# LA-UR-24-26605

**Approved for public release; distribution is unlimited.**

**Title:** 2024 NMDC Ambassador Training Materials

**Author(s):** Kelliher, Julia Mae; Rodriguez, Francisca Ester; Johnson, Leah Young Davenport; Prime, Kaelan Jantira; Roux, Simon; Eloe-Fadrosh, Emiley A.; Smith, Montana; Clum, Alicia; De Santiago, Alejandro; Gajigan, Andrian; Maher, Rebecca; Hanson, Buck Timothy; Robinson, Christopher; Betacurt-Anzola, Daniela; Skoog, Emilie; Cirolia, Giana Teresa; Skeen, Heather; Oduwole, Iyanu; Kajihara, Kacie; Pham, Kent; Camuy-Vélez, Lennel; et al.

**Intended for:** NMDC Ambassador training materials that will receive a DOI
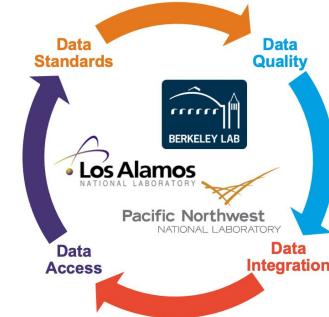
**Issued:** 2024-07-02

## Los Alamos
### NATIONAL LABORATORY

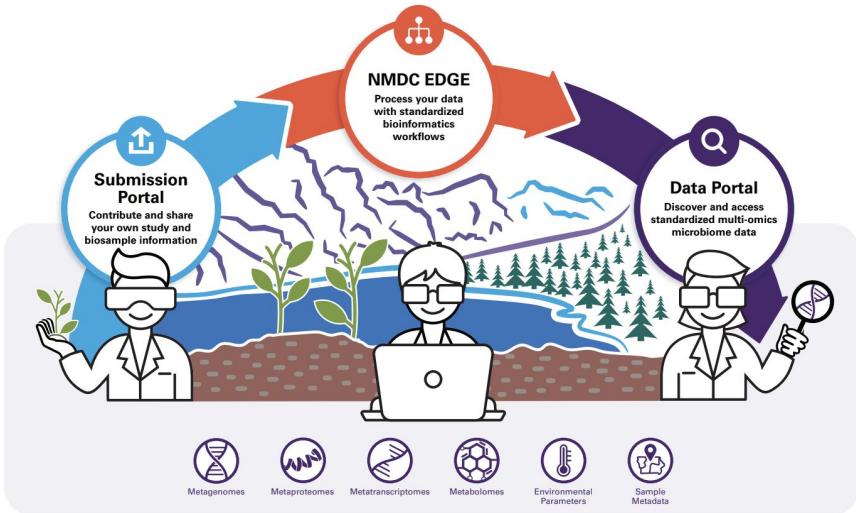# What is the NMDC?

2024 Ambassador Cohort

# What is the NMDC?



The NMDC is a *sustainable data discovery platform* that promotes open science and shared-ownership across a broad and diverse community of researchers, funders, publishers, societies, and other collaborators. The NMDC aims to enable multi-omic microbiome research to accelerate scientific discovery. The NMDC is a Department of Energy funded program that is a collaboration between 3 National Laboratories: Lawrence Berkeley National Laboratory (LBNL), Los Alamos National Laboratory (LANL), and Pacific Northwest National Laboratory (PNNL)

**"Enabling inclusive and interdisciplinary environmental microbiome science by connecting data, people, and ideas"**

# Why we need standards

# Vision and Mission

**nmdc**
National Microbiome
Data Collaborative

**Vision**
To **connect data**, **people**, and **ideas** to advance microbiome innovation and discovery

**Mission**
To support a FAIR microbiome data sharing network, through **infrastructure**, **data standards**, and **community building**, that addresses pressing challenges in environmental sciences
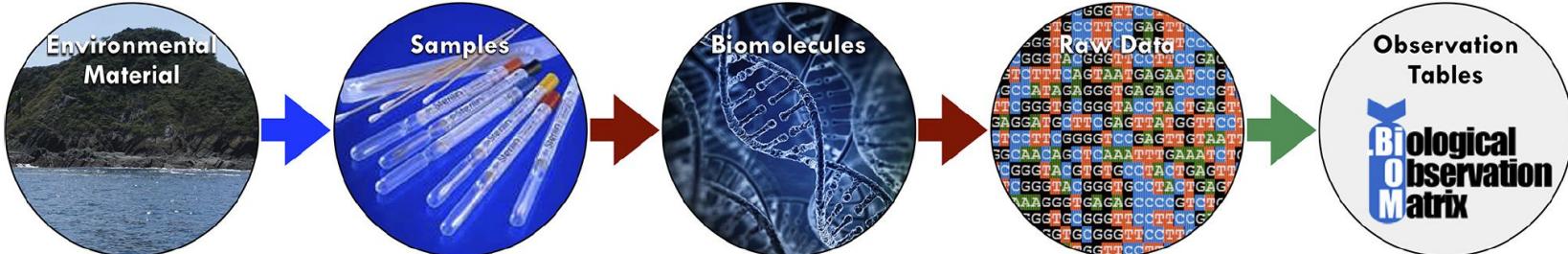
# Metadata Standards

**Adopting standards for reporting makes data human and machine readable.**

Metadata (data about data):

- Contextual data about your data
- Vital for data
  - Preservation
  - Discovery
  - Access
  - Reuse

**Sample metadata** includes information about:

- ***When*** it was collected
- ***Where*** it was collected
- ***What*** kind of sample is it
- ***Treatment*** applied during experimentation
- ***Environmental Properties*** from which the sample was taken

# NMDC Submission Portal



Upload your data

**Submission Portal**

Lower barriers to collect study and biosample data

Validate your submission against pre-made MIxS templates
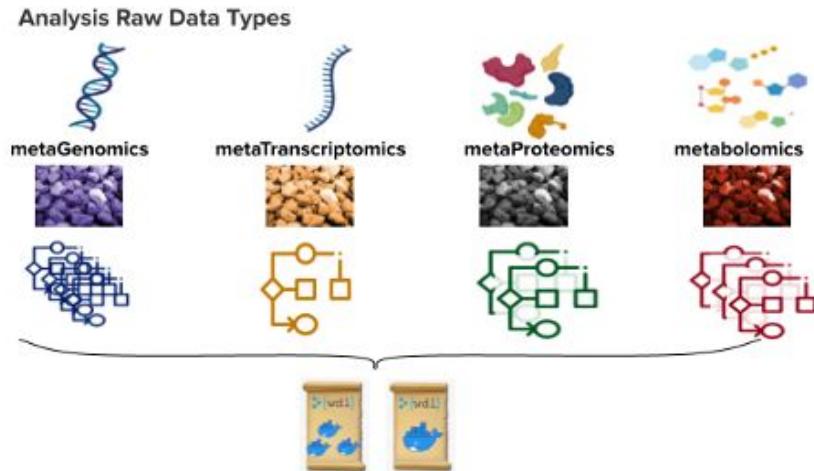
Receive guidance on how to meet standards

# Standardized Bioinformatic Workflows



The NMDC has integrated state of the art open-source bioinformatics tools into standardized workflows for processing raw multi-omics data to produce **interoperable** and **reusable** annotated data products

- **NMDC Workflows:**
  - Metagenome data
    - ReadsQC
    - Read-based taxonomy classification
    - Assembly
    - Annotation
    - Metagenome assembled genomes (MAGs)
  - Metatranscriptome data
  - Natural organic matter data
  - Metabolome data
  - Metaproteome data
  - Viruses & Plasmids

# NMDC EDGE



**NMDC EDGE**

Streamline multi-omics data processing

# Data Stewardship + FAIR Data

- NMDC is committed to FAIR data principles of making microbiome data Findable, Accessible, Interoperable, and Reusable
- Both raw and processed data should be FAIR
- Processing data in a standardized way makes it increasingly interoperable and reusable

