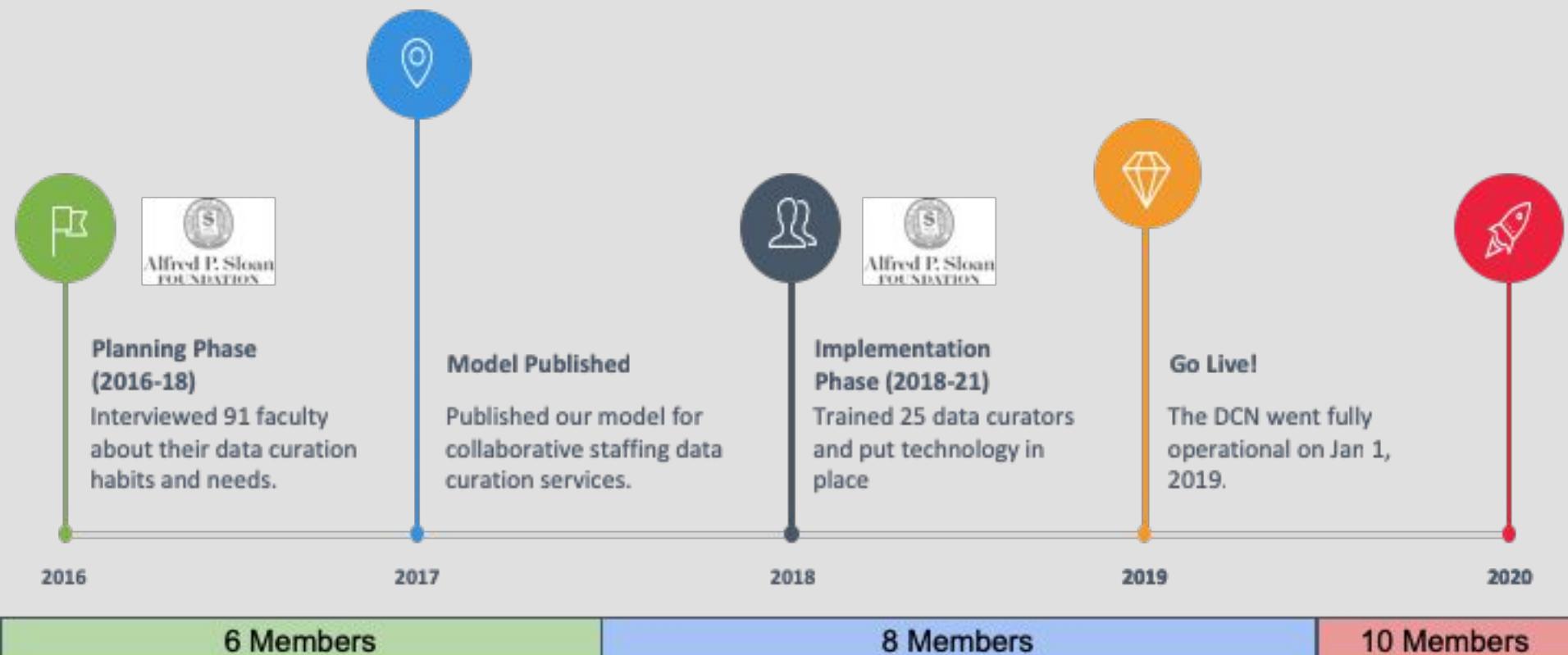
Cornell University
Lead: Wendy Kozlowski
Curators: Sarah Wright, Henrik Spoon

Johns Hopkins University
Lead: Mara Blake, Co-PI
Curators: Chen Chiu, Dave Fearon, Marley Kalt

Dryad Digital Repository
Lead: Elizabeth Hull
Curators: Erin Clary, Debra Fagan, Rich Yaxley

New York University
Lead: Katie Wissel
Curator: Andrew Battista

DATA CURATION NETWORK

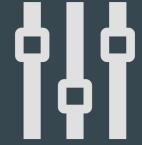




DCN Curation



DCN Education



DCN Resources



DCN R&D



DCN Sustainability

**Vision for the
Data Curation
Network**



DCN Curation



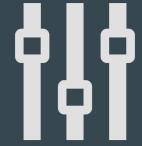
Provide expert data curation services for network partners



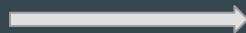
DCN Education



Offer professional development opportunities for an emerging data curator professional community



DCN Resources



Create and openly share data curation best practices



DCN R&D



Demonstrate that curated datasets are measurably of greater reuse value than non-curated data



DCN Sustainability



Expand into a sustainable entity that grows beyond our initial partner institutions

Data Curation

The DCN provides the training, coordination, and technical infrastructure to seamlessly connect expert data curators across the network with all types of data sets for robust curation.

