# Research Data Management

## Data Skills: Planning for Research Success

# Instructors

----

**Julie Goldman**

Research Data Services Librarian
Countway Library of Medicine
Julie_Goldman@hms.harvard.edu

**Meghan Kerr**

Archivist and Records Manager
Center for the History of Medicine
Meghan_Kerr@hms.harvard.edu

Slides: bit.ly/rdm2018

Harvard Medical School | Data Management Working Group

Countway Library of Medicine
An Alliance of the Harvard Medical School and Boston Medical Library

Center for the History of Medicine

Harvard Chan Bioinformatics Core

hms | hsdm
office for postdoctoral fellows

HARVARD MEDICAL SCHOOL | OFFICE FOR Academic and Research Integrity

Department of Systems Biology

HARVARD MEDICAL SCHOOL
Research Information Technology Solutions - RITS

HMS Information Technology

ICCB-Longwood Screening Facility

DRSC/TRiP Functional Genomics

The Neurobiology Imaging Facility
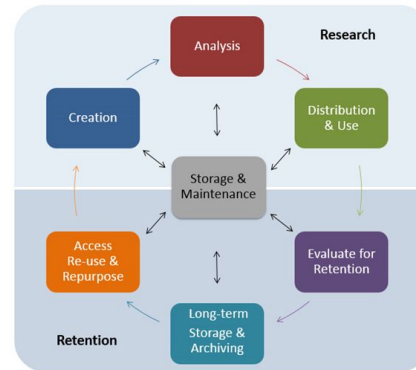in the Neurobiology Department of Harvard Medical School

Hi+S
Harvard Program in Therapeutic Science

# Harvard Biomedical Data Management Website

https://datamanagement.hms.harvard.edu

# Research Data Management Checklist

https://datamanagement.hms.harvard.edu/hms-data-lifecycle

# Introduce Yourself!

– – – –

**Name**

**School / Department**

**Most common data format**

*(Text, Excel, SPSS, Google Docs, etc.)*

# Data Lifecycle for Biomedical Data

# Storage affects the whole cycle

# Why Manage Data?

----

- Easier to analyze organized, documented data

- Avoid accusations of fraud & misconduct

- Don't lose data

- Find data more easily

- Get credit for your data

- Don't drown in irrelevant data



Data Sharing and Management Snafu in 3 Short Acts

https://datamanagement.hms.harvard.edu/biomedical-data-management-planning

# Data Management Plan

————

A data management plan (DMP) is a written document that describes the data you expect to acquire or generate during the course of a research project, how you will manage, describe, analyze, and store those data, and what mechanisms you will use at the end of your project to share and preserve your data.



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

https://datamanagement.hms.harvard.edu/biomedical-data-management-planning-0

# 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.



https://dmptool.org

https://datamanagement.hms.harvard.edu/biomedical-data-management-planning-0

# Data
____

**Raw data**: What is being measured or observed? This is the data that is being generated during the research project.

**Processed data**: How can the raw data be made useful- able to be manipulated?

**Analyzed data**: What does the data tell us? Is it significant? How so?

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

https://datamanagement.hms.harvard.edu/metadata-overview

## Creation

✓ Raw data

✓ Working files

## Analysis

✓ Analytical methods

✓ Analysis results

# Metadata

----

"Good metadata is standardized, consistent and interoperable, and facilitates discovery, preservation and archiving of data."



https://www.ersa.edu.au/understanding-metadata

https://datamanagement.hms.harvard.edu/metadata-overview

# On-Your-Own Exercise
# Documentation

For your most common data type, make a list of the most important information to record for each dataset.

# File Conventions

----

## Versioning

- For analyzed data use version numbers
- Save files often to a new version
- Label the final version FINAL
- For code, consider GIT or SVN

## Organization

- Any system is better than none
- One project, one folder
- Separate folders for data or project stages
- Date-based folders (pairs well with lab notebook)

# File Conventions

----

**Files with naming conventions:**

20161104_ProjectA_Ex1Test1_SmithE_v1.xlsx

20180204-ProjectA-Report-SmithE-v5-FINAL.docx



"FINAL".doc

FINAL.doc!

FINAL_rev.2.doc

FINAL_rev.6.COMMENTS.doc

FINAL_rev.8.comments5.CORRECTIONS.doc

track changes

FINAL_rev.18.comments7.corrections9.MORE.30.doc

FINAL_rev.22.comments49.corrections.10.#@$%WHYDID ICOMETOGRADSCHOOL????.doc

WWW.PHDCOMICS.COM

http://phdcomics.com/comics/archive.php?comicid=1531

# Document Your Conventions

––––

**No point to have a system without documentation!**

- README.txt (use .txt over .doc because it's more durable)

- Front cover of research notebook

- A printout by the computer

# README File

____

- Basic project information

- Title, Contributions, Grant Info

- Contact information

- All locations of where data live,
  including backups

- Useful information about the files
  and how they are organized

- Explain file naming conventions
  and abbreviations

https://datamanagement.hms.harvard.edu/readme-files



Example Template: http://data.research.cornell.edu/content/readme

**On-Your-Own Exercise**

# Conventions

Develop a file naming convention for your most common data type.