**Findable**
Ensure all data registered within NMDC are human and machine readable

**Accessible**
Identify data sets that are available, including any authentication and authorization requirements

**Interoperable**
Provide provenance, metadata, and uniformly processed data, we are lowering the barriers to making data interoperable

**Reusable**
Enable download of data, data products, and workflows for external reprocessing
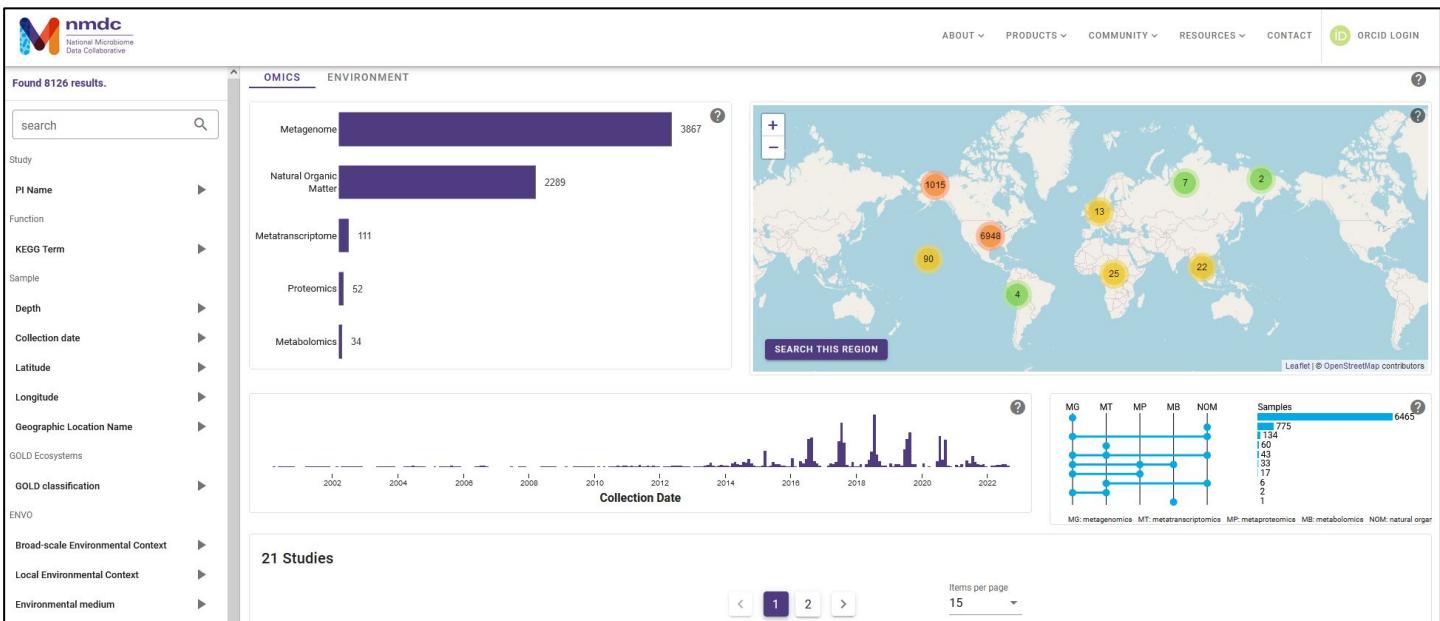
# NMDC Data Portal



**Data Portal & API**

Access and discovery of microbiome information

# NMDC Engagement



**Goal 1: Recognize and support the diverse research needs and perspectives of the microbiome research community**

**Goal 2: Promote best practices across the microbiome community, from researchers to funders**

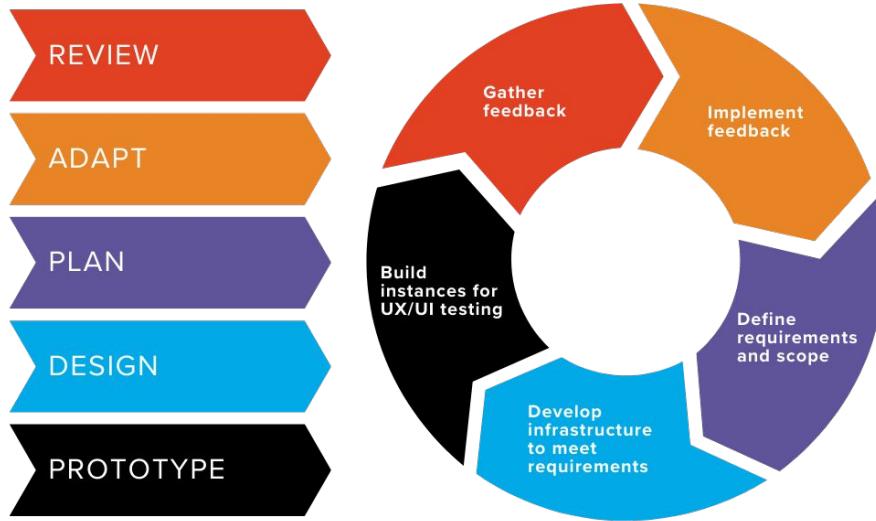**Goal 3: Build a microbiome ecosystem that enables scientific discovery and innovation**

| | Community | Collaborator | Champion | Ambassador |
|---|---|---|---|---|
| Works with microbiome data | ✓ | ✓ | ✓ | ✓ |
| Contributed or cited NMDC data | + | ✓ | + | + |
| Advocates for FAIR microbiome data | + | + | ✓ | ✓ |
| Steward of well-curated data | + | + | ✓ | ✓ |
| Speaks on behalf of the NMDC | | | | ✓ |
| Hosts community events | | | | ✓ |

+ Optional

# Partnerships

# Community-centered design process



REVIEW

ADAPT

PLAN

DESIGN

PROTOTYPE

Gather feedback

Implement feedback

Define requirements and scope

Develop infrastructure to meet requirements

Build instances for UX/UI testing

**2023-2024 User Research**

| 34 | 321 | 120 |
|---|---|---|
| User research participants | Insights from user interviews | Action items from insights |

- The NMDC program utilizes a community-centered design approach to its product development process
  - User interviews, usability testing, beta-testing
  - Understand the needs of the community and how researchers use the products
- Workshops and event feedback is critical too
- Want to develop products that are as useful to the community as possible

14

# User Research



**Examples of feedback:**
- Updates to interface design, tutorials, and help guidance
- Feature implementation
  - Batch processing in NMDC EDGE
  - Data download updates in the Data Portal
  - Template accessibility updates in the Submission Portal

**Volunteer for our user research program!**

https://microbiomedata.org/user-research/

# Annual Reports





### Advancing microbiome science for the benefit of all

In January of this year, the White House Office of Science and Technology Policy launched the Year of Open Science with new actions to advance open and equitable science policies across the federal government. The National Microbiome Data Collaborative (NMDC) is committed to open and equitable research in microbiome science. This commitment is the foundation for our infrastructure development activities and has been a prominent driver this past year across our three products: the Submission Portal, NMDC EDGE, and the Data Portal. It has been an exciting and busy time for our team as we work to serve the scientific community in a way that enables microbiome innovation and discovery. I am exceptionally proud of the progress the NMDC team has made this past year in fostering strong community partnerships and advancing our powerful products into tools that drive scientific impact.

We launched the NMDC persistent identifier service in January 2023 and deployed programmatic access to all NMDC data through a public application programming interface (API). Together, these efforts support a larger findable, accessible, interoperable, and reusable (FAIR) data ecosystem to programmatically exchange and link data across resources. The Data Portal now hosts over 7,700 biosamples with a collective nearly 90 TB of multi-

omics microbiome data and links across complementary data platforms, including the Integrated Microbial Genomes and Microbiomes (IMG/M) and Genomes OnLine Database (GOLD) of the Joint Genome Institute (JGI), the Department of Energy (DOE) Systems Biology Knowledgebase (KBase), the National Center for Biotechnology Information (NCBI), the Mass Spectrometry Interactive Virtual Environment (MassIVE), the Environmental System Science Data Infrastructure for a Virtual Ecosystem (ESS-DIVE), and the National

Emiley Eloe-Fadrosh
National Microbiome
Data Collaborative Lead

### A community-driven data infrastructure

The NMDC is tackling existing gaps in microbiome research by using proven approaches and new innovations in distributed data infrastructure and linked data technologies. Our three products — the Submission Portal, NMDC EDGE, and the Data Portal — are driven by community needs. They support data, information, knowledge sharing, and access. This past year, we worked closely with the research community to strengthen our existing infrastructure in ways that will catalyze new research. This included launching a new persistent identifier service, supporting programmatic access to NMDC data, expanding the amount of available high-quality data and workflows, and contributing to major updates of community data standards.

