



Intermediate Curation on the Dataverse Supported Repository









Code of Conduct



Our Pledge

We as members, contributors, and leaders pledge to make participation in our community a harassment-free experience for everyone, regardless of age, body size, visible or invisible disability, ethnicity, sex characteristics, gender identity and expression, level of experience, education, socio-economic status, nationality, personal appearance, race, caste, color, religion, or sexual identity and orientation.

Our Standards

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

- Demonstrating empathy and kindness toward other people
- Being respectful of differing viewpoints and experiences
- Being respectful of differing opinions, viewpoints, and experiences
- Giving and gracefully accepting constructive feedback
- Accepting responsibility and apologizing to those affected by our mistakes, and learning from the experience
- Focusing on what is best not just for us as individuals, but for the overall community

Examples of unacceptable behavior by participants include:

- The use of sexualized language or imagery, and sexual attention or advances of any kind
- Trolling, insulting or derogatory comments, and personal or political attacks
- Public or private harassment
- Publishing others' private information, such as a physical or email address, without their explicit permission
- Other conduct which could reasonably be considered inappropriate in a professional setting





Sonia Barbosa

Associate Director of Dataverse Support, Data Curation, and The Murray Archive



Agenda

Intermediate Data Curation

Review of Fundamentals:

Recap of Data Curation Fundamentals

Curation Features and Tools on the Dataverse Platform

- Features/Tools: Collections, Datasets, Files
- Use Cases and Examples
- Q&A Session



Review of Fundamentals in Data Curation



Introduction to Data Curation

The Importance of Curation:

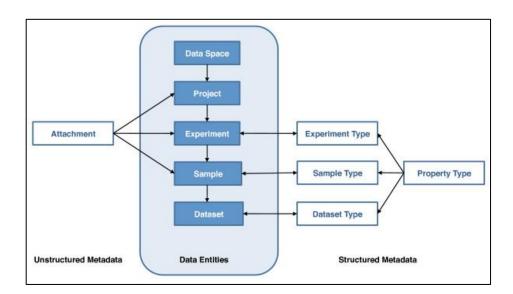
- Enhance Data Quality
- Efficient Organization
- Advanced Analysis
- Collaborative Workflows
- Ethical Compliance
- Future-proofing Data

Basic Data Curation Principles:

- Data Identification
- Data Collection
- Data Cleaning
- Data Organization
- Data Storage
- Metadata Management
- Data Preservation
- Data Sharing
- Documentation
- Compliance and Ethics

When to Document Data

https://managing-qualitative-data.org/modules/2/a/ © Social Science Research Council and licensed under a Creative Commons, Attribution Share-Alike (CC-BY-SA) license.



OpenBIS: A flexible framework for managing and analyzing complex data in biology research - Scientific Figure on ResearchGate. Available from:

https://www.researchgate.net/figure/Data-organization-and-metadata-Data-are-organized-using-entities-and-relations-that-are_fig7_51861855 [accessed 3 Dec, 2023]



Understanding Data Quality:

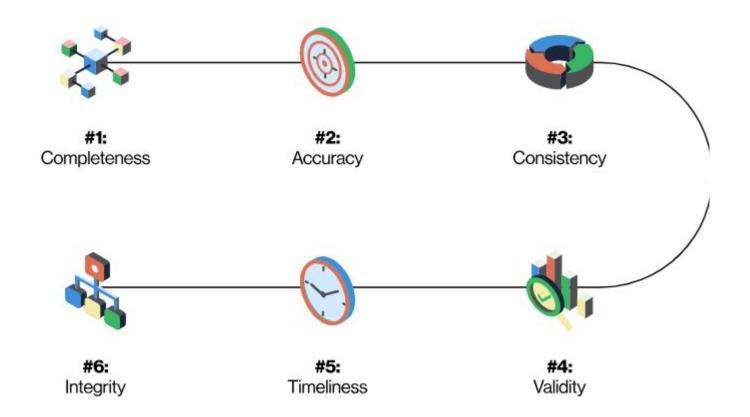
Assessing Data Quality, Metrics:

- Data Accuracy
- Data Completeness
- Data Consistency
- Data Validity
- Timeliness
- Relevance
- Usability
- Duplication
- Privacy/Security

Data Cleaning Techniques:

- Handle Missing Values
- Detect outliers and Correct
- Normalization and Standardization
- Deduplication of Data
- Parsing and Formatting
- Error Correction
- Handle Inconsistencies





https://www.cloverdx.com/blog/6-data-quality-metrics-you-cant-ignore



| a. Emple | oyees | | b. Emple | oyees | |
|----------|-------|--------|----------------------------------|-----------|-------------|
| name | age | salary | name | age | salary |
| Paul | 1978 | NULL | Paul | ? | 29,000 |
| Paul | NULL | 29,000 | Melanie | 1990 | NULL |
| Paul | 1979 | NULL | Bob | 1977 | 37,000 |
| Melanie | 1990 | NULL | Charlie | 1978 | 32,000 |
| Bob | NULL | 37,000 | | | |
| Bob | 1977 | NULL | Paul.ag | je = 1978 | 8 OR 1979 ? |
| Charlie | 1978 | 32,000 | Taxon socrator socrator socrator | | |

Figure 3 from "++Spicy: an Open-Source Tool for Second-Generation Schema Mapping and Data Exchange".



