

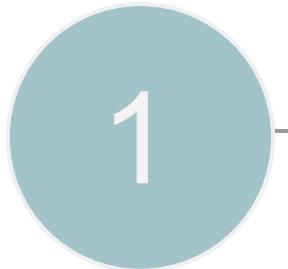
DATAVERSE

101: For Data Discovery Group

Elizabeth Quigley & Eleni Castro
Countway Library of Medicine - May 8, 2014

Agenda

What is Dataverse?



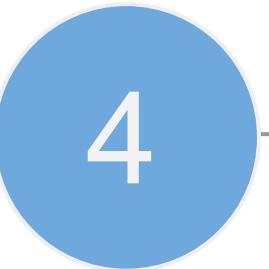
Why use Dataverse?



Biomedical
Metadata



Demo of 4.0
& DataTags



Questions?

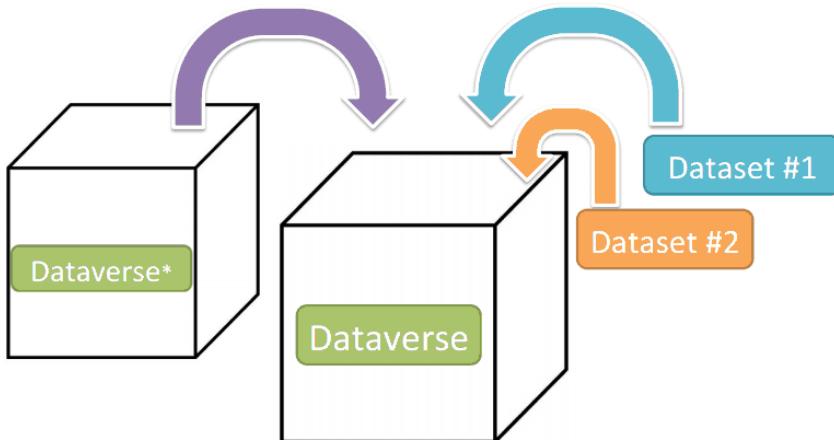


What is Dataverse?

A repository for sharing, citing, analyzing, and preserving research data.

What is a dataverse?

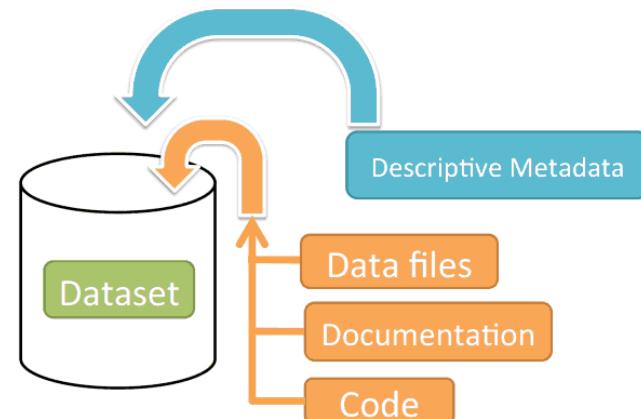
Schematic Diagram of a **Dataverse** in Dataverse 4.0



Container for your **Datasets** and/or **Dataverses***

* Dataverses can now contain other Dataverses (this replaces Collections & Subnetworks)

Schematic Diagram of a **Dataset** in Dataverse 4.0



Container for your data, documentation, and code.

Why use Dataverse?

Upload datasets and files:

- Receive a formal data citation for a dataset and its files, with a persistent identifier (DOI).
- Provide as much descriptive metadata as possible:
 - for others to **discover** a dataset and its files
 - for others to **reuse** a dataset for their own research

Why use Dataverse?

Find and use data:

- Search for dataverses, datasets, and files (incl. variables)
- Browse dataverses, datasets, and files
- Download data
- Subset & Analysis within Dataverse using Zelig statistical analysis software (R)

Why use Dataverse?