Sevilleta Long Term Ecological Research (LTER)
Site in New Mexico. Credit: Buck Hanson

# NMDC Resources

## nmdc
### National Microbiome Data Collaborative

**Website:** https://microbiomedata.org/
**Data Portal:** https://data.microbiomedata.org/
**Submission Portal:** https://data.microbiomedata.org/submission/home
**NMDC EDGE:** https://nmdc-edge.org/home
**Github:** https://github.com/microbiomedata
**Docker Hub:** https://hub.docker.com/u/microbiomedata
**Documentation:**
https://nmdc-documentation.readthedocs.io/en/latest/overview/nmdc_overview.html
**YouTube:** https://www.youtube.com/channel/UCyBqKc46NQZ_YgZlKGYegIw/featured

### Get involved!

**Sign up for our newsletter**
microbiomedata.org

**Become a NMDC Champion**
bit.ly/championsapp

**Find us on X/Twitter**
@microbiomedata

**Find us on LinkedIn**
https://bit.ly/NMDC_LinkedIn

**Find us on Instagram**
@microbiomedata

## Read more about the NMDC

Kelliher JM *et al.* Cohort-based learning for microbiome research community standards. *Nat Microbiol* (2023). doi.org/10.1038/s41564-023-01361-7.

Hu B, Canon S, Eloe-Fadrosh EA, et al.. Challenges in Bioinformatics Workflows for Processing Microbiome Omics Data at Scale. *Front Bioinform*. 1:826370. (2022) doi: 10.3389/fbinf.2021.826370.

Eloe-Fadrosh EA *et al.* The National Microbiome Data Collaborative Data Portal: an integrated multi-omics microbiome data resource. *Nucleic Acids Res.* 7;60(D1):D828–D836. (2022) doi: 10.1093/nar/gkab990.

Wood-Charlson, E.M., Anubhav, Auberry, D. *et al.* The National Microbiome Data Collaborative: enabling microbiome science. *Nat Rev Microbiol* **18,** 313–314 (2020). doi.org/10.1038/s41579-020-0377-0

Vangay, P *et al.* Microbiome metadata standards: Report of the National Microbiome Data Collaborative's workshop and follow-on activities. *mSystems* **6**, e01194-20 (2021). doi.org/10.1128/mSystems.01194-20

# 'One Slide on the NMDC' for Presentations

# National Microbiome Data Collaborative

## Mission

To support a FAIR microbiome data sharing network, through **infrastructure**, **data standards**, and **community building**, that addresses pressing challenges in environmental sciences



Findable
Accessible
Interoperable
Reusable

## The NMDC offers 3 Products:

### Submission Portal

Lower barriers to collect study and biosample data

### NMDC EDGE

Streamline multi-omics data processing

### Data Portal & API

Access and discovery of microbiome information



Submission Portal



NMDC EDGE



NMDC Data Portal

## Get involved!

🌐 **Website/newsletter**
microbiomedata.org

𝕏 **Find us on X/Twitter**
@microbiomedata

💬 **Become a Champion**
bit.ly/champions-program

in **Find us on LinkedIn**
https://bit.ly/NMDC_LinkedIn

📷 **Find us on Instagram**
@microbiomedata

## Engagement

**User Facilities**

JGI JOINT GENOME INSTITUTE DEPARTMENT OF ENERGY

EMSL

**Individuals**

The NMDC Champions Program

The NMDC Ambassador Program

**Strategic Partners**

NSF | neon Operated by Battelle

NASA

GeneLab

# Purpose of this training

- Provide an overview of the benefits of standardized bioinformatics workflows and how they promote FAIR data
- Discuss the importance and benefits of using standardized workflows for the future of microbiome data
- Provide template slides that anyone can use in their own events
- Introduce audience to NMDC EDGE
  - The hands-on activity can be used (or modified and used) for use in your own events too! You can upload or link your own data

# FAIR Data

- NMDC is committed to FAIR data principles of making microbiome data Findable, Accessible, Interoperable, and Reusable
- Raw and **processed** data should be FAIR
- NMDC prioritizes standardization and user-centered design to achieve FAIR and open data

**0101** **Findable**
Ensure all data registered within NMDC are human and machine readable

**Accessible**
Identify data sets that are available, including any authentication and authorization requirements

**Interoperable**
Provide provenance, metadata, and uniformly processed data, we are lowering the barriers to making data interoperable

**Reusable**
Enable download of data, data products, and workflows for external reprocessing

# Intro to multi-omics techniques

-   Multi-omics: the integration of data from multiple omics techniques, such as genomics, transcriptomics, proteomics, and metabolomics, to provide a comprehensive assessment of a biological system

-   By analyzing various layers of molecular information simultaneously, taxonomic, genetic, and functional information can be used to unravel information about an organism and/or a community

# Multi-omics data processing

- Multi-omics data collection is rapidly becoming one of the most effective ways to interrogate microbiomes
  - The infrastructure surrounding these data is not keeping up
    - Large amounts of compute power needed to store and process these data
    - Data not FAIR (not made publicly accessible, not utilizing standards, etc.)
    - Different omics data not connected, data not comparable within an omics type
- Processing these omics data streams is currently very difficult to do in a way that allows for effective cross-study comparisons and data re-use



Metagenomes    Metaproteomes    Metatranscriptomes    Metabolomes

# Bioinformatics explosion



The massive explosion in the number of bioinformatics tools and workflows has led to data being processed in many ways, thus **limiting between-study comparability**

These are *just* mapping tools for genomic data developed between 2001-2018.

Timeline of NGS read aligners. Image from Nuno Fonseca https://www.ecseq.com/support/ngs/what-is-the-best-ngs-alignment-software

# Benefits of standardized workflows

- Help to make data and data products FAIR
- Reproducibility within a lab and between labs
- Better keep track of what was done to data
- Allows for integrations and comparisons with other datasets from other studies
- No need to put together own workflows
  - Existing workflows can include the best tools for omics data processing
  - Saves research, implementation, and testing time

# NMDC Bioinformatics Workflows

The NMDC has integrated state of the art open-source bioinformatics tools into standardized workflows for processing raw multi-omics data to produce **interoperable** and **reusable** annotated data products.

**NMDC Workflows:**
- Metagenome data
  - ReadsQC
  - Read-based taxonomy classification
  - Assembly
  - Annotation
  - Metagenome assembled genomes (MAGs)
- Metatranscriptome data
- Natural organic matter data
- Metabolome data
- Metaproteome data
- Viruses & Plasmids

# NMDC Workflows

# Why the NMDC workflows?

Benefits of using the NMDC workflows:
- Tools were carefully researched, selected, and modified for optimal performance
  - Many tools are the tools are production quality, regularly used to process thousands of datasets from DOE user facilities
- Workflows have been extensively tested on data from dozens of institutions and sample types
- Users are able to run these workflows through shared computing resources
  - Users don't need to download these tools or databases, nor have access to their own computing clusters
- The workflows are offered in a user-friendly interface for users with any level of bioinformatics experience
- Open source platform, extensive documentation

# NMDC Data Portal



Processed datasets available on the NMDC Data Portal have all been run through the NMDC workflows allowing for direct comparisons between this data and data processed in NMDC EDGE

# Reads QC

Performs quality control on raw metagenome Illumina reads to trim/filter low quality data and to remove artifacts, linkers, adapters, spike-in reads and reads mapping to several hosts and common microbial contaminants.