The second image shows the same x-ray image with all sample PII and PHI metadata removed or obscured



"After de-identification, all image metadata is removed, and all text burned into the image is obscured with an opaque rectangle. This configuration of de-identification is useful for when you need only the image pixel data for further analysis, machine learning (ML) model training, or inference."

https://cloud.google.com/architecture/de-identification-of-medical-images-through-the-cloud-healthcare-api



Data Organization:

Taxonomies and Ontologies:

Taxonomies

Helpful in organizing data (organisms) into hierarchical structures, easing navigation and retrieval of data

Example:

Field - Biology

Species Taxonomy:

Kingdoms

Phyla

Classes

Orders

Families

Genera

Species

Ontologies

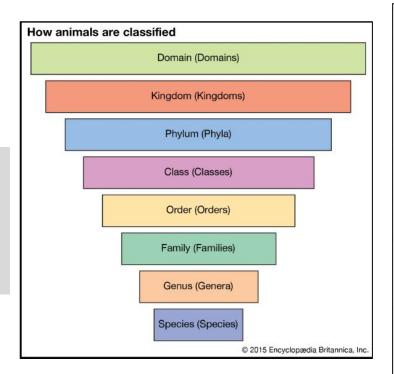
Enable a complex representation of knowledge. Allows inferences and reasoning with a domain

Example:

Field: Medical Field

Diabetes Mellitus Diagnosis Ontology (DDO) is an ontology for diagnosis of diabetes, containing diabetes-related complications, symptoms, drugs, lab tests, etc. The Adverse Event Reporting Ontology (AERO) is aimed at increasing quality and accuracy of reported adverse clinical events





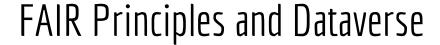
Taxonomies and Ontologies



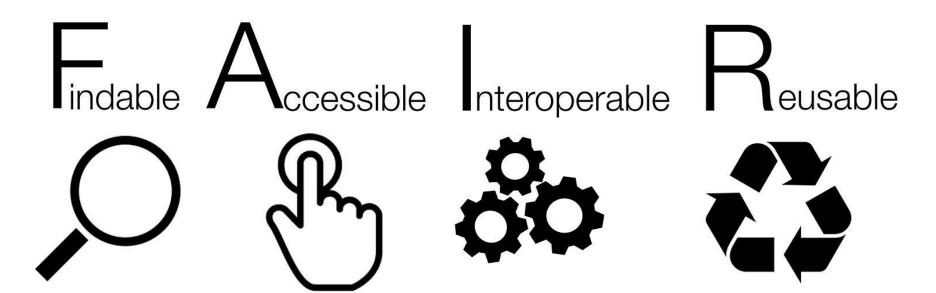
https://www.arabidopsis.org/about/nrg1295.pdf



Curation Features and Tools on the Dataverse Platform







https://en.wikipedia.org/wiki/File:FAIR_data_principles.jpg



Collection-level



Current Features



Dataverse Collections

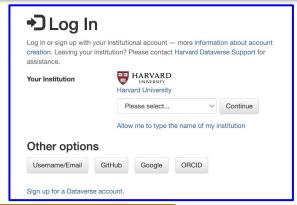
- Own administration
- Own branding (and can be embedded anywhere)

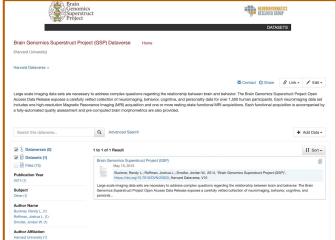
Datasets

- Citation
- Metadata
- Versioning
- Private URL
- Custom Terms/Multiple License/Permissions
- Guestbooks
- Publishing Workflows



- Citation
- Ingest
- Preview/Explore
- Metadata
- Versioning
- Permissions/Embargo/Re strictions





| 8 | Brain Genomics Superstruct Project (GSP) Dataverse (Harvard University) GSP |
|-------|---|
| Gene | ral Information |
| Them | e + Widgets |
| Permi | ssions |
| Group | os |
| Datas | et Templates |
| Datas | et Guestbooks |
| Featu | red Dataverses |



Dataset-level



Current Features



Dataverse Collections

- Own administration
- Own branding (and can be embedded anywhere)

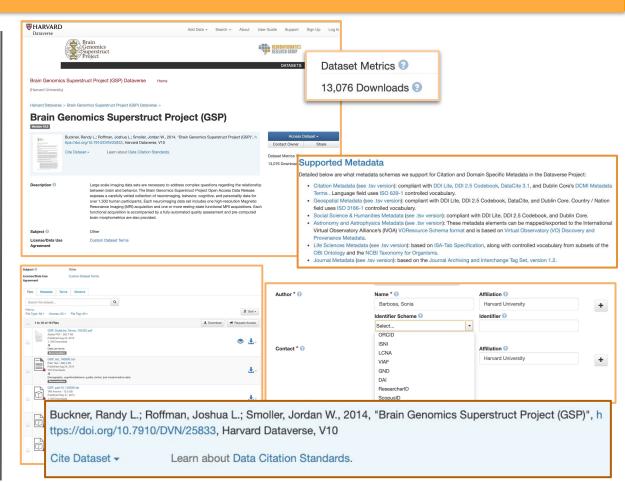


Datasets

- Citation
- Metadata
- Versioning
- Private URL
- Custom Terms/Multiple License/Permissions
- Guestbooks
- Publishing Workflows



- Citation
- Ingest
- Preview/Explore
- Metadata
- Versioning
- Permissions/Embargo/Re strictions





File-level



Current Features



Dataverse Collections

- Own administration
- Own branding (and can be embedded anywhere)

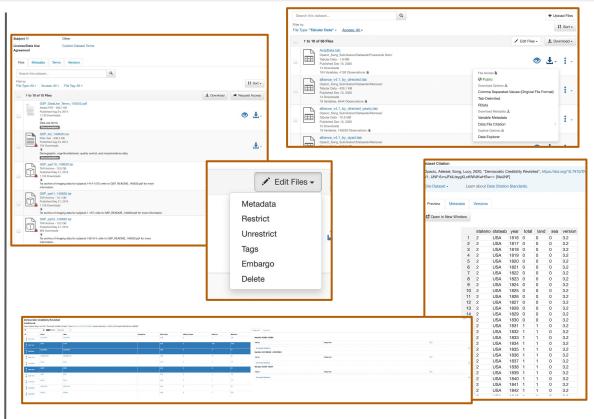
Datasets



- Citation
- Metadata
- Versioning
- Private URL/*Anonymous Peer Review
- Custom Terms/*Multiple Licenses/Permissions
- Guestbooks
- Publishing Workflows



- Citation
- Ingest
- Preview/Explore
- Metadata/Provenance
- Versioning
- Permissions/Embargo/Re strictions

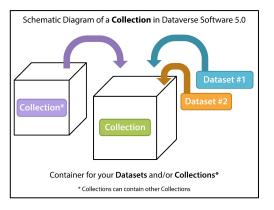


Review of Basics in Dataverse Curation

Best practices:

- Request a <u>COLLECTION</u>
- Deposit <u>DATASETS</u> within COLLECTION
- Upload <u>FILES</u> (data, code, documentation) to DATASETS

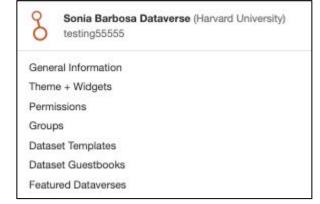




Collections

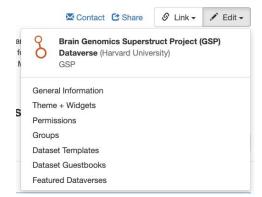
Best practices for "COLLECTION" creation:

- Collection name
- Identifier
- Category
- Collection level contact
- Collection Description
- Metadata selection
 - Custom metadata?
- Facet selections



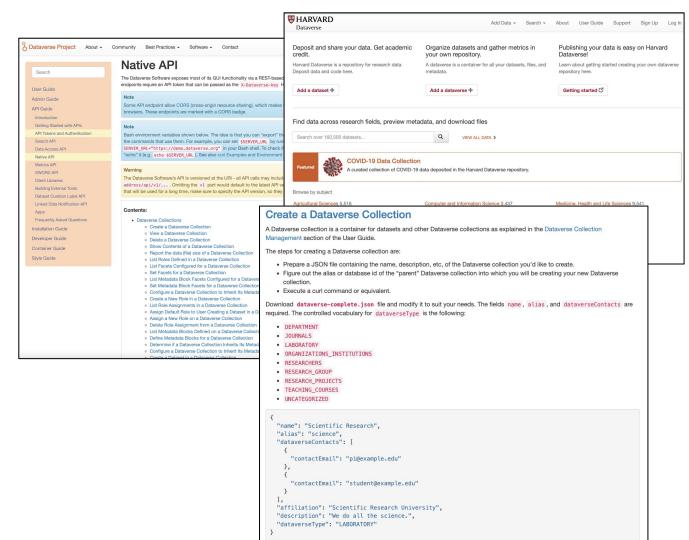
While some of this is required and can't be skipped, not enabling or using these other options results in underutilization of the features available to you

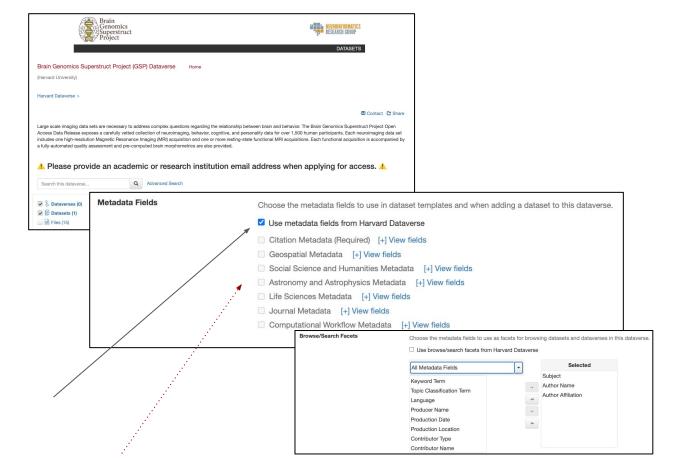




Example collection edit options from <u>Brain</u> <u>Genomics</u> Collection, Harvard Dataverse

Intermediate Options in:



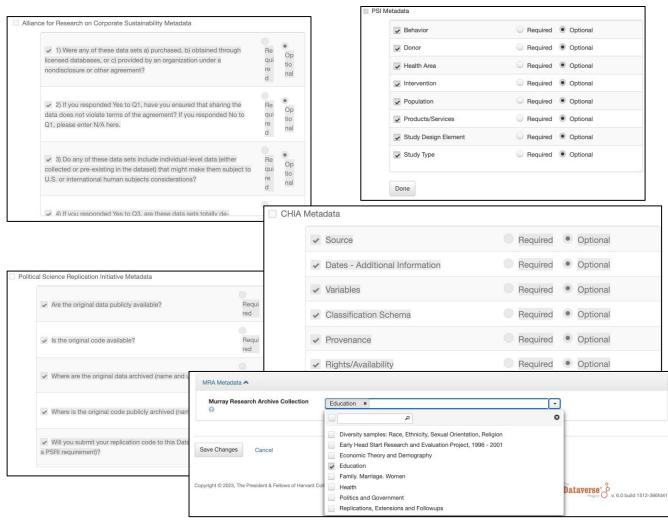






(As seen in repository)

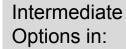
Intermediate Options in:

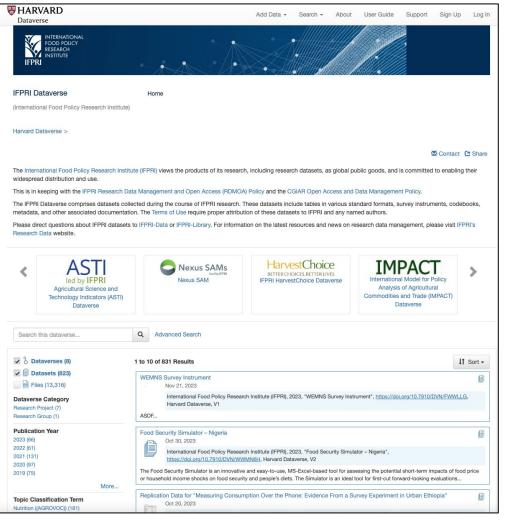




Many sols on smallholder farms in Malawi have poor soil organic matter content. This results in poor make productivity when sufficient mineral

ligars are not added. Building soil organic matter requires improving both cereal and legume-crops primary productivity through mi





Datasets

Best practices for "Dataset" creation:

- Gather all data and documentation to be shared
- Know what license you will use
- Data are deidentified
- All pre sharing requirements have been fulfilled (see data quality above)
- Determine what metadata you will populate to describe the dataset, some are required by default

While some of this is required and can't be skipped, not enabling or using these other options results in underutilization of the features available to you



https://guides.dataverse.org/en/latest/api/ native-api.html#create-a-dataset-in-a-dat averse-collection

Intermediate Options in:

Datasets

Create a Dataset in a Dataverse Collection

A dataset is a container for files as explained in the Dataset + File Management section of the User Guide.

