# The Dataverse Project Is Also A Home For Life Sciences Data



Sonia Barbosa & Eleni Castro, Harvard University
bioCADDIE webinar: June 8, 2016
http://dataverse.org

# The Dataverse Project

Open source research data repository software

Share, preserve, cite, explore, & analyze data

# Collaborations

- The Institute for Quantitative Social Science (IQSS)

- the Harvard University Library

- Harvard University Information Technology

- The [Open Data Assistance Program at Harvard](#) (a collaboration with Harvard Library, the Office for Scholarly Communication and IQSS)

- The Library Technology Services at HUIT provides hosting and backup support

## Researchers

Enjoy full control over your data. Receive *web visibility, academic credit,* and *increased citation counts.* A personal dataverse is easy to set up, allows you to display your data on your personal website, can be branded uniquely as your research program, makes your data more discoverable to the research community, and satisfies data management plans. Want to set up your personal dataverse?

## Journals

Seamlessly manage the submission, review, and publication of data associated with published articles. Establish an *unbreakable link* between *articles in your journal* and *associated data.* Participate in the open data movement by using Dataverse as part of your journal data policy or list of repository recommendations. Want to find out more about journal dataverses?

## Developers

Participate in a vibrant and growing community that is helping to drive the norms for sharing, preserving, citing, exploring, and analyzing research data. Contribute code extensions, documentation, testing, and/or standards. *Integrate research analysis, visualization* and *exploration tools,* or other research and data archival systems with Dataverse. Want to contribute?

## Institutions

Establish a research data management solution for your community. Federate with a growing list of Dataverse repositories worldwide for increased discoverability of your community's data. Participate in the drive to set norms for sharing, preserving, citing, exploring, and analyzing research data. Want to install a Dataverse repository?

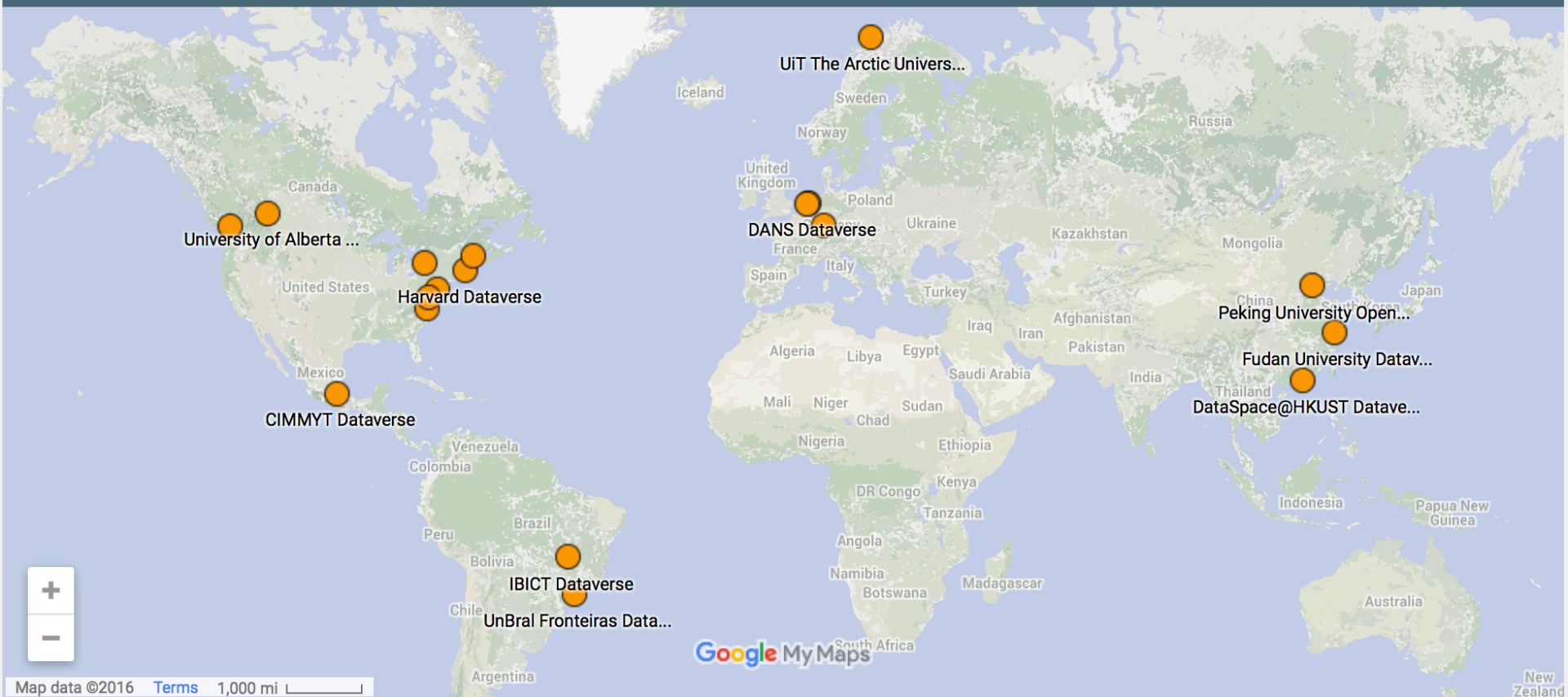# DATAVERSE REPOSITORIES - A WORLD VIEW

| 17 Installations | 1,500+ Dataverses | 65,000+ Datasets | 1,700,000+ Downloads |
|---|---|---|---|

## Dataverse Repositories



UiT The Arctic Univers...