**Input:** Raw Illumina data

**Output:** File of cleaned reads and QC statistics

# Read-based Taxonomy Classification



Takes in Illumina sequencing files and profiles the reads using 3 taxonomic classification tools (GOTTCHA2, Kraken2, and Centrifuge) with a range of sensitivity and specificity

**Input:** Illumina data: it is highly recommended to input clean reads from the ReadsQC workflow

**Output:** Results for each tool at three taxonomic levels (Species, Genus, and Family). Interactive Krona plots are also generated

# Metagenome Assembly

Takes in Illumina data, runs error correction, assembly, and assembly validation

**Input:** Illumina data: recommended input is the output from the ReadsQC workflow

**Output:** File of assembled contigs; assembly statistics

# Metagenome Annotation

Takes in assembled metagenomes and generates structural and functional annotations

**Input:** Assembled contigs: recommended input is the output from the Metagenome Assembly workflow

**Output:** Structural annotation file, functional annotation file, several summary files

# Metagenome Assembled Genomes

**nmdc**
National Microbiome
Data Collaborative

Classifies contigs into bins, bins are refined using functional annotation file, bins are evaluated for completeness and contamination. Quality of bins is determined and lineage is assigned

**Input:** Assembled contigs, read mapping file from the assembly, functional annotation of the assembly

**Output:** Summary statistics, file of high quality (HQ) and medium quality (MQ) bins

# Additional -omics Workflows

# Metatranscriptome Workflow



Takes in raw metatranscriptome data, filters data for quality, removes rRNA reads, assembles and annotates the transcripts. Data is mapped back to the genomic features in the transcripts and RPKMs (Reads Per Kilobase of transcript per Million mapped reads) are calculated for each feature in the functional annotation file.

**Input:** Illumina data

**Output:** Assembled transcripts, annotated features file, annotation files

# Natural Organic Matter Workflow

Takes mass spectrometry data collected from organic extracts to determine the molecular formulas of natural organic biomolecules in the input sample.

**Input:** The output from a mass spec experiment; a calibration file of molecular formula references is also required when running via command line

**Output:** Primary output file is the Molecular Formula Data Table

# Metabolome Workflow

The GC-MS based metabolomics workflow leverages PNNL's CoreMS software framework.

**Input:** Raw GC-MS data

**Output:** Metabolites data table

# Metaproteome Workflow

This workflow is an end-to-end data processing workflow for protein identification and characterization using MS/MS data

→ **Input:** Raw LC-MS/MS data and an associated metagenome file

→ **Output:** Protein crosstab; QC plots

# Viruses & Plasmids Workflow

**geNomad** is a tool that identifies virus and plasmid genomes from nucleotide sequences. It provides state-of-the-art classification performance and can be used to quickly find mobile genetic elements from genomes, metagenomes, or metatranscriptomes.

| Speed | Taxonomic assignment | Functional annotation |
|---|---|---|
| geNomad is significantly faster than similar tools and can be used to process large datasets. | The identified viruses are assigned to taxonomic lineages that follow the latest ICTV taxonomy release. | Genes encoded by viruses and plasmids are functionally annotated using geNomad's marker database. |

# Workflow Availability

All workflows can be downloaded and run locally on your own computing resources, or can be run through NMDC EDGE

GitHub: **https://github.com/microbiomedata**

Docker Hub: **https://hub.docker.com/u/microbiomedata**

Questions?

# NMDC EDGE


National Microbiome Data Collaborative

The workflows can be run in a user-friendly interface:
**https://nmdc-edge.org/home**

Designed for bioinformaticians of every level of expertise (including novices)

# NMDC EDGE Walkthrough

# Tutorials, User guides, Information

# Public Projects

# Upload Files

# Retrieve SRA Data

# Run Metagenome Pipeline

# Metatranscriptomics Workflow

# Natural Organic Matter Workflow

# Viruses & Plasmids Workflow

# Metaproteomics Workflow

# My Projects

# My Projects: Possible Statuses

# Example Results: Metagenome

# Example Results: Metagenome

# Example Results: Metagenome

# Example Results: Metagenome



**Metagenome_pipeline_test**

**Project Summary:**

**Description:** This is a test of the NMDC metagenome pipeline for training purposes
**Owner:** jkelliher@lanl.gov
**Submission Time:** Mon Aug 23 2021 09:25:01 GMT-0600
**Status:** Complete
**Type:** Metagenome Pipeline

expand | close sections

General

ReadsQC Result

Read-based Taxonomy Classification Result

Metagenome Assembly Result

Metagenome Annotation Result

Metagenome MAGs Result

Browser/Download Outputs

## Metagenome Assembly Result

| Name | Status |
| --- | --- |
| scaffolds | 25,324 |
| contigs | 25,726 |
| scaf_bp | 52,206,897 |
| contig_bp | 52,201,077 |
| gap_pct | 0.011 |
| scaf_N50 | 691 |
| scaf_L50 | 4,103 |
| ctg_N50 | 724 |
| ctg_L50 | 3,971 |
| scaf_N90 | 14,186 |
| scaf_L90 | 726 |
| ctg_N90 | 14,473 |
| ctg_L90 | 716 |
| scaf_logsum | 645,093 |
| scaf_powsum | 120,098 |
| ctg_logsum | 638,015 |
| ctg_powsum | 116,432 |
| asm_score | 33.765 |
| scaf_max | 1,491,105 |
| ctg_max | 859,644 |

# Example Results: Metagenome



**Metagenome_pipeline_test**

**Project Summary:**

**Description:** This is a test of the NMDC metagenome pipeline for training purposes
**Owner:** jkelliher@lanl.gov
**Submission Time:** Mon Aug 23 2021 09:25:01 GMT-0600
**Status:** Complete
**Type:** Metagenome Pipeline

expand | close sections

- General
- ReadsQC Result
- Read-based Taxonomy Classification Result
- Metagenome Assembly Result
- Metagenome Annotation Result
- Metagenome MAGs Result
- Browser/Download Outputs

### Metagenome Annotation Result

#### Processed Sequences Statistics

| Data type | Number of seqs | Number of bps | Median length | Average length | Length shortest seq | Length longest seq | Standard deviation |
|---|---|---|---|---|---|---|---|
| final_fasta | 25,726 | 52,201,077 | 818.5 | 2,029.118 | 200 | 859,644 | 16,939.403 |
| sequences_with_genes | 24,248 | 51,497,305 | 865 | 2,123.775 | 200 | 859,644 | 17,443.493 |
| sequences_without_genes | 1,478 | 703,772 | 404 | 476.165 | 203 | 1,918 | 217.554 |

#### Predicted Genes Statistics

| Feature type | Prediction method | Number of seqs | Number of bps | Median length | Average length | Length shortest seq | Length longest seq | Standard deviation | Number of predicted features |
|---|---|---|---|---|---|---|---|---|---|
| CDS | Prodigal v2.6.3 | 12,478 | 3,694,932 | 180 | 228.831 | 75 | 1,935 | 156.372 | 16,147 |
| CDS | GeneMark.hmm-2 v1.05 | 18,576 | 35,352,681 | 480 | 669.267 | 90 | 16,545 | 616.622 | 52,823 |
| tRNA | tRNAscan-SE v.2.0.7 (Oct 2020) | 451 | 67,404 | 76 | 79.486 | 56 | 146 | 10.062 | 848 |
| misc_feature | INFERNAL 1.1.3 (Nov 2019) | 4 | 1,454 | 366.5 | 363.5 | 349 | 372 | 10.408 | 4 |
| regulatory | INFERNAL 1.1.3 (Nov 2019) | 4 | 1,454 | 366.5 | 363.5 | 349 | 372 | 10.408 | 4 |
| ncRNA | INFERNAL 1.1.3 (Nov 2019) | 4 | 1,454 | 366.5 | 363.5 | 349 | 372 | 10.408 | 4 |
| rRNA | INFERNAL 1.1.3 (Nov 2019) | 4 | 1,454 | 366.5 | 363.5 | 349 | 372 | 10.408 | 4 |