<http://datacurationnetwork.org>

DCN CURATE Steps

DCN Curators will take **CURATE** steps for each data set, that

C **Check** data files and read documentation

U **Understand** the data (try to), if not...

R **Request** missing information or changes

A **Augment** the submission with metadata for findability

T **Transform** file formats for reuse and long-term preservation

E **Evaluate** and rate the overall submission for FAIRness.

Table A1. Draft checklist of DCN CURATE steps and FAIRness scorecard

CURATE Actions	Curation Checklist
Check data files and read documentation <ul style="list-style-type: none">▪ Review the content of the data files (e.g., open and run the files or code).▪ Verify all metadata provided by the author and review the available documentation.	<input type="checkbox"/> Files open as expected <ul style="list-style-type: none"><input type="checkbox"/> Issues _____ <input type="checkbox"/> Code runs as expected <ul style="list-style-type: none"><input type="checkbox"/> Produces minor errors<input type="checkbox"/> Does not run and/or produces many errors <input type="checkbox"/> Metadata quality is rich, accurate, and complete <ul style="list-style-type: none"><input type="checkbox"/> Metadata has issues _____ <input type="checkbox"/> Documentation Type (<i>circle</i>) Readme / Codebook / Data Dictionary / Other: _____ <ul style="list-style-type: none"><input type="checkbox"/> Missing/None<input type="checkbox"/> Needs work
Understand the data (or try to) <ul style="list-style-type: none">▪ Check for quality assurance and usability issues such as missing	<i>Varies based on file formats and subject domain. For example....</i>

DCN Workflow

Uncurated Data

Presenting scale and expertise challenges to individual institutions



Curated Data

at scale and with great efficiency through shared Data Curation Network

DCN Workflow

Uncurated Data

Presenting scale and expertise challenges to individual institutions



Curated Data
at scale and with great efficiency through shared Data Curation Network



