

# From Clinic to Analysis: Why Setting Up and Understanding Your Metadata Matters

---

Presented by:

Nicole Jimenez, PhD

Department of Obstetrics and  
Gynecology Postdoc

2022- 2023 Data Science Fellow

03-24-2023



# Agenda

- Temperature Check
- Data Science/ Open Science
- Data Science Initiatives
- Metadata
- Study Design for Downstream Analysis
- Metadata Examples
- Public Databases
- Inspecting Data
- Tools, Tips, Resources -> I'll send slides so you have these!

# Gauge: Where does the audience background?

Basic Science Researcher (1)



Translational Data Scientist (3)



Clinician/Clinical Team Researcher (2)



Software Developer/  
Informatics Guru (4)



# Gauge: Where does the audience stand on DS?

No background (1)



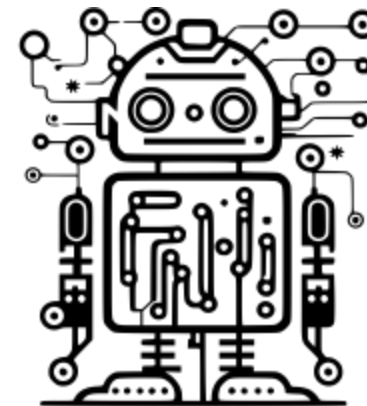
More than Average/ Still learning (3)



Some/ Minimal (2)



I dream in Code and Data (4)





# Gauge: Do you have Datasets to analyze or will be analyzing

Learning/Planning Project (1)



Utilizing Public Data/ Data Organization (3)



Recruitment/Data Collection(2)

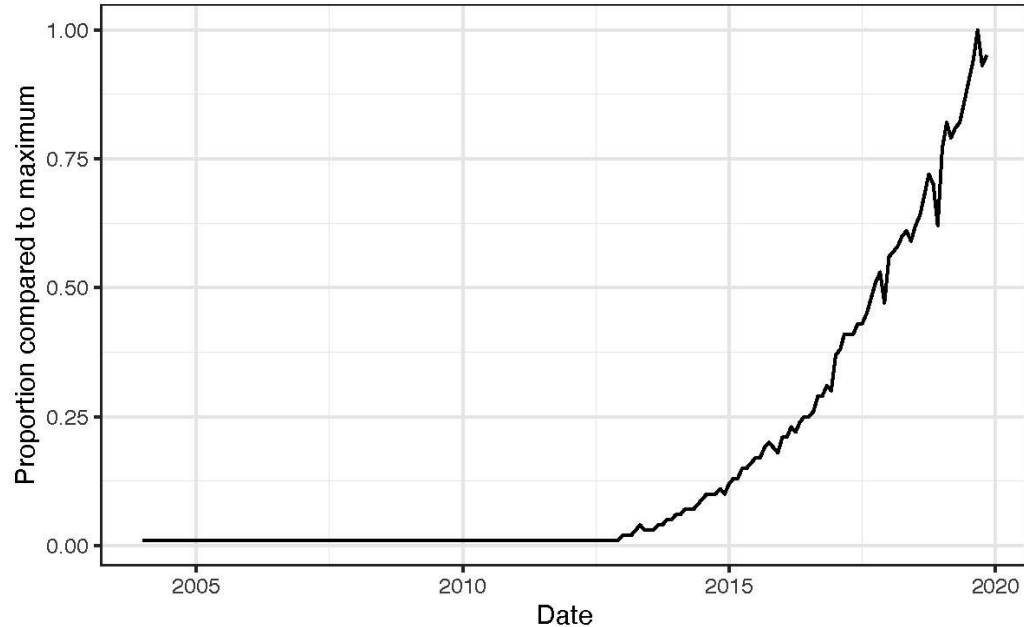


I have Tons of Data and Need to Analyze ASAP (4)



# What is Data Science?

Trends in Google searches of Data Science



Creation of a data-driven workforce; the world is data!

‘data science is an umbrella term to describe the entire complex and multistep processes used to extract value from data.’ - Jeannette Wing , 2019

## Backend data science\*:

- Hardware
- Efficient computing
- Data storage infrastructure
- Data engineering

## Frontend data science\*:

- Data analysis
- Machine learning engineers
- Deep learning developers

\*may depend on field

## Biomedical Research Data:

- Fundamental research using model organisms (such as mice, fruit flies, and zebrafish)
- Clinical studies (including medical images)
- Observational and epidemiological studies (including data from electronic health records and wearable devices).

# Open Science



Open Science (OS) is the movement to make scientific research, data and their dissemination available to any member of an inquiring society, from professionals to citizens. - **Open Responsible research and Innovation** to further **Outstanding Knowledge (ORION)**

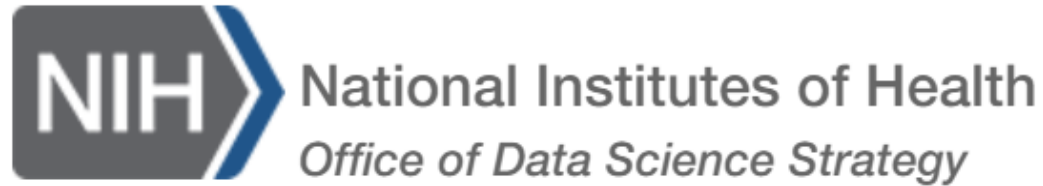
## What is FAIR?

Biomedical research data should adhere to FAIR principles, meaning that it should be Findable, Accessible, Interoperable, and Reusable.

- To be **Findable**, data must have unique identifiers, effectively labeling it within searchable resources.
- To be **Accessible**, data must be easily retrievable via open systems and effective and secure authentication and authorization procedures.
- To be **Interoperable**, data should “use and speak the same language” via use of standardized vocabularies.
- To be **Reusable**, data must be adequately described to a new user, have clear information about data-usage licenses, and have a traceable “owner’s manual,” or provenance.



# Data Science Initiatives in Biomedical Research



Open Knowledge Roadmap

Pathways to Enable Open-Source Ecosystems

Data Infrastructure	Modernized Data Ecosystem	Data Management, Analytics, and Tools	Workforce Development	Stewardship and Sustainability
<ul style="list-style-type: none"><li>•Optimize data storage and security</li><li>•Connect NIH data systems</li></ul>	<ul style="list-style-type: none"><li>•Modernize data repository ecosystem</li><li>•Support storage and sharing of individual datasets</li><li>•Better integrate clinical and observational data into biomedical data science</li></ul>	<ul style="list-style-type: none"><li>•Support useful, generalizable, and accessible tools and workflows</li><li>•Broaden utility of and access to specialized tools</li><li>•Improve discovery and cataloging resources</li></ul>	<ul style="list-style-type: none"><li>•Enhance the NIH data-science workforce</li><li>•Expand the national research workforce</li><li>•Engage a broader community</li></ul>	<ul style="list-style-type: none"><li>•Develop policies for a FAIR data ecosystem</li><li>•Enhance stewardship</li></ul>

