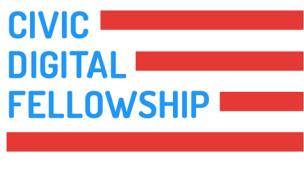


LADDER: Data Management at NIAMS

Scientific Information Technology Branch

Christine Winchester – Intramural Research Project Liaison



RALPH BUAN
Santa Monica College
Interaction Design

MARS IKEDA
Smith College
Statistical & Data Sciences

LADDER

- Laboratories
- Advancing
- Dependable
- Data
- Environments &
- Reuse



DATA MANAGEMENT MATTERS

- **Modernizing infrastructure** at NIAMS requires knowing how much data you have now so you can estimate what you'll need in the future
- Scientific instruments at NIAMS are replaced every few years. **Newer instruments create larger data files**
- NIAMS/NIH attracts researchers from all over the world. High turnover can lead to more **orphaned data (data created by people who are no longer at NIAMS)**

NIH & DATA POLICY

- January 25, 2023: Final NIH Policy for Data Management & Sharing: provides institutional support for Open Science
- **Open Science** is a new approach to the scientific research based on:
 - **cooperative** work
 - **diffusing knowledge** using digital technologies and collaborative tools*

*from FosterOpenScience.uk



OPEN SCIENCE & DATA

- Open Science requires **researchers to create and share high quality data** so it can be found and reused by other researchers
- **Data that is high quality** means:
 - **F**indable
 - **A**ccessible
 - **I**nteroperable
 - **R**Reusable
- You can't make high quality, FAIR data without strong data management at the institute and researcher levels

OPEN SCIENCE BENEFITS

- **Efficiency**
 - Reduce duplication & cost
 - More research from the same data
- **Quality & Integrity**
 - Wider evaluation & scrutiny of data & ideas
 - Early identification of scientific malpractice
- **Innovation**
 - Less delay in reuse of results of research (e.g., articles, data sets) could mean a shorter window between research and new products & services

LADDER ACTIONS

- **Review** previous fellows' findings and next steps
- **Research** data management practices available via NIH resources
- **Analyze** existing data storage utilization at NIAMS
- **Understand** partners' (researchers & technologists) needs. Survey took <5 mins., interviews up to 1 hour
- **Develop & share** solutions

Data Management Survey

Data Organization & Management

15 Responses 03:56 Average time to complete Active Status

[View results](#) [Open in Excel](#)

1. Lab/Group
[More Details](#)

15 Responses

Latest Responses
"Laboratory of Skin Biology"
"SAB/LCTU"
"Translational Genetics and Genomics Section"

- 15 total responses
- Behaviors, pain points and roadblocks, and awareness of data management practices
- Preparation for stakeholder interviews

Stakeholder Interviews

- Talked to **Principal Investigators** and **Core Facilities**
- Follow up on survey answers
- Understand current challenges around data management
- Validate potential solutions and practices



Data Usage Reports

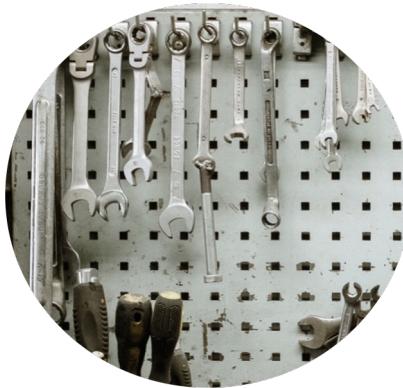
IRP Home - Isilon				
Size	Username	Trimmed User	SurName	Given Name
216G	BONELIMM	BONELIMM	Bonelli	MICHAEL
366K	DAVISFP	DAVISFP	Davi	FRED
933G	KANNOY	KANNOY	Kanno	YUKA
622G	MIKAMIY	MIKAMIY	Mikami	YOHEI
832K	MORIMOTOY	MORIMOTOY	Morimoto	YUKI
243G	MORRISONT/MORRISONT/	Morrison	TASHA	
1.4G	NAGASHIMAI	NAGASHIMAI	Nagashima	HIROYUKI
2.7T	VILLARINOAV	VILLARINOAV	Villarin	ALEJANDRO
236G	YAOC2	YAOC2	Yao	CHEN
104K	dadaho	dadaho	Dada	HANNAH
15G	IRISP	IRISP	Pratt	IRIS
491G	OSHEAOJ	OSHEAOJ	Oshea	JOHN
190G	PLOTZP	PLOTZP	Plotz	PAUL
5.6TB				
HPC Home				
Size	Username			
544G	kannoy			
799G	mikar			
32K	naga			
624K	yaoc2			
1.3TB				
IRP/Data2				
6.7T	ARBD			
907G	ARBD			
7.6TB				
33	Current Data		14.5TB	
34	Orphan Data		3.8TB	
35	Total:		18.3TB	
36				
33	Current Data		14.5TB	
34	Orphan Data		3.8TB	
35	Total:		18.3TB	