# Example Results: Metagenome

# Example Results: Metagenome



**Metagenome_pipeline_test**

**Project Summary:**

**Description:** This is a test of the NMDC metagenome pipeline for training purposes
**Owner:** jkelliher@lanl.gov
**Submission Time:** Mon Aug 23 2021 09:25:01 GMT-0600
**Status:** Complete
**Type:** Metagenome Pipeline
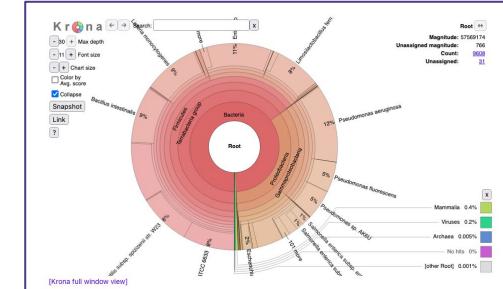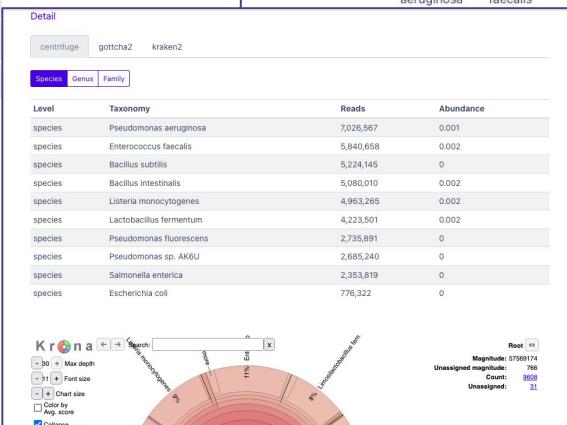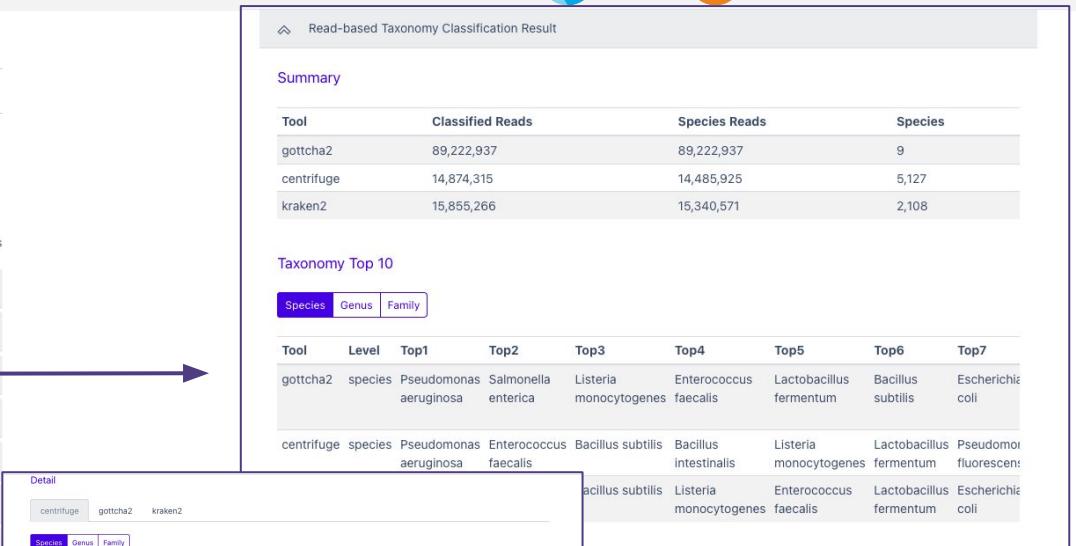
expand | close  sections

General

ReadsQC Result

Read-based Taxonomy Classification Result

Metagenome Assembly Result

Metagenome Annotation Result

Metagenome MAGs Result

Browser/Download Outputs

Browser/Download Outputs

| File | Size | Last Modified |
|------|------|---------------|
| MetagenomeAnnotation | | |
| MetagenomeAssembly | | |
| MetagenomeMAGs | | |
| ReadbasedAnalysis | | |
| centrifuge | | |
| gottcha2 | | |
| kraken2 | | |
| activity.json | 919 B | over 1 year ago |
| data_objects.json | 3 kB | over 1 year ago |
| Metagenome_pipeline_test.json | 2.11 MB | over 1 year ago |
| ReadsQC | | |

# Questions?

# Hands-on Activity

- Navigate to [https://nmdc-edge.org/home](https://nmdc-edge.org/home)

- Log in using your ORCiD

- Select Viruses & Plasmids workflow
  - Select 'Run a single workflow'

- Enter a Project/Run name
  - Example: LastName_geNomad_Test

- Optional: enter a project description

- Select button to the right of 'Select a File' for the Input Assembled Fasta File

- Go to public data ➜ virus_plasmid ➜ you can choose any of the test files here (we recommend one of the ones that says _over5k)

- Run Option: Default

- Select Submit

# Results

- Go to Public Projects; Select button on corresponding



Viruses_Plasmids_Training_Test          virus_plasmid          Complete          1/17/2023, 10:42:44 AM

- General: provides run information
- Explore the virus_plasmid results!
- Browser/Download outputs: Provides downloadable files

# Other Runs

- You can run any of the workflows at any time!
  - Can upload and run your own data
- Feel free to explore other public project results
  - Metagenome_pipeline_test
- Send comments or issues to: nmdc-edge@lanl.gov

# Questions?

# Community-centered design process





## 2023-2024 User Research



34
User research participants

321
Insights from user interviews

120
Action items from insights

- The NMDC program utilizes a community-centered design approach to its product development process
  - User interviews, usability testing, beta-testing
  - Understand the needs of the community and how researchers use the products
- Workshops and event feedback is critical too
- Want to develop products that are as useful to the community as possible

# Community-centered design process



## 2023-2024 User Research

| 34 | 321 | 120 |
|---|---|---|
| User research participants | Insights from user interviews | Action items from insights |

NMDC EDGE Feedback examples
- Data visualization updates
- Support for batch processing
- Updates to how results are displayed and accessed
- Import SRA data easily with SRA accession numbers

# NMDC EDGE Beta Testing

NMDC EDGE feedback areas:
- Running the workflows in NMDC EDGE
- NMDC EDGE user interface
- NMDC workflow training materials

Interested in becoming a beta-tester? We will give you credit on the website and through ORCiD!

- Volunteer at https://microbiomedata.org/user-research/
- Select "NMDC EDGE" for product interest

Implementation of NMDC EDGE beta-tester feedback



4.4%

15.6%

80%

- Implemented suggestions
- Working to implement suggestions
- Cannot implement suggestions

# Reporting Mechanisms

- Beta tester form (NMDC EDGE homepage)
  - Specific workflows that you ran
  - Information about file sizes, issues with jobs (killed or errors), etc
- Email: nmdc-edge@lanl.gov
  - Can get troubleshooting help with NMDC EDGE workflows
- General reporting form
  - Did you come across any issues when using an NMDC product?
  - Features feedback
  - https://forms.gle/yxu9gkbufPigtbrB8

nmdc

National Microbiome
Data Collaborative

Questions?

# NMDC Resources



**Website:** https://microbiomedata.org/
**Data Portal:** https://data.microbiomedata.org/
**Submission Portal:** https://data.microbiomedata.org/submission/home
**NMDC EDGE:** https://nmdc-edge.org/home
**Github:** https://github.com/microbiomedata
**Docker Hub:** https://hub.docker.com/u/microbiomedata
**Documentation:**
https://nmdc-documentation.readthedocs.io/en/latest/overview/nmdc_overview.html
**YouTube:** https://www.youtube.com/channel/UCyBqKc46NQZ_YgZlKGYeglw/featured

## Get involved!

🌐 **Sign up for our newsletter**
microbiomedata.org

💬 **Become a NMDC Champion**
bit.ly/champions-program

𝕏 **Find us on X/Twitter**
@microbiomedata

in **Find us on LinkedIn**
https://bit.ly/NMDC_LinkedIn

📷 **Find us on Instagram**
@microbiomedata

## Read more about the NMDC

*nature microbiology* — Kelliher JM *et al*. Cohort-based learning for microbiome research community standards. *Nat Microbiol* (2023). doi.org/10.1038/s41564-023-01361-7.

