

# Research Data Management

## Getting Started with Data Management Plans

# Instructors

— — —



**Julie Goldman**

Research Data Services Librarian  
Countway Library of Medicine  
[Julie\\_Goldman@hms.harvard.edu](mailto:Julie_Goldman@hms.harvard.edu)



**Meghan Kerr**

Archivist and Records Manager  
Center for the History of Medicine  
[Meghan\\_Kerr@hms.harvard.edu](mailto:Meghan_Kerr@hms.harvard.edu)



Slides: <https://hlrdm.library.harvard.edu/presentation-slides>

A reference guide with information and resources to help you manage your research data

## DATA LIFECYCLE

## Planning Data Management

How can I best manage my data throughout the lifecycle of my research to save time and money in the future?

- Data Management Plans (DMPs)
- DMP requirements and tools
- What research objects should be tracked and documented

## Data Acquisition and Collection

How can I acquire data in an efficient and ethical way, and how can I ensure that my data is used appropriately?

- Data Use Agreements (DUAs)
- Institutional Review Boards (IRB and IACUC)
- Subscription data

## Storage, Security, and Analysis

What are my options for effectively organizing, storing, securing, computing, and analyzing my research data?

- Data security
- Computing, research methods, data science, and viz support
- Electronic Lab Notebooks

## Dissemination and Preservation

Why is it worthwhile to share my data? What do funders and journals require? Can I get help with data curation?

- Data repositories
- Open Access
- Data citation, FAIR principles
- Data disposal

## LATEST NEWS



## New Data Use Agreement Guidance for Harvard Researchers

January 4, 2019

As part of a larger initiative to improve transparency of research compliance processes, the University has introduced a tool for submission, review, and management of Data Use Agreements (DUAs). In support of recently published guidance by the Office for the Use of Research Data (OURD) on the management of DUAs, researchers sharing data can use the system to request a DUA review, correspond with the DUA reviewer, track the status of review, and manage active DUAs (including amendments and extensions...

[Read more](#)

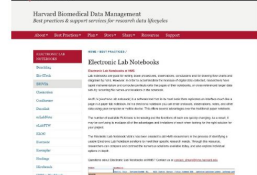
## TWEETS FROM HARVARD LIBRARY @HLRMP



**@scstatelibrary** Don't forget about your data for #PreservationWeek! One of the most important things you can do to prevent disaster with your files is back up, back up, back up! Remember the 3-2-1 rule for backing up: 3 copies, 2 different formats, 1 off-site location back up. #preswk1 c0PBNKXKX3



**@Hayleiman** Harvard (@g0lts2 & others at @HUSCountway) I c0P78W8WZq4 t c0P5GZ23W8B 7H2 t c0W8E8W8W8H



## DATA SHARES

## One Step Closer to the "Paper of the Future"

March 13, 2019



**Catherine Zucker** is a fourth-year PhD student in astronomy at Harvard. Her work focuses on understanding the structure of our Milky Way Galaxy, through the combination of observations, numerical simulations, statistics, and data visualization. She is advised by Professors Alyssa Goodman and Douglas Finkbeiner. **Alyssa Goodman** is the Robert Wheeler Wilson Professor of Applied Astronomy, a Co-Founder of the Initiative

Innovative Computing and a member of the Harvard Data Science Initiative steering committee, whose research spans astronomy, data visualization and online systems for research and education. **Douglas Finkbeiner** holds joint professorships in the departments of Astronomy and Physics, whose research spans high-energy astrophysics, dark-matter annihilation, Galactic structure, interstellar dust, and large photometric surveys.

