

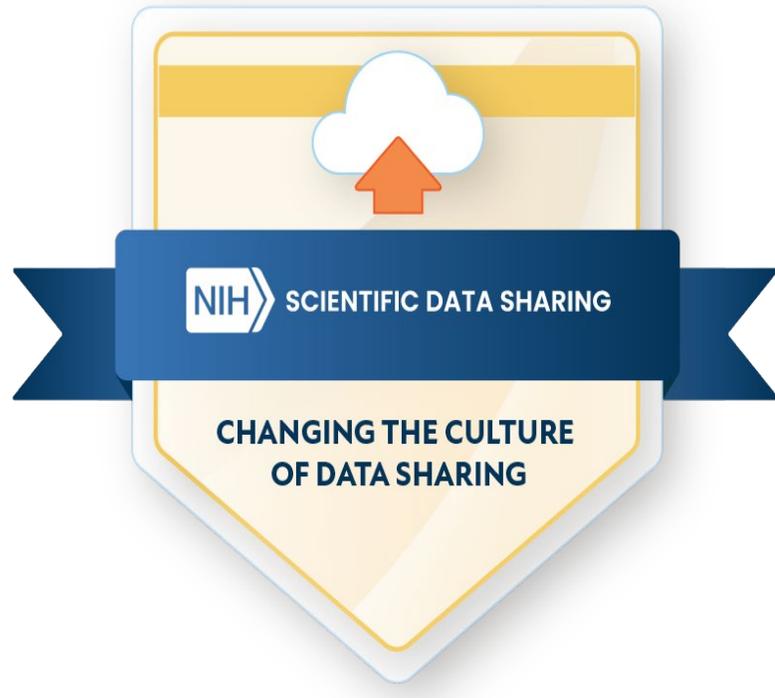
# Metadata and Data Standards – What and Why

Matthew Schu, PhD

27 June 2023



# Maximizing the value of scientific data



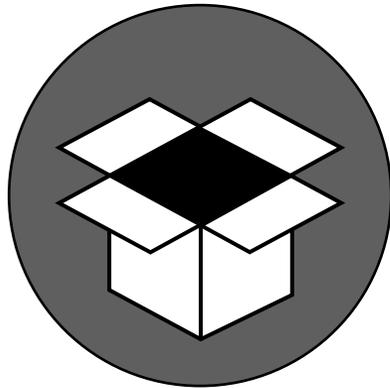
- Test validity of research findings
- Strengthen analysis through combined data sets
- ➔ • Allow for reuse of hard to generate data
- Open new frontiers of discovery
- Foster trust in publicly funded research

# Metadata - adding value to scientific data

Metadata is “data about data” and is the information needed to discover, use, and understand data and describes the



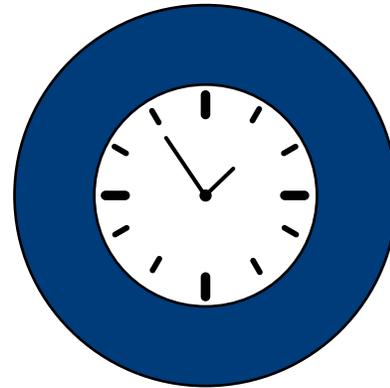
Who



What



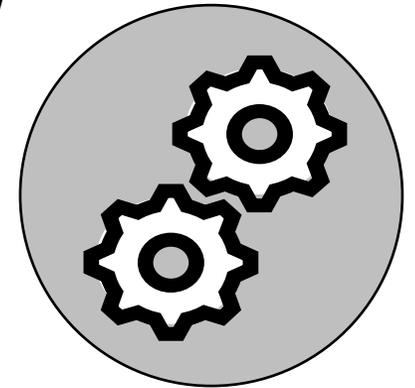
Where



When



Why



How

Reconstruct the context so users know what can and cannot be done with the data

# What can you do with data?



What can you do with this data?

- Enjoy it (it's cute)
- Share it with a friend

When all you have is the data, there are a limitations on what you can do with it.

