

Curating Big Data for Reuse

Harvard Dataverse

Open Repositories 2023

Katie Mika
Data Services Librarian
Harvard Library & IQSS

Gustavo Durand
Tech Lead / Architect
IQSS

Leonid Andreev
Senior Software Developer
IQSS

Challenges

Volume

Large files

Large number of files

Variety

Diverse data types

Velocity

Instrument data

Streaming data

Veracity

Data quality

Repositories & Big Data

1. Review several projects from Harvard Dataverse that addressed “big data” in some fashion. What challenges did they present & how did we address them?
2. A few updates from Dataverse on how we’re thinking about addressing several of these
3. Lingering questions for the community about how repositories and technical infrastructure could meet these challenges?

Levy Collection

1 to 10 of 28,167 Results

- 60k files
- Researcher wanted unique DOIs
- Each file would be accessed individually by external image viewer
- Emphasized importance of communicating best practices to users

HARVARD Dataverse

Photographs from the Leon Levy Expedition to Ashkelon Collection.

1985 Photographs 1986 Photographs 1987 Photographs 1988 Photographs

Search this dataverse... Advanced Search

1 to 10 of 28,167 Results

Dataverses (28)
Datasets (28,167)
Files (28,167)

Publication Year
2022 (12,050)
2021 (16,117)

Subject
Arts and Humanities (28,167)

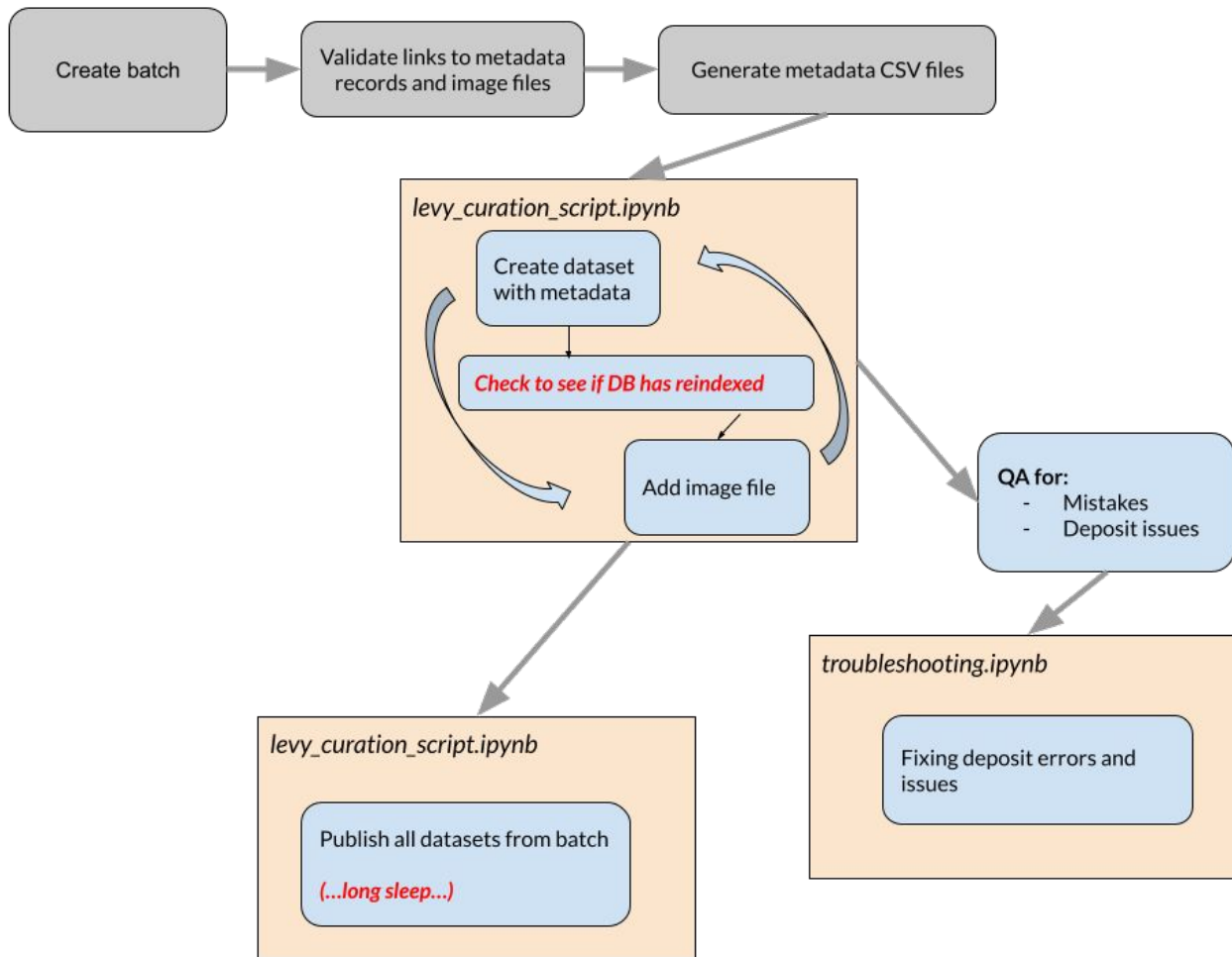
Author Name
Master, Daniel M. (28,167)
Stager, Lawrence E. (28,167)