DCN Workflow

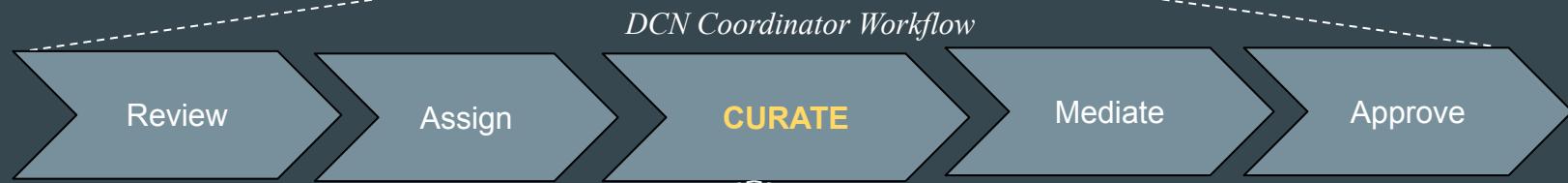
Uncurated Data

Presenting scale and expertise challenges to individual institutions

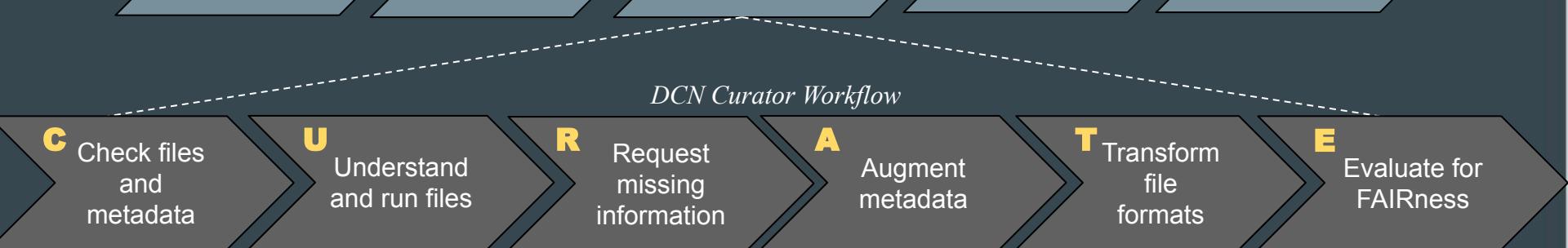


Curated Data

at scale and with great efficiency through shared Data Curation Network



DCN Coordinator Workflow



DCN Curator Workflow

Tools we use to run the Network

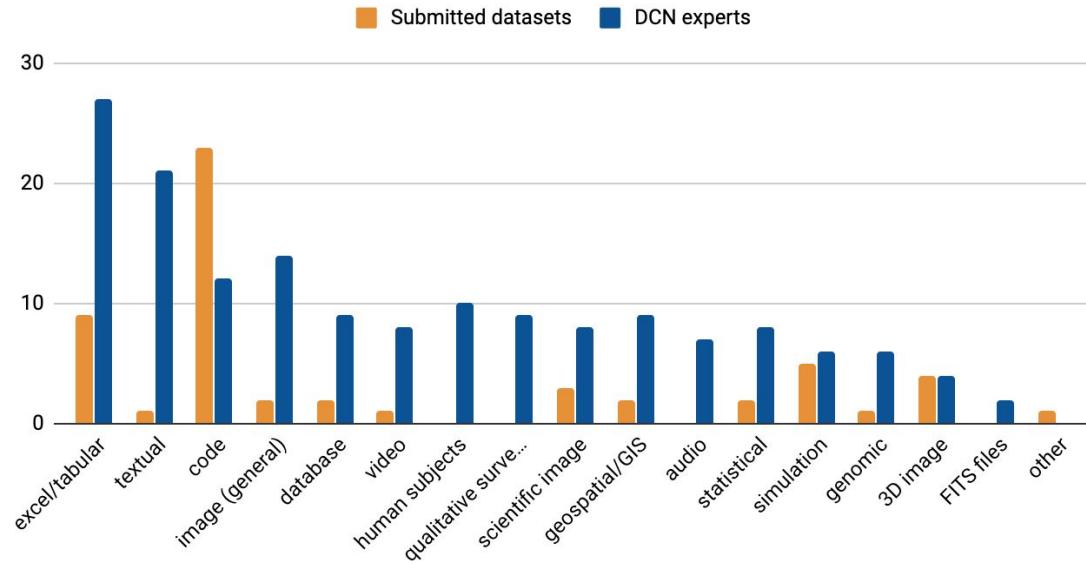
- Jira time-tracking software¹
- Survey to capture curator expertise
- Annual training and networking opportunities
- Slack, listserv for ongoing community



1. Kozlowski, Wendy, Elizabeth Coburn and Mara Blake. Walk it Like you Talk it: Jira as a tool for documenting the curation process. RDA 13th Plenary Meeting; 2019 April 2-4; Philadelphia, PA.

Measures of Success

Data types of datasets submitted to the DCN relative to data type expertise within the DCN



Data curation stats (viz)