# Storage

----

**Storage, backup, and security are interrelated**



## Storage & Maintenance

✓ Store on appropriate tier, with proper security

✓ Store locally on servers or in the cloud
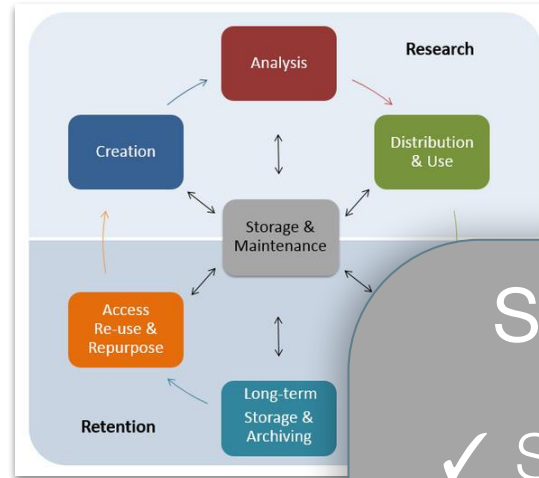
✓ Plan to maintain system

# Security

————

## Access

Limiting the availability of your data

## Systems

Protecting your hardware and software

## Data Integrity

Ensure that your data is not manipulated in an unauthorized way

| 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 |

# Electronic Lab Notebook Matrix

https://datamanagement.hms.harvard.edu/electronic-lab-notebooks

## On-Your-Own Exercise
# Storage

1. Conduct a quick inventory of your data:

   - *What datasets do you have?*
   - *How big are they?*

2. Inventory where your files are currently stored, including backups:

   - *How safe are your data?*

3. Do you have any PHI or HRCI data?

   - *What do you need to ensure their security?*

# Ownership

————

Do you know who owns your data or the dataset you are using?



Who owns your data? (Hint: It's not you)

# Data Sharing

----

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.

### Distribution & Use

✓ Share data with collaborators

✓ Annotate datasets & upload to public repositories

✓ Include in relevant publications & reports

# Citation & Attribution

**give credit
get credit
cite data**

————

Acknowledgement of the use of someone else's information or work is a long-accepted practice in scholarly communication.

The following elements are generally considered the core elements of a data citation:

- *Author/Creator(s): creators of the data; can be one or more people or organizations*

- *Title: title of the data set*

- *Version: exact version or edition of the data set used*

- *Publication Date: date when the data set was published or released*

- *Publisher/Archive: data center or repository that is archiving and distributing*

- *Identifier/Locator: URL or other linkable locator for the data; a persistent, permanent URL such as a DOI (Digital Object Identifier) or a handle is preferred*

**FORCE11 Joint Declaration of Data Citation Principles:** https://doi.org/10.25490/a97f-egyk

# Unique Identifiers

----

**Digital Object Identifier**

Permanently assigned to an object to provide a resolvable persistent network link to current information about that object, including where the object, or information about it, can be found on the Internet

**Open Researcher and Contributor ID**

Provides a persistent digital identifier that distinguishes you from every other researcher and supports automated linkages between you and your professional activities ensuring that your work is recognized

# On-Your-Own Exercise
## ORCiD

Don't have an ORCID?

Create one now!

It's free, easy and will last throughout your professional career!

https://orcid.org

# Data Repository Comparison Matrix
https://datamanagement.hms.harvard.edu/overview-data-repositories

**On-Your-Own Exercise**
# Repositories

Consider your grant funding or project goals:

- Are you required to deposit your data in a repository?

- What repository(ies) will work for your dataset?

# Research Records

## Four Types of Records



Active → Inactive → Destroyed / Archived

# Retention

----

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*

Evaluate for Retention

✓ Identify and retain <u>essential</u> research records

✓ Organize and annotate appropriately

# Appraisal & Archiving

----

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*

**Long-term Storage & Archiving**

✓ In compliance with HMS & federal policy

✓ As requested by investigators

https://datamanagement.hms.harvard.edu/data-evaluation-appraisal

# Questions?

# Open Online Course via Canvas
http://bit.ly/HMS-RDM-MOOC

# Upcoming Seminars

----

**Getting Started with Data Management Plans**

Monday, July 23
12:30 - 1:20 pm
HSPH FXB Building Room G12

**Register:** http://bit.ly/RDM-7-23

**Tips and Tools for Data Storage at Harvard**

Wednesday, August 8
12:30 - 1:20 pm
HSPH FXB Building Room G12

**Register:** http://bit.ly/RDM-8-8

bit.ly/rdm-survey

# Key Resources

————

**Harvard Biomedical Data Management**
http://datamanagement.hms.harvard.edu

**Center for the History of Medicine | Archives and Records Management**
https://www.countway.harvard.edu/chom/archives-and-records-management

**Research Information Technology Solutions**
http://rits.hms.harvard.edu

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

**Harvard Catalyst | The Harvard Clinical and Translational Science Center**
http://catalyst.harvard.edu

**Office for Scholarly Communications**
https://osc.hul.harvard.edu/policies