[Read more](#)

## UPCOMING EVENTS

- 2019 MAY 07** Getting Started with Data Management Plans 12:00pm to 1:00pm
- 2019 MAY 16** Webinar: What's New at ORCID? 10:00am to 11:00am
- 2019 MAY 17** Webinar: Research Data Management and Workflows in Universities - Transforming Data to Knowledge 10:00am to 11:30am

# Research Data Management @Harvard Website

<https://researchdatamanagement.harvard.edu>

# Introduce Yourself!

— — —



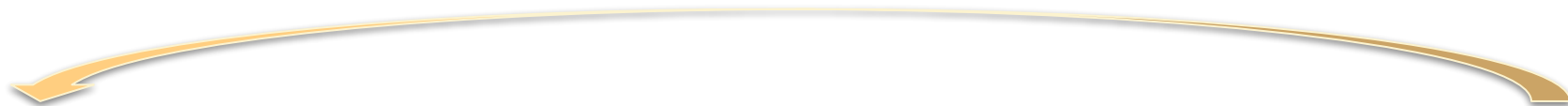
**Name**

**School / Department**

**Have you written/followed a DMP before?**

*(for a grant, class research project, etc.)*

# RDM Workflow



PLANNING DATA MANAGEMENT	ACQUISITION AND COLLECTION	STORAGE, SECURITY, AND ANALYSIS	DISSEMINATION AND PRESERVATION
<b>Plan for research data needs</b>  Research object documentation  Data Management Plans (DMPs)  DMPTool	<b>Find, acquire &amp; collect data</b>  Instruments, Researchers, Vendors  Data use agreements (DAU)  Institutional Review Boards (IACUC, IRB, IBC)	<b>Organize, store &amp; process data</b>  R, Python, OpenRefine  Statistical software  File systems, Asset management  Data security	<b>Share data in repository</b>  Data repository  Data curation; Appraise for enduring value  Data citations, DOIs  Migrate data to preservation repository

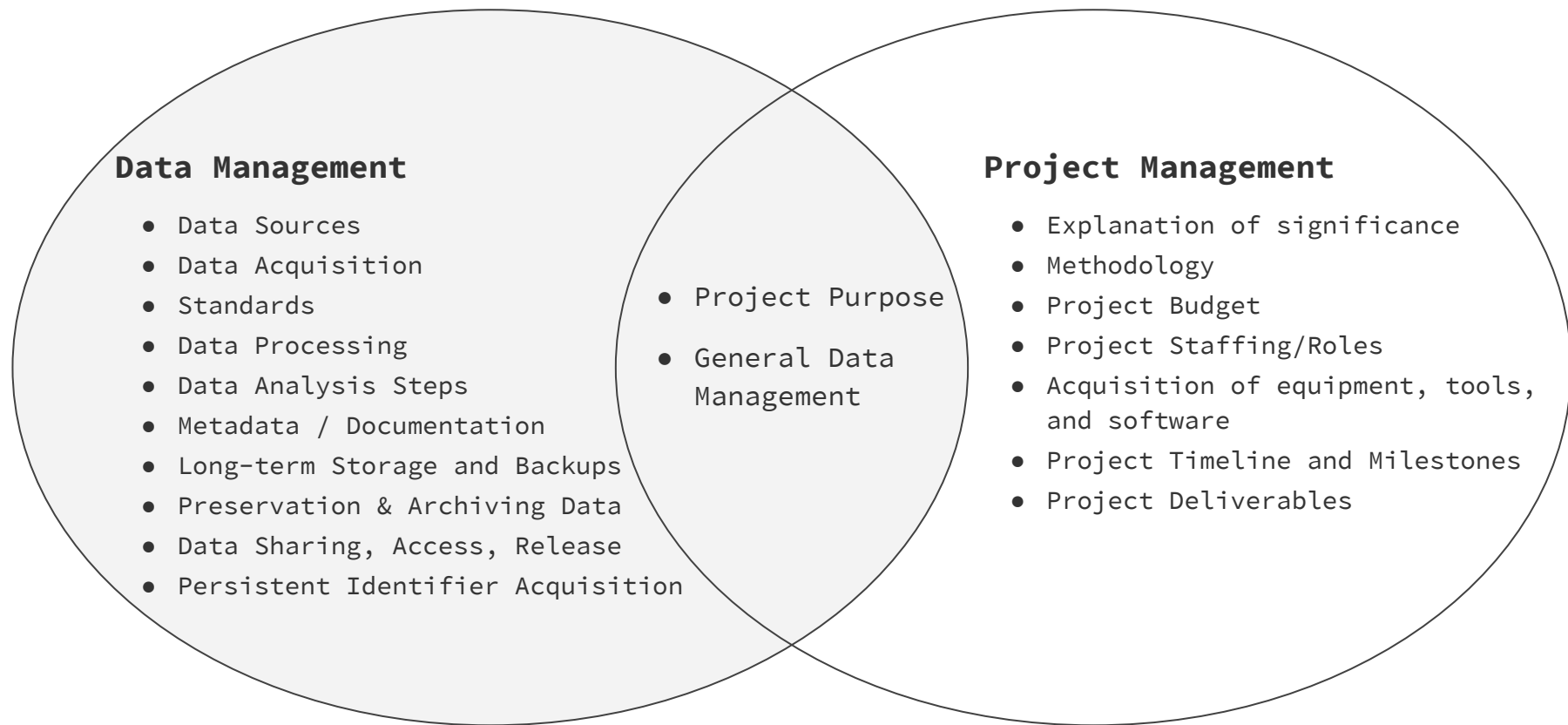
# Why Manage Data?