# What can you do with data + metadata?

I'm looking for a pet and this puppy is cute. What breed of dog is this?

- Pet Lover

What were the camera settings to focus on the puppy and blur the background?

- Photographer

I want to make this into a poster. Are there any license considerations? What is the format of the image.

- Home Owner



## Metadata:

Image Title: Puppy with leaves

Subject: Golden retriever puppy

Subject Age: 8 weeks

Subject Sex: Female

Photographer: Cam Jansen

Camera Manufacturer: Canon

Camera Model: EOS 5D Mark IV

Camera Lens: EF 24-105

Zoom Level: 4x

Aperture: f1.4

Focal Length: 60mm

Image Date: 9-16-2021

Image Format: PNG

Image Resolution: 300 dpi

Image Processing Software: Lightroom

Licensing: Shared under CC agreement

# Metadata Levels

Study or Dataset	Variable	File
<p><b>Information describing the purpose of collecting the data</b></p> <ul style="list-style-type: none"><li>• Study name</li><li>• Study URL or Digital Object Identifier (DOI)</li><li>• Funding information</li><li>• Investigator contact information</li><li>• Associated publications</li><li>• Release version</li><li>• Study duration and dates active</li><li>• Collection protocol</li><li>• Sample size and population description</li></ul>	<p><b>Information about measured variables that belong to a study or dataset</b></p> <ul style="list-style-type: none"><li>• Variable name</li><li>• Variable description</li><li>• Variable data type (e.g., numeric, character, Boolean)</li><li>• Data format (e.g., field length)</li><li>• Dataset that contains the variable</li><li>• Generation method</li></ul> <p>Often stored in tabular format (data dictionary or code book)</p>	<p><b>Information about files that have been produced during the course of the study</b></p> <ul style="list-style-type: none"><li>• File name</li><li>• File URL</li><li>• File description</li><li>• File format (e.g., .csv, .xlsx, .png)</li><li>• Access control information</li><li>• License</li></ul>

You're probably already collecting metadata, often required for manuscript or repository submission

# Metadata Standards

Data associated with rich metadata is more useful, but what level of metadata would allow for appropriate data reuse?

Data Type	Metadata Standards
Biochemical	Minimum Information About a Bioactive Entity ( <a href="#">MIABE</a> )
Clinical	Clinical Data Interchange Standards Consortium (CDISC) <a href="#">Analysis Data Model</a>
Genomics	Minimum Information about a (Meta)Genome Sequence ( <a href="#">MIGS/MIMS</a> )
Transcriptomics	Minimal Information about a high throughput SEQuencing Experiment ( <a href="#">MINSEQE</a> )
Imaging	Minimum Information about Tissue Imaging ( <a href="#">MITI</a> )
Metabolomics	Core Information for Metabolomics Reporting ( <a href="#">CIMR</a> )
Proteomics	The Minimum Information About a Proteomics Experiment ( <a href="#">MIAPE</a> )

Metadata standards for various biological disciplines are available at [fairsharing.org](http://fairsharing.org)

# Nutrition for Precision Health Powered by *All of Us* (NPH)



## Nutrition for Precision Health

Powered by the *All of Us* Research Program,  
part of the National Institutes of Health

[About the Research](#) | [Eligibility](#) | [Study Activities](#)

[GET INVOLVED](#)

## Can a personalized diet help improve health and prevent chronic diseases?

The Nutrition for Precision Health study is trying to answer this question by studying how individual people respond to different foods. Nutrition for Precision Health is a partner of the *All of Us* Research Program. This is a large effort to speed up health research.

