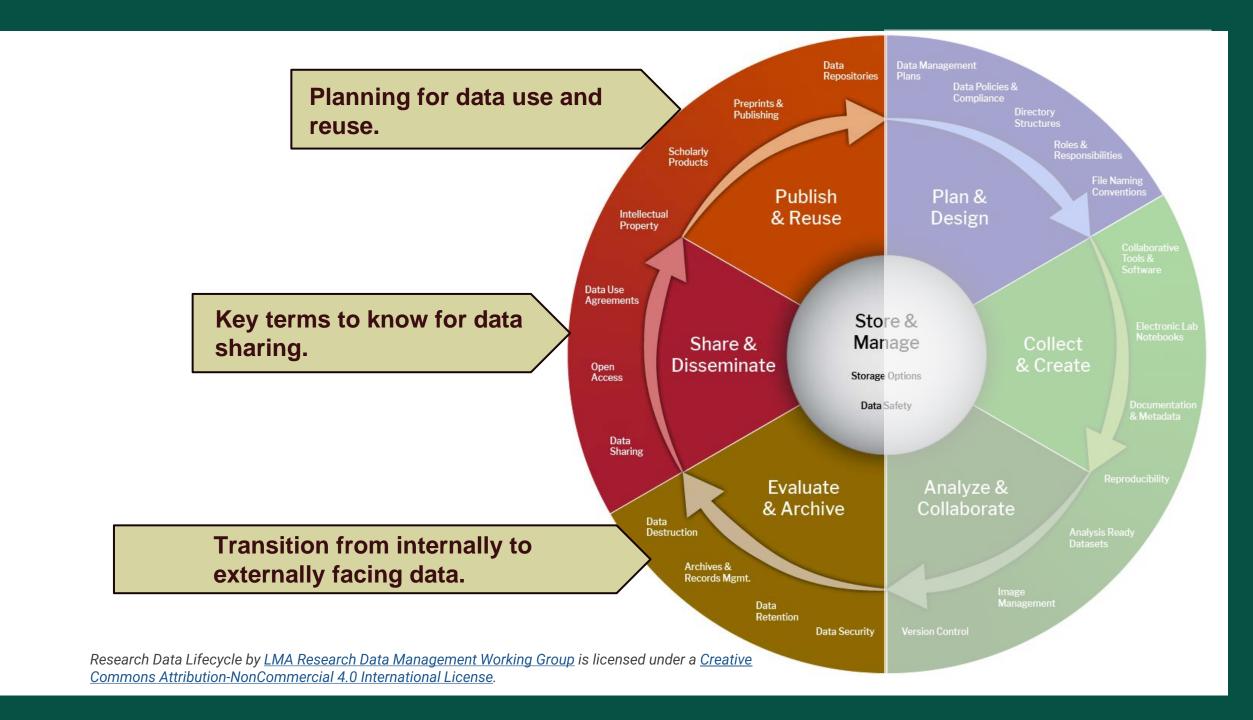
Part II

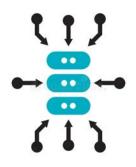
Best Practices: Publishing, Disseminating, and Maintaining Data



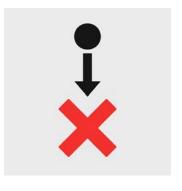




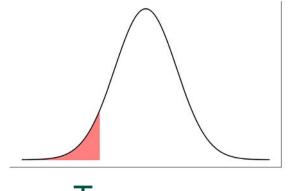
Internal to External: Human Subjects Considerations



Aggregation



Suppression



Top or Bottom Coding



Restricted
Data
Release

Internal to External: Data Destruction and Retention



- Know the timeline for data destruction and retention.
- Automate data destruction if possible.
- Know how to destroy and document the destruction.
- Many options. No one solution. Once the data is external, it cannot truly be destroyed.

Internal to External: Documentation

- Avoid internal jargon!
- Ensure any cross-tabulations or values reported use the public version of the data.
- Provide examples for suggested use:
 - Consider examples across a wide range of software packages.
 - Choose examples for essential data manipulation and best usage practice.

```
%>% group_by(var1 = haven::as_factor(var1))
     %>% summarize(var1 = weighted.mean(var1, PERWT)
STATA
 svyset cluster [pweight=perwt], strata(strata)
  svy: mean var1
 proc sort data = datasetname;by strata cluster;
  proc surveymeans data = datasetname;
  weight perwt;
  strata strata;
 cluster cluster;
 var var1;
```

Source:

https://usa.ipums.org/usa/complex_survey_vars/userNotes_variance.shtml

Internal to External: Project Branding



National Neighborhood Data Archive





- Consider a project name and identity.
- Identify ways to connect to users.
 - A project website
 - Social media accounts
 - Email list serve

Key Terms for Public Release: The Data Management and Sharing Plan (DMSP)

- The externally facing version of a data management plan.
- Their structure, content and length are determined by the funder.
 - NSF Guidelines
 - NIH Guidelines
- Unlike the DMP, the DMSP is a static document. This will be submitted with a grant application.