- - -
- Easier to analyze organized, documented data
- Find data more easily
- Don't lose data
- Don't drown in irrelevant data
- Get credit for your data
- Avoid accusations of misconduct



Data Sharing and Management Snafu in 3 Short Acts

# Data Management vs Project Management

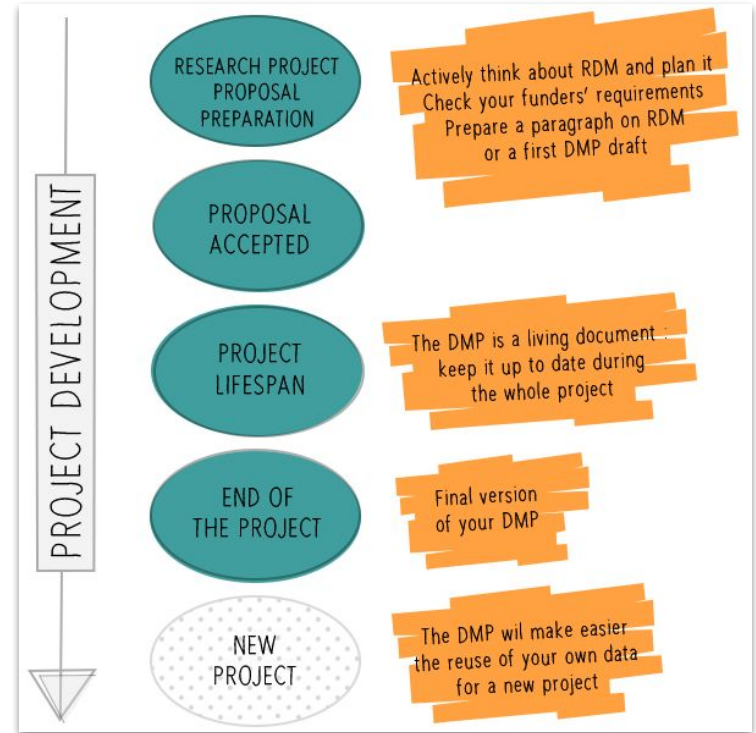




# Data Management Plan

Short (2pg) document that describes what you will do with your data. DMPs now required by all major federal funders & many private funders. Part of your grant approval & reporting.