*frontiers* — Hu B, Canon S, Eloe-Fadrosh EA, et al.. Challenges in Bioinformatics Workflows for Processing Microbiome Omics Data at Scale. *Front Bioinform*. 1:826370. (2022) doi: 10.3389/fbinf.2021.826370.

*Nucleic Acids Research* — Eloe-Fadrosh EA *et al*. The National Microbiome Data Collaborative Data Portal: an integrated multi-omics microbiome data resource. *Nucleic Acids Res*. 7;60(D1):D828–D836. (2022) doi: 10.1093/nar/gkab990.

*nature REVIEWS MICROBIOLOGY* — Wood-Charlson, E.M., Anubhav, Auberry, D. *et al*. The National Microbiome Data Collaborative: enabling microbiome science. *Nat Rev Microbiol* **18**, 313–314 (2020). doi.org/10.1038/s41579-020-0377-0

*mSystems* — Vangay, P *et al*. Microbiome metadata standards: Report of the National Microbiome Data Collaborative's workshop and follow-on activities. *mSystems* **6**, e01194-20 (2021). doi.org/10.1128/mSystems.01194-20
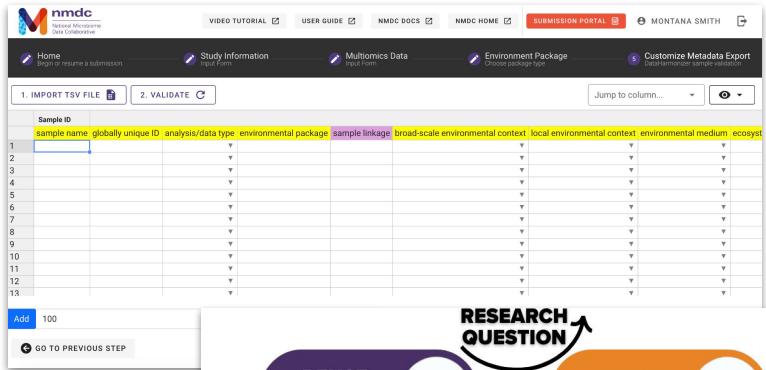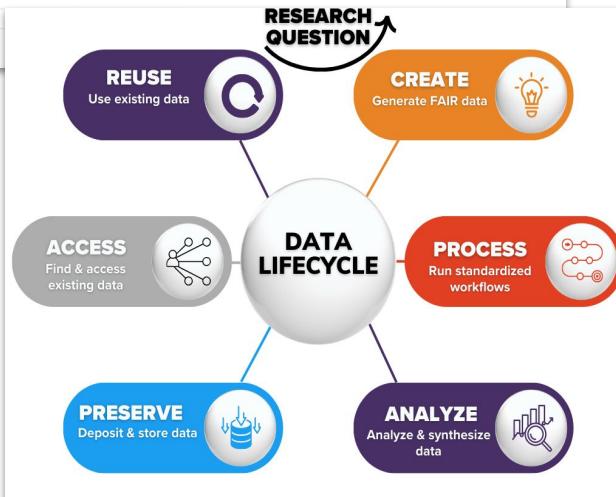
# Metadata Standards and Submission Portal
2024 Ambassador Cohort

# Purpose of this training

- Provide an overview of the benefits of metadata standardization and how this promotes interoperability
- Discuss the importance of metadata standards to enable data reuse
- Provide template slides to use for events
- Introduce audience to the NMDC Submission Portal
    - The hands-on activity can be used (or modified and used) for use in your events too!

What are metadata?

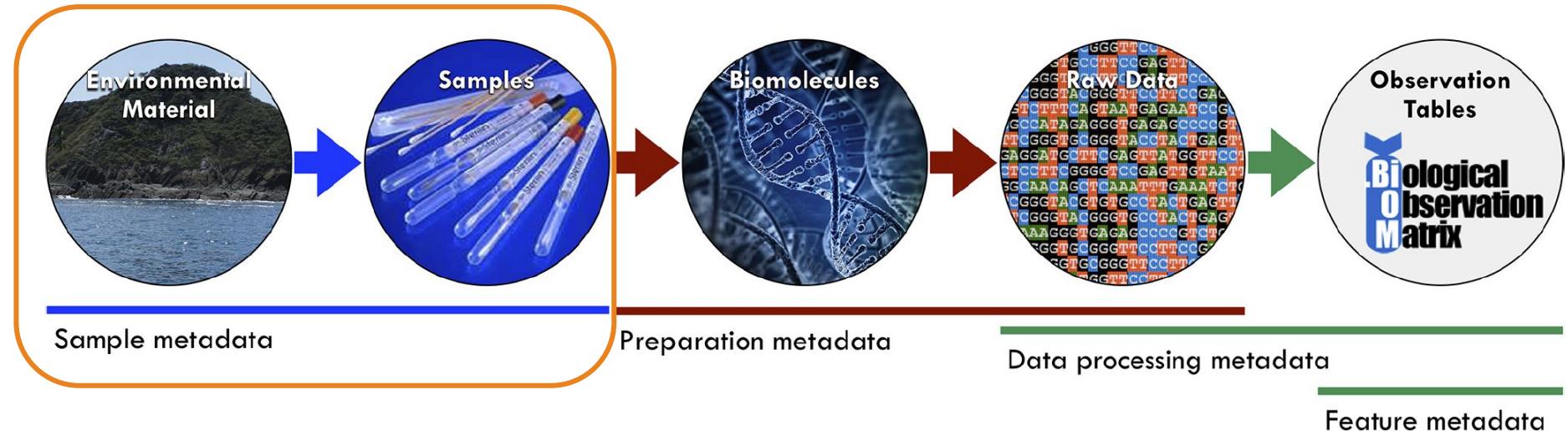# What are Metadata?



Metadata are …

- Contextual data about your data
- Vital for data
  - Publication & deposition
  - Preservation
  - Discovery
  - Access
  - Reuse

# Sample Metadata

## Microbiome Environmental & Sample Metadata



From: **Introduction to Metadata and Ontologies:** Everything You Always Wanted to Know About Metadata and Ontologies (But Were Afraid to Ask)  DOI: 10.25979/1607365