Satisfaction surveys

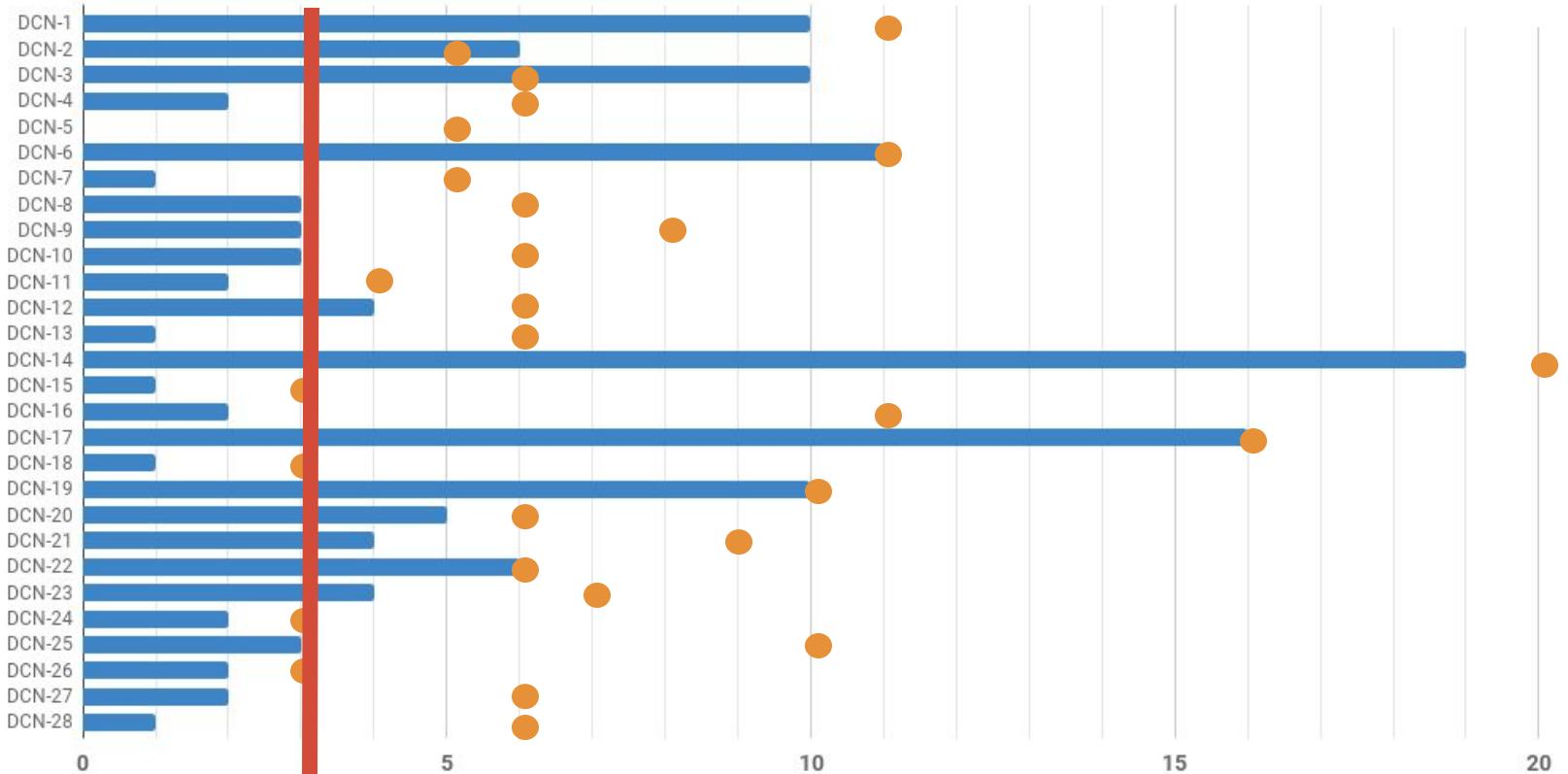
Efficiencies gained

- 58 data sets curated
- Averaging 2.5 hours/dataset
- Assignments made within 24 hours

DCN Dataset Turnaround Time in Working Days (M-F)

Number of Working Days from Submission to Curation Completion

Dataset Identifier



Median Turnaround Time = 3 Days

○ = Due Date (in Working Days)

Ashley Hetrick



Assistant Director for Research Data Engagement and Education
University of Illinois

Hetrick leads the development and delivery of data management work, [Follow](#) ..., well as data management plan (DMP) reviews, for the [Research Data Service \(RDS\)](#). Prior to this position, Hetrick worked in various roles within Illinois' central Information Technology (IT) group, including social media analytics, IT communications, and leading an AV/IT help desk service.

Data sets curated by Ashley

[DCN-7: Forest Resources Database](#)



DATA CURATION NETWORK

[Home](#) [About](#) [Our Curators](#) [Resources](#) [News](#) [Events](#) [Contact](#)

DCN-7: Forest Resources Database

Data set citation

"Cloquet Forestry Center Continuous Forest Inventory (1959-2014)" available at the Data Repository for the University of Minnesota, <https://doi.org/10.13020/096z-kg59>.

Curated by [Lisa Johnston](#) at the Data Repository for the University of Minnesota and [Ashley Hetrick](#) at the Illinois Data Bank.