**Figure 2.** NIH Strategic Plan for Data Science: Overview of Goals and Objectives

Experiential Learning for Emerging and Novel Technologies (ExLENT) program





# Data Science Initiatives in Biomedical Research



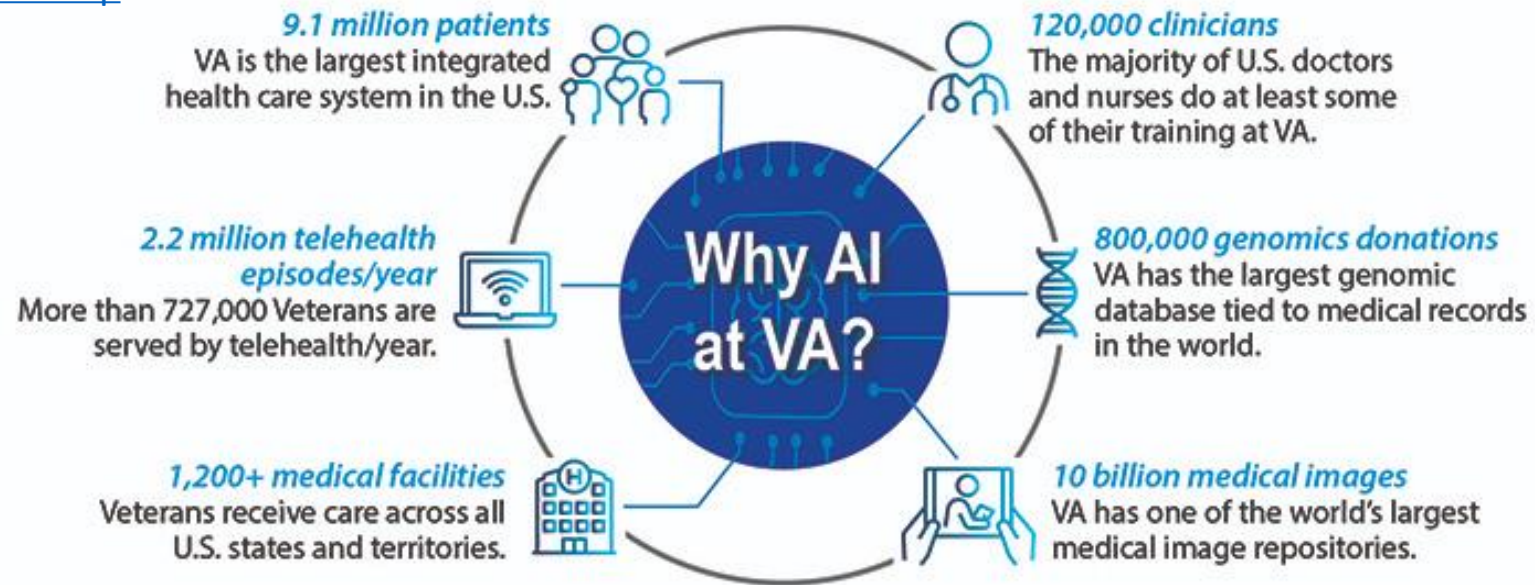
**The Division of Clinical Analytics and Decision Support (CDADS)** provides timely, targeted, meaningful data analytics and clinical decision support systems (CDSS) for the purposes of research, clinical care and quality improvement at the University of Arizona College of Medicine – Phoenix and associated institutions. <https://phoenixmed.arizona.edu/clinicaldata>

**Biomedical Informatics - Clinical Fellowship Program** is designed to effectively prepare fellows for a career in clinical informatics by providing them with the knowledge, skills and attitudes to support informatics-enabled improvement of clinical services.

<https://phoenixmed.arizona.edu/clinical-informatics-fellowship>



**Data Science Institute** aims to foster the next generation of data-driven research by encouraging university-wide interdisciplinary collaboration, gaining external visibility, developing industry alliances, and increasing funding for research at the University of Arizona (UA). <https://datascience.arizona.edu/>



<https://www.research.va.gov/naii/default.cfm>

<https://www.research.va.gov/naii/BD-STEP/>

# What is Metadata?

Metadata, “**data about data**,” provides information such as data content, context, and structure, which is also valuable to the biomedical research community as it affects the ability of data to be found and used. - NIH

One definition commonly applied to the concept of metadata is the simple phrase “**data about data**.” This simplistic definition, however, belies the significance and complexity of the nature of metadata. - American Health Information Management Association

What comes to mind when  
you hear the word  
metadata?

<https://library.ahima.org/doc?oid=106378#.ZBy0PHbMK5d>


[NIH Data Science Strategic Plan](#)

# Types of Metadata?

- **Administrative metadata:** data about a project or resource that are relevant for managing it; E.g. project/resource owner, principal investigator, project collaborators, funder, project period, etc. They are usually assigned to the data, before you collect or create them.
- **Descriptive or citation metadata:** data about a dataset or resource that allow people to discover and identify it; E.g. authors, title, abstract, keywords, persistent identifier, related publications, etc.
- **Structural metadata:** data about how a dataset or resource came about, but also how it is internally structured. E.g. the unit of analysis, collection method, sampling procedure, sample size, categories, variables, etc. Structural metadata have to be gathered by the researchers according to best practice in their research community and will be published together with the data.

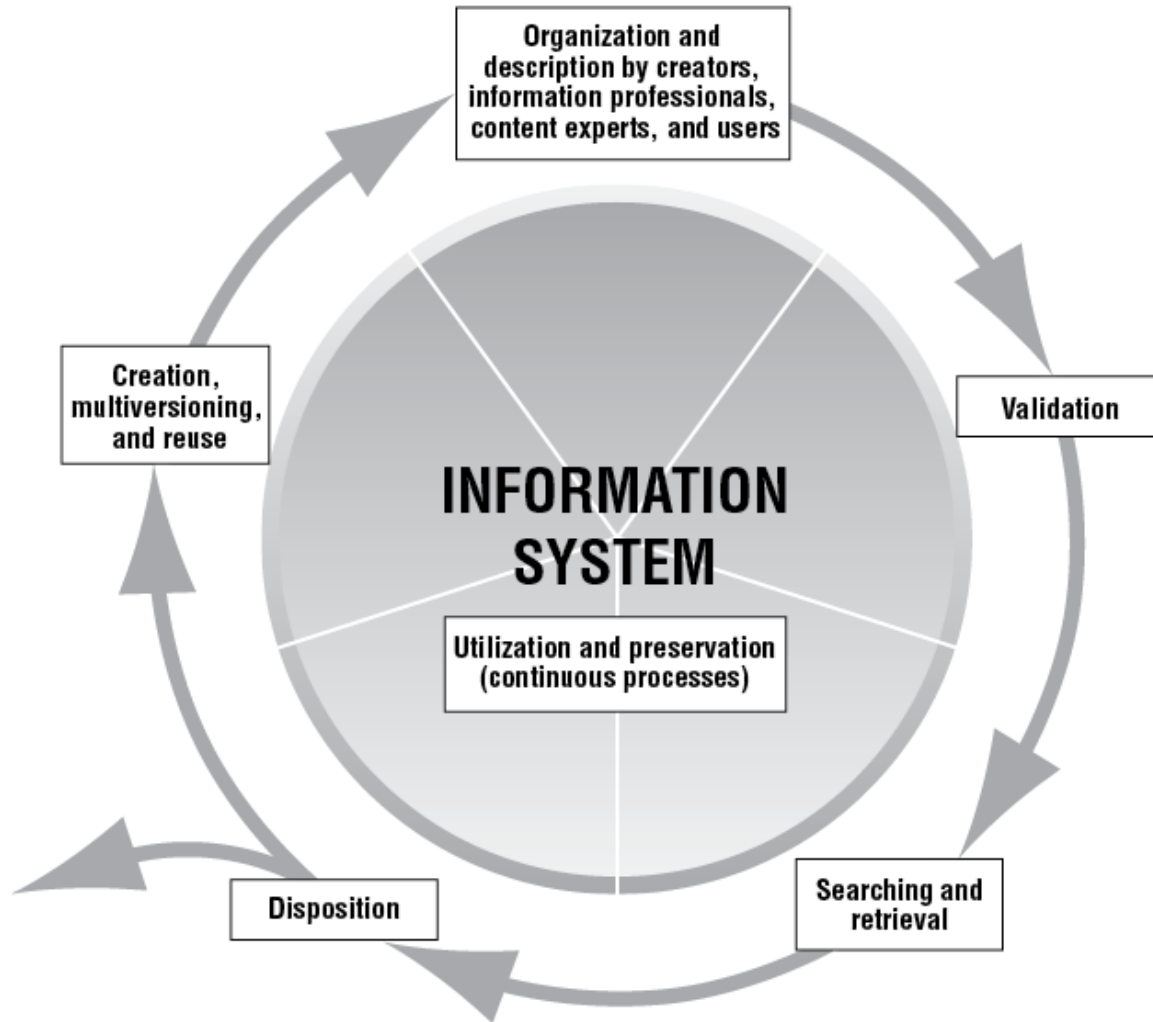
# Types of Metadata? - Exercise

What are the types of metadata observed?