1. Project, experiment, and data description
2. Documentation, organization, and storage
3. Access, sharing, and re-use
4. Archiving



<https://researchdata.epfl.ch/plan-fund/dmp>





# Research Data

## Data Through the Lifecycle

**Raw:** What is being measured or observed?

**Processed:** How can the raw data be manipulated?

**Analyzed:** What does the data tell us?

**Finalized/Published:** How does the data support your research question?

*Consider: type of data, formats, size & complexity*

### Acquisition & Creation

- ✓ Raw data
- ✓ Working files

### Analysis

- ✓ Analytical methods
- ✓ Analysis results

# Example: *Research Data Description*

---

- ❑ Data types will include plain text files and PDFs, ready for Libra deposit and distributed version control using git. (3)
- ❑ Primary experimental Data
  - a. Voltage data...data are initially acquired and stored using LabChart Pro and then converted to HDF5 using a custom converter.
  - b. High speed video recordings...stored as uncompressed AVI files or as HDF5 files.
  - c. Laboratory notebooks and other notes. These are stored electronically using the LabArchives software. (8)

# Metadata

---

**Data documentation** provides the information necessary to fully understand and interpret the data

**Metadata** should be standardized, consistent and interoperable, and facilitates discovery, preservation and archiving of data

***Consider: templates & standards, project vs data level***

<https://researchdatamanagement.harvard.edu/best-practices-organizing-documenting-research-data>



Andy Warhol, *Big Torn Campbell's Soup Can (Pepper Pot)*, 1962 The Andy Warhol Museum, Pittsburgh Founding Collection, Contribution The Andy Warhol Foundation for the Visual Arts, Inc.

# Example: *Metadata and Documentation System*

— — —

- ❑ Metadata will be provided. The project will document information about the context, content, quality, provenance, and/or accessibility of the data used. This will also include information embedded in the raw FID files. Additionally, the project will seek to document information about authors, dates and brief descriptions for scanned PDFs, notebooks and lab work. (9)
- ❑ Metadata will be stored using the TEI XML encoding. Metadata will be stored in English and in compliance with ISO 639-2 in order to make these data more easily readable by machines. (2)



# Storage, Backup, and Security

## Storage, Security & Maintenance

- ✓ Store on appropriate tier, with proper security
- ✓ Store locally on servers or in the cloud
- ✓ Plan to maintain system

**Consider: storage type, backup location**

<https://researchdatamanagement.harvard.edu/data-security-0>

LEVEL 1	Public information	► Level 1 Data Types
LEVEL 2	Level 2 is information the University has chosen to keep confidential but the disclosure of which would not cause material harm.	► Level 2 Data Types
LEVEL 3	Level 3 information could cause risk of material harm to individuals or the University if disclosed.	► Level 3 Data Types
LEVEL 4	Level 4 information would likely cause serious harm to individuals or the University if disclosed.	► Level 4 Data Types
LEVEL 5	Level 5 information would cause severe harm to individuals or the University if disclosed.	► Level 5 Data Types

# Example: *Storage and Security Plan*

— — —

- ❑ All of the project data will be maintained on servers, local computers, and hard drives maintained by the project director. The costs of data management are projected to be minimal, and will be borne by the project director. (4)
- ❑ Data security and confidentiality are protected by using Microsoft Active Directory authentication, and the storage is backed up to LT0-4 tape on a daily and weekly basis and stored offsite at Iron Mountain facilities. (9)



# Protection and Privacy

— — —

**Access:** Limiting the availability of your data

**Systems:** Protecting your hardware and software

**Data Integrity:** Ensure your data is not manipulated in an unauthorized way

**Ethics:** Consider the wider consequences of your research

**Personal Data:** Remove data which are not used; ensure subject confidentiality

***Consider: consult ethics committees, anonymize data to protect privacy of your participants***

# Example: *Provisions for Data Privacy and Access*

— — —

- ❑ Research records will be kept confidential, and access will be limited to the PI, primary research team members, and project participants. Data will be housed on a local server controlled by the PI, and will be accessible via SSH and VPN. Data containing identifiable information, or information covered by an NDA, will be held in an encrypted format. (6)
- ❑ The website that presents the BPS tool-kit has a standard UC Berkeley privacy policy that is linked from every page. It notes that while information may be collected to run the services, personal information will not be disclosed without a user's consent, except for "certain explicit circumstances in which disclosure is required by law." (1)



# Policies for Re-use

---

When establishing data sharing and access policies and provisions, consider *whom* you will share your data with, *how* it will be shared, and *when* in the research process you will share it.