IRP Home Folders					
Size	Username	Trimmed User	SurName	Given Name	Email
79G	ACEVEDON	ACEVEDON	Acevedo Luna	NATALIA	natalia.acev@NIAMS.IRP
5.4G	CARETTIG	CARETTIG	Caretti	GIUSEPPINA	giuseppina.c@NIAMS.IRP
36G	CIUFFOLIV2	CIUFFOLIV2	Ciuffoli	VERONICA	veronica.ciul@NIAMS.IRP
1.0T	FENGX4	FENGX4	Feng	XUESONG	fengx4@mai@NIAMS.IRP
267M	KABAMA	KABAMA	Kaba	MAMA	mama.kaba@NIAMS.IRP
192G	KHATEBMN	KHATEBMN	Khateb	MAMDOH	mamduh.kh@NIAMS.IRP
7.7T	KOK3	KOK3	Ko	KYUNG DAE	kyungdae.ko@NIAMS.IRP
234G	RIPARINIG2	RIPARINIG2	Riparini	GIULIA	giulia.riparir@NIAMS.IRP
5.5G	SABBAGHSE	SABBAGHSE	Sabbagh	SARA	sara.sabbag@NIAMS.IRP
99G	SARTOREV	SARTOREV	Sartorelli	VITTORIO	sartorev@m.NIAMS.IRP
6.5G	WANGHONGJ	WANGHONGJ	Wang	HONG JUN	wanghongj@NIAMS.IRP
3.6G	lopezchm	lopezchm	Lopez	CHRISTOPHER	christopher.l@NIAMS.IRP
3.5G	suzawam2	suzawam2	Suzawa	MASATAKA	masataka.su@NIAMS.IRP
9.3TB					
30					
31	Current		32.6TB		
32	Orphan		4.3TB		
33	Total:		36.9TB		
34					
35					
31	Current	32.6TB			
32	Orphan	4.3TB			
33	Total:	36.9TB			
34					
35					

Insights & Solutions



Clear the Clutter
Campaign



Toolbox



CORES Data
Portal

Clear the Clutter Campaign: Insights

Researchers' Blockers	IT's Solutions
Lack of resources: Staff & Time	Divide work into incremental steps so work can be done in stages
Apprehension: Taking the first step	IT deletes data declared of no value at beginning of project
Apprehension: “Making things worse”	Logical structure from easiest to manage data (e.g., oldest) to most difficult (e.g., most recent)

Clear the Clutter Campaign

Stage 1

Files 10+ years old

- Deleted

Stage 2

Last accessed 7-10 years ago

- Assessed
- Deleted
- Archived

Stage 3

All remaining files created 3-7 years ago

- Assessed
- Deleted
- Archived

Stage 4

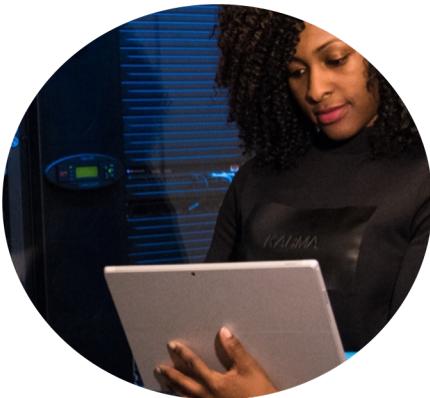
Orphan Data

- Claimed
- Deleted
- Archived

LADDER Toolbox



Exercises and
Workshops



Data Flow
Diagrams



Best Practices
and Resources

LADDER TOOLBOX

Data Management Plans

Effective January 25, 2023, all research funded or conducted by the NIH that result in the generation of scientific data will require the submission of a data management and sharing plan. It should outline how scientific data will be managed and shared throughout the lifecycle of a research project.