To create a dataset, you must supply a JSON file that contains at least the following required metadata fields:

- Title
- Author Name
- · Point of Contact Email
- Description Text
- Subject

Submit Incomplete Dataset

Note: This feature requires dataverse.api.allow-incomplete-metadata to be enabled and your statasetValid field.

Providing a .../datasets?doNotValidate=true query parameter turns off the validation of Name" is required. For example, a minimal JSON file would look like this:

```
"datasetVersion": {
  "metadataBlocks": {
    "citation": {
      "fields": [
          "value": [
              "authorName": {
                "value": "Finch, Fiona",
                "typeClass": "primitive",
                "multiple": false.
                "typeName": "authorName"
          "typeClass": "compound".
          "multiple": true,
          "typeName": "author"
      "displayName": "Citation Metadata"
```

The following is an example HTTP call with deactivated validation:

Note: You may learn about an instance's support for deposition of incomplete datasets via Show Support Of Incomplete Metadata Deposition.

Submit Dataset

As a starting point, you can download dataset-finchl.json and modify it to meet your needs. (dataset-finchl_fr.json is a variant of this file that includes setting the metadata language (see :Metadatal.anguages) to French (fit). In addition to this minimal example, you can download dataset-create-new-all-default-fields.json which populates all of the metadata fields that ship with a Dataverse installation.)

The curl command below assumes you have kept the name "dataset-finch1.json" and that this file is in your current working directory.

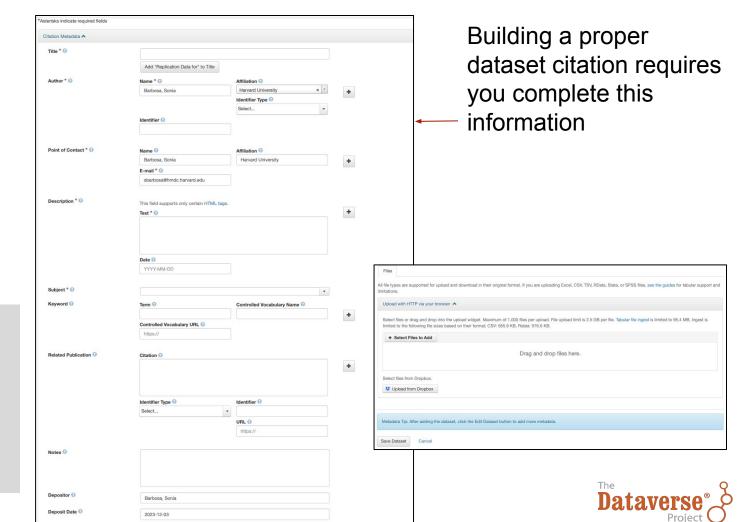
Next you need to figure out the alias or database id of the "parent" Dataverse collection into which you will be creating your new dataset. Out of the box the top level Dataverse collection has an alias of "root" and a database id of "1" but your installation may vary. The easiest way to determine the alias of your root Dataverse collection is to click "Advanced Search" and look at the URL. You may also choose a parent Dataverse collection under the root Dataverse collection.

Note

See curl Examples and Environment Variables if you are unfamiliar with the use of export below.

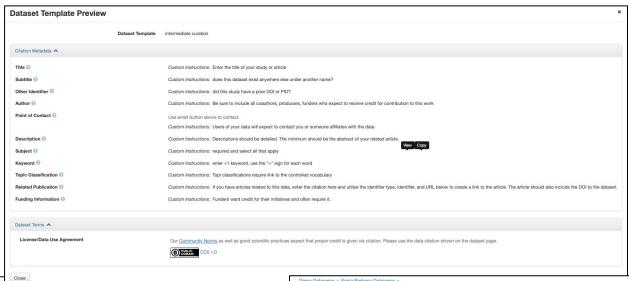


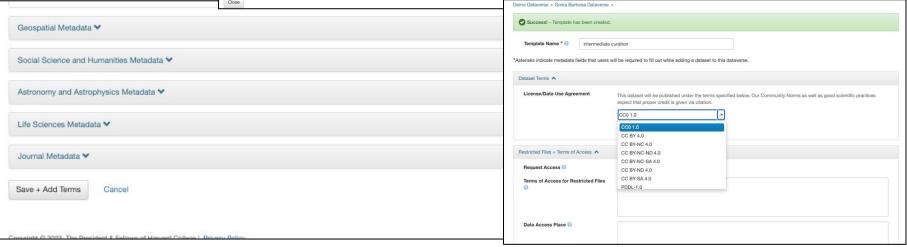
The following is an example HTTP call with deactivated validation:

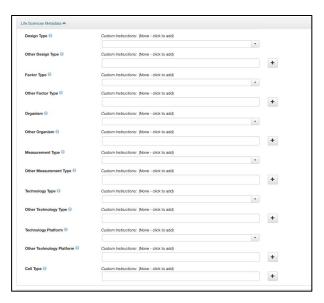


Datasets

Intermediate Options in: Datasets







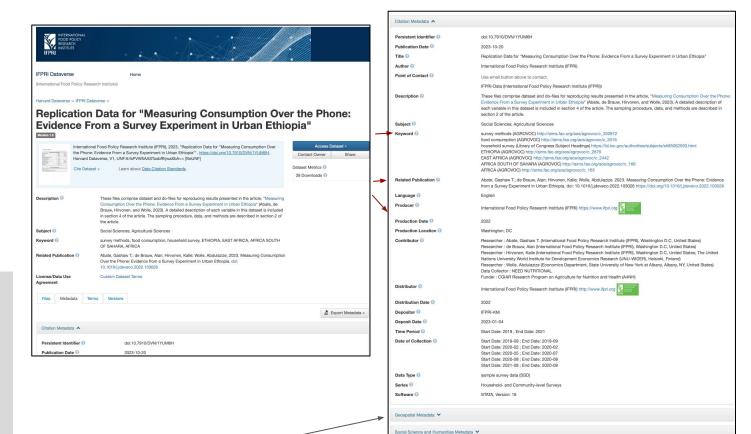
| Unit of Analysis 🕤 | Custom Instructions: (None - click to add) | | | | |
|------------------------------------|--|------------------------------------|---|--------------------------------|---|
| | | | + | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| Universe ① | Custom Instructions: (None - click to add) | | | | |
| | | Astronomy and Astrophysics Metadal | 10 | | |
| | | Type O | Custom Instructions: (None - click to add | | |
| | | | | | |
| | | Facility () | Custom Instructions: (Nove - click to add | | |
| | | | | | + |
| Time Method (2) | Custom Instructions: (None - click to add) | Instrument 0 | Custom instructions: (None - click to add | | |
| | | | | | + |
| | No. of the same of | Object © | Custom Instructions: (None - click to add | | |
| Data Collector O | Custom Instructions: (None - click to add) | | | | + |
| | FamilyName, GivenName or Organization | Spatial Resolution © | Custom instructions: (None - click to add | | |
| Collector Training () | Custom Instructions: (None - click to add) | | | | |
| Collector training 0 | Castori instruccions. (None - Circk to abd) | Spectral Resolution () | Custom instructions: (None - click to add | | |
| | | | | | |
| Frequency () | Custom Instructions: (None - click to add) | Time Resolution @ | Custom instructions: (None - click to add | | |
| | | | | | |
| | | Bandpass () | Custom instructions: (None - click to add | | |
| Sampling Procedure | Custom Instructions: (None - click to add) | | | | + |
| | | Central Wavelength (m) © | Custom instructions: (None - click to add | | |
| | | | Enter a floating-point number. | | + |
| | | Wavelength Range () | Custom Instructions: (None - click to add | | |
| | | mareningui narige U | Minimum (m) (iii | Maximum (m) O | |
| | | | Enter a floating-point number. | Enter a floating-point number. | 4 |
| Target Sample Size ① | Custom Instructions: (None - click to add) | | | | |
| | Actual © Fe | Dataset Date Range () | Custom instructions: (None - click to add Start () | End ① | |
| | Enter an integer | | YYYY-MM-DD | YYYY-MM-00 | 4 |
| | | | | | |
| Major Deviations for Sample Design | Custom Instructions: (None - click to add) | Sky Coverage O | Custom instructions: (None - click to add | | |
| 0 | | | | | + |
| | | Depth Coverage 0 | Custom Instructions: (None - click to add | | |
| Collection Mode | Custom Instructions: (None - click to add) | | Enter a floating-point number. | | |
| | | | + | | |

Intermediate Options in: Datasets

| Journal 0 | Custom Instructions: (None - click to add) | |
|--------------------|--|---------|
| | Volume 🕤 | Issue 🕢 |
| | Publication Date | |
| | YYYY or YYYY-MM or YYYY-MM-DD | |
| Type of Article () | Custom Instructions: (None - click to add) | |
| 7,000 | Select | |

| Geographic Coverage 🕢 | Custom Instructions: (None - click to add) | | |
|-------------------------|---|------------------|---|
| | Country / Nation 🕟 | State / Province | |
| | Select • | | + |
| | | City 0 | |
| | Other ① | | |
| Geographic Unit 💿 | Custom Instructions: (None - click to add) | | |
| | | | + |
| Geographic Bounding Box | Custom Instructions: (None - click to add) West Longitude | East Longitude | |
| | | | + |
| | | | |
| | North Latitude | South Latitude ① | |
| | | | + |

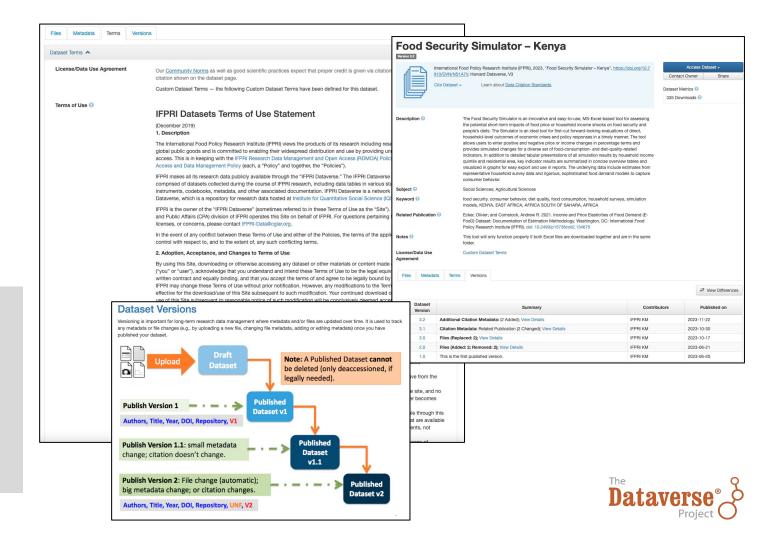
Datasets

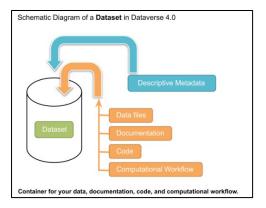


Multiple metadata blocks



Datasets





Files

Best practices for "Files" creation:

- File types and sizes
 - o <u>Compressed</u>
 - <u>Tabular</u>
- Know what license to use
- Data are deidentified
- All pre sharing requirements have been fulfilled (see data quality above)
- File level descriptions
- Know what additional features are available to support multi file uploads, folders

While some of this is required and can't be skipped, not enabling or using these other options results in underutilization of the features available to you

- File Handling
 - File Previews
 - Tabular Data Files
 - Research Code
 - Computational Workflow
 - Computational Workflow Definition
 - FAIR Computational Workflow
 - How to Create a Computational Workflow
 - How to Upload Your Computational Workflow
 - How to Describe Your
 Computational
 Workflow
 - How to Search for Computational Workflows
 - Astronomy (FITS)
 - GeoJSON
 - GeoTIFF
 - Shapefile
 - NetCDF and HDF5
 - H5Web Previewer
 - NcML
 - Geospatial Bounding Box
 - Compressed Files
 - Other File Types

When adding a file to a dataset, you can optionally specify the following:

- A description of the file.
- The "File Path" of the file, indicating which folder the file should be uploaded to within the dataset.
- Whether or not the file is restricted.
- Whether or not the file skips tabular ingest. If the tabIngest parameter is not specified, it defaults to true.