Digital Object Identifier



Open Access: free & unrestricted



creative commons

Creative Commons Licenses

**Consider: access categories (open, registered, limited, embargo) & licensing (CC)**

## Example: *Data Re-use and Copyright Statement*

— — —

- ❑ The researchers associated with this study are not aware of any reasons that might prohibit the sharing of the data to be generated under this project for public use and potential secondary uses, assuming data is handled consisted with IRB and NDA guidelines. The principal investigators retain the right for first use of the data. (6)



- ✓ Share data with collaborators

- ✓ Annotate datasets & upload to public repositories
- ✓ Include in relevant publications & reports

**Consider: timing, data papers, consulting an expert**

	Yes	No					
	<b>Page last updated July 2, 2018</b>						
Requirement	<a href="#">Dataverse</a>	<a href="#">Dryad</a>	<a href="#">figshare</a>	<a href="#">Zenodo</a>	<a href="#">GigaScience</a>	<a href="#">Scientific Data</a>	
<b>Data Size and Format</b>							
Hosting of common file formats (e.g. csv, tsv, xls,.xlsx, doc, pdf)						-	
Hosting of proprietary file formats (e.g. raw image files)						-	
Unlimited size per file						-	
Unlimited total dataset size						-	
<b>Data Licensing</b>							
CC0 waiver†	recommended	required	recommended	available	required	-	
<b>Data Attribution and Citation Tools</b>							
Assignment of dataset DOIs						-	
<b>User Access Controls</b>							
Tiered access (e.g. administrator-level, collaborator-level, curator-level)						-	
Journal-integrated, anonymous access (for peer review pre-publication)						-	
Optional embargo to data release following publication						-	
<b>Data Access Tools</b>							
Comprehensive data and metadata search tools							
Data access via direct download						-	
Data downloading via API						-	
Built-in tools for reading proprietary file formats						-	
Integrated data analysis tools							
<b>Cost</b>							
Data deposition fees	none	tiered	none	none	none	-	
Data maintenance fees	none	none	none	none	none	-	

## Example: *Data Sharing Plan*

— — —

- ❑ Data will be made available for sharing to qualified parties by the Co-PIs, so long as such a request does not compromise intellectual property interests, interfere with publication, invade subject privacy, betray confidentiality, or precede data curation. (7)



# Archiving and Preservation

— — —

Data retention requirements are put in place by funding agencies and sponsoring institutions for a number of reasons:

- promote the reuse of data within and across disciplines
- protect intellectual property rights
- make research findings available
- support open data initiatives

Appraisal process for evaluating research records and data:

- ***Inventory of the records:*** volume, data types, formats, metadata, other relevant information
- ***Interview about the project:*** impact of the project, significance of the research or researcher, basic information about the grant

***Consider: does your dataset have reuse potential, is your dataset reusable***

## Example: *Long-term Preservation of Data*

— — —

- ❑ While UVa's Records Management protocol specifies a 5-year retention period for all grant-related material, the Library and UVa Information Technology Services plan to preserve content deposited in Libra is anticipated indefinitely. (3)
- ❑ Storing LOGAR records in TEI XML provides assurance that the project's data will be available for long-term scholarly research. Storing GeoPACHA data using its open data standards assures long-term support. (5)

# Activity

# DMP Bingo



## How to Play:

1. The *BINGO* card squares describe various types of data management decisions & choices
2. Players will try to find matches between their card's squares and the DMP assigned
3. The *BINGO* cards have both "good" and "bad" DMP attributes which should be taken into consideration
4. Groups are encouraged to discuss and evaluate their DMP together as all the cards have the same criteria (in different places)
5. Players should mark the squares that match their DMP in some way
6. A player gets *BINGO* when a straight line of 5 matching squares are marked!