Keyword Term
Archaeology (28,167)

Geographic Coverage City
Ashkelon (28,167)

Search Results:

- A00_12858.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12858.JPG", <https://doi.org/10.7910/DVN/HNE1R1>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/8f41364b-18bd-9e15-f3e1-7979e41d6fc0>
- A00_12859.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12859.JPG", <https://doi.org/10.7910/DVN/RQDSH2>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/d7ef688d-0db8-f6a1-38e1-badf553bcb6c>
- A00_12860.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12860.JPG", <https://doi.org/10.7910/DVN/E8VIMR>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/392f82ad-f6e2-1c8a-4891-bf22aea2d9b5>
- A00_12861.JPG



Levy Collection

1 to 10 of 28,167 Results

- 60k files
- Researcher wanted unique DOIs
- Each file would be accessed individually by external image viewer
- **Emphasized importance of communicating best practices to users**

HARVARD Dataverse

Photographs from the Leon Levy Expedition to Ashkelon Collection.

1985 Photographs 1986 Photographs 1987 Photographs 1988 Photographs

Search this dataverse... Advanced Search

1 to 10 of 28,167 Results

Datasets (28) Files (28,167)

Publication Year: 2022 (12,050), 2021 (16,117)

Subject: Arts and Humanities (28,167)

Author Name: Master, Daniel M. (28,167), Stager, Lawrence E. (28,167)

Keyword Term: Archaeology (28,167)

Geographic Coverage City: Ashkelon (28,167)

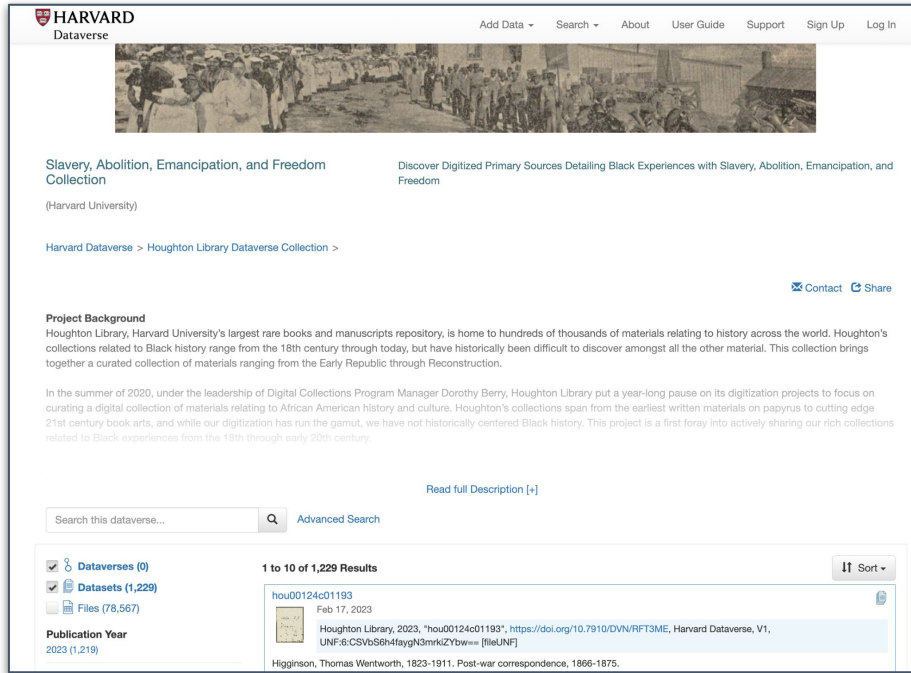
A00_12858.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12858.JPG", <https://doi.org/10.7910/DVN/HNE1R1>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/8f41364b-18bd-9e15-f3c1-7979e41d6fc0>