Curation actions

[DISCOVERY SERVICES](#)

[DOCUMENTATION](#)

[EVALUATE FAIRNESS](#)

[Follow](#) ...

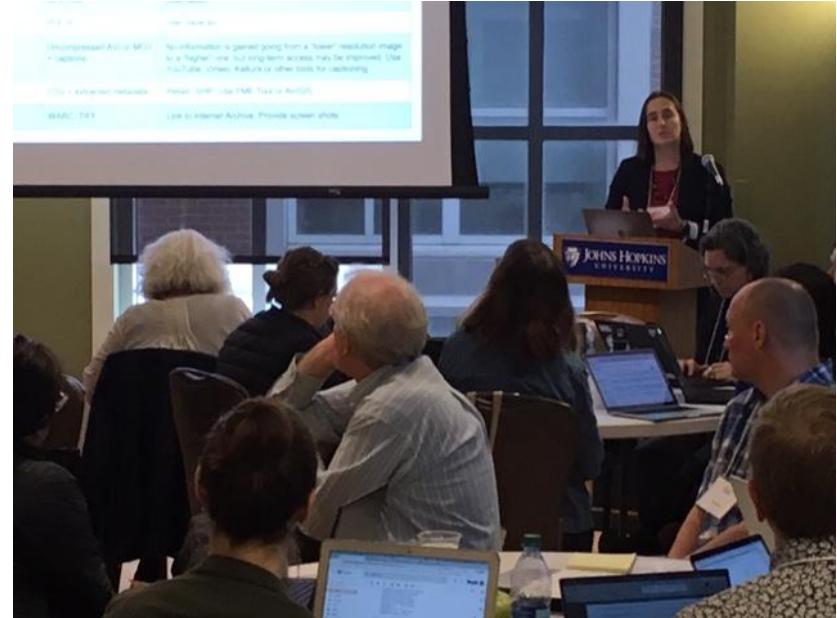
DCN Education

We offer professional development opportunities
for an emerging data curator professional
community

<https://sites.psu.edu/dcnworkshops>

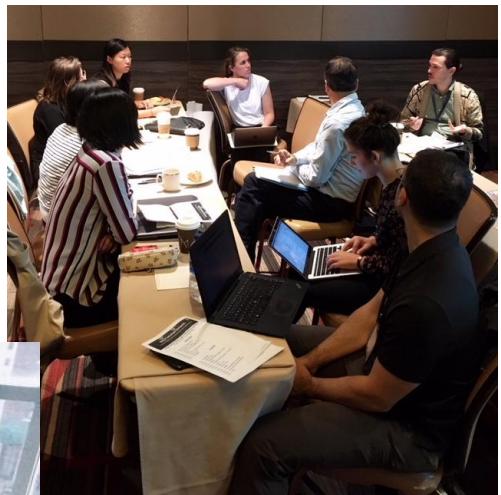
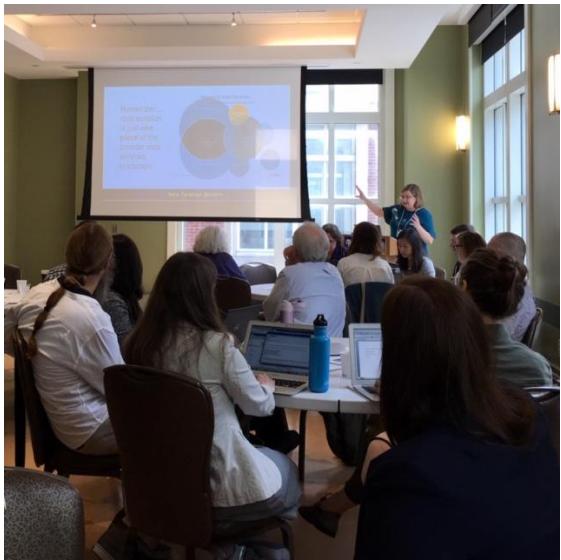
Enhancing Expertise throughout the Broader Community

Cynthia Hudson Vitale, PI



Specialized Data Curation Workshop @JHU 2019

DCN Education



<https://sites.psu.edu/dcnworkshops/>

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 for curation nationwide.
3. Meet like-minded colleagues who are interested in building and extending curation practices.