[LEARN MORE](#)



**Nutrition is not one-size-fits-all.**

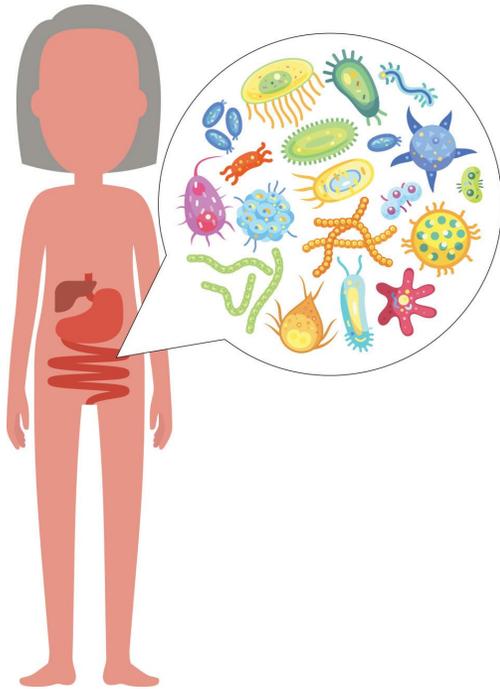
# All of Us Researcher Workbench

The screenshot shows the 'Data Browser' page of the All of Us Researcher Workbench. At the top, there is a navigation bar with the 'All of Us Research Hub' logo, the NIH logo, and menu items for 'ABOUT', 'DATA & TOOLS', 'DISCOVER', and 'SUPPORT'. Below the navigation bar, the page title 'Data Browser' is centered. A paragraph explains that the data is aggregate-level, derived from multiple sources, and that personal identifiers have been removed for privacy. A search bar labeled 'Search Across Data Types' with a 'Keyword Search' input field is present. Below the search bar, it states 'Data includes 409,420 participants as of 2/15/2023.' To the right of the search bar are three circular icons: a yellow one with a question mark and speech bubble labeled 'FAQs', a green one with a question mark and document labeled 'Introductory Videos', and a blue one with a question mark and database icon labeled 'User Guide'. Below these icons is a section titled 'EHR Domains' with four cards: 'Conditions' (25,638 medical concepts, 254,700 participants), 'Drug Exposures' (29,865 medical concepts, 239,740 participants), 'Labs & Measurements' (16,216 medical concepts, 252,980 participants), and 'Procedures' (30,328 medical concepts, 242,580 participants). Each card has a 'View' link at the bottom.

- Data from NPH will be made available via the *All of Us* Researcher workbench

- AI analytic tools will be brought to the data to help uncover patterns within the high dimensional data captured
- Other researchers will be able to access the data from NPH and potentially combine these results with those from other studies in the Researcher Workbench

# Scientific metadata example



## Microbiome metadata

1.sample\_name

2.host\_subject\_id

3.sample\_type

4.sample\_quality [Bristol stool scale]