A00_12859.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12859.JPG", <https://doi.org/10.7910/DVN/RQDSH2>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/d7ef688d-0db8-f6a1-38e1-badf553bcb6c>

A00_12860.JPG
Oct 13, 2021 - 2000 Photographs
Master, Daniel M.; Stager, Lawrence E., 2021, "A00_12860.JPG", <https://doi.org/10.7910/DVN/E8VIMR>, Harvard Dataverse, V1
Link to OCHRE database: <http://p.lib.uchicago.edu/1001/org/ochre/392f82ad-f6e2-1c8a-4891-bf22aea2d9b5>

A00_12861.JPG

SAEF & Historic Datasets



The screenshot shows the Harvard Dataverse interface for the 'Slavery, Abolition, and Freedom Collection'. At the top, there is a navigation bar with 'Add Data', 'Search', 'About', 'User Guide', 'Support', 'Sign Up', and 'Log In'. Below the navigation bar is a large historical photograph of a group of people, likely enslaved individuals, in a field. The main heading is 'Slavery, Abolition, and Freedom Collection' with a sub-heading 'Discover Digitized Primary Sources Detailing Black Experiences with Slavery, Abolition, Emancipation, and Freedom'. The page includes a 'Project Background' section with text about Houghton Library's digitization efforts and a search bar at the bottom with a search button and a 'Sort' dropdown menu. The search results show 1,229 datasets and 78,567 files, with a specific result for 'hou00124c01193' by Houghton Library, 2023.

- Large number of files
- Diversity of content: images, XML, json, text, csv, documentation, code
- How to curate data and aux files for diverse reuse cases?
- How do you describe and make files findable? Browseable?

Dataset Metadata

Citation Metadata

Dataset Persistent ID doi:10.70122/FK2/9XXPIX

Title hou00201c00195

Author Houghton Library (Harvard University)

Contact Use email button above to contact.
Jacobson, Christine (Houghton Library)

Description [Unidentified author]. Rev. Mr. Broadwater (says) Every child is allowed to go to the government school, but the colored people can have separate schools [first line] : Ms (in unidentified hand), [no place], 1863?

Subject Arts and Humanities

Kind of Data Institutional papers

Origin of Sources URN-3:FHCL.HOUGH:100528159

Geospatial Metadata

Geographic Coverage Tuskegee
Alabama
Canada

SAEF Metadata

MMD ID 990091469160203941

Created 1863

Person/Org Tags Broadwater, Rev. [?]

Theme Education; Freedmen

Genre Government documents

Custom metadata block supports discovery by project metadata

File Metadata & Tags

 [hou00201c00009_0003.jpg](#)
 JPEG Image - 776.7 KB
 Deposited Jun 9, 2022
 MD5: a40...e00

File associated with: Balch, F. V. (Francis Vergnies), 1839-1898. Letters to the Commission, Washington, D.C., 1864 Jan. 11 and 14. Origin of source: <https://nrs.harvard.edu/URN-3:FHCL.HOUGH:100526440>

City:Washington **Country:United States** **Created:1864** **Data** **Genre:Correspondence**

Person/Org:Balch, F. V. (Francis Vergnies) **Physical Format:Institutional papers**

State:District of Columbia **Theme:Freedmen** **UID:h00201c00009**

Auto-Generated Documentation: Relationship Files

	filename_source	filename_target	relationships	source_file_format	target_file_format	relationships
1	hou00201c00009_0001.jpg	hou00201c00009_mets.xml	belongs_to	JPEG 2000 JP2	Extensible Markup Language	
2	hou00201c00009_mets.xml	hou00201c00009_0001.jpg		Extensible Markup Language	JPEG 2000 JP2	contains
3	hou00201c00009_0002.jpg	hou00201c00009_mets.xml	belongs_to	JPEG 2000 JP2	Extensible Markup Language	
4	hou00201c00009_mets.xml	hou00201c00009_0002.jpg		Extensible Markup Language	JPEG 2000 JP2	contains
5	hou00201c00009_0003.jpg	hou00201c00009_mets.xml	belongs_to	JPEG 2000 JP2	Extensible Markup Language	
6	hou00201c00009_mets.xml	hou00201c00009_0003.jpg		Extensible Markup Language	JPEG 2000 JP2	contains
7	hou00201c00009_0004.jpg	hou00201c00009_mets.xml	belongs_to	JPEG 2000 JP2	Extensible Markup Language	
8	hou00201c00009_mets.xml	hou00201c00009_0004.jpg		Extensible Markup Language	JPEG 2000 JP2	contains

