

IMPROVING DATA CURATION SERVICES WITHIN AN INSTITUTIONAL REPOSITORY

UNC Chapel Hill Libraries
Julie Rudder - Repository Program Librarian
Rebekah Kati - Institutional Repository Librarian



Overview

- Introduction
- Levels of Data Curation Activities Exercise
- Using the Levels of Data Curation at UNC
- Hyrax Audit for Data
- Next Steps



UNIVERSITY
LIBRARIES

[Advanced Search](#)

Carolina Digital Repository

[Home](#)[Collections](#)[Departments](#)[Deposit](#)

Deposit, share, and promote your
scholarly and creative work.

Deposit Your Work



Articles



Student Papers



Datasets



Other Deposits



Hyrax

Hyrax: An open-source, Samvera-powered repository front-end

[Learn more](#)



Community Support

Hyrax is maintained and supported by the [Samvera](#) community. Samvera is an open-source repository solution built



Flexible Workflow

Hyrax supports multiple workflows including the ability to [define custom workflows](#), allowing implementations to



Fully Featured Solution

Hyrax offers a repository solution that can meet the needs of institutional/data repositories and digital object

Why research data curation, why now?

- Support FAIR data, Open Science
- Support faculty goals with grants and data sharing
- OA policy (resources and renewed interest in data)
- New staff

Sure, we do
data curation!



Specific & shared
understanding

Triangle Research Libraries Network (TRLN) Institute

- <https://osf.io/preprints/lissa/zj5pq/>
- Data Curation Team at Duke University
 - **Digital Repository Content Analysts:** Moira Downey, Susan Ivey
 - **Senior Research Data Management Consultants:** Sophia Lafferty-Hess, Jennifer Darragh

Goals for TRLN Institute

- provide more specificity around data curation within our individual contexts
- determine a method to discuss our service model
- identify gaps we would like to fill
- determine what is currently out of scope for our repositories.

ARL SPEC KIT 354

“respondents conflated data curation activities with research data management services...this indicates that a common understanding of data curation is not widespread or ubiquitous”

(Hudson-Vitale et al., 2017)

A Definition for Research Data Curation

“the encompassing work and actions taken by curators of a data repository in order to provide meaningful and enduring access to data.”

Date Curation Network(Johnston et al., 2016).

For full definitions of data curation activities see the [Data Curation Network: Data Curation Terms and Activities](#).

Level 1 Curation	Level 2 Curation	Level 3 Curation
<p>Ingest</p> <ul style="list-style-type: none"> • Authentication • Chain of Custody • Deposit Agreement • Documentation • File Validation • Metadata <p>Appraisal</p> <ul style="list-style-type: none"> • <i>Rights Management (licenses)</i> <p>Curate</p> <ul style="list-style-type: none"> • Arrangement & Description • File Inventory or Manifest • Indexing • Persistent ID <p>Access</p> <ul style="list-style-type: none"> • Contact Information • Data Citation • Discovery Services • Embargo • File Download • Full Text-Indexing • Metadata Brokerage • Terms of Use • Use Analytics • <i>Restricted Access (system automated)</i> <p>Preservation</p> <ul style="list-style-type: none"> • File Audit • Migration • Secure Storage • Versioning • Succession Planning • Tech/Monitoring Refresh • Cease Data Curation 	<p>Appraisal</p> <ul style="list-style-type: none"> • <i>Rights Management (DUAs)</i> • <i>Risk Management (file review)</i> • Selection <p>Curate</p> <ul style="list-style-type: none"> • Contextualize • Curation Log • File Format Transformations • File Renaming • <i>Quality Assurance</i> • Restructure <p>Access</p> <ul style="list-style-type: none"> • <i>Restricted Access (mediated requests)</i> <p>Preservation</p> <ul style="list-style-type: none"> • Repository Certification 	<p>Appraisal</p> <ul style="list-style-type: none"> • <i>Risk Management (remediation)</i> <p>Curate</p> <ul style="list-style-type: none"> • Code Review • Conversion (Analog) • Data Cleaning • De-Identification • Interoperability • Peer Review • <i>Quality Assurance</i> <p>Access</p> <ul style="list-style-type: none"> • Data Visualization

Level 1 (UNC) - Systems and Policies

Ingest: Authenticity, Chain of Custody, Deposit Agreement, Documentation, File Validation, Metadata

Appraisal: *Rights Management (licenses only)

Curate: Arrangement & Description, File Inventory or Manifest, Indexing, Persistent ID, Transcoding

Access: Contact information, Data Citation, Discovery Services, Embargo, File Download, Full Text-Indexing, Metadata Brokerage, Terms of Use, Use Analytics, *Restricted Access (system automated)

Preservation: File Audit, Migration, Secure Storage, Succession Planning, Tech/Monitoring Refresh, Versioning, Cease Data Curation

Level 2 (Duke) - Human intervention, data knowledge

Appraise/Accept: *Rights Management (DUAs), *Risk Management (file review), Selection

Curate: Contextualize, Curation Log, File Format Transformations, File Renaming, *Quality Assurance, Restructure

Access: *Restricted Access (mediated requests)

Preserve: Repository Certification

Level 3: Human intervention, domain-specific, data knowledge

Appraise/Accept: *Risk Management (remediation)

Curate: *Code Review Conversion (Analog), Data Cleaning, De-Identification, Interoperability, Peer Review, Quality Assurance, Software Registry

Access: Data Visualization

Preserve: Emulation

Using the Levels in the CDR

- Feasibility of FAIR with current staffing levels and generalist subject knowledge?
- What is our sweet spot? Where can we add value?
- Is our new system Hyrax good for data? Capabilities?
- How do we get institutional support?
- What training do we need?