Note that when a Dataverse installation is configured to use S3 storage with direct upload enabled, there is API support to send a file directly to S3. This is more complex and is described in the Direct DataFile Upload/Replace API guide.

Intermediate Options in:



Add a File to a Dataset

When adding a file to a dataset, you can optionally specify the following:

- · A description of the file.
- . The "File Path" of the file, indicating which folder the file should be uploaded to within the dataset.
- · Whether or not the file is restricted.
- . Whether or not the file skips tabular ingest. If the tabIngest parameter is not specified, it defaults to true.

Note that when a Dataverse installation is configured to use S3 storage with direct upload enabled, there is API support to send a file directly to S3. This is more complex and is described in the Direct DataFile Upload/Replace API quide.

In the curl example below, all of the above are specified but they are optional,

Note

See curl Examples and Environment Variables if you are unfamiliar with the use of export below.

curl -H "X-Dataverse-key:\$API_TOKEN" -X POST -F "file=@SFILENME" -F 'jsonData=("description":"My descrip
tion.","directoryLabel":"data;subdir1","categories":["Data"], "restrict":"false", "tablngest":"false")'
"\$SERVER URL/api/datasets/:persistentId/add?persistentId=\$PERSISTENT ID"

The fully expanded example above (without environment variables) looks like this:

You should expect a 201 ("CREATED") response and JSON indicating the database id that has been assigned to your newly uploaded file

Please note that it's possible to "trick" a Dataverse installation into giving a file a content type (MIME type) of your choosing. For example, you can make a text file be treated like a video file with - f' file=@REAME.txt; type=video/mpeg4', for example. If the Dataverse installation does not properly detect a file type, specifying the content type via API like this a potential workaround.

The curl syntax above to upload a file is tricky and a Python version is provided below. (Please note that it depends on libraries such as "requests" that you may need to install but this task is out of scope for this guide.) Here are some parameters you can set in the script:

- dataverse_server e.g. https://demo.dataverse.org
- api_key See the top of this document for a description
- persistentId Example: doi:10.5072/FK2/6XACVA
- · dataset id Database id of the dataset

In practice, you only need one the dataset_id or the persistentId. The example below shows both uses.





Files

Uploading Files

To prepare, log in to Dataverse and:

- · find the DOI for the dataset you wish to add files to, and
- find or generate an API key for yourself in the Dataverse instance you are using (from the popup menu under your profile).

The simplest way to run the DVUploader is to place the jar file into the directory containing a subdirectory with the files intended for upload. (The DVUploader can be placed anywhere on disk and can upload files from any directory, but this requires adding these paths to the command lime and/or configuration of Java's classpath.)

REQUIRED: Run the jar with the following command line:

java -jar DVUploader-v1.2.0beta3.jar -key=<api key> -did=<dataset doi> -server=<server URL> <dir or file names> where

<apikey> is replaced with the API Key generated by the user in Dataverse

<dataset doi> is replaced with the DOI of the target Dataset

<serverURL> is replaced by the URL of the Dataverse server being used (with no trailing '/' and do not include any path to a specific Dataverse on the server), and

<dir or file names> is replaced by the name of a directory and/or a list of individual files to upload.

These four arguments are always required. There are additional options listed below. **Note: **For a first test, adding - listonly is useful - it will make the DVUploader list what it would do, but will not perform any uploads.

For example, Java – Jar DWgloader-v1.10.4 Jar – key=5599882-559-496r-4232-2380886765c – did-dis-118.9372/FZ/TUNBWF = serve-in-tips://dastwers-cs.tol.org testifur would upload all of the files in the 'testifur' directory (relative to the current directory where the java command is run) to the Dataset at https://dataverse.dull.org/dataset.html/presistentied-150.5072/FX/TUNBWF (if it existed: the dataset in this example is not

The output from the DVUploader looks like:

Dataverse Mode: Uploading files to a Dataverse instance Using apiKey: 8599b802-659e-49ef-823c-20abd8efc05c
Adding content to: doi:10.5072/FK2/TUNNVE Using server: https://dataverse.tdl.org Request to upload:
testdir

PROCESSING(C): testdir Found as: doi:10.5072/FK2/TUNNVE

File Upload





- Folder Upload
- rsync + SSH Upload
 - File Upload Script

Dataverse >

- Command-line <u>DVUploader</u>
 - Usage
- Integrations Dashboard Uploader
- Duplicate Files
- BagIt Support

BagIt Support

Bagit is a set of hierarchical file system conventions designed to support disk-based storage and network transfer of arbitrary digital content. It offers several benefits such as integration with digital libraries, easy implementation, and transfer validation. See the Wikipedia article for more information.

If the Dataverse installation you are using has enabled Bagit file handling, when upic checksum values listed in each Bagit's manifest file against the uploaded files and g repository will identify a certain number of errors, such as the first five errors in each

| Select files or | r drag and drop into | he upload widget | . Maximum of 1,000 files per upload. |
|-----------------|----------------------|------------------|--------------------------------------|
| + Select | Files to Add | | |
| | | | Drag and drop files here |
| | | | |

The manifest declared a file, "basic/data/invalid-file-08.rb", that is not found in the Bagit package
 The manifest declared a file, "basic/data/invalid-file-01.rb", that is not found in the Bagit package
 The manifest declared a file, "basic/data/invalid-file-06.rb", that is not found in the Bagit package

Supported File Formats

Tabular Data ingest supports the following file formats:

| ile format | Versions supported | | | |
|------------------------------|----------------------------------|--|--|--|
| SPSS (POR and SAV formats) | 7 to 22 | | | |
| STATA | 4 to 15 | | | |
| R | up to 3 | | | |
| Excel | XLSX only (XLS is NOT supported) | | | |
| CSV (comma-separated values) | (limited support) | | | |

See the subsections in the left sidebar for more information on each of these supported formats.