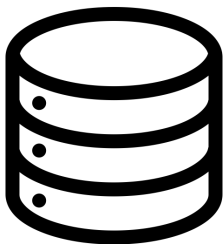
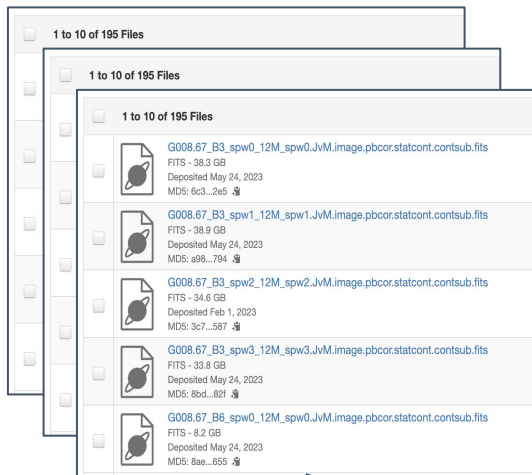
HMS Reich Lab

	Dataset Version	Summary	Contributors	Published on
<input type="checkbox"/>	8.0	Citation Metadata: Description (1 Changed); Files (Added: 13; Removed: 12); View Details	Katherine Mika, Swapan Mallick	2023-04-05
<input type="checkbox"/>	7.0	Files (Added: 12; Removed: 12); View Details	Swapan Mallick	2023-04-05
<input type="checkbox"/>	6.0	Citation Metadata: Description (1 Changed); Files (Added: 12; Removed: 12); View Details	Swapan Mallick	2023-04-05
<input type="checkbox"/>	5.0	Files (Added: 12; Removed: 12); View Details	Swapan Mallick	2023-04-05
<input type="checkbox"/>	4.0	Files (Added: 12; Removed: 13); View Details	Swapan Mallick	2023-04-05
<input type="checkbox"/>	3.0	Citation Metadata: Author (1 Added); Title (Changed); Files (Added: 13; Removed: 9); View Details	Katherine Mika, Swapan Mallick	2023-04-05
<input type="checkbox"/>	2.0	Files (Added: 9; Removed: 9); View Details	Katherine Mika, Swapan Mallick	2023-04-05
<input type="checkbox"/>	1.0	This is the first published version.	Swapan Mallick, David Reich, Katherine Mika	2023-04-01

server

- Hack versioning?

ALMA-IMF Large Program dataset



- Secured approval for very large files, but will require post-publication updates
- Versioning retains old files & explodes storage
- Do we keep files from old versions?

Solutions!