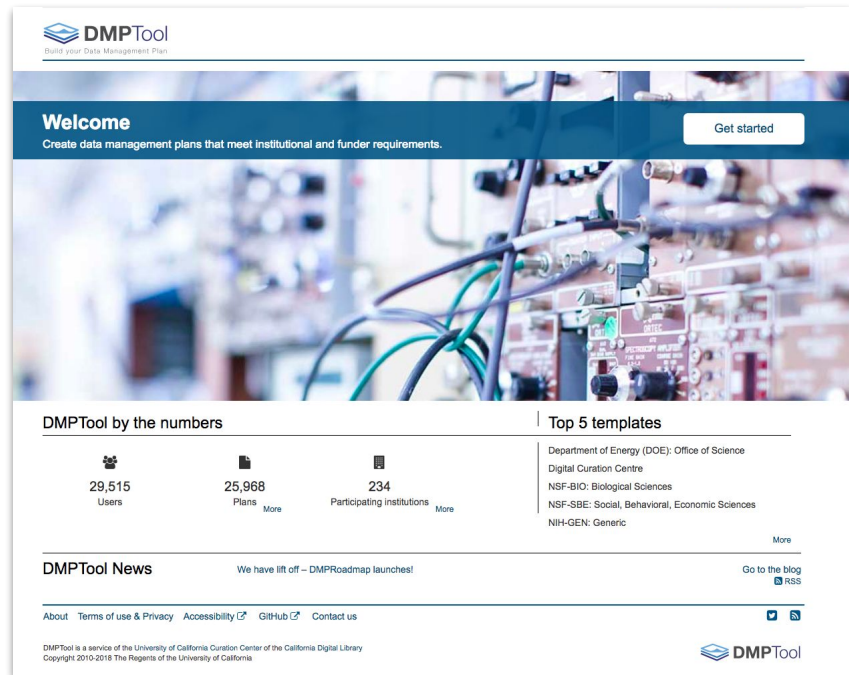
**\*\*CAVEAT: *BINGO* is not guaranteed\*\***

O'Donnell, Megan (2016): DMP Bingo - the good, the bad, the ugly (v.2).  
figshare. <https://doi.org/10.6084/m9.figshare.1564825.v2>




# DMPTool

The DMPTool is an online tool that includes data management plan templates for many of the large funding agencies that require them.

Harvard is an affiliated partner institution. You can login as a user from your institution with your HarvardKey. By being affiliated Harvard, you will be presented with institution-specific guidance to help you complete your plan.



The screenshot shows the DMPTool website homepage. At the top is the DMPTool logo with the tagline "Build your Data Management Plan". Below the logo is a "Welcome" banner with the text "Create data management plans that meet institutional and funder requirements." and a "Get started" button. The main content area features a large image of electronic equipment. Below this, there are two sections: "DMPTool by the numbers" and "Top 5 templates".

DMPTool by the numbers			Top 5 templates
 29,515 Users	 25,968 Plans <a href="#">More</a>	 234 Participating institutions <a href="#">More</a>	Department of Energy (DOE): Office of Science Digital Curation Centre NSF-BIO: Biological Sciences NSF-SBE: Social, Behavioral, Economic Sciences NIH-GEN: Generic <a href="#">More</a>

Below the statistics, there is a "DMPTool News" section with the headline "We have lift off – DMPRoadmap launches!" and a link "Go to the blog". At the bottom, there is a footer with links for "About", "Terms of use & Privacy", "Accessibility", "GitHub", and "Contact us", along with social media icons and the DMPTool logo.