DATA CURATION NETWORK

Specialized Data Curation Workshop Agenda

April 17th & 18th ♦ Johns Hopkins University ♦ Baltimore, Maryland

Wednesday

9:00	Welcome & Breakfast
9:30	The Value of Curation
10:00	Curation Deep Dive #1: C Step
10:30	Break
10:45	Curation Deep Dive #1: U Step
12pm	Lunch
1:00	Primer Timer → pitch idea of primer topics
1:30	Curation Deep Dive #2: R & A Steps
2:30	Break
3:00	Curation Deep Dive #2: R & A Steps continued
4:00	End of Day One
5:30	Reception

Thursday

9:00	Breakfast
9:30	Coffee with Data
10:15	Review Day 1
10:30	Curation Deep Dive #3: T Step
11:30	Lunch
12:15	Curation Deep Dive #3: E & D Step
1:15	Primer Time 2
2:00	Group feedback on primers
2:15	Wrap up
2:30	Everyone Disperses

Check files
Understand or try to
Request missing information
Augment the submission
Transform the format
Evaluate for FAIRness
Document throughout

www.datacurationnetwork.org



Pictured: Group activity at the DCN Specialized Data Curation Workshop, co-located at the DLF Forum on October 17-18, 2018.

Our curriculum engages attendees with lectures, group activities and demonstrations.

Hands-on data curation activities

Data Curation Assignment: Images (Penn State)

 Survey Data

 Tabular Data

 Code

 Image Data

 Geospatial Data



Title: S'Urachi Site-Based Archaeological Survey
2015

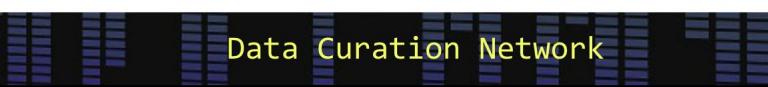
Author: Victor T. Hail

Discipline: Archeology

Date: 2015

Access: Public

Reason for deposit: Connect to published article and report

 Data Curation Network

DCN Workshops



INSTITUTE of
Museum and Library
SERVICES

Workshop #1
Las Vegas, NV
Oct 2018 (DLF)
for 22 accepted
applicants

- Geodatabases
- netCDF files
- Wordpress
- SPSS
- Microsoft Excel
- Jupyter Notebooks
- Microsoft Access

Workshop #2
Baltimore, MD
April 2019 (JHU)
for 27 accepted
applicants

- Atlas.ti
 - Confocal microscopy
 - GeoJSON
 - Google Docs
 - Lidar Point Clouds
 - NVivo
- PDF
 - R
 - STL files
 - Tableau
 - Text/character encoding

FUTURE
Workshop #3
St Louis, MO
(Wash U)

- 31 attendees
- Primers to be published in 2020

Data Curation Resources

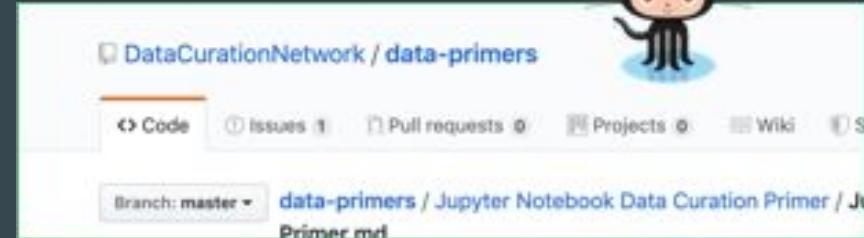
Create and openly share data curation best practices

<https://github.com/DataCurationNetwork/data-primers>

DCN Resources

Data Curation Checklists
format-specific CURATE(D) checklists

Data Curation Primers
concise, actionable resources meant to assist data curation in adding value to a dataset.



Primer creation process:

- Attendees select topic (2-3 people)
- Working session (Roadmap)
- DCN mentor
- Meet monthly (6 months)
- Peer review; co-occurring with internal webinar
- Revisions
- Publish to Github / IR



Our Peer Review Process

Primers are formally peer reviewed half-way through the six month process, or additionally if needed.

1. What steps in the curation process were you able to complete using this primer?
2. What sections of the primer did you find most useful?
3. Do you have suggestions for how content may be revised or enhanced?