<input type="checkbox"/> Assembly	GenBank	Scientific name ↑	Size (...)	Date	BioProject	Genes
<input type="checkbox"/> ASM1888363v1	GCA_018883635.1	Candidatus Lactobacillus ...	1.748	Jul, 2021	PRJNA543206	1,856
<input type="checkbox"/> BRZ_IL__bin99	GCA_944327185.1	Candidatus Lactobacillus ...	1.546	Jul, 2022	PRJEB53581	
<input type="checkbox"/> ASM1888367v1	GCA_018883675.1	Candidatus Paralactobaci...	1.161	Jul, 2021	PRJNA543206	1,178
<input type="checkbox"/> ASM883148v1 	GCA_008831485.1	Lactobacillus acetotolera...	1.684	Nov, 2019	PRJNA566216	1,658



# Metadata Lifecycle



**Table 2: Metadata in the Lifecycle**

Information Lifecycle Phase	Example Metadata
Creation/Generation	Source, Date created, Time entered, Author, Version number
Classify/Index	File name, Document name
Store/Maintain	Last date accessed, Date archived
Search/View/Share	Patient name, Record number
Secure/Disclosed	Date disclosed, Party disclosed, Status of disclosure
Retain/Preserve/Dispose	Retention date, Disposition of data, Status of record (hold/active/inactive)

# Standardization

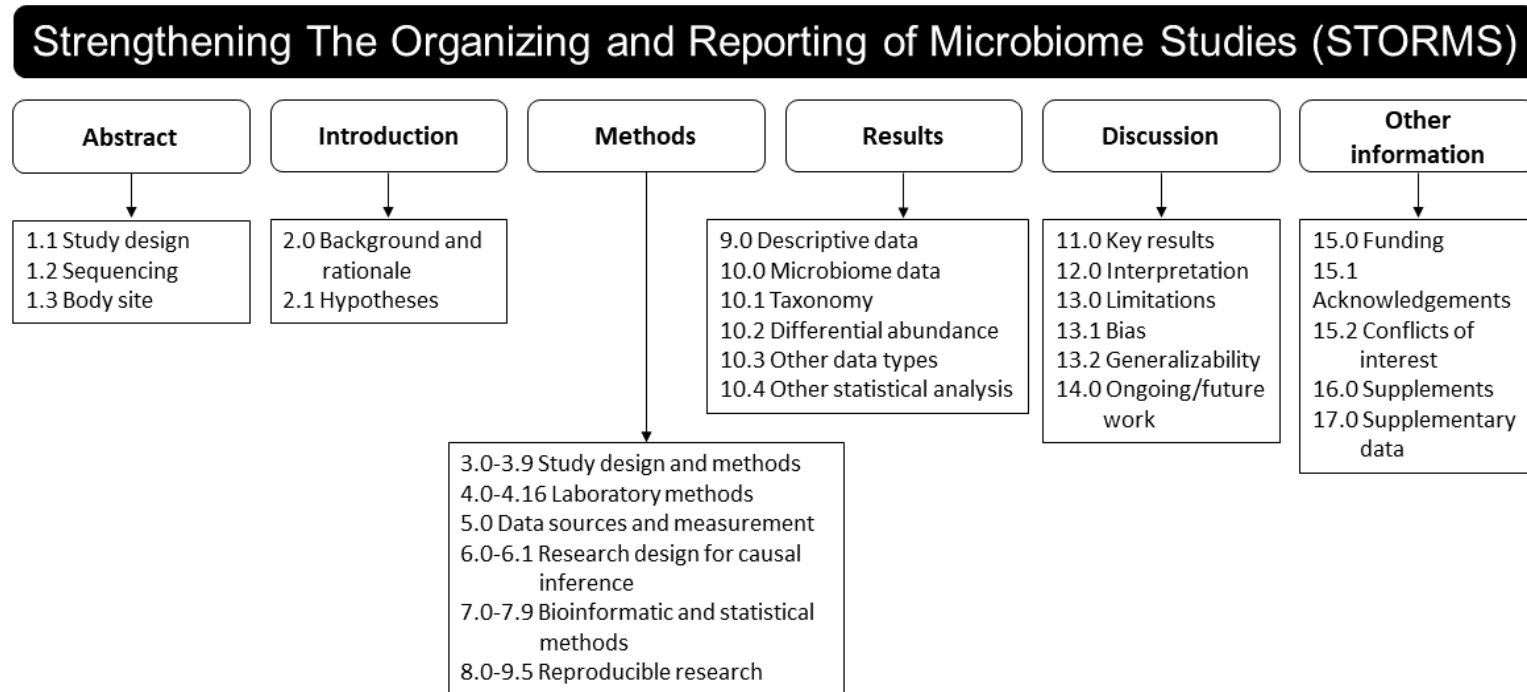
The **minimum information standard** is a set of guidelines for reporting data derived by relevant methods in biosciences. If followed, it ensures that the data can be easily verified, analyzed and clearly interpreted by the wider scientific community. Keeping with these recommendations also facilitates the foundation of structuralized databases, public repositories and development of data analysis tools. Individual minimum information standards are brought by the communities of cross-disciplinary specialists focused on issues of the specific method used in experimental biology.

Standardization of methods, data collection techniques or reporting will ultimately help with better metadata for easier cross-comparison studies or to identify datasets that best meet your research questions.

<https://www.dcc.ac.uk/guidance/standards/metadata>

<https://carpentries-incubator.github.io/fair-bio-practice/05-intro-to-metadata/index.html>

# Standardization



What standards on  
data  
collection/metadata  
reporting does your  
field have?

<https://www.stormsmicrobiome.org/>

<https://datamanagement.hms.harvard.edu/collect-analyze/documentation-metadata>

# Plan Ahead -> "Data Management"

- What is your study question?
- Will there be other factors you need to account for?
- How are you collecting your data?
- Will it be interpretable by computers? -> Interoperable to others
- How will you store your data?

All will help you organize metadata and data!

**Try not throw the kitchen sink at your research cohort or data collection**, this will ultimately lead to more time cleaning the data and take more thought on how to utilize, create stewardship for this data, and best report the data.



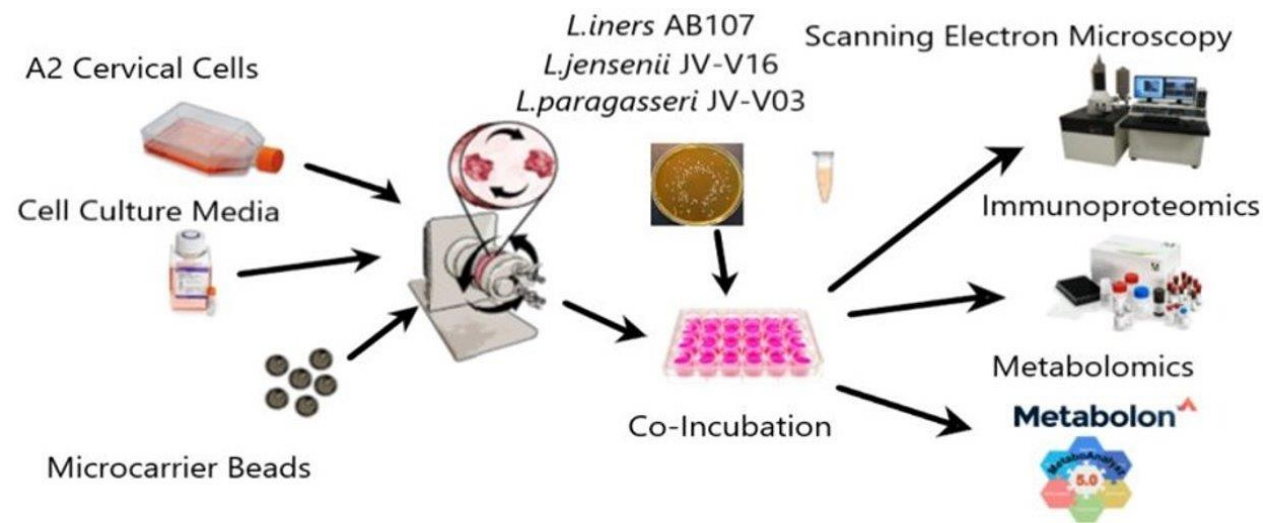
# Plan Ahead -> Example

- **What is your study question?**
  - Does a particular species of bacteria associate with a diagnosis (cervical cancer)?
- **Will there be other factors you need to account for?**
  - Do the patients utilize cigarettes, douching hygiene, have history of BV? Are the patients pregnant or postmenopausal? Do we need to collect HPV status or genotypes?
- **How are you collecting your data?**
  - Clinical Team coordination at clinic, IRB consent, colposcopy biopsy to derive diagnosis group
- **Will it be interpretable by computers? -> Interoperable to others**
  - Yes, and working through it for future use ( two other studies have utilized already)
- **How will you store your data?**
  - Stored at NIH Sequence Read Archive and data use agreement through Women's Health Office due to IRB consent form. Code at Github.

How many of you work in wet  
lab/field research?

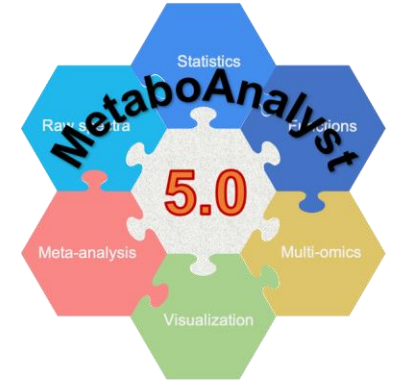
What type of metadata do you encounter?