DATA MANAGEMENT AND SHARING POLICY

GENOMIC DATA SHARING POLICY

PUBLIC ACCESS POLICY

Home > Data Management and Sharing Policy > Planning & Budgeting for Data Management and Sharing > Writing a Data Ma

Writing a Data Management & Sharing Plan

Learn what NIH expects Data Management & Sharing Plans to address, as well as how to submit your Plan.

Applications for Receipt Dates
BEFORE Jan 25 2023

Applications for Receipt Dates ON/AFTER Jan 25 2023

Applications for Receipt Dates ON/AFTER Jan 25 2023

ON THIS PAGE:

- Submitting Data Management and Sharing Plans
- Data Management and Sharing Plan Format
- Sample Plan
- Revising Data Management and Sharing Plans
- Additional Considerations

DMSP Templates

Example: NIH Data Management and Sharing Plan

OMB No. 0925-0001 and 0925-0002 (Rev. 07/2022 Approved Through 01/31/2026)

DATA MANAGEMENT AND SHARING PLAN

If any of the proposed research in the application involves the generation of scientific data, this application is subject to the NIH Policy for Data Management and Sharing and requires submission of a Data Management and Sharing Plan. If the proposed research in the application will generate large-scale genomic data, the Genomic Data Sharing Policy also applies and should be addressed in this Plan. Refer to the detailed instructions in the application guide for developing this plan as well as to additional guidance on sharing nih.gov. The Plan is recommended not to exceed two pages. Text in italics should be deleted. There is no "form page" for the Data Management and Sharing Plan. The DMS Plan may be provided in the format shown below.

Public reporting burden for this collection of information is estimated to average 2 hours per response, including the time for reviewing instructions, searching existing data sources, gathering, and maintaining the data needed, and completing and reviewing the collection of information. An agency may not conduct or sponsor, and a person is not required to respond to, a collection of information unless it displays a currently valid OMB control number. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to: NIH, Project Clearance Branch, 6705 Rockledge Drive, MSC 7974, Bethesda, MD 20892-7974, ATTN: PRA (0925-0001 and 0925-0002). Do not return the completed form to this address.

Element 1: Data Type

- A. Types and amount of scientific data expected to be generated in the project: Summarize the types and estimated amount of scientific data expected to be generated in the project,
- B. Scientific data that will be preserved and shared, and the rationale for doing so:

 Describe which scientific data from the project will be preserved and shared and provide the rationale for this decision.
- C. Metadata, other relevant data, and associated documentation: Briefly list the metadata, other relevant data, and any associated documentation (e.g., study protocols and data collection instruments) that will be made accessible to facilitate interpretation of the scientific data.

Example: NSF Data Management and Sharing Plan

DATA MANAGEMENT AND SHARING PLAN¹

Given that most proposed research involves the generation of scientific data, proposals are subject to the NSF Data Sharing Policy and require submission of a Data Management and Sharing Plan. Data should be findable, accessible, interoperable, and reusable (FAIR). Refer to the guidance on https://new.nsf.gov/funding/data-management-plan for developing this plan. The Plan should not exceed two pages. The Plan may be provided in the format shown below. Only the sections appropriate for the project need to be completed. A valid Data Management and Sharing Plan may include only the statement that no detailed plan is needed if the statement is accompanied by a clear justification. The Plan is not intended to circumvent the number of pages that are allowed for the Project Description. Text in italic with a yellow background should be deleted before submission.

Element 1: Data Types

A. Types and amount of scientific data expected to be generated in the project: Summarize the types and estimated amount of scientific data expected to be generated in the project.

B. Of the generated scientific data in 1.A, the scientific data that will be preserved and shared, and the associated rationale:

Describe which types of scientific data from the project will be preserved and shared and briefly provide the rationale for this decision.

C. Metadata and associated documentation:

Briefly list the metadata (provides information about the scientific data) and any associated documentation (e.g., study protocols, research methods, and data-collection instruments) that will be made accessible to facilitate interpretation of the preserved and shared scientific data.