Publication

<https://github.com/DataCurationNetwork/data-primers>

The screenshot shows a GitHub repository page. At the top, it displays the repository name "DataCurationNetwork / data-primers". Below the repository name are standard GitHub navigation links: "Unwatch" (2), "Star" (7), "Fork" (2), "Code", "Issues 1", "Pull requests 0", "Projects 0", "Wiki", "Security", "Insights", and "Settings". A dropdown menu indicates the branch is "master". The main content area shows a single file, "Primer.md", which is a Jupyter Notebook Data Curation Primer. The file was added by "cynhusdon" via upload on May 30 at commit 566c52a. The file has 253 lines (186 sloc) and is 17.2 KB. Below the file information, there is a large preview window showing the content of "Primer.md". The preview title is "Jupyter Notebooks: A Primer for Data Curators". Underneath the title, the word "Participants:" is listed, followed by a bullet point: "Daina Bouquin, Center for Astrophysics. Harvard & Smithsonian. (daina.bouquin@cfa.harvard.edu)".

- Published on GitHub
- Primers are expected to grow from their original version
- The community may suggest revisions

SPSS

Authors: Joshua Dull, Sai Deng, Shahira Khair & Jeanine Finn

DCN Mentor: Sophia Lafferty-Hess

<https://github.com/DataCurationNetwork/data-primers>

Key Curatorial Considerations:

- Preservation actions
 - Save as .por? To ASCII or not to ASCII?
 - Preservation recommendations
 - ICPSR, LOC and others
 - Suggested software for converting & reviewing SPSS files
- Further considerations
 - SPSS Version
 - Researcher feedback
 - Which files do researchers save?
- Other highlights
 - SPSS Tutorials
 - Bibliography for more curation resources

Microsoft Access

Author: Fernando Rios &
Dave Fearon

DCN Mentor: Dave Fearon

<https://github.com/DataCurationNetwork/data-primers>

Key Curatorial Considerations:

What is the complexity of the database?

- Simple DBs (few tables, no forms, queries, macros) could be curated like a spreadsheet

As a base level for preservation:

- Keep original files + export tables to flat CSVs
- Screenshot the Relationships Diagram
- Run the Database Documenter and save the report alongside the DB
- Check for linked tables
- Other objects (SQL, forms, VB)?

Need help from creator

- Table relations, meaning of column names, how data is to be queried

Microsoft Excel

Authors: Ho Jung Yoo, Sandra Sawchuk & Greg Janée

DCN Mentor: Wendy Kozlowski

<https://github.com/DataCurationNetwork/data-primer>

Key Curatorial Considerations:

There are no metadata standards for Microsoft Excel, so detailed documentation from the depositor is encouraged. Documentation should contain info about:

- Context of the original study
- Description of each file
- Description of each worksheet (ideally one table per worksheet)
- Revisions of the data
- Description of each variable in the files

Jupyter Notebooks

Authors: Daina Bouquin,
Matthew Benzing, Sophie
Hou & Lee Wilson

DCN Mentor: Susan Borda

<https://github.com/DataCurationNetwork/data-primers>

Code needs curation too!

Jupyter notebooks contain code, incorporate data, and require different considerations

Different metadata for different situations

- Minimal deposit
- Runnable deposit
- Comprehensive deposit

} Consider repository suitability

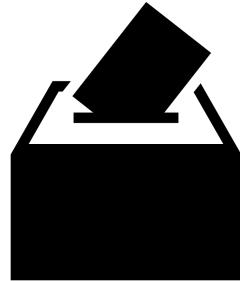
Primer Topic Preview (coming in January 2020)

- Atlas.ti
- Confocal microscopy
- GeoJSON
- Google Docs
- Lidar Point Clouds
- NVivo
- PDF
- R
- .STL files
- Tableau
- Text/character encoding



Workshop #2 at Johns Hopkins University

Share your expertise



Community Authored Data Curation Primers

<https://github.com/DataCurationNetwork/data-primers>

Contribute to these community resources via Github

Thanks to the DataOne Community for making this presentation possible

**DATA
CURATION
NETWORK**

Data Curation R&D

We create and openly share data curation procedures and best practices.

<http://datacurationnetwork.org>

Research Topics of Interest

Importance of Data Curation?