# From Experiment to Basic Science Metadata Management



[Jimenez et al, Msphere. 2023.](#)

Sample	Treatment	(14 or 15)-	(S)-3-hydr	1,2-diolec
PBS_116	PBS	405,887	76,009	206,658
PBS_118	PBS	847,138	108,050	177,241
PBS_120	PBS	511,835	58,643	274,290
PBS_123	PBS	667,936	64,681	259,518
PBS_124	PBS	392,964	106,190	239,745
PBS_125	PBS	485,556	59,092	110,821
L. paragasseri	L.paragasseri	700,587	94,376	321,208
L. paragasseri	L.paragasseri	985,133	75,839	207,689
L. paragasseri	L.paragasseri	660,800	90,472	410,015
L. jensenii	L.jensenii	527,266	67,944	398,375
L. jensenii	L.jensenii	381,775	90,894	205,589
L. jensenii	L.jensenii	463,027	78,480	219,223
L. jensenii	L.jensenii	247,194	88,932	350,906
L. iners_1	L.iners	506,928	98,945	267,280
L. iners_1	L.iners	492,555	56,428	138,259
L. iners_1	L.iners	591,220		840,044



# From Recruitment to Clinic Metadata Management

How many of you work in clinical research?

What type of metadata do you encounter?



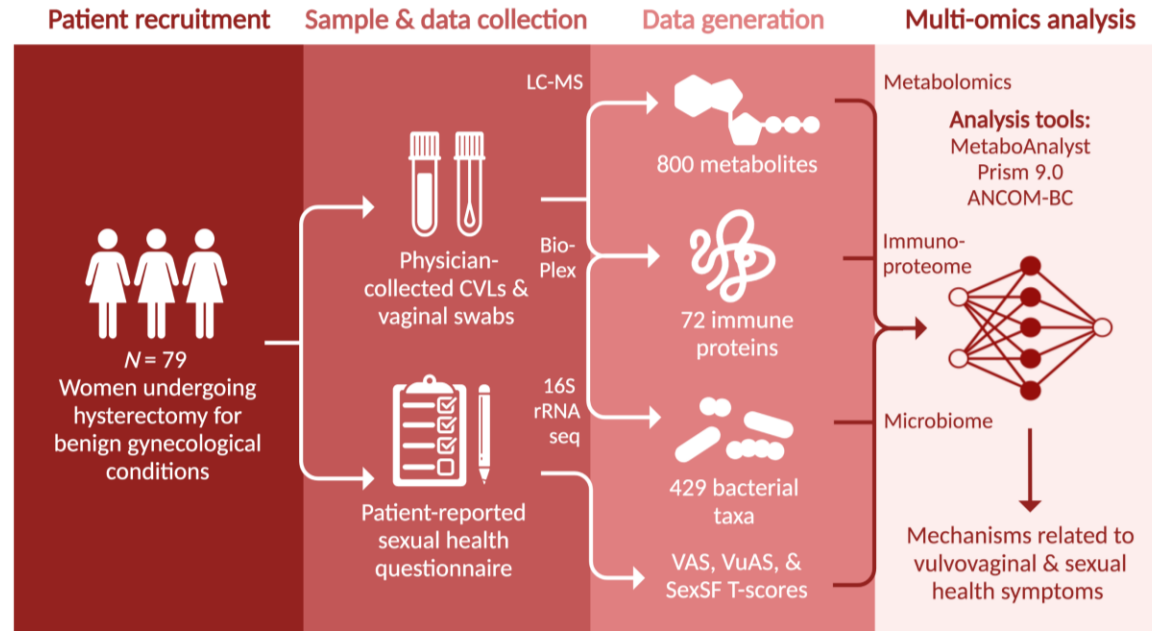
# From Recruitment to Clinic Metadata Management

## Medical/Clinic data:

Diagnosis for study:  
EMC, Hyperplasia,  
Benign condition

Benign condition  
diagnosis: fibroids,  
adenomyosis,  
endometriosis,  
abnormal uterine  
bleeding, etc

Grade, size, location,  
invasion, molecular tests



Crossley et al, unpublished, 2023

## Laboratory Data:

HPLC metabolite  
concentrations

Microbiome 16S read counts

Milliplex – immune marker  
concentrations

## Validated Survey Data:

Vaginal Assessment: EMC, Hyperplasia, Benign condition

Vulvar assessment: fibroids, adenomyosis, endometriosis, abnormal uterine  
bleeding, etc

Sexual health and Satisfaction questionnaire:

# Clinical Data Importance and Data Stewardship

## Clinical Data and Information Security (NIH, but applicable elsewhere):

- Proper handling of the vast domain of clinical data that is being continually generated from a range of data producers is a challenge for biomedical research community.
- Patient-related data is both quantitative and qualitative and can arise from a wide array of sources, including specialized research projects and trials; epidemiology; genomic analyses; clinical-care processes; imaging assessments; patient-reported outcomes; environmental-exposure records; and a host of social indicators now linked to health such as educational records, employment history, and genealogical records.
  - **Key take-homes:**
    - Practice robust and proactive information-security approaches to ensure appropriate stewardship of patient data
    - Curation of authentic, trusted data sources for future research
    - Protect against patient data compromise or loss.

# Clinical Data Importance



<https://www.researchallofus.org/>



<https://www.gida-global.org/care>

# Electronic Health Record (EHR)

## Medical/Clinic data:

Diagnosis for study: EMC,  
Hyperplasia, Benign condition

Benign condition diagnosis:  
fibroids, adenomyosis,  
endometriosis, abnormal uterine  
bleeding, etc

Grade, size, location, invasion,  
molecular tests

TABLE 1. SUGGESTED ELECTRONIC HEALTH RECORD  
CONTENTS BY PROVIDER

<i>Data provider</i>	<i>Kinds of data constituting the electronic health record</i>
Pharmacist	Medicine prescribed (linked to information about chemical composition, interactions, side effects, etc.)
Physician	Dosage Comments about the patient's observable physical condition at time of examination (coded) Diagnosis Treatment plan
Radiologist	Diagnosis (coded) after reading computed tomography scans, magnetic resonance imaging scans, X-rays, etc.
Raw test results	Electronic versions of image-based tests Unanalyzed content from monitoring devices
Clerical staff	Demographics (coded), such as sex, age, blood type, and family history
Nurse	Vital signs at time of examination (coded), including blood pressure and temperature Self-reported symptoms

<https://www.liebertpub.com/doi/10.1089/big.2013.0023>

<https://library.ahima.org/doc?oid=106378#.ZBywGnbMK5c>



# Survey

## Validated Survey Data:

Vaginal Assessment: EMC, Hyperplasia, Benign condition

Vulvar assessment: fibroids, adenomyosis, endometriosis, abnormal uterine bleeding, etc

Sexual health and Satisfaction questionnaire:

How would you set up  
this metadata and  
data?

In last month, XYZ?

<b>Never</b>	<b>Almost Never</b>	<b>Sometimes</b>	<b>Fairly Often</b>	<b>Very Often</b>
<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Series of  
questions then  
total score  
calculated based  
on survey  
requirements

# Recording Metadata

- **README:** A README File is a text file located in a project-related folder that describes the contents and structure of the folder and/or a dataset so that a researcher can locate the information they need.
- **Data Dictionary:** Also known as a *codebook*, a data dictionary defines and describes the elements of a dataset so that it can be understood and used at a later date.
- **Protocol:** A protocol describes the procedure(s) or method(s) used in the implementation of a research project or experiment. If you need to maintain protocols, we strongly recommend a tool like protocols.io.
- **Lab Notebook:** For research groups that use them, Electronic Lab Notebooks offer several advantages over traditional paper notebooks.

How do you currently record metadata?

# Planning ahead with ReDCaP

The University of Arizona Center for Biomedical Informatics and Biostatistics (CB2) is pleased to be a partner in the REDCap Consortium. It supports a secure web application (REDCap) designed exclusively to support data capture for research studies and institutional projects. **The REDCap application allows users to build and manage online surveys and databases quickly and securely.**

The screenshot shows the REDCap form builder interface. On the left, there's a sidebar with 'View & Stats' and 'Feature' sections. The main area displays a form titled 'Demographics Information' with several fields: 'Study ID', 'First Name', 'Date of Birth', 'Last Name', 'Gender', and 'Street, City, State, ZIP'. Each field has an 'Add Field Here' button. The 'Date of Birth' field has a calendar icon and a 'Today' button. The 'Gender' field has a dropdown menu. The 'Street, City, State, ZIP' field has an 'Expand' button.