The [NIH Final Policy for Data Management and Sharing](#) emphasizes the importance of good data management practices, promotes accountability and transparency, and establishes the expectation for sharing scientific data from NIH-funded or conducted research.

The NIH recommends addressing all of the elements below in your plan within 2 pages or less.

Data type

- Describe the types and estimated amount/size of scientific data to be generated and/or used in the study.
- Decide and indicate which scientific data will be preserved and shared based on legal, ethical and technical factors that may affect the extent to which scientific data will be preserved and shared.
- List metadata, other relevant data, and any associated documentation.

Tools, software, or code

- Indicate which specialized tools are needed to access or manipulate shared scientific data to support replication or reuse.
- Specify how needed tools can be accessed and, if known, whether such tools are likely to be available for as long as the scientific data remain available.

Standards

- Indicate what standards will be applied to the scientific data and associated metadata.

Preservation, access and timelines

- List the name(s) of the repository(ies) where scientific data and metadata arising from the project will be archived.
- Specify how the scientific data will be findable and identifiable.
- Identify any differences in timelines for different subsets of scientific data to be shared.

Access, distribution, and reuse considerations

- Describe any applicable factors affecting access, distribution, or reuse of scientific data related to: informed consent, privacy and confidentiality protections, whether access to scientific data derived from humans will be controlled, and any legal restrictions.

Oversight of data management or sharing

- Indicate how compliance with the plan will be monitored and managed, frequency of oversight, and by whom.

File Naming Conventions

Establish a naming convention for your research files that is short, consistent, and descriptive to help you stay organized and identify your files when you need them. Do this before you start generating or collecting files and write it down so others in your lab can follow this standard.

Best Practices

- Avoid special characters or spaces
- Use capitals and underscores instead of periods, spaces, or slashes
- Use versioning
- Write down naming convention for you and your team
- Put the most important information at the beginning of the file name
- For date format, use YYYYMMDD or YYMMDD

Elements to consider using for your naming convention

- Date of creation or data range of experiment
- Location/lab name
- Project name or acronym
- Sample ID or experiment number
- Version number
- Researcher name or initials
- File Type
- Conditions
- Instrumentation

Steps for making your naming convention:

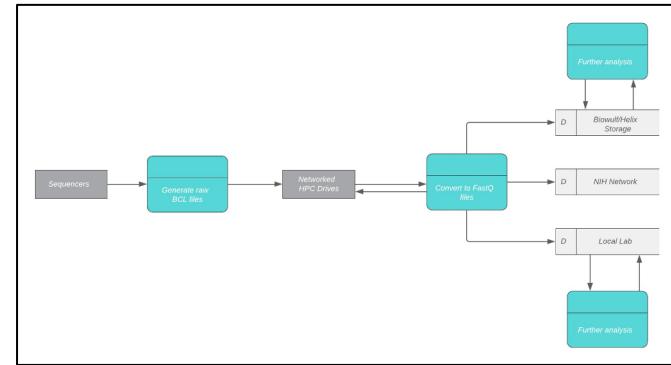
1. Choose a set of files that you'll be making this naming convention for. This could be for an experiment, a particular file type, an instrument(s), etc.
2. Pick 3-5 elements from the above list to include in your convention. Choose elements that will help you differentiate and identify these files.
3. Arrange these elements in a way that makes sense to you. Start with the most important element first.
4. Write it down so that everyone in your team can follow this standard.

Sample naming convention:

Date_ProjectName_FileType_Initials.ext

File name using the above convention:

20210719_DNA_InterviewNotes_RB.docx



FAIR Guiding Principles

The FAIR Data Principles highlight the need to embrace good practice by defining essential characteristics of data objects to ensure that data are reusable by humans and machines.



Findable

Metadata and data should be easy to find by both humans and computers.



Accessible

Once found, users need to know how the data can be accessed.



Interoperable

Metadata use an accessible and standardized language so it can be read and used by a broad range of applications.