Reasons to share data:

- Fulfill data management plan requirements (e.g., PLOS, NIH, NSF, etc).
- Get recognition and credit via data citations.
- Allow collaborators to contribute to your Dataverse (e.g., Journal, Project).
- Restrict data to your team until ready to publish
- Facilitate discovery and reuse of your data through extensive metadata.
- Enable reproducible research, which contributes to the validation and verification of science.
- ...



photo: [Flickr Commons](#)

Biomedical Metadata

Initially working with Stem Cell Commons research metadata/ontologies for Proof of Concept.

Design Type

- Case Control
- Cross Sectional
- Not Specified
- Parallel Group Design
- Perturbation Design

Factor Type

- Age
- Biomarkers
- Developmental Stage
- Cell Surface Markers
- Cell Type/Cell Line

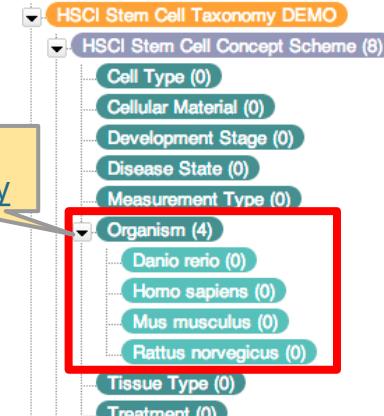
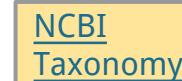
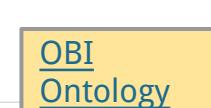
Measurement Type

- DNA Methylation Profiling (Bisulfite-Seq)
- DNA Methylation Profiling (MeDIP-Seq)
- Histone Modification (ChIP-Seq)
- Protein-RNA Binding (RIP-Seq)
- Transcription Factor Binding (ChIP-Seq)

Organism

- Danio rerio
- Homo sapiens
- Mus musculus
- Rattus norvegicus

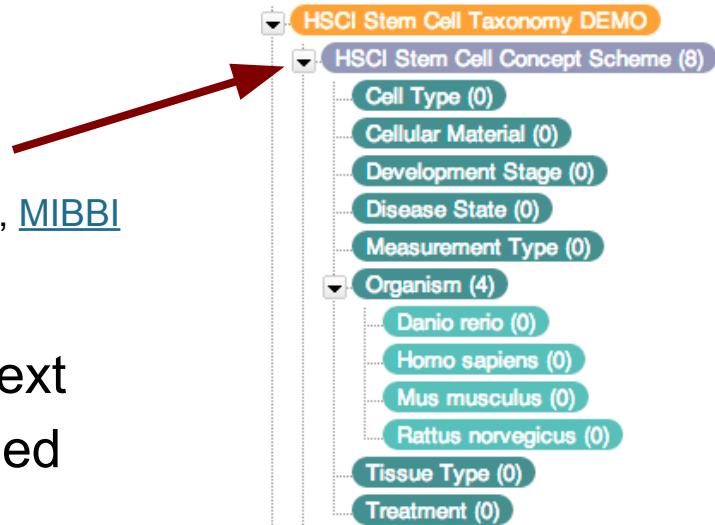
Cell Type

+


Bio Metadata Challenges

How feasible is it to support a one-size-fits-all solution for biomedical metadata?

1. Template Approach: work w/
different groups to support their
specific needs. (\uparrow Time)
(e.g., HSCI [Stem Cell Commons](#), [FAIRPORT](#), [MIBBI](#)
(e.g., [fMRI](#))).
2. For broader support allow free text
(\downarrow Quality) instead of just controlled
vocabularies (e.g., [OBO](#))