The screenshot shows a Microsoft Excel spreadsheet with a table of demographic data. The table has columns for Study ID, Date subject signed consent, Last Name, First Name, Street, City, State, ZIP, Phone number, Second phone number, E-mail, Gender, and Has the subject given consent. The data is organized into rows, with the first row being the header. The table contains 11 rows of data.

Study ID	Date subject signed consent	Last Name	First Name	Street, City, State, ZIP	Phone number	Second phone number	E-mail	Gender	Has the subject given consent
1	8/3/1992	Duck	Donald	123 Main St. Orlando, FL 12345	(415) 555-1212	(456) 545-1252	just_ducky@vrdw.com	Male	
2	8/4/1992	House	Hansel	124 Main St. Orlando, FL 12345	(415) 555-1213	(456) 545-1253	lost1@woods.org	Male	
3	8/5/1992	Forest	Gretel	125 Main St. Orlando, FL 12345	(415) 555-1214	(456) 545-1254	lost2@woods.org	Female	No
4	8/6/1992	Mouse	Minnie	126 Main St. Orlando, FL 12345	(415) 555-1215	(456) 545-1255	minnieme@wdw.com	Female	Yes
5	8/7/1992	Mouse	Mickey	127 Main St. Orlando, FL 12345	(415) 555-1216	(456) 545-1256	franchise@wdw.com	Male	
6	8/8/1992	Dog	Goofy	128 Main St. Orlando, FL 12345	(415) 555-1217	(456) 545-1257	maxdad@hyuck.biz	Male	
7	8/9/1992	Dog	Pluto	129 Main St. Orlando, FL 12345	(415) 555-1218	(456) 545-1258	woofn@hotmail.com	Male	
8	8/10/1992	Duck	Daisy	130 Main St. Orlando, FL 12345	(415) 555-1219	(456) 545-1259	daisy@wdw.com	Female	Yes
9	8/11/1992	Castle	Cinderella	131 Main St. Orlando, FL 12345	(415) 555-1220	(456) 545-1260	princess@castle.biz	Female	Yes
10	8/12/1992	Tiger	Jasmine	132 Main St. Orlando, FL 12345	(415) 555-1221	(456) 545-1261	jt314@agrabah.edu	Female	No
11	9/5/2011	Long		122 Main St. Here, TN 37443				Male	

<https://cb2.uahs.arizona.edu/services-tools/surveys-clinical-databases-redcap/redcap-training>

<https://cb2.uahs.arizona.edu/services-tools/surveys-clinical-databases-redcap>

# NIH Repositories

DOMAIN-SPECIFIC REPOSITORIES

GENERALIST REPOSITORIES

DOWNLOAD(.csv)


Domain-Specific Repositories

CLEAR ALL

25 PER PAGE ^

Displaying 1 - 25 of 133 results

NAME/DESCRIPTION	ICO	SUBJECT AREA	MODEL SYSTEM	ACCESS TYPE	PROPERTIES	REPOSITORY LINKS
search name & description	All	All	All	All	All	
<b>Federal Interagency Traumatic Brain Injury Research (FITBIR) Informatics System</b> The Federal Interagency Traumatic Brain Injury Research (FITBIR) informatics system was developed to share data across the entire TBI research field ..More	CIT NINDS	Clinical research Imaging Neuroscience	human	controlled registered	<div>Open data submission</div> <div>Open timeframe for data deposit</div> <div>NIH funding support</div> <div>Sustained support</div>	<div>DATA ACCESS</div> <div>DATA SUBMISSION</div>
<b>Metabolomics Workbench</b> The NIH Common Fund's National Metabolomics Data Repository (NMDR) is now accepting metabolomics data for small and large studies on cells, tissues ..More	Common Fund	Clinical research Computational biology Other	human non-human	open	<div>Open data submission</div> <div>Open timeframe for data deposit</div> <div>NIH funding support</div> <div>Sustained support</div>	<div>DATA ACCESS</div> <div>DATA SUBMISSION</div>
<b>exRNA Atlas</b> Includes exRNA profiles derived from various biofluids and conditions and currently stores data profiled from small RNA sequencing assays.	Common Fund	Clinical research Neuroscience Sequence biology	human non-human	registered open	<div>Open data submission</div> <div>Open timeframe for data deposit</div> <div>NIH funding support</div> <div>Sustained support</div>	<div>DATA ACCESS</div> <div>DATA SUBMISSION</div>
<b>Illuminating Druggable Genome</b> The Pharos interface provides facile access to most data types around proteins collected and harmonized by the project. Pharos integrates data from a ..More	Common Fund	Sequence biology	human	open	<div>Open data submission</div> <div>Open timeframe for data deposit</div> <div>NIH funding support</div> <div>Sustained support</div>	<div>DATA ACCESS</div> <div>DATA SUBMISSION</div>



[HOME](#)
[Taxonomy](#)
[Genomes](#)
[Proteomes](#)
[16S rRNA Microbiome](#)
[Animal Microbiomes](#)
[Downloads](#)
[Resources](#)
[HOMD-v2](#)
[\[page-maps\]](#)

[16S rRNA RefSeq V15.22](#)
[Genomic RefSeq V16.1](#)

## Download HOMD Data

[\[top\]](#)  
[Taxonomy data](#)  
[Abundance data](#)  
[Genomic data](#)  
[Phylo data](#)  
[16S rRNA seqs](#)  
[Database Schema](#)

### Taxonomic Data: Batch Downloads

Type	Formats		
<a href="#">Taxon Table [page]</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
<a href="#">Taxonomic Hierarchy [page]</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
<a href="#">Taxonomic Level [page]</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>

### Abundance Data: Batch Downloads [\[page\]](#)

Type	Formats		
<a href="#">Eren (v1v3)</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
<a href="#">Eren (v3v5)</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
<a href="#">DeWalt</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
<a href="#">Segata</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>

### Genomic Meta Information: Batch Downloads

Type	Formats		
<a href="#">Sequence Meta Information [page]</a>	<a href="#">Tab Delimited Text (View in browser)</a>	<a href="#">Tab Delimited Text (Save to file)</a>	<a href="#">MS Excel Format (Save to file)</a>
NCBI Genome Annotations			
		<a href="#">[FTP Site for download]</a>	
PROKKA Genome Annotations			
		<a href="#">[FTP Site for download]</a>	
Genomic Trees:			
Conserved Protein, Ribosomal Protein and 16S rRNA		<a href="#">[FTP Site for download]</a>	

<https://hombd.org/>

[https://www.nlm.nih.gov/NIHbmic/domain\\_specific\\_repositories.html](https://www.nlm.nih.gov/NIHbmic/domain_specific_repositories.html)

# Public Databases are not Created Equal

- Data Repository submission: What you put in will be what you get out
  - Nicole's example as a comparative genomics Grad student:

Identifiers	BioSample: SAMD00290000; SRA: DRS179169																											
Organism	<a href="#">Bifidobacterium bifidum</a> cellular organisms; Bacteria; Terrabacteria group; Actinomycetota; Actinomycetes; Bifidobacteriales; Bifidobacteriaceae; Bifidobacterium																											
Package	<a href="#">MIGS: cultured bacteria/archaea; version 6.0</a>																											
Attributes	<table><tr><td><b>sample name</b></td><td>B. bifidum BI-28</td></tr><tr><td><b>collection date</b></td><td>2017-02-17</td></tr><tr><td><b>broad-scale environmental context</b></td><td>missing</td></tr><tr><td><b>local-scale environmental context</b></td><td>human-associated habitat</td></tr><tr><td><b>environmental medium</b></td><td>feces</td></tr><tr><td><b>geographic location</b></td><td><a href="#">Japan</a></td></tr><tr><td><b>host</b></td><td>Homo sapiens</td></tr><tr><td><b>isolation and growth condition</b></td><td><a href="https://doi.org/10.1038/s41396-021-00937-7">doi.org/10.1038/s41396-021-00937-7</a></td></tr><tr><td><b>latitude and longitude</b></td><td>missing</td></tr><tr><td><b>number of replicons</b></td><td>1</td></tr><tr><td><b>project name</b></td><td>Infant gut SCFA and microbiota relationship</td></tr><tr><td><b>reference for biomaterial</b></td><td>missing</td></tr><tr><td><b>strain</b></td><td>BI-28</td></tr></table>		<b>sample name</b>	B. bifidum BI-28	<b>collection date</b>	2017-02-17	<b>broad-scale environmental context</b>	missing	<b>local-scale environmental context</b>	human-associated habitat	<b>environmental medium</b>	feces	<b>geographic location</b>	<a href="#">Japan</a>	<b>host</b>	Homo sapiens	<b>isolation and growth condition</b>	<a href="https://doi.org/10.1038/s41396-021-00937-7">doi.org/10.1038/s41396-021-00937-7</a>	<b>latitude and longitude</b>	missing	<b>number of replicons</b>	1	<b>project name</b>	Infant gut SCFA and microbiota relationship	<b>reference for biomaterial</b>	missing	<b>strain</b>	BI-28
<b>sample name</b>	B. bifidum BI-28																											
<b>collection date</b>	2017-02-17																											
<b>broad-scale environmental context</b>	missing																											
<b>local-scale environmental context</b>	human-associated habitat																											
<b>environmental medium</b>	feces																											
<b>geographic location</b>	<a href="#">Japan</a>																											
<b>host</b>	Homo sapiens																											
<b>isolation and growth condition</b>	<a href="https://doi.org/10.1038/s41396-021-00937-7">doi.org/10.1038/s41396-021-00937-7</a>																											
<b>latitude and longitude</b>	missing																											
<b>number of replicons</b>	1																											
<b>project name</b>	Infant gut SCFA and microbiota relationship																											
<b>reference for biomaterial</b>	missing																											
<b>strain</b>	BI-28																											