<u>Files</u>

File Previews

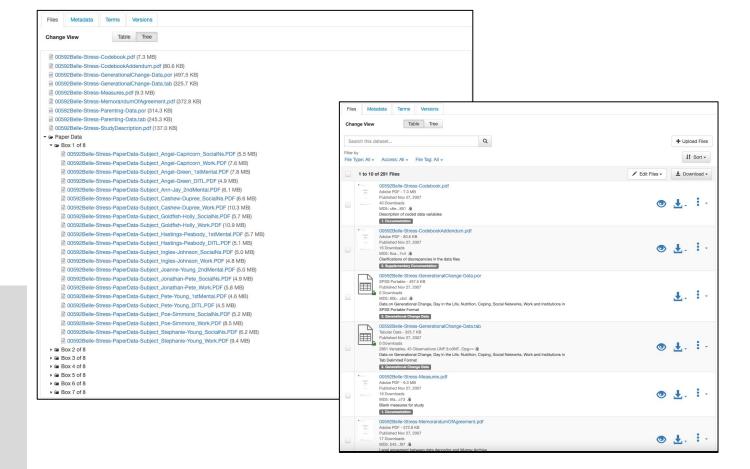
Dataverse installations can add previewers for common file types uploaded by their research communities. The previews appear on the file page. If a preview tool for a specific file type is available, the preview will be created and will display automatically, after terms have been agreed to or a guestbook entry has been made, if necessary. File previews are not available for restricted files unless they are being accessed using a Private URL. See also Private URL to Review Unpublished Dataset.

Previewers are available for the following file types:

- Text
- PDF
- Tabular (CSV, Excel, etc., see Tabular Data File Ingest)
- Code (R, etc.)
- Images (PNG, GIF, JPG)
- Audio (MP3, MPEG, WAV, OGG, M4A)
- Video (MP4, OGG, Quicktime)
- Zip (preview and extract/download)
- HTML
- GeoJSON
- GeoTIFF
- Shapefile
- NetCDF/HDF5
- Hypothes.is

Additional file types will be added to the dataverse-previewers repo before they are listed above so please check there for the latest information or to request (or contribute!) an additional file previewer.







File Citation European Patent Office, 2023, "PATSTAT Biblio Spring 2021", https://doi.org/10.7910/DVN/TAHKXA, Harvard Dataverse, V1; scripts.zip [fileName] Cite Data File -Learn about Data Citation Standards. **Dataset Citation** European Patent Office, 2023, "PATSTAT Biblio Spring 2021", https://doi.org/10.7910/DVN/TAHKXA, Harvard Dataverse, V1 Cite Dataset -Learn about Data Citation Standards. Preview Metadata Versions C Open in New Window

Intermediate Options in: Files

To download the complete zip file, please use the Access File button above.

| | ntation_Scripts | | |
|--------|--|--------------|----------|
| Creat | | | |
| | reateDatabaseScript | | |
| | createDatabaseSQLserver.sql | 2.31 kB | ₹ |
| | reateIndexScripts | | |
| Ŀ | allIindexes.sql | 10.39 kB | . ± |
| Ŀ | tls201_appln.IX_tls201_appln_date.sql | 545.00 Bytes | 1 |
| L. | tls201_appln.IX_tls201_appln_internat.sql | 557.00 Bytes | ±. |
| i i | tls201_appln.IX_tls201_appln_nr.sql | 566.00 Bytes | ₫. |
| | tls204_appln_prior.IX_tls204_prior_appln_id.sql | 572.00 Bytes | ±. |
| - | tls206_person.IX_tls206_person_cty.sql | 547.00 Bytes | 1 |
| lì. | tls207_pers_appln.IX_tls207_pers_appln_id.sql | 561.00 Bytes | .±. |
| là. | tls207_pers_appln.IX_tls207_pers_appln_pers_id.sql | 577.00 Bytes | |
| L. | tls209_appln_ipc.tls209_appl_ipc_XLS209M1.sql | 569.00 Bytes | |
| - | tls211_pat_publn.tls211_pat_publn_XLS211M2.sql | 604.00 Bytes | 4 |
| lì. | tls211_pat_publn.tls211_pat_publn_XLS211M3.sql | 564.00 Bytes | |
| à | tls211_pat_publn.tls211_pat_publn_XLS211M4.sql | 566.00 Bytes | 4 |
| L. | tls212_citation.tls212_citation_XLS212C2.sql | 568.00 Bytes | |
| L. | tls212_citation.tls212_citation_XLS212C3.sql | 586.00 Bytes | |
| li li | tls222_appln_jp_class.tls222_appln_jp_class_XLS222C1.sql | 626.00 Bytes | <u>+</u> |
| li | tls223_appln_docus.tls223_appln_docus_XLS223C1.sql | 586.00 Bytes | |
| L. | tls226_person_orig.IX_tls226_person_id.sql | 553.00 Bytes | |
| L | tls227_pers_publn.IX_tls227_pers_publn_id.sql | 565.00 Bytes | |
| l' | tls227_pers_publn.IX_tls227_pers_publn_pers_id.sql | 577.00 Bytes | 1 |
| i i | tls231_inpadoc_legal_event_UK_appln_id_seq_nr.sql | 654.00 Bytes | 4 |
| ∨ 🚈 Cr | reateTableScripts | | |
| - | alltablescreate.sql | 962.00 Bytes | 4 |
| L | createtic?01annin col | 1 75 kB | 1 |