What did we do?

- Training
- Created policies
- Created documentation
- Created new workflow
- Assessed existing deposits
- Evaluated Hyrax for data

Training

- Data management online course
- Literature review
- Review of existing services
- Conference attendance
- Workshops

Why a new data policy?

- Formalized and documented service activities
 - Defines activities in Levels
 - Transparency in services
- Large deposit approval took a long time
- Data definition was vague
- Administrative and library leadership buy-in for service and staff time

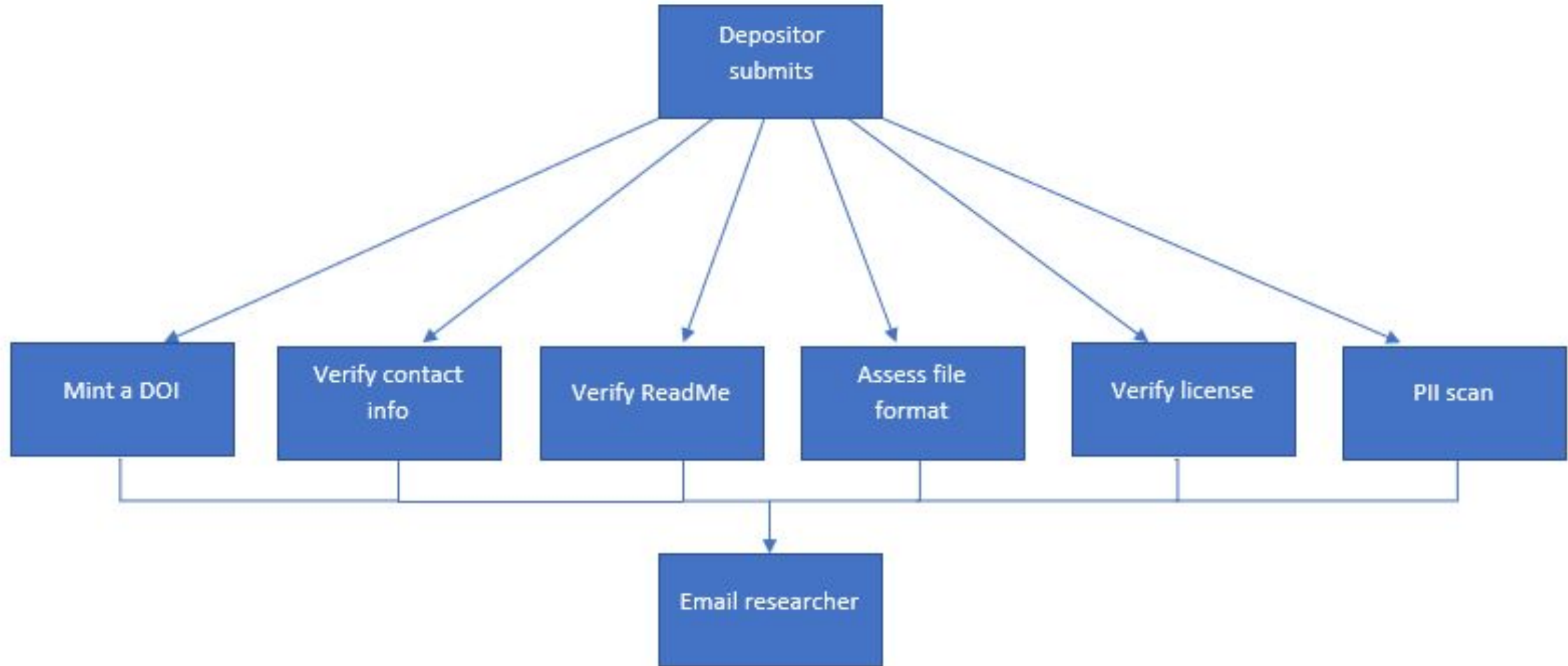
What is in the new data policy?

- Data definition
- Review of large deposit requests is now reviewed by data librarians
- 10 year retention review
- Tombstone record
- Submission review by CDR staff
- ReadMe or other documentation is now required

Documentation

- CDR staff documentation
 - Mediated deposit questions
 - Workflow for self and mediated submission
 - Comparison document contrasting old practice with new policy
- Depositor documentation
 - Data deposit FAQs
 - ReadMe template

Workflow



Existing Data Deposit Assessment Project Goals

- Determine compliance with new policy
- Test scope of service
- Test new workflow
- Target areas for improvement

Results

- Compliance with new policy
 - ReadMe needed
 - Open formats needed
 - DOIs and licenses are good!
- Scope of the service
 - Covers self and mediated deposit, not supplemental data
- Test of new workflow
- Further questions

Where Do We Go Next?