## Metagenome or environmental sample from human gut metagenome

Identifiers	BioSample: SAMN27962527; Sample name: PRL2010-MUC-C; SRA: SRS12856508													
Organism	<a href="#">human gut metagenome</a> unclassified entries; unclassified sequences; metagenomes; organismal metagenomes													
Package	<a href="#">Metagenome or environmental; version 1.0</a>													
Attributes	<table><tr><td><b>host</b></td><td>not applicable</td></tr><tr><td><b>isolation source</b></td><td>not applicable</td></tr><tr><td><b>collection date</b></td><td>2022-03-08</td></tr><tr><td><b>geographic location</b></td><td><a href="#">Italy</a></td></tr><tr><td><b>latitude and longitude</b></td><td>not collected</td></tr><tr><td><b>source_organism</b></td><td>Bifidobacterium bifidum</td></tr></table>		<b>host</b>	not applicable	<b>isolation source</b>	not applicable	<b>collection date</b>	2022-03-08	<b>geographic location</b>	<a href="#">Italy</a>	<b>latitude and longitude</b>	not collected	<b>source_organism</b>	Bifidobacterium bifidum
<b>host</b>	not applicable													
<b>isolation source</b>	not applicable													
<b>collection date</b>	2022-03-08													
<b>geographic location</b>	<a href="#">Italy</a>													
<b>latitude and longitude</b>	not collected													
<b>source_organism</b>	Bifidobacterium bifidum													
BioProject	<a href="#">PRJNA833139</a> Retrieve <a href="#">all samples</a> from this project													

<https://www.nature.com/articles/sdata201921>



# Public Databases are not Created Equal

- Data Repository submission: What you put in will be what you get out

Figure 3: Metadata submissions to NCBI BioSample from 2009–2017.

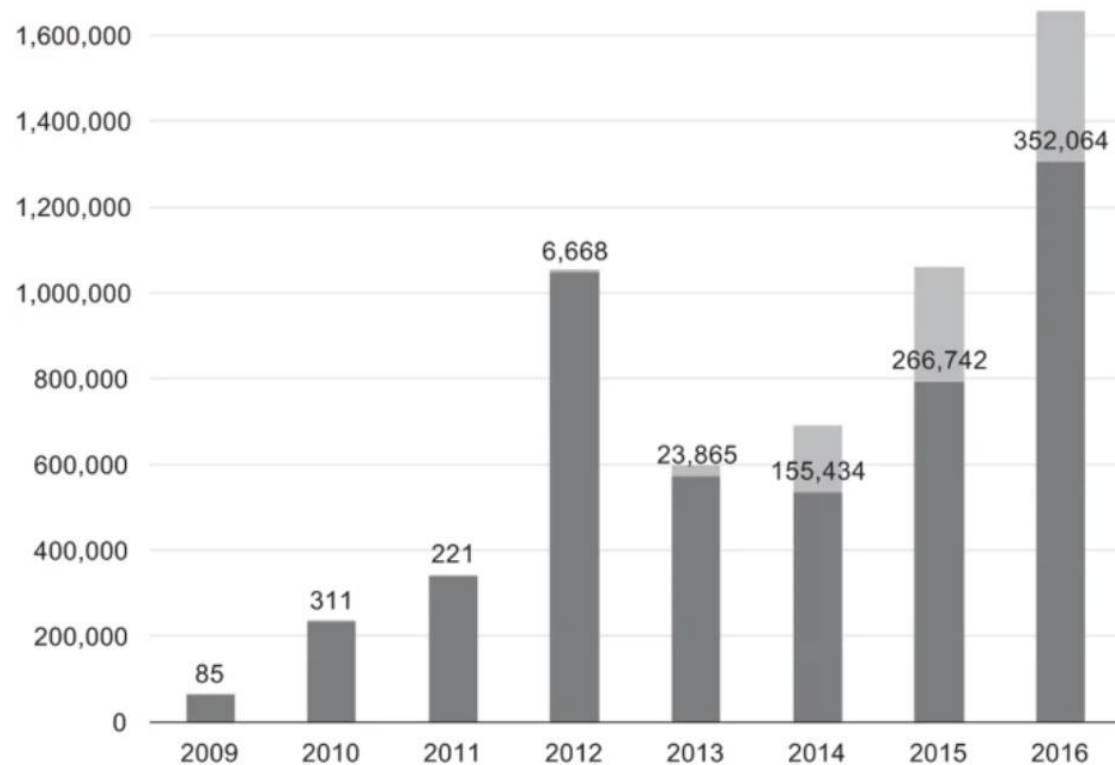
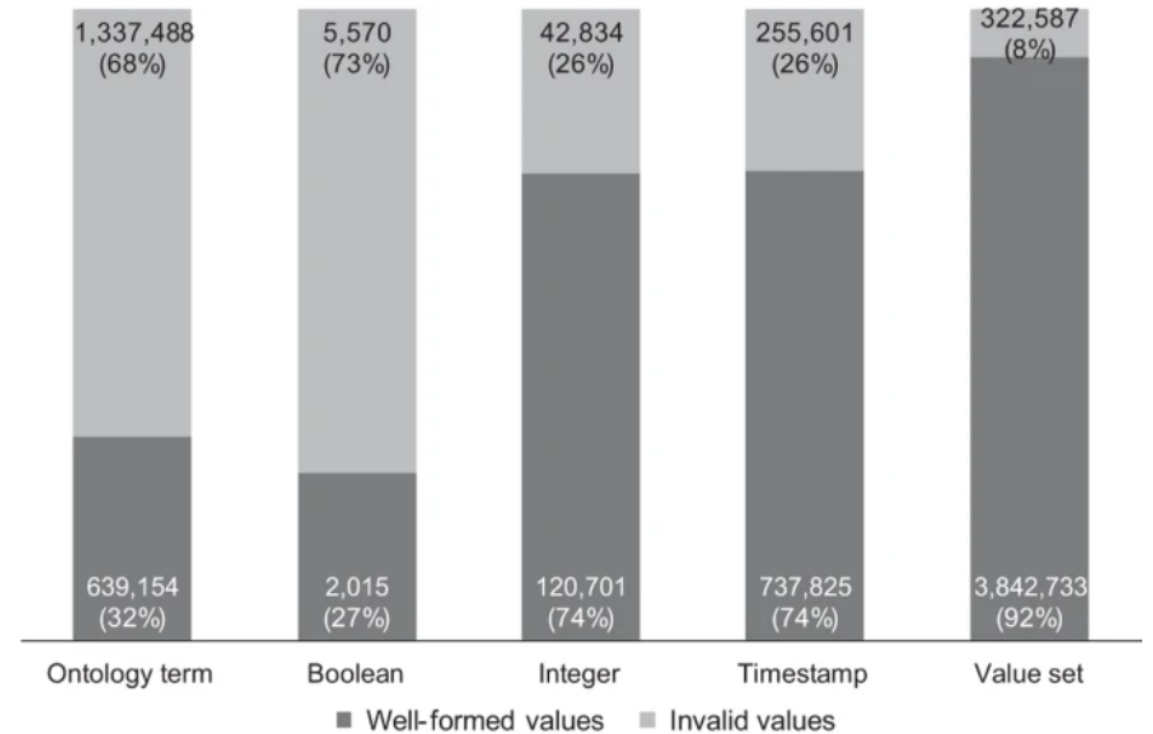


Figure 4: Quality of dictionary attributes in NCBI BioSample according to their type.