Transition from 3.6 to 4.0



The screenshot shows the Harvard Dataverse Network homepage. At the top, it features the IQSS logo, the Harvard Library logo, and the text "Share, Cite, Reuse, Archive Research Data" followed by "Scientific data for reproducible research". Below this, the title "Harvard Dataverse Network" is displayed, along with a search bar and links for "Create Account", "Harvard Affiliate", and "Log in". A note below the title states: "The Harvard Dataverse Network is open to all scientific data from all disciplines worldwide. It includes the world's largest collection of social science research data. Learn more about the Dataverse Network." The page is divided into several sections: "Dataverses" (with 574 entries), "Studies" (with 5250 entries), "RECENTLY RELEASED DATAVERSES" (listing "Elections Project", "Blanca de Lizaur", "Research of Archivaisip", "Code Metrics", and "Timothy Ryan"), "RECENTLY RELEASED STUDIES" (listing "Block Group-Level Data, 1970: San Diego, California", "Census of Population and Housing, 1970: [San Diego, California]: Summary Statistic File 4A: Housing by NIA", "HMS and MS Data Summary by Kerdpaipoj, Prad", "Testing for near I(2) trends when the signal to noise ratio is small", and "Consultants of the Fear of Success in Black High School Women, 1974"), and "POPULATION SERVICES INTERNATIONAL (PSI) Dataverses" (with 11 entries). A sidebar on the left shows the PSI logo and a list of recent releases. At the bottom, there are "View More" links for each section.

Dataverse v3.6

Transition from 3.6 to 4.0

What's changing?

- User interface and simpler workflows
- Support for more metadata
- Account page redeveloped with new functionality
- New interactive data exploration and analysis tool



Let's check out Dataverse 4.0 Beta.....

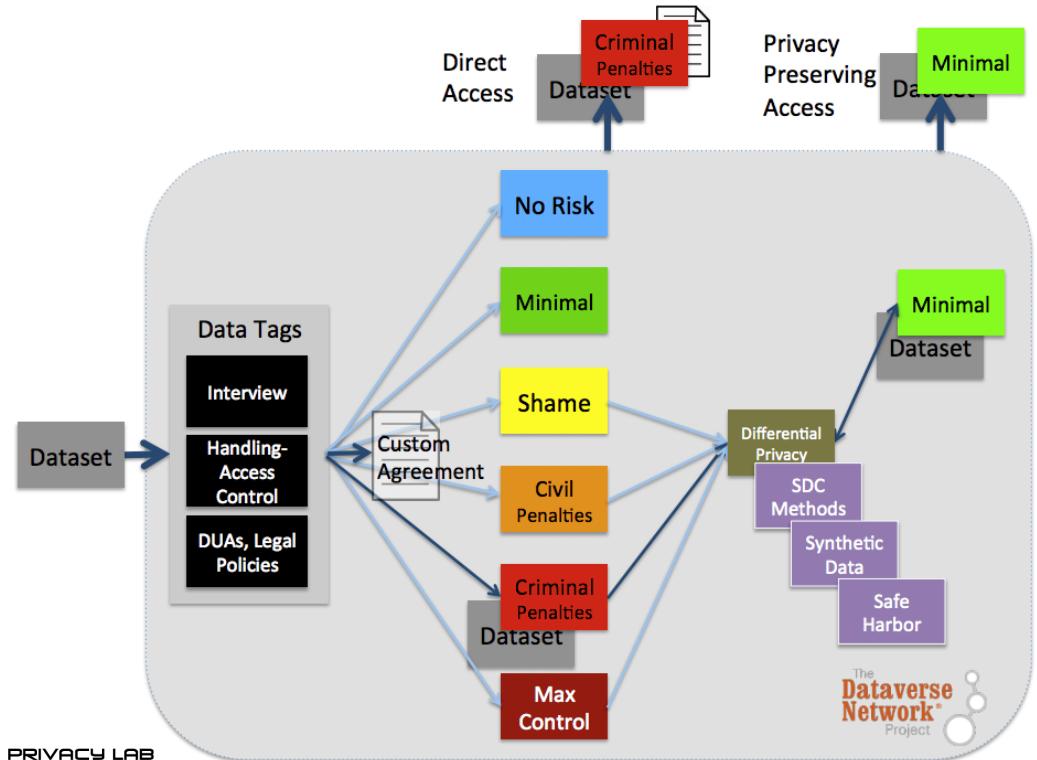
dataverse-demo.iq.harvard.edu

Sharing Privacy Sensitive Data

Coming soon:

Tools for sharing
privacy sensitive
research data.

→ DataTags (demo)





Questions?

To learn more about Dataverse and the Data Science Team at IQSS, visit datascience.iq.harvard.edu

Elizabeth Quigley: equigley@iq.harvard.edu & Eleni Castro: ecastro@fas.harvard.edu