- Locally
 - Launch Hyrax
 - Evaluate level of data support as submissions increase
 - Investigate subject area data curation support
 - Assess Level Two activities to plan future service expansion
- Community
 - Audit
 - Hyrax work to support data

Hyrax Audit for FAIR/Level One

Repository Data Audit Community Template

<https://docs.google.com/spreadsheets/d/1NXyELgm1YVnE4mLtNTyuhPWp3-aSR Cz5SyWnicX8i2M/edit#gid=0>

- Hyrax Capability
- Local Solution Needed
- Policy Needed

[+](#)
[☰](#)
[README](#)
[Repository Data Audit Template](#)
[6 Hyrax-Audit](#)

Repository Data Audit Template

File Edit View Insert Format Data Tools Add-ons Help

All changes saved in Drive

100% £ % .0 .00 123 Arial 10 B I S A

fx

	A	B	C	D	E	F	G	H	I	J	K	
1	Data Curation Network Activity - Level One	FAIR Principle	Hyrax 2.1.x	Local Solution Needed	Policy Needed	Notes						
2	Authentication			x		Shibboleth config docu:						
3	Chain of Custody					This could use a good look.						
4	Deposit Agreement		Deposit Aggrement Included		x							
5	Documentation		Supports Supplemental Files		x							
6	File Validation		Hyrax uses FITS for file characterization									
7	Metadata		Custom WorkTypes; RDF; GeoNames	Data work type creation								
8	Rights Management (licenses)		CC Licenses available	Additional License can be added	x							
9	Arrangement & Description		File ordering, Work Nesting, Collection Building									
10	File Inventory or Manifest		Supports Supplemental Files									
11	Indexing		SOLR; formatted date fields, ??									
12	Persistent Identifier		ActiveFedora::Noid gem to mint Noid-style identifiers	DOIs need to be locally implemented	x							
13	Transcoding		Supported via FFmpeg			May need to be tweaked from default settings						
14	Contact Information		Depositor and Creator information	Customize Worktype to collect additional information	x							
15	Data Citation		MLA, APA, Chicago suggested formatting									
16	Discovery Services		Blacklight, SOLR, SEO									
17	Embargo		Embargo supported		x							
18	File Download		File Download Supported		x							
19	Full Text-Indexing		Runs agains PDFs and office docs uses SOLR (needs to be configured)	could extend service to more types		https://github.com/samvera/hyrax-derivatives/blob/3434afecb720b163415be78d8b371cd20b514a7b/lib/hyrax/derivatives/processes/full_text.rb	https://github.com/samvera/hyrax/blob/21117ad865f62a15668ad59f474764856a56dbdf/app/services/hyrax/file_set_derivatives_service.rb					
20	Metadata Brokerage		OAI-PMH; ResourceSync	Would need to register for any specific registries		OAI-PMH via Blacklight OAI Provider						
21	Restricted Access (system automated)		Robust Access Controls	Any HIPPA or PII restrictions	x							
22	Terms of Use		Licenses are displayed on the view pages		x							
			Hyrax provides support for capturing usage									

+ ☰

README Repository Data Audit Template Hyrax-Audit

Works Cited

FORCE11. *The FAIR Data Principles*. Retrieved June 1, 2018 from: <https://www.force11.org/group/fairgroup/fairprinciples>

Hudson-Vitale, Cynthia; Imker, Heidi; Johnson, Lisa R.; Carlson, Jake; Kozlowski, Wendy; Olendorf, Robert; and Stewart, Claire. (2017). *SPEC Kit 354 Data Curation*. Washington DC: Association of Research Libraries. <https://doi.org/10.29242/spec.354>

Johnson, Lisa R.; Carlson, Jake; Hudson-Vitale, Cynthia; Imker, Heidi; Kozlowski, Wendy; Olendorf, Robert; and Stewart, Claire. (2016, October 23). *Data Curation Network: Data Curation Terms and Activities*. Retrieved from: <http://hdl.handle.net/11299/188638>

Johnson, Lisa R.; Carlson, Jacob; Hudson-Vitale, Cynthia; Imker, Heidi; Kozlowski, Wendy; Olendorf, Robert; and Stewart, Claire. (2018). How Important Are Data Curation Activities to Researchers? Gaps and Opportunities for Academic Libraries. *Journal of Librarianship and Scholarly Communication*, 6(General Issue), eP2198. <https://doi.org/10.7710/2162-3309.2198>

Lafferty-Hess, Sophia; Rudder, Julie; Downey, Moira; Ivey, Susan; and Darragh, Jen. (2018). *Conceptualizing Data Curation Activities Within Two Academic Libraries*. Retrieved from: <https://share.osf.io/preprint/460F1-F92-9F9>