# What are Metadata?
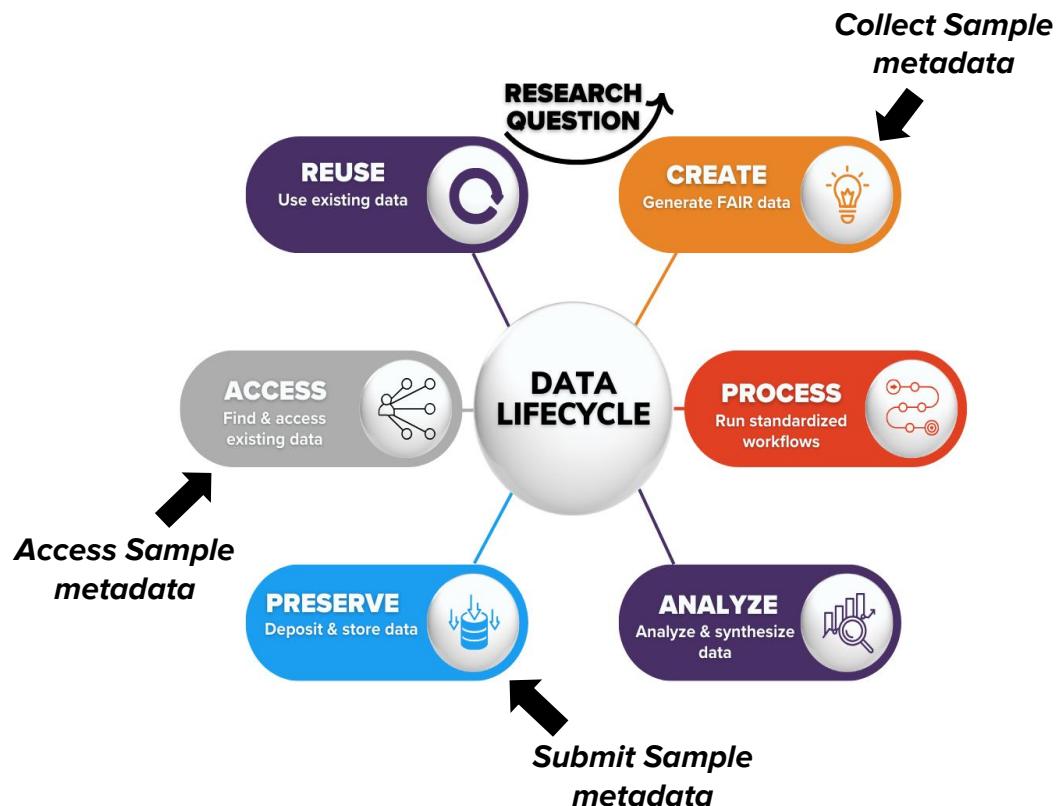


**Sample metadata** includes information about:

- **When** it was collected
- **Where** it was collected
- **How** it was collected
- **What** kind of sample is it
- **Treatment** applied during experimentation
- **Environmental Properties** from which the sample was taken

# Data Lifecycle and Metadata



**During experimental design:**

- Plan for every stage: be intentional about the process of collecting, storing, processing, and protecting (meta)data

- Detail a data management plan that includes metadata for your samples

# Breadth of information in microbiome science



**Sequencing method**

Nucleic acid extraction
Sequencing method

**Multiple different biogeochemical and environmental parameters**

climate
elevation

veg type
plant part

lat/lon
biome
type
material
...

pH
temperature
carbon
nitrogen
...

depth
soil type
land use

**Data Types**

Metagenomes

Metaproteomes

Metatranscriptomes

Metabolomes

**Study information**

PI & Contributors
experimental design

**Different sample Treatments and preparations**

Sample collection device
Sample processing

**Analysis outputs**

assembly statistics
gene function
metabolite and peptide counts
taxon abundance

# Valid and complete metadata

- Lots of samples and studies exist with little to no metadata
- Metadata can be incorrect (e.g. pH of 100) or not usable
- For your own studies, important that you collect valid and complete metadata for data to be FAIR

| Sampling date | 04.06.2023 |
|---|---|
| pH | 100 |
| Geographic location | n/a |
| Host | Plant |
| Elevation | 10,000 m |
| Sequencing technology | |

**Real example from NCBI SRA**

| collection date | not applicable |
|---|---|
| broad-scale environmental context | not applicable |
| local-scale environmental context | not applicable |
| environmental medium | not applicable |
| geographic location | not applicable |
| latitude and longitude | not applicable |

# Valid and complete metadata



| Sampling date | 04.06.2023 |
|---|---|
| pH | 100 |
| Geographic location | n/a |
| Host | Plant |
| Elevation | 10,000 m |
| Sequencing technology | |

→ Format not standardized, not ISO compliant

→ Invalid entry, not in range

→ Field not filled out

→ Entry not very specific

→ Likely incorrect units

→ Missing fields

**nmdc**
National Microbiome Data Collaborative

# Why metadata standards?



Datasets can be difficult to compare and reuse if they lack consistent language and formatting

Why aren't these datasets directly comparable? What isn't standardized?

| idNumber | material | sample depth | temperature |
|----------|----------|--------------|-------------|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| sampleNum | substance | sample depth | temp |
|-----------|-----------|--------------|------|
| 8725 | dirt | 45    cm | 21.1 |
| 2312 | ground liquid | 105  cm | 7 |

# Why metadata standards?



Metadata fields and titles can be inconsistent



| idNumber | material | sample depth | temperature |
|----------|----------|--------------|-------------|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| sampleNum | substance | sample depth | temp |
|-----------|-----------|--------------|------|
| 8725 | dirt | 45    cm | 21.1 |
| 2312 | ground liquid | 105  cm | 7 |

# Why metadata standards?



## The terms used to describe metadata can be different

| idNumber | material | sample depth | temperature |
|---|---|---|---|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| sampleNum | substance | sample depth | temp |
|---|---|---|---|
| 8725 | dirt | 45   cm | 21.1 |
| 2312 | ground liquid | 105  cm | 7 |

# Why metadata standards?



The units used can be inconsistent and not directly comparable

| idNumber | material | sample depth | temperature |
|---|---|---|---|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| sampleNum | substance | sample depth | temp |
|---|---|---|---|
| 8725 | dirt | 45    cm | 21.1 |
| 2312 | ground liquid | 105  cm | 7 |

# Why metadata standards?

National Microbiome Data Collaborative

Some metadata may lack units or descriptive information

| idNumber | material | sample depth | temperature |
|----------|----------|--------------|-------------|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| sampleNum | substance | sample depth | temp |
|-----------|-----------|--------------|------|
| 8725 | dirt | 45    cm | 21.1 |
| 2312 | ground liquid | 105  cm | 7 |

# Why metadata standards?



Adopting standards for reporting makes data human and machine readable.



| idNumber | material | sample depth | temperature |
|----------|----------|--------------|-------------|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| idNumber | material | sample depth | temperature |
|----------|----------|--------------|-------------|
| 8725 | soil | .45 m | 21.1 °C |
| 2312 | groundwater | 1.05 m | 7 °C |

# Why metadata standards?



Adopting standards for reporting makes data human and machine readable.



| idNumber | material | sample depth | temperature |
|---|---|---|---|
| 3928 | soil | 0.03 m | 23.2 °C |
| 3234 | groundwater | 1 m | 9.02 °C |

| idNumber | material | sample depth | temperature |
|---|---|---|---|
| 8725 | soil | .45 m | 21.1 °C |
| 2312 | groundwater | 1.05 m | 7 °C |

# Consequences of not standardizing metadata



I SENT YOU THE DATA.

THANKS!

...THIS IS A WORD DOCUMENT CONTAINING AN EMBEDDED PHOTO YOU TOOK OF YOUR SCREEN WITH THE SPREADSHEET OPEN.

YEAH? DOES YOUR COMPUTER NOT SUPPORT .NORM FILES? MAYBE YOU NEED TO UPDATE.

SINCE EVERYONE SENDS STUFF THIS WAY ANYWAY, WE SHOULD JUST FORMALIZE IT AS A STANDARD.

https://xkcd.com/2116/

- Can miss critical contextual information when performing analyses
  - May miss other confounding variables
- Data may not be able to be published or deposited into certain repositories
- Limits data comparisons and reuse within a group and beyond
- Prevents reproducibility
- Difficult to search for, difficult to compare with other datasets

# The standards gap



Metadata standards usage in SRA

**Microbiome Metadata Standards: Report of the National Microbiome Data Collaborative's Workshop and Follow-On Activities**

Pajau Vangay [a], Josephine Burgin [b], Anjanette Johnston[c], Kristen L. Beck [d], Daniel C. Berrios [e], Kai Blumberg [f], Shane Canon[a], Patrick Chain[g], John-Marc Chandonia [a], Danielle Christianson[a], Sylvain V. Costes[e], Joan Damerow[a], William D. Duncan[a], Jose Pablo Dundore-Arias [h], Kjiersten Fagnan[a], Jonathan M. Galazka [e], Sean M. Gibbons [i,j], David Hays[a], Judson Hervey [k], Bin Hu [g], Bonnie L. Hurwitz [f], Pankaj Jaiswal [l], Marcin P. Joachimiak[a], Linda Kinkel[m], Joshua Ladau[a], Stanton L. Martin[n], Lee Ann McCue [o], Kayd Miller [a], Nigel Mouncey[a], Chris Mungall[a], Evangelos Pafilis [p], T. B. K. Reddy [a], Lorna Richardson [b], Simon Roux [q], Lynn M. Schriml[w], Justin P. Shaffer [r], Jagadish Chandrabose Sundaramurthi [a], Luke R. Thompson [s,t], Ruth E. Timme [u], Jie Zheng [v], Elisha M. Wood-Charlson [a], Emiley A. Eloe-Fadrosh [a]