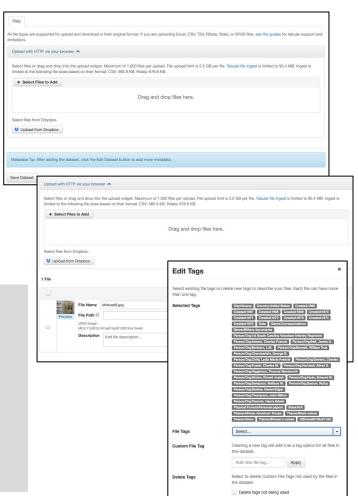
Access File ▼

Share

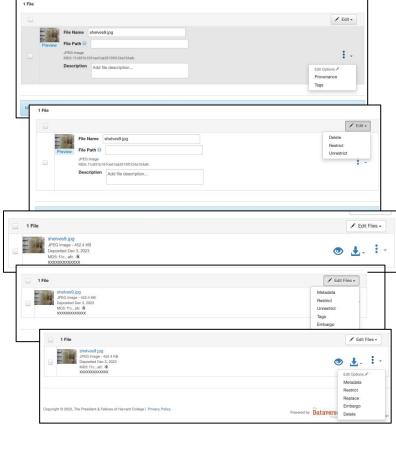
Contact Owner

File Metrics 🕣

1 Download @



Save Changes Cancel



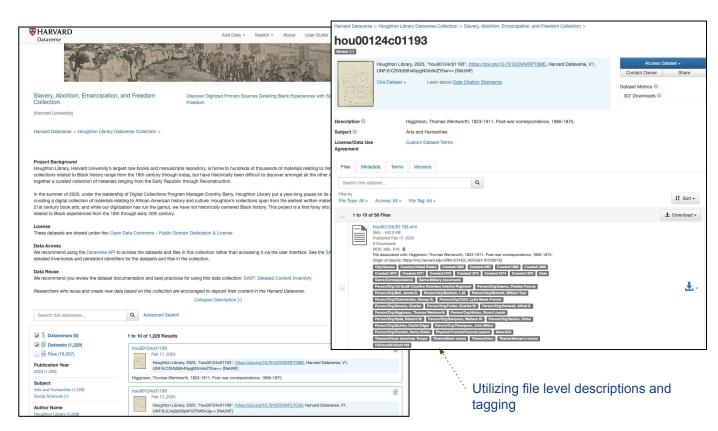
Intermediate Options in: Files

SAEF Collection

SAEF on GitHub

Utilizing a "collection" with detailed descriptions and links to related materials as well as information on utilizing APIs for downloading this many-files dataset

Intermediate & Advanced Settings in:





SAEF Collection

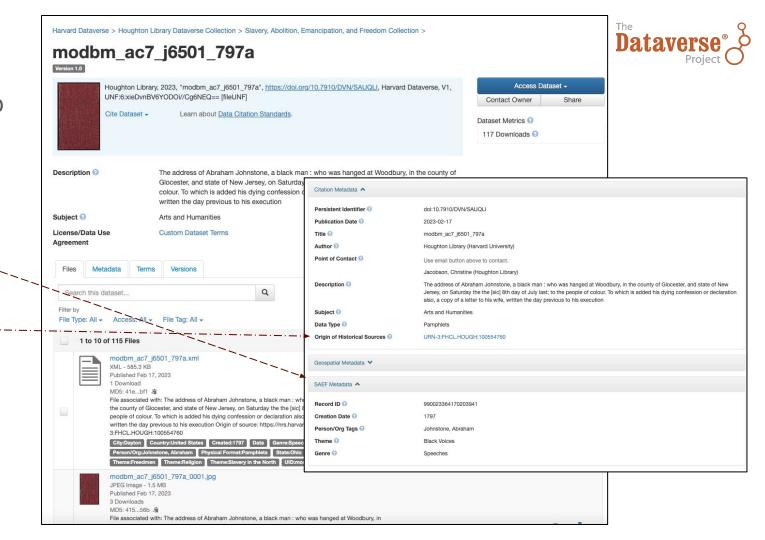
SAEF on GitHub

Utilizes a
"Custom metadata
block" as well as
multiple metadata
blocks

Utilizes linking to external data sources

Intermediate & Advanced Settings in:

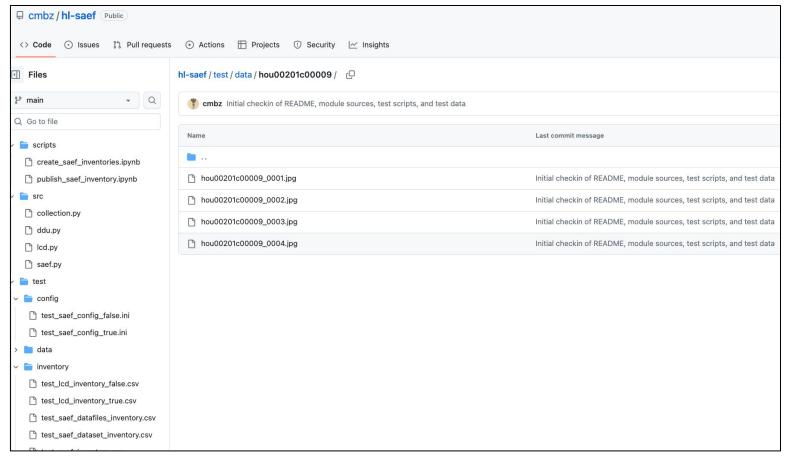
<u>Files</u>



SAEF Collection

SAEF on GitHub

Intermediate & Advanced Settings in:





SAEF Collection

SAEF on GitHub

Intermediate & Advanced Settings in:

<u>Files</u>

```
# create the pyDataverse datafile
datafile = Datafile()
# set the metadata on the datafile
datafile.set({'pid': pid, 'filename': filepath, 'restrict':restrict,
              'description':description, 'categories':categories})
# upload the datafile via the api
response = api.upload_datafile(pid, filepath, datafile.json(), is_pid=True)
status = int(response.status code)
if (not ((status >= 200) and
   (status < 300))):
   msg = 'Upload failed: {}'.format(response.status_code)
    print('SAEFDataset::upload_datafiles: Error - failed to upload datafile {}. {}'.format(filepath, msg))
    # log the event
    message = '{} - filename: {} - {}'.format(self._object_osn, filepath, msg)
    self.log_api_message('SAEF::upload_datafiles', 'api.upload_datafile: {}'.format(filepath), status, message)
else:
    # log the successful event
    msq = '{} - {}'.format(self._object_osn, filepath)
    self.log_api_message('SAEF::upload_datafiles', 'api.upload_datafile', status, msg)
```





Summary

Before data sharing:

- Explore data cleaning and validation options and techniques for handling missing data, outliers, and anomalies
- Ensure data consistency and integrity
- Employ data normalization and transformation techniques
- Examine strategies for data integration and metadata management
- Know the importance of versioning and data provenance

When sharing data (Dataverse):

- Utilize APIs and tools for "many" files deposits
- Utilize "collection" creation
 - Select additional metadata blocks outside of the default block
 - Select searchable facets to enhance discoverability of your datasets
- When creating the "dataset," complete as many metadata fields as relevant to your dataset, and include useful detailed information for title, description, keywords, and related articles to increase the Findability, Accessibility, and ReUsability of your dataset
- Organize Files in supported formats (tabular, compressed) to take advantage of the preview tools and archival functionality that supports Interoperability and ReUse

Thank you!