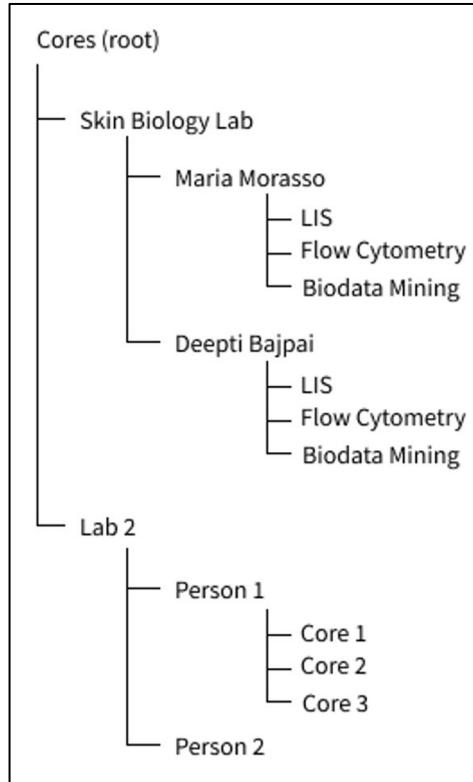
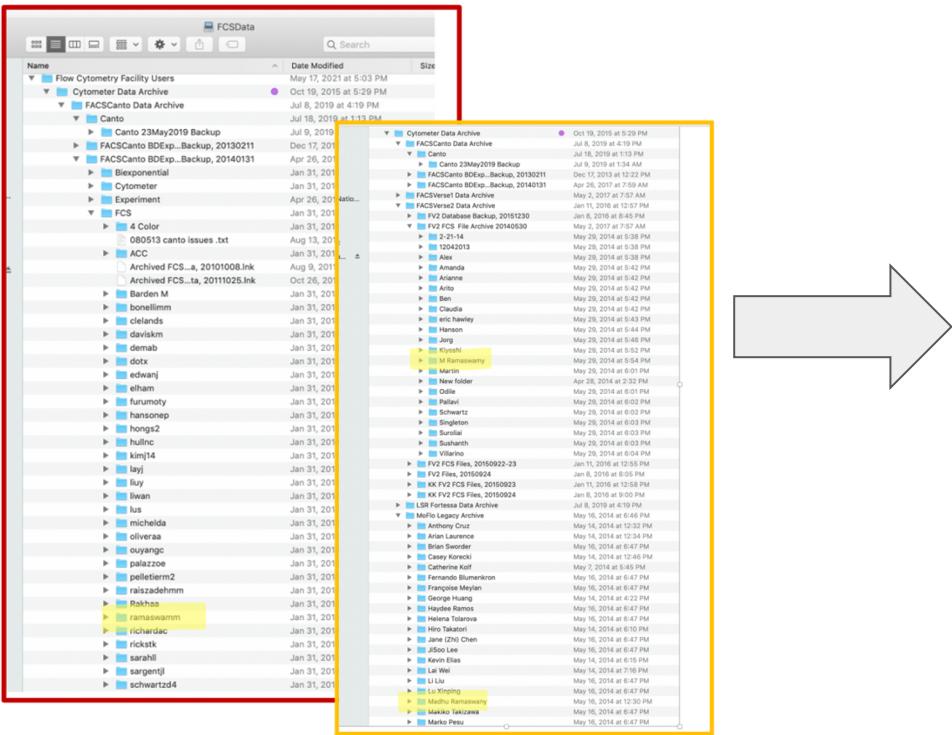
Reusable

Data is well-described and provides clear usage of licenses.

CORES Data Portal: Insights

Researchers' & Technologists' Blockers	IT's Solutions
Don't know who owns data	Organize folders by owner
Processes are hectic & patchwork	All CORES handle data the same way
CORES have low staff	Fewer places to look for data will put less stress on current staff of CORES & researchers

CORES Data Portal



Next Steps

- Start the **Clear the Clutter Campaign**
- Make the **Toolbox a living document & resource**, perhaps stored on Github
- **Build out CORES Data Portal** with appropriate security permissions
- As the 2023 Data Management policy nears, **educate researchers on importance of version control & Open Science**
- Consider forming a Data Management advisory board to help keep lines of communication open

Thank You!

- **Christine Winchester, La'Tanya Burton, Diana Mungai**, and SITB
- **Rachel Dodell** and **Ariana Soto** from Coding it Forward
- **Jacqueline Cattell, Erin Walker**, and **Allissa Dillman**
- **Martyn Green** and **Christine Cutillo**
- **NIAMS Intramural Research Program** PIs and Core Facilities

Questions