Iceland

Sweden

Norway

Russia

United Kingdom

DANS Dataverse

Poland

France

Italy

Kazakhstan

Mongolia

Spain

Turkey

China

Japan

Peking University Open...

Canada

United States

Afghanistan

Pakistan

Iraq

Iran

India

Fudan University Datav...

University of Alberta ...

Algeria

Libya

Egypt

Saudi Arabia

Thailand

DataSpace@HKUST Datave...

Harvard Dataverse

Mali

Niger

Chad

Sudan

Mexico

Nigeria

Ethiopia

CIMMYT Dataverse

Venezuela

Colombia

DR Congo

Kenya

Tanzania

Indonesia

Papua New Guinea

Peru

Brazil

Bolivia

Angola

Namibia

Botswana

Madagascar

IBICT Dataverse

Chile

UnBral Fronteiras Data...

South Africa

Argentina

Australia

New Zealand

Google My Maps

Map data ©2016    Terms    1,000 mi

# Features

**Data Citation**
automatically generated

**Multiple Publishing Workflows**
dataset in draft, in review, and then published

**Terms of Use + Guestbook**
CC0 waiver default, custom terms of use, and download metrics

**Account + Data Notifications**
access request, roles granted, and when data is published to name a few

**Faceted Search**
metadata fields based facets

**Pull header metadata from Astronomy (FITS) files**

**APIs for interoperability**
search API, data deposit API

**Three Levels of Metadata**
description/citation, domain-specific or custom fields, file metadata

**Access Control Support**
pre-defined and custom roles

**Restricted Files + Ability to request access to restricted files**
allow anyone, certain people, or no one to be able to download files

**Customization of dataverses**
branding, metadata based facets, sub-dataverses, featured dataverses

**Re-format, Summary Statistics, and Analysis for Tabular Files**
integration with TwoRavens

**Mapping of Geospatial files**
integration with WorldMap

# Features

- Standard Citation: Title, DOI, UNF, Versioning, Repository (following FORCE11 Joint Declaration of Data Citation Principles)

- File level support: MD5, UNF, Tabular data, multiple download options, tags, descriptions, zip extraction, audio, video, PDF preview, image files w/preview, unlimited files, all file types

- Metadata support

- Terms: CC0, additional terms, restricted/open, application forms

- Versioning

# Features…

- dataverse or dataset
- Themes and widgets
- Permissions
- Groups
- Guestbook
- Templates
- Featured dataverses

# Next Releases

| NEXT RELEASES | CURRENT PROJECTS | PAST RELEASES |
|---|---|---|

**Version 4.4, June 16 2016:**
This release can be tracked here: https://github.com/IQSS/dataverse/milestones/4.4

- Updates to widgets for personal websites
- Support for remote authentication with Shibboleth
- Guestbook feature bug fixes

**Version 4.5, End of June, 2016:**
- Metadata Harvesting and Export Metadata in standard formats
- Private URL for reviewing unpublished datasets

# Current Projects

| NEXT RELEASES | CURRENT PROJECTS | PAST RELEASES |
|---|---|---|

These projects will be integrated into the Dataverse in 2016:

## Summer 2016
- Handles
- Internationalization
- File-level metadata, file-level landing page and provenance metadata

## Fall 2016
- Support for sensitive data
- Support for large-scale data

# Past Releases

Previous Dataverse 4.x releases can be found here: **https://github.com/IQSS/dataverse/releases**
Each release includes release notes outlining what features or functionality have been added as well as the bugs fixed.

## Version 4.3, March 21, 2016:
This release code can be found be here: https://github.com/IQSS/dataverse/releases/tag/v4.3

- DataCite API support (extension from current support of DOIs from EZID)
- Ability to add custom text to the dataset publishing pop up (only available for Dataverse installations)
- Ability to log in using your email address
- Ongoing bug fixes

# Dataverse is working on being FAIR

# Life Sciences Metadata

- ISA-Tab with Scientific Data flavor (see next slide)
- Various ontologies from bioportal (ex. OBI)
- NCBI Taxonomy
- Plan to support export (Fall 2016) – discoverability with NIH's discovery index

# Life Sciences Metadata ⌃

**Design Type**
- ☐ Case Control
- ☐ Cross Sectional
- ☐ Not Specified
- ☐ Parallel Group Design
- ☐ Perturbation Design
- ☐ Technological Design

**Factor Type**
- ☐ Age
- ☐ Biomarkers
- ☐ Developmental Stage
- ☐ Cell Surface Markers
- ☐ Cell Type/Cell Line
- ☐ Disease State