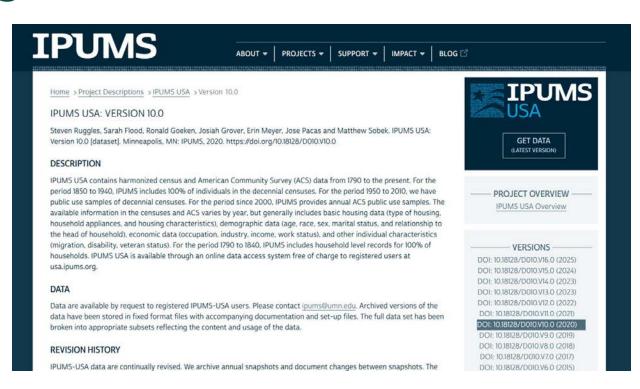
Key Terms for Public Release: Metadata

- Metadata: 'information about our data collections that help others discover, understand, and use them' (Source: ICPSR)
- Should follow specific guidelines such as the <u>Data Documentation Initiative</u> (DDI) and the <u>Dublin Core</u>.
- Why are metadata important?
 - Metadata help your data be FAIR (Findable, Accessible, Interoperable, and Reusable) by helping digital systems classify and identify your data and its key elements

Example of Metadata

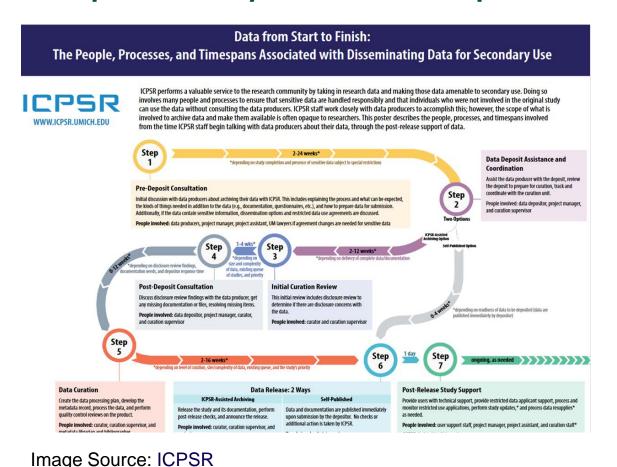
Key Terms for Public Release: Unique Persistent Identifiers

- For data, digital object identifiers (DOIs) are often used. Regulated by the <u>Digital</u> <u>Object Identifier Foundation</u>.
- DOIs provide a persistent and unique way to identify a resource even as technology changes.
- DOIs also can serve as public facing identifiers and project version control.
- Broadly share the doi not the data!



https://doi.org/10.18128/Do10.V10.0

Key Terms for Public Release: Data Repository, Data Deposit, and Data Curation



- Data Repository: 'a centralized place to hold data, share data publicly, and organize data in a logical manner.' Harvard Dataverse
- Data Deposit: the process for submitting data for dissemination through a data depository.
- Data curation: the processing and maintaining of project data.

What is Shared with a Repository?

ICPSR Depositor Checklist

Depends on the repository!

At a minimum:

- Documentation
- Data
- Metadata

Depositor Checklist

Depositors should review this list before depositing. If you have questions or concerns, please email ICPSR-help@umich.edu and we will work with you to resolve any issues. While not all of these items are mandatory, providing more complete information increases the chances of the materials being easily replicated and reused.

Ease of Understanding:

- ☐ File and folder structure named with unique, descriptive titles; if providing multiple datasets, include dataset to documentation legend or name correspondingly
- Any publication or promotion deadlines related to the plan for the release of the data and materials
- ☐ Intended access type; public or restricted access
- If restricted access is needed, plan to discuss dissemination with the archive and collaborate on terms of use agreement
- Have you already distributed data with us? If so, is this deposit an update to the data you supplied in the past? Please provide ICPSR study number and dataset relation.

Complete Documentation:

- □ Original questionnaire
- ☐ Codebook; with unweighted frequencies is best for data comparison
- ☐ Identify and describe computed and derived variables
- ☐ Interviewer instructions
- ☐ User Guide
- ☐ Final Report
- ☐ Citations of related publications the data were used in
- Question text, DDI, or another text-source for the question
- Variable groupings (Especially for large datasets, it is useful to categorize variables into conceptual groupings)

Selecting a Repository: Some Popular Options

ICPSR and OpenICPSR

Harvard Dataverse

Open Science Foundation

Github

PubMed Central

Project or Department Website (Not Recommended!)