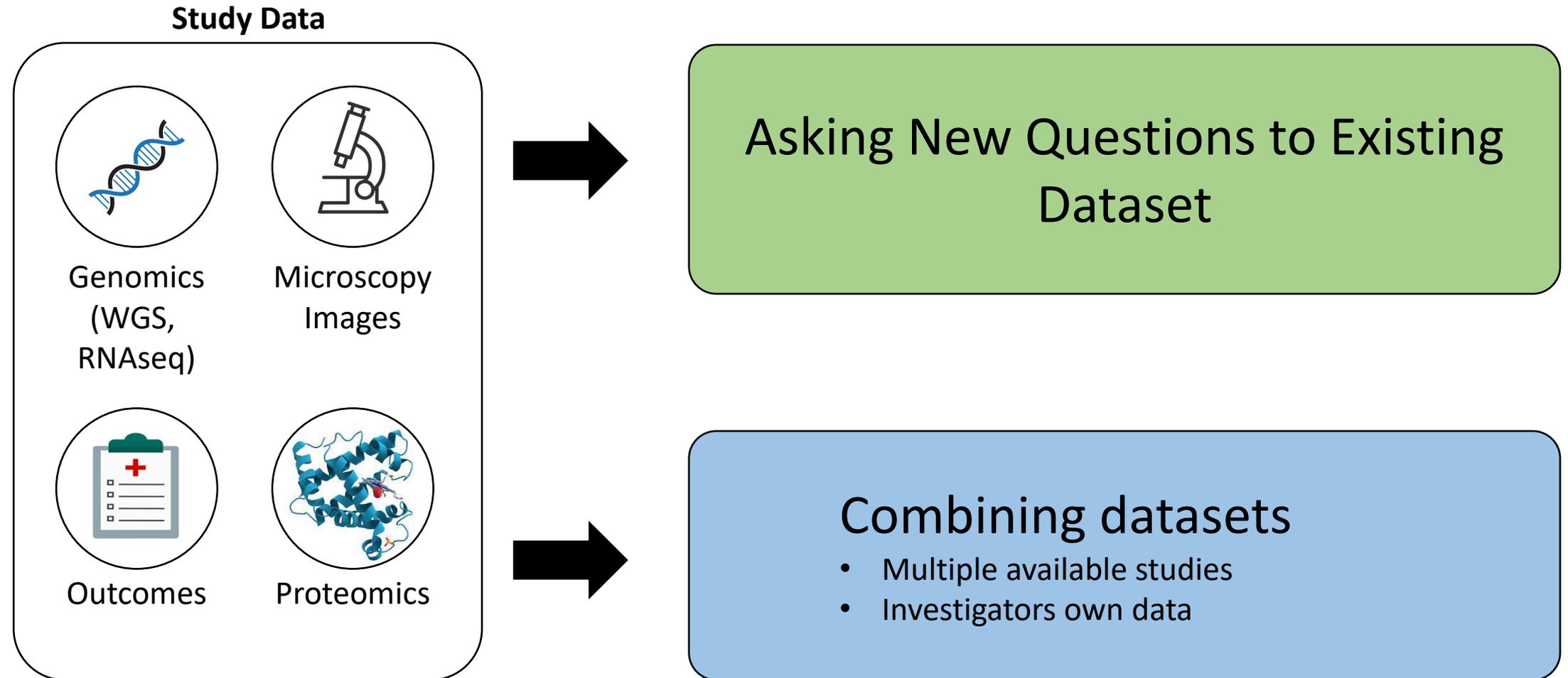
5.physical\_specimen\_location

6.collection\_date

7.country

- Type 1: Separate hard lumps, like nuts (difficult to pass)
- Type 2: Sausage-shaped, but lumpy
- Type 3: Like a sausage but with cracks on its surface
- Type 4: Like a sausage or snake, smooth and soft (average stool)
- Type 5: Soft blobs with clear cut edges
- Type 6: Fluffy pieces with ragged edges, a mushy stool (diarrhea)
- Type 7: Watery, no solid pieces, entirely liquid (diarrhea)