<https://dmptool.org>

<https://researchdatamanagement.harvard.edu/data-management-plans>



# Questions?

## Research Data Management @Harvard

[Home](#) [Vision](#) [Data Lifecycle ▾](#) [Policies](#) [Resources](#) [Contacts](#)

[HOME](#) /

### Contact

Office of the Vice Provost for Research:

- [Mercè Crosas](#), Harvard University Research Data Management Officer
- [Rachel Talentino](#), Research Compliance Officer

Research Data Management Library Contacts:

- Harvard Business School - Baker Library: [Katherine McNeill](#), Research Data and Collections Librarian
- Harvard Graduate School of Education - Gutman Library: [Alex Hodges](#), Faculty Director, Gutman Library, & HGSE Librarian
- Harvard Library: [Ceilyn Boyd](#), Research Data Program Manager
- Longwood - Countway Library: [Julie Goldman](#), Countway Research Data Services Librarian

Data Use Agreements - Negotiating Offices:

- [Office for Sponsored Programs: dua@harvard.edu](#)
- [HMS Office of Research Administration: SPAContracts@hms.harvard.edu](#)
- [HSPH Office of Research Administration: duahelp@harvard.edu](#)

Human Subjects Research:

- [Cambridge and Allston IRB Contact Page](#) or email [cuhs@harvard.edu](#) if your department is not listed
- [Longwood-Area IRB Contact Page](#)

Harvard Dataverse Repository:

- [support@dataverse.harvard.edu](#)

Information Security:

- Find your [School-specific Information Security Officers](#), or email the Information Security Office at [itsec-ec@harvard.edu](#)

Don't know who to contact? Contact [Mercè Crosas](#)

**[bit.ly/rdm-survey](https://bit.ly/rdm-survey)**

# Key Resources

— — —

**Research Data Management @Harvard**  
[researchdatamanagement.harvard.edu](https://researchdatamanagement.harvard.edu)

**Harvard University Archives**  
[library.harvard.edu/libraries/harvard-university-archives](https://library.harvard.edu/libraries/harvard-university-archives)

**Office of the Vice Provost for Research | Research Data Security & Management**  
[vpr.harvard.edu/pages/research-data-security-and-management](https://vpr.harvard.edu/pages/research-data-security-and-management)

**Harvard Catalyst | The Harvard Clinical and Translational Science Center**  
[catalyst.harvard.edu](https://catalyst.harvard.edu)

**Office for Scholarly Communications**  
[osc.hul.harvard.edu/policies](https://osc.hul.harvard.edu/policies)

# Sources: DMP Examples

— — —

1. HK-50161-14. University of California, Berkeley. Berkeley Prosopography Services: Implementing the Tool-Kit. *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
2. HD-228971-15. CUNY Research Foundation, Graduate School and University Center. DH Box: A Digital Humanities Laboratory in the Cloud. *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
3. HD-51674-13. University of Virginia. “Are We Speaking in Code?” (Voicing the Craft & Tacit Understandings of Digital Humanities Software Development). *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
4. HD-228966-15. Ohio State University. Automatic Music Performance Analysis and Comparison Toolkit (AMPACT). *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>

# Sources: DMP Examples

— — —

5. HD-229071-15. Vanderbilt University. Deep Mapping the Reduccion: Building a Platform for Spatial Humanities Collaboration on the General Resettlement of Indians. *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
6. HD-229062-15. Georgia State University Research Foundation, Inc. Notoriously Toxic: Understanding the Language and Costs of Hate and Harassment in Online Games. *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
7. HD-229002-15. University of Utah. Poemage Prototype. *Data Management Plans From Successful Grant Applications (2011 - 2014)* <https://www.neh.gov/about/foia/library>
8. Example Data Management Plan: Biology (2). New England Collaborative Data Management Curriculum. Editor: Lamar Soutter Library, University of Massachusetts Medical School. <https://library.umassmed.edu/resources/necdmc/dmp>
9. Example Data Management Plan: Chemistry. New England Collaborative Data Management Curriculum. Editor: Lamar Soutter Library, University of Massachusetts Medical School. <https://library.umassmed.edu/resources/necdmc/dmp>