Surveyed researchers at 6 institutions and ARL members about the perceived importance of curation activities

Value of Curation?

Are curated data more usable than non-curated?
Apply a validated metadata quality schema to DCN...

Maturity Models?

(Future) What is the state of research data services among varying types of Libraries?

Sustainability

We will expand into a sustainable entity that grows beyond our initial partner institutions.

<http://datacurationnetwork.org>

Growing the Data Curation Network

- Slow, intentional growth
- Y1-2 Recruit 4 new partners
 - 2019: 2 new partners!
 - 2020: Future
- Expand more broadly in 2021

	Year 1	Year 2	Year 3	Year 4	Year 5	Year 6
Support	Grant Funded (Y1-Y2) transition to partnership model (Y3)		Curation-as-service (Y4-6)			
Timing	2017-19		2020-22		2022-2023	
Phase	Implementation		Transition		Sustaining	
Partners	8 initial partners + 4 more incrementally			Recruit new partners as use and demand dictate		

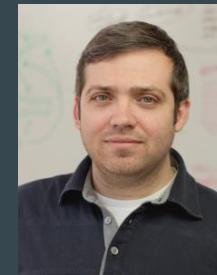
Sustainability Planning

- 2019 Advisory Panel
- Consultant RFP process
- Lyrasis (June 2019-Dec 2019)
 - Market analysis (focus groups, interviews).
 - Administrative Structures (legal support, not-for-profit)
 - Financial Models (in-kind, membership, fee-for-service, or a hybrid)
 - Community Engagement case studies

DCN 2019 Advisory Panel



Yasmeen Shorish,
James Madison



Jeff Spies, 221B



Limor Peer,
Yale University



John Chodacki,
Univ California



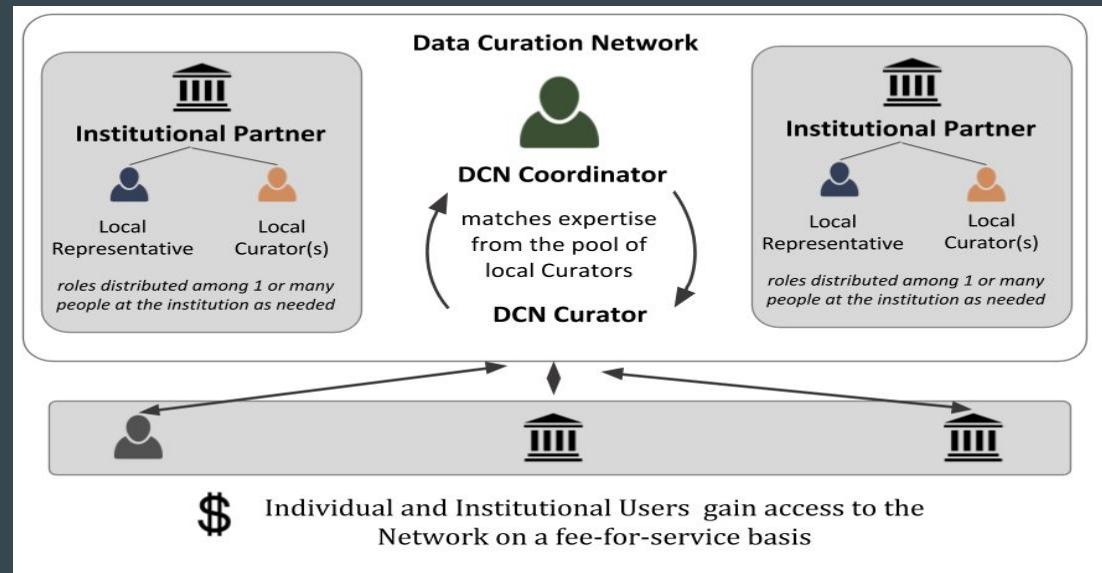
Mike Roy,
Middlebury College



Jay Brodeur,
McMaster College

DRAFT : Hybrid Model

- Opt 1: In Kind
- Opt 2: Fee for service



What is next?

Future directions

- Advocacy?
- Consultation?
- Domain repositories?
- Professional curator community?



We are happy to work together! Data curation without borders!



Thank you

Contact us!

dcn-team@googlegroups.com

<http://datacurationnetwork.org>



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