<https://www.nature.com/articles/sdata201921>

# Public Database Inspection

## Considerations:

- Does this dataset fit my needs? answer my research question?
- Do I trust the data from this dataset? What are variables or metadata that might key me in on the validity?
- What is the format? Data type?
- Is it in a format that I can run basic statistics or run exploratory analyses to investigate the data set? Is the data complete?

**There are many resources here at UA and on GitHub to help you inspect your data, Check out the calendar and office hours of DSI to figure out best ways to investigate your data!**

# Preparation of Data/Metadata

**Data also exist in a wide variety of formats**, which complicates the ability of researchers to find and use biomedical research data generated by others and creates the need for extensive data “cleaning.”

According to a 2016 survey, data scientists across a wide array of fields said they spend most of their **work time (about 80 percent) doing what they least like to do: collecting existing datasets and organizing data.**

That leaves less than **20 percent of their time for creative tasks like mining data for patterns** that lead to new research discoveries.

# Data Cleaning

## Box 1. Terms Related to Data Cleaning

**Data cleaning:** Process of detecting, diagnosing, and editing faulty data.

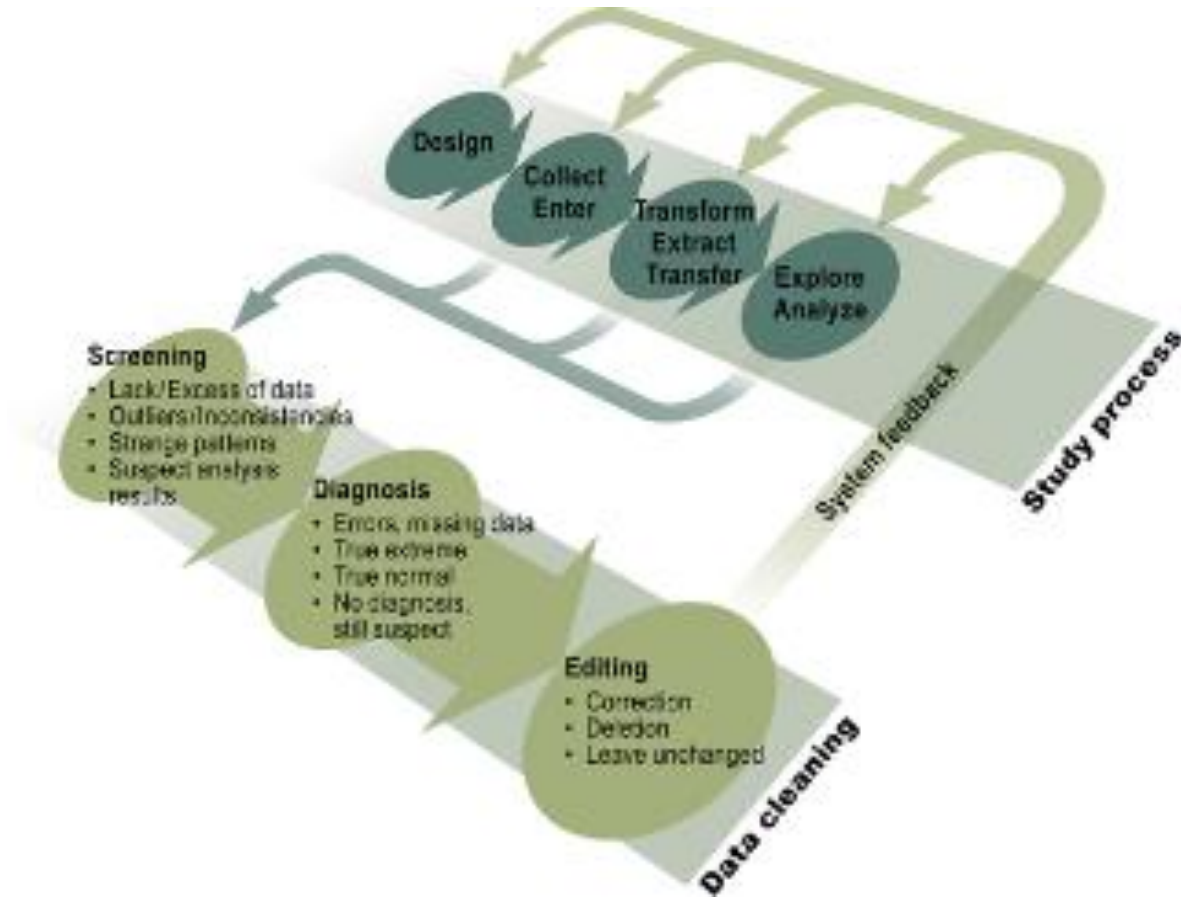
**Data editing:** Changing the value of data shown to be incorrect.

**Data flow:** Passage of recorded information through successive information carriers.

**Inlier:** Data value falling within the expected range.

**Outlier:** Data value falling outside the expected range.

**Robust estimation:** Estimation of statistical parameters, using methods that are less sensitive to the effect of outliers than more conventional methods.



<https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.0020267>

<http://varianceexplained.org/r/tidy-genomics/>

<https://careerfoundry.com/en/blog/data-analytics/what-is-data-cleaning/>

# Formatting for Analysis

Excel is easy to use, flexible and powerful, however, it often gives us too much freedom which leads to bad practices and difficult to re-use data and metadata.

## Items to Keep in the Computer's Mind:

- Using multiple tables in a sheet
- Using multiple tabs
- Not filling in true zeros
- Using problematic null values
- Using formatting(highlighting, bolding, merging columns) to convey information
- Placing comments or units in cells
- Entering more than one piece of information in a cell
- Inconsistency in used values
- Using problematic field names( beware the space)
- Using special characters in data
- Inspect Data name for misnumbering or dates
- Etc....

**Table 1.** Commonly used null values, limitations, compatibility with common software and a recommendation regarding whether or not it is a good option. Null values are indicated as compatible with specific software if they work consistently and correctly with that software. For example, the null value "NULL" works correctly for certain applications in R, but does not work in others, so it is not presented in the table as R compatible.

Null values	Problems	Compatibility	Recommendation
0	Indistinguishable from a true zero		Never use
Blank	Hard to distinguish values that are missing from those overlooked on entry. Hard to distinguish blanks from spaces, which behave differently.	R, Python, SQL	Best option
-999, 999	Not recognized as null by many programs without user input. Can be inadvertently entered into calculations.		Avoid
NA, na	Can also be an abbreviation (e.g., North America), can cause problems with data type (turn a numerical column into a text column). NA is more commonly recognized than na.	R	Good option
N/A	An alternate form of NA, but often not compatible with software		Avoid
NULL	Can cause problems with data type	SQL	Good option
None	Uncommon. Can cause problems with data type	Python	Avoid
No data	Uncommon. Can cause problems with data type, contains a space		Avoid
Missing	Uncommon. Can cause problems with data type		Avoid
-,+,,	Uncommon. Can cause problems with data type		Avoid



# Formatting for Analysis

## Check your File type and your software dependencies!!!

- Plain text files like comma or tab separated values (.csv, .tsv) can be accessed without any special software.
- If you analyse your data with R or Python, or you know that your data are meant to be processed that way you should be using text formats whenever possible, and as soon as you capture your data.
- if you only use Excel and so does your community, just keep using it. Just be aware of the possible pitfalls discussed, especially when working with gene or protein names and accession numbers.

If utilizing a Graphical User Interface or specialized software, check the sample data files this will save time and help you figure out what formats may be required if not explicitly stated

# Overview Metadata? - Exercise

- What contextual details (metadata) are needed to make your data meaningful?
- What form or format will the metadata describing your data take? Which metadata standards will you use? If there is no applicable standard, how will you describe your data in a way that will make them accessible to others?
- How will metadata files be generated for each of the data sets that you produce? Who will do the work of data description?
- Who on your team will be responsible for ensuring that metadata standards are followed and are correctly applied to the corresponding data sets?

# Data Science Fellowship or Ambassador Programs



Each semester, a new cohort featuring up to 12 fellows will be offered. Fellows are expected to attend and participate in twice weekly virtual training activities. **Fall 2023 semester** participants will be awarded a \$7,000 stipend for successful completion of the program.

<https://datascience.arizona.edu/education/data-science-fellows>

## 2022-2023 Ambassadors



**Chen Chen**

College of  
Humanities

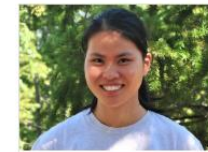
[chenc33@arizon...](mailto:chenc33@arizona...)



**Rongbo Jin**

School of  
Government and  
Public Policy

[rongbojin@arizon...](mailto:rongbojin@arizon...)