Considerations for Selecting a Repository



Funder requirements



Popularity with your target user base



Institutional partnerships for reducing data curation and archival costs



Specialized repository features



Costs

Information to Get from a Funder

Desirable Characteristics for All Data Repositories

When choosing a repository to manage and share data resulting from Federally fu

- Unique Persistent Identifiers: Assigns datasets a citable, unique persistent ic support data discovery, reporting, and research assessment. The identifier po de-accessioned or no longer available.
- Long-Term Sustainability: Has a plan for long-term management of data, inc
 on a stable technical infrastructure and funding plans; and having contingenc
 unforeseen events.
- Metadata: Ensures datasets are accompanied by metadata to enable discove ideally widely used across, the community(ies) the repository serves. Domain generalist repositories.
- Curation and Quality Assurance: Provides, or has a mechanism for others to integrity of datasets and metadata.
- Free and Easy Access: Provides broad, equitable, and maximally open access submission, consistent with legal and ethical limits required to maintain priva data.
- Broad and Measured Reuse: Makes datasets and their metadata available w attribution, citation, and reuse of data (i.e., through assignment of adequate m
- Clear Use Guidance: Provides accompanying documentation describing term data use committee).
- Security and Integrity: Has documented measures in place to meet generally release of data, with levels of security that are appropriate to the sensitivity c
- Confidentiality: Has documented capabilities for ensuring that administrative

NIH Supported Data Repositories



Repositories for Sharing Scientific Data

Detailed NSF Information



Questions to Ask a Repository

- How is metadata handled? Who creates the metadata?
- Does the repository allow an embargo period? How about a limited release?
- How does the repository manage persistent identifiers and version control?
- Does the repository do any data de-identification review or data curation?
- Does the repository support restricted data releases? Who approves researcher access?
- What are the costs associated with the deposit? Are there size limitations?

- Will users have any costs to access the data? use and reuse agreements? Are there user terms of use?
- How is the data stored? What is the long term plan for data preservation?
- What is your user support policy?

A Comparison of Generalist Repositories

Generalist Repository Comparison Chart

of data type, format, content, or disciplinary focus. For this chart, we included a repository available to all searchers specific to dinical trials (Vivil) to bring awareness to those in this field.

TOPIC	HARVARD DATAYERSE. BEPOSITORY	DEYAD	FIGSHARE	MENDELEY DATA	ose	MALI	ZENODO
Brief Description	Historial Debawence Historial points in historial data repository gaper to all researchers, from any disciplient, both made and solubide of the Historial community, where you can share, within a change, and appliant research data.	Doyad is an open data publishing platform and community committed to the open awaitability and routine re-ose of all research data. Dryad fully constess all data and metadata and publishes exclusively under a Creative Common's Vollac Domain License (CCCO).	If igniture is a freely available open data publishing platform for all peasanches where they can share and get credit for all types of achieful polytomer including any file hype from any recommendation of achieful peasanches aftering of larger rates etc.	Mendaley Data is a free repository apecialmed for research data. Search more than 201 million debaseds indexed from 1000 for old data repositories and collect and shore datasets with the research controvally following the IAIR data principles.	OSI to a free and open source project management too linat supports researches throughout their entire project likespoods respect to a per project likespood in open science heat practices.	Will is an independent, non-profit organization that has developed a global data-sharing and analytics patients. Our focus is on sharing individual participant-level state from correplated of initial trails to serve the international research community.	Powering Open Science Itsult on Open Science Built by reservatives for the CERN data carries, whose purposes a long farm preservation of digital objects. CERN matches one of the larged warnful for datasets in the world in 19th-energy physics.
Size limits	No type size limit per clataset. Herward Datoward Repost currently sets a file size limit of 2 MGB.	3000B per chinaset through inconaer sachranakon system and up to 110 with assistance in turn helpilytatashyad org	2008 to their supplier experience of the supplier experience of their aborage in their bagginning at 10008 age to 1018 in particular and System limit of 5 fB this.	10GS per datasat	Projects and childhous projects currently have a 50 GB storage him if it they are public, and 5 GB limit if they are public, and 5 GB limit if they are public, and 5 GB limit if they are provide upload from the rate of 5 GB storage. There is not limit imposed by CBS storage, manual of all or public and content and co	If more than 1 B or along take, much part to use it exposed gives from a consumption of the control of the cont	50GB per defeated, contact us via https:// seroducing/support for ingles i mits
Storage space per researcher	178	Noimt	Notest	Nolimit	Noimt	No imt	Noimt

Questions? Discussion

Heide Jackson Associate Research Professor Maryland Population Research Center heidej@umd.edu