DVUploader – batch file uploads

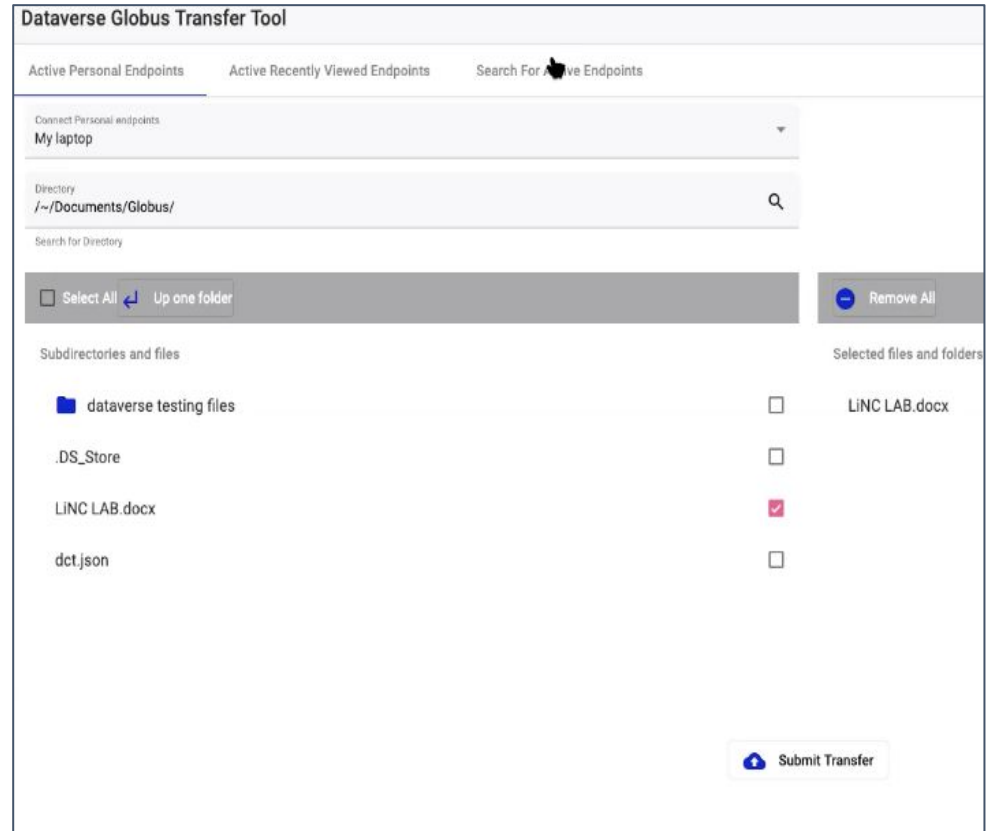
A Java program that uses the Dataverse API to upload a file(s)/a directory tree of files to Dataverse:

```
java -jar DVUploader-v1.1.0.jar -key=<api key> -did=<dataset doi> -server=<server URL> <dir/file(s)>
```

- Many options including:
 - -recurse – go into subdirectories
 - -ex=<pat> - ignore files with the specified pattern
 - -verify – check fixity
 - -limit=<x> – max number to process
 - -directupload – use new direct-upload-to-S3-capability
- Useful For:
 - Automating uploads
 - Many files
 - Incremental uploads (will upload new files in a dir)
- Upcoming
 - proof-of-concept code to re-create a dataset from an archival bag

Globus Integration

- Ability to add files (or just file metadata, leaving file at source) to Dataverse
- Developed as an external tool to be integrated into the Dataverse upload workflow



TRSA

- Trusted Remote Storage Agents
 - Agent - Dataverse can communicate with this
 - Storage - especially for sensitive or big data
 - Remote - Dataverse does not control access
 - Trusted - service agreement guarantees

Signed URLs for secure API access from External Tools

The Dataverse External Tool mechanism has been sending the user's API key to the tool

- Allows the tool to read/write on the user's behalf, but
- API key is visible in browser and using tools on draft datasets shouldn't be done on public machines
- Tool is trusted to only do the intended operations and only access the intended dataset/datafile

Signed URLs - a much more secure alternative allowing the tool to make specific API calls for a limited time

- URLs are cryptographically signed and validated when used - will be rejected if modified or used after the configured time period
- Allowed URLs and valid durations (usually seconds to hours) are defined as part of configuring the tool (i.e. admins control what's allowed)
- URLs are also limited by the permissions of the user running the tool
- Tools can be configured to avoid having any URL visible in the browser, or be limited to having one very short-lived URL (e.g. 30 seconds)

SPA Re-architecture Project

- Single Page Application (frontend)
 - **modernize** the application
 - **separate backend and frontend** to increase **interoperability**
 - Dataverse backend becomes an **API-first application**
 - **extend modularization** of backend and frontend
 - **speed up development** and implementation of new UI/UX ideas
 - Native **accessibility** (A11y) and **internationalization** (i18n) support
- **empower the community**

*Persistent problems
require community
solutions & reimagining
curation best practices*

Community frameworks

1. What is a realistic, sustainable storage limit for institutional repositories? What level/tier do we set limits at?
2. How long should we preserve Big Data?
3. What does stricter appraisal look like?
4. What are the ongoing stewardship requirements?
5. How does data management planning differ?
6. How to communicate better with researchers about the realities of sharing and preserving/archiving Big Data

Thanks!

katherine_mika@harvard.edu

gdurand@iq.harvard.edu