**Lia Ossanna**

School of Natural  
Resources and the  
Environment

[lossanna@arizona...](mailto:lossanna@arizona...)



**Salena Torres  
Ashton**

School of  
Information

[salena@arizona.e...](mailto:salena@arizona.e...)



**Zisu Wang**

Eller College of  
Management

[zisuwang@arizon...](mailto:zisuwang@arizon...)



**Sarah Yates**

College of Medicine

[syates@arizona.e...](mailto:syates@arizona.e...)

In short, DSAs are champions for data science literacy in your college!

<https://datascience.arizona.edu/dsa>

# Data Science on Phoenix Campus

5th Annual Research Symposium

## reimagine HEALTH

Translating Big Science and Big Data into Better Health Care

January 27, 2023



  
**Workshop - Data Science on the Phoenix Biomedical campus!**

Let's talk Data Science on the UA Phoenix Biomedical Campus!



**REGISTER FOR THIS WORKSHOP**  
**November 8th and December 7th**  
**5:30-6 p.m.**  
HSEB room A446, UA Health Sciences Library - Phoenix. Enter via the library main entrance on the 3rd floor and then take the stairs or the elevator within the library to reach room A446.  
Or join virtually with the zoom link that you will receive upon registration.

**YOU ARE NOT ALONE!**  
Join other like-minded researchers and Luisa Rojas, a Data Science Institute Data Science Fellow and CTS doctorate student, as she hosts a series of engaging workshops dedicated to open science principles: Data Management Plans and Open Data Science tools for the lab - right here at the UA Phoenix Biomedical Campus.

Bring your questions and suggestions and, of course, feel free to bring a friend!



## Future Programming focused on DS Initiatives

- 1 Computational and Experimental Research
- 2 Clinical and Translational Research
- 3 Epidemiology and Population-Based Research
- 4 Research Using Data Generated Outside a Research Context

Booz Allen identified five competencies necessary for making data FAIR and AI-ready:

Competency	Description
Dataset Documentation	<ul style="list-style-type: none"><li>Documentation of study design and data elements, including methodology, study population description, inclusion/exclusion criteria, sampling procedures, and expected relationships between variables through use of conceptual modeling or knowledge graphs.</li><li>Ethical, legal, and social implications of poor or incomplete dataset documentation on the analysis of current studies and potential secondary usage.</li></ul>
Ontology Usage & Data Encoding	<ul style="list-style-type: none"><li>Domain-specific standardized ontologies and when to request the addition of new content, attachment of granular and accurate metadata to datasets through use of secure data/metadata entry software (such as REDCap).</li><li>Self-identification in human-subject data collection (i.e., free response of race/ethnicity, sex/gender, and disability status) and compliance with federal data collection standards).</li></ul>
Data Cleaning & Formatting	<ul style="list-style-type: none"><li>Importing of structured and unstructured data into a coding environment, removal of identifiable information from datasets before sharing, and transformation of data into a structured, machine-readable format that is cleaned, accurate, and consistent.</li><li>Ethics of handling self-identification in human-subject data collection including 'roll-up' to federal classification standards and interpretation of classification categories in findings, as well as an understanding that these are evolving concepts that may not be accurately/consistently reflected in older datasets.</li></ul>
Data Governance	<ul style="list-style-type: none"><li>Creation of long-term management plan for data, including curation, sharing, access, reuse, and archiving/preservation.</li><li>Human-subject data storage and sharing standards, including the need to de-identify or perturb data to remove PII or PHI prior to sharing or release; and a sufficient understanding of study context and design to identify and document potential threats to inference.</li></ul>
Data Collaboration	<ul style="list-style-type: none"><li>Facilitation of interdisciplinary collaboration across biomedical and data science research teams during iterations of study design, data collection, and model development/analysis, which requires sufficient knowledge of AI concepts to assist biomedical researchers with algorithmic data needs and sufficient knowledge of biomedical concepts to convey to AI developers relevant and precise biomedical contextual information, to the degree that biomedical and AI researchers can generate and assess clinically meaningful results.</li><li>Avoidance of or mitigation for damaging ethical implications in deployed health models, such as (a) poor data collection and aggregation methods, namely misappropriation of data collected in non-research efforts without consideration of intended collection context and/or misapplication to target populations without accounting for differences between study and target populations; (b) inadequately documented and/or labeled historical or discriminatory bias in input data, and (c) insufficiently rigorous explorations of the effects that findings from algorithmic deployments could have on patients if incorporated into clinical workflows, including the difference between statistically and clinically significant findings and future model retraining due to increasing noncomparability between pre- and post-treatment datasets.</li></ul>

We want to hear from you, what do you want to see on the PBC!

VentureCafe: Hackathon  
April 6th 11:30 am to 5 pm  
[PBC information](#)

Viewing/Panel discussion in the works!



# Data Science Institute Calendar

🏠 [Home](#) > Calendar

## Calendar

MAR 22	<b>Classical Machine Learning (ML)</b> 3 to 4 p.m., March 22, 2023
MAR 23	<b>Data Science Tapas</b> 3 to 4 p.m., March 23, 2023
MAR 24	<b>Learning through Open Science Contributions</b> 10 a.m. to noon, March 24, 2023
MAR 24	<b>Data Science Workshop Series on the Phoenix Biomedical Campus</b> 11:30 a.m. to 1 p.m., March 24, 2023
MAR 28	<b>Data Wrangling in R with the Tidyverse</b> 9 to 11 a.m., March 28, 2023
MAR 28	<b>R Workshop</b> 9 to 11 a.m., March 28, 2023
MAR 28	<b>Data &amp; Viz Drop-in</b> 9 to 11 a.m., March 28, 2023

TODAY
THIS WEEK
THIS MONTH
◀ March 2023 ▶
S M T W T F S
26 27 28 1 2 3 4
5 6 7 8 9 10 11
12 13 14 15 16 17 18
19 20 21 22 23 24 25
26 27 28 29 30 31 1

Category

- Any - ▼

**Need help with your data?**  
Check out Data Science Institute Calendar, there is likely a workshop or office hours to assist in your analysis journey!

<https://datascience.arizona.edu/calendar>



# Data Science Resources at University of Arizona

## [Foundational Open Science Skills](#)

[Open Science Framework](#)

[Roots for Resilience](#)

## [Data Management](#)

[Soteria](#)

[Cyverse](#)

[ReData](#)

## [Learn to Code](#)

[rstudio-connect](#)

[PlanetMicrobe](#)

[University of Arizona - High Performance Computing](#)

## [UA Data Science Slack Channel](#)

[ResBaz](#)

[Women in Data Science - Tucson](#)

# Acknowledgements



RESEARCH, INNOVATION & IMPACT  
**Data Science Institute**

- **Funding**

- Data Science Institute – Data Science Fellows
- BIO5 Institute – Postdoctoral Fellowship
- Community Foundation for Sothern Arizona



- **Herbst-Kralovetz Lab Members and Clinical Team**

- **College of Medicine-Phoenix - Department of Obstetrics and Gynecology**

- **Early Investigator Support**

- Postdoctoral Affairs Office
- UA-CoM-P Research Office,
- Native American Cancer Prevention – GUIDeS: U54CA143924, U54CA143925



# References/resources

- <https://www.getty.edu/publications/intrometadata/setting-the-stage/>
- <https://guides.lib.umich.edu/c.php?g=283277&p=1888478>
- <https://www.ncbi.nlm.nih.gov/sra/docs/submitmeta/>
- <https://www.ncbi.nlm.nih.gov/grc>
- <https://metadata-wizard-tutorial.readthedocs.io/en/latest/>
- <https://databrowser.researchallofus.org/>
- <https://datacarpentry.org/>
- [https://cran.r-project.org/web/packages/eatGADS/vignettes/meta\\_data.html](https://cran.r-project.org/web/packages/eatGADS/vignettes/meta_data.html)
- [https://cran.r-project.org/doc/contrib/de\\_Jonge+van\\_der\\_Loo-Introduction\\_to\\_data\\_cleaning\\_with\\_R.pdf](https://cran.r-project.org/doc/contrib/de_Jonge+van_der_Loo-Introduction_to_data_cleaning_with_R.pdf)
- <https://allofus.nih.gov/>
- <https://hdsr.mitpress.mit.edu/pub/gg6swfqh/release/2>

# Hands-On:

There are many ways to clean/inspect data:

- Python, R, etc
- For today's example we will utilize a tool recommended by Data Carpentries called OpenRefine: <https://openrefine.org/>
- Activity: <https://datacarpentry.org/OpenRefine-ecology-lesson/>
- Follow the tutorial or try it out with your data

OpenRefine