ISA-Tab metadata in Dataverse 4.3

**Organism**
- ☐ Arabidopsis thaliana
- ☐ Bos taurus
- ☐ Caenorhabditis elegans
- ☐ Chlamydomonas reinhardtii
- ☐ Danio rerio (zebrafish)
- ☐ Dictyostelium discoideum

**Other Organism**
[                                                    ] [ **+** ]

**Measurement Type**
- ☐ genome sequencing
- ☐ cell sorting
- ☐ transcription factor binding site identification
- ☐ hematology
- ☐ cell counting
- ☐ DNA methylation profiling

**Other Measurement Type**
[                                                    ] [ **+** ]

**Technology Type**
- ☐ nucleotide sequencing
- ☐ flow cytometry
- ☐ DNA microarray
- ☐ mass spectrometry
- ☐ gel electrophoresis
- ☐ protein microarray

**Technology Platform**
- ☐ 210-MS GC Ion Trap (Varian)
- ☐ 220-MS GC Ion Trap (Varian)
- ☐ 225-MS GC Ion Trap (Varian)
- ☐ 240-MS GC Ion Trap (Varian)
- ☐ 300-MS quadrupole GC/MS (Varian)
- ☐ 320-MS LC/MS (Varian)

**Cell Type**
[                                                    ] [ **+** ]

# Dataverse in isaexplorer

# Structural Biology Data + Dataverse



https://figshare.com/articles/SBGRid_DB_Poster_Force2016_pdf/3175417

# Journals + Dataverse

**DATA IN BRIEF**

Data in Brief (DiB) Dataverse (Elsevier)   Home

✉   ⤴

# High resolution 3D laboratory x-ray tomography data of femora from young, 1-14 day old C57BL/6 mice

Bortel, Emely L; Duda, Georg N; Mundlos, Stefan; Willie, Bettina M; Fratzl, Peter; Zaslansky, Paul, 2015, "High resolution 3D laboratory x-ray tomography data of femora from young, 1-14 day old C57BL/6 mice", http://dx.doi.org/10.7910/DVN/29628, Harvard Dataverse, V1

≡ Download Citation ▾

If you use these data, please add this citation to your scholarly resources. Learn about Data Citation Standards.

| | |
|---|---|
| **Description** | This data article contains high resolution (1.2 μm effective pixel size) lab-based micro-computed tomography (μCT) reconstructed volume data of young C57BL/6 mouse femur bone midshafts. This data formed the basis for the analyses of bone structure development in healthy mice, including closed and open porosity as reported in http://dx.doi.org/10.1016/j.actbio.2015.03.027. The data reveal changes seen in bone material and porosity distribution in young animals aged 1 to 14 days old.The mouse bones transform from porous scaffolds into solid structures during normal organogenesis. The large data may be freely used by others and in all research areas. |
| **Keyword** | C57BL/6 growth |
| **Related Publication** | Bortel EL, Duda GN, Mundlos S, Willie BM, Fratzl P, Zaslansky P. Long bone maturation is driven by pore closing: A quantitative tomography investigation of structural formation in young C57BL/6 mice. Acta Biomater (2015) |

Files   Metadata   Terms   Versions

Search this dataset...   Q Find

**19 Files**   ⬇ Download

☐

☐   day10_sample1.zip
ZIP Archive - 1.3 GB - Apr 6, 2015 - 16 Downloads
MD5: 2f03c776e509e1dac5cf212efd78fbf3;   ⬇ Download

day10_sample2.zip

open health data

UPmetajournals          ]u[ ubiquity press open access

**Open Health Data Dataverse** (Ubiquity Press)          Home Page

📊 Metrics          49 Downloads

# WHO Mortality database

de Roos, Albert, 2015, "WHO Mortality database", http://dx.doi.org/10.7910/DVN/28948, Harvard Dataverse, V1

≣ Download Citation ▾

If you use these data, please add this citation to your scholarly resources. Learn about Data Citation Standards.

### Description

The WHO mortality data was transformed into a corpus of mortality data in a standard relational database format that allows for easy data mining. The set includes corresponding population data, calculated mortality rates and an ICD-code reference table encompassing all years of ICD registration. The database can be downloaded and imported into a relational database or be combined with other epidemiological or demographic data. The easy of access of these data for researchers may be of great benefit to the research into global trends and causes of death.

| Files | Metadata | Terms | Versions |
| --- | --- | --- | --- |

🔍 Search this dataset...          🔍 Find

**5 Files**          ⬇ Download

☐

**Dump20141228.zip**
ZIP Archive - 271.6 MB - Feb 2, 2015 - 14 Downloads
MD5: 9dd3a6e773dd257be65b085c33502922;
All the tables in the mortality database in SQL dump format, to be imported into a
relational database such as MySQL
**SQL database dump**

⬇ Download

# Future Work w/ bioCADDIE

Will be attending bioCADDIE workshop in late June to learn more from you!

# Thank you!

Contact: support@dataverse.org

Twitter: @dataverseorg

Web: http://dataverse.org