# Scientific Data Reuse

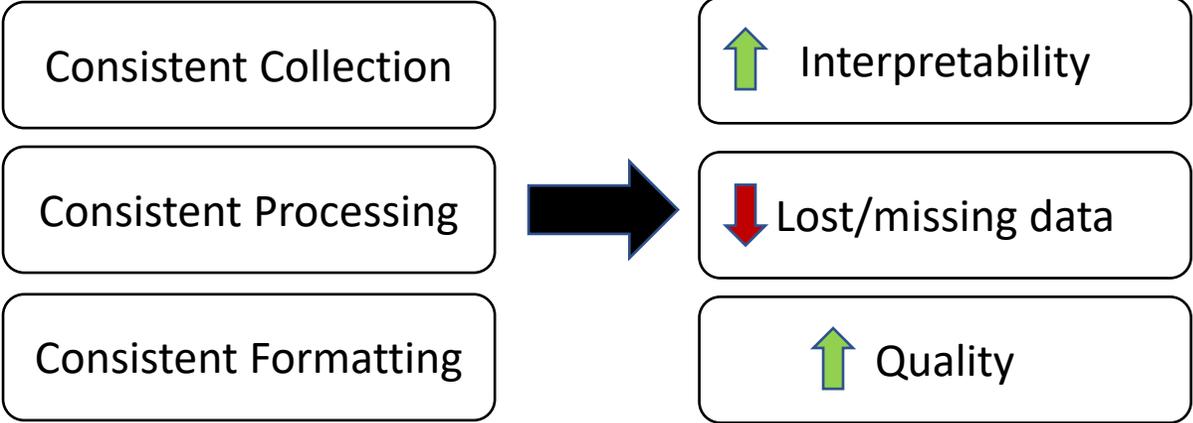


# Data Standards and Harmonization

## Data Harmonization

	Study 1	Study 2
Question	How many packs a day do you <u>currently</u> smoke?	What is your smoking status?
Responses	1 – 0 2 – 1-2 packs 3 – 3-5 packs 4 – >5 packs	1 – Never smoker 2 – Past smoker 3 – Current smoker

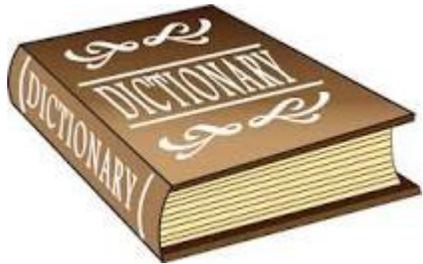
How to combine the smoking status across the two studies?



Your own data will be readily compatible with all other studies that utilize the same or comparable standards

# Data Standards

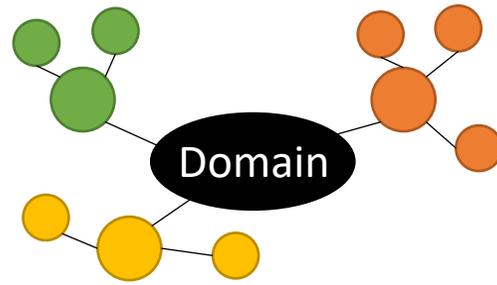
## Vocabularies / Terminologies



Clearly communicating the definition of terms within the study context

- [UMLS](#)
- [SNOMED](#)
- [LOINC](#)
- [RxNorm](#)

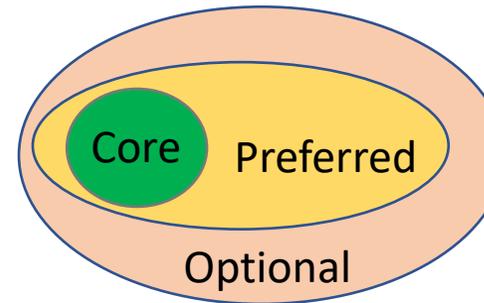
## Ontologies



Describes the relationship between concepts within a scientific domain

- [BioPortal](#)
- [Ontology of Precision Medicine and Investigation](#)

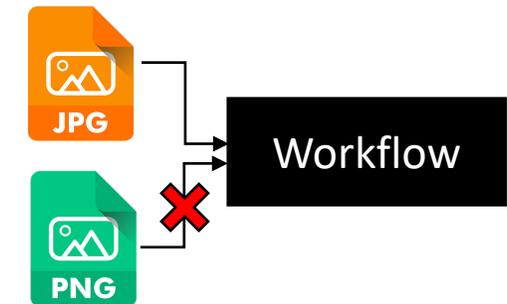
## Common Data Elements



Instrument with questions and responses for consistent collecting of study data

- [NIH CDE repository](#)
- [PhenX Toolkit](#)

## File Formats



Common file formats allow for consistent data processing and tool utilization

# Example Data Standards for NPH (microbiome)

- Example from QIIME data standards

## Sample information file

The *sample information file* will define the biological context of each sample, with categories like `sample_type`, `treatment`, etc. The `sample_name` defined in this file is used to relate each sample in the preparation file with the biological sample.

## Required fields for Qiita

Note that Qiita require to have at least two columns, including `sample_name`, for a sample information file to be added to the system:

Field name	Format	Description
<code>sample_name</code>	free text with restrictions	Identifies a sample. It is the primary key and must be unique. Allowed characters are alphabetic <code>[A-Za-z]</code> , numeric <code>[0-9]</code> , and periods <code>.</code>

# Questions