[a]Lawrence Berkeley National Laboratory, Berkeley, California, USA
[b]European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK
[c]National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA
[d]IBM Almaden Research Center, San Jose, California, USA
[e]NASA Ames Research Center, Moffett Field, California, USA

What we found:
- Lack of awareness of metadata standards, inconsistent usage, lack of training
- Metadata managed as Excel spreadsheets
- Not machine-actionable
- Not following FAIR

# Community Metadata Standards



**The NMDC metadata standards utilize and enhance**
**_existing_ community-driven standards**

1. MIxS: Minimum Information about any (x) Sequence
   Genomic Standards Consortium (GSC)

2. GOLD: Genomes OnLine Database
   Joint Genome Institute (JGI)

4. EnvO: Environment Ontology
   Open Biological and Biomedical Ontology (OBO) Foundry

# Standards for describing sample context



**GSC MIxS**

- Minimal Information about any (x) Sequence (MIxS)
- ~600 standardized fields for reporting
  - Different environmental packages
    - Soil
    - Water
    - Sediment
    - Plant-associated
    - ...

# The NMDC + Community Standards

**NMDC derives _required metadata_ from community standards**

**MIxS**

Examples for MIxS:

Sample Identifiers

Growth Facility

Geographic Location (latitude longitude)

Geographic Region (country and/or sea, region)

Collection Date

Sample Material Processing

Storage Conditions & Temperature

**EnvO**

Broad-scale Environmental Context

Local Environmental Context

Environmental Medium

**GOLD**

GOLD Environment Path

# MIxS Environmental Extensions



17 extensions (*currently available in NMDC)
New extensions added with community input

| MIxS Environmental Extensions | |
|---|---|
| *air | *built environment |
| *host-associated | human-associated |
| human-gut | human-oral |
| human-skin | human-vaginal |
| *hydrocarbon resources-cores | *hydrocarbon resources-fluids/swabs |
| *microbial mat/biofilm | *miscellaneous natural or artificial environment |
| *plant associated | *sediment |
| *soil | wastewater/sludge |
| *water | agriculture |
| *Pending:* food, parasite, etc. | |

MIxS

EnvO

GOLD

# EnvO - Sample Metadata Standards

**MIxS**

**EnvO**

**GOLD**

## EnvO (Environment Ontology)

- Dynamic, community resource
- Hierarchical classification of samples by environment
- Mandated by MIxS for environment fields

# MIxS requires environment terms

MIxS descriptors specify the *sample environment* with the Environment Ontology (EnvO)

**broad-scale environmental context**



**local-scale environmental context**



**environmental medium**



http://obofoundry.org/ontology/envo

# The NMDC + Community Standards



| MIxS / EnvO | | |
|---|---|---|
| Broad-scale environment | Local-scale environment | Environmental Medium |
| Freshwater lake biome | Lake Shore | Sediment |
| Freshwater lake biome | Lake | Algal bloom |

EnvO & GOLD terms together gives us improved environmental context to the microbiome biosamples!

| GOLD Ecosystem classification | | | | |
|---|---|---|---|---|
| Ecosystem | Ecosystem Category | Ecosystem Type | Specific Ecosystem | Ecosystem Tree |
| Environment | Aquatic | Freshwater | Lake | Sediment |
| Environment | Aquatic | Freshwater | Lake | Algal bloom |

MIxS

EnvO

GOLD

# GOLD - Sample Metadata Standards

**MIxS**

**EnvO**

**GOLD**

## The Genomes OnLine Database (GOLD)

- Manually curated collection of genome projects and their metadata
- Metadata fields: ~600
- Controlled Vocabulary fields: 76 (3,873 terms)
- Currently contains hundreds of thousands of microbiome biosamples

# GOLD - Five-level ecosystem path



**GOLD Ecosystem classification**  —  **Example: Lake Sediment**

| | |
|---|---|
| **Ecosystem** | Environmental |
| **Ecosystem Category** | Aquatic |
| **Ecosystem Type** | Freshwater |
| **Ecosystem Subtype** | Lake |
| **Specific Ecosystem** | Sediment |

nmdc
National Microbiome Data Collaborative

nmdc

National Microbiome
Data Collaborative

**Activity: Explore the GOLD ecosystem tree viewer to identify your sample metadata path**

**https://gold.jgi.doe.gov/ecosystemtree**

**https://data-sandbox.microbiomedata.org/**
For workshops, exploration; submissions deleted every 7 days on Sundays

or

**https://data.microbiomedata.org/**
Navigate to Products ➜ Submission Portal

# The NMDC Submission Portal

# The NMDC Submission Portal

# The NMDC Submission Portal

# The NMDC Submission Portal

# The NMDC Submission Portal

# The NMDC Submission Portal



Submission Context
- Indicate if data already exists
- Are you using NMDC to complete metadata for samples going to a DOE user facility?

# The NMDC Submission Portal



## Study Information ⓘ

A study summarizes the overall goal of a research initiative and outlines the key objective of its underlying projects.

**Study Name ***
Name is required

**Principal Investigator Name**
The Principal Investigator who led the study and/or generated the data.

**Principal Investigator Email ***
E-mail is required

**Principal Investigator ORCID**
ORCID iD of the Principal Investigator.

**Webpage Links** ▼
Link to the Principal Investigator's research lab webpage or the study webpage associated with this collection of samples. Multiple links can be provided.

**Study Description**
Provide a description of your study. This should include some general context of your research goals and study design. For examples, please see existing study landing pages on the data portal.

**Optional Notes**
Add any additional notes or comments about this study.

# The NMDC Submission Portal



**Contributor Roles & Access Permissions**

CRediT Roles : https://credit.niso.org/

Permission Levels

- Author/Owner, Editor, Metadata Contributor, Viewer

# The NMDC Submission Portal

# The NMDC Submission Portal

| Home | Submission Context | Study Information | Multi-omics Data | Environment Package | Customize Metadata Export |
|------|-------------------|------------------|-----------------|--------------------|--------------------------|
| Begin or resume a submission. | Input form | Input Form | Input Form | 5 Choose package type | 6 DataHarmonizer sample validation |

## Environment Package

Choose environment package for your data.

- ◯ air
- ◯ built environment
- ◯ host-associated
- ◯ hydrocarbon resources - cores
- ◯ hydrocarbon resources - fluids swabs
- ◯ microbial mat_biofilm
- ◯ miscellaneous natural or artificial environment
- ◯ plant-associated
- ◯ sediment
- ⦿ soil
- ◯ water

Under development

- ◯ human-associated
- ◯ human - gut
- ◯ human - oral
- ◯ human - skin
- ◯ human - vaginal

# The NMDC Submission Portal

# The NMDC Submission Portal



Templates are downloadable if users prefer to work in Excel

# The NMDC Submission Portal



Users can then upload excel files with completed metadata fields for validation

# The NMDC Submission Portal



A color key assists researchers in knowing types of errors and required and recommended fields

# The NMDC Submission Portal



The validate button allows for live validation and will point out invalid and empty cells

# The NMDC Submission Portal



The Column Help window can help users address errors

# The NMDC Submission Portal



Users can show/hide certain columns

# The NMDC Submission Portal



Users can also search for and jump to certain columns

NMDC Data Portal

# NMDC Data Portal



**Data Portal & API**

Access and discovery of microbiome information

Filter datasets by standardized metadata terms

Questions?

Activities

# Submission Portal Activity

- Select '**Create New Submission**'
- For Submission Context, select 'No' for if the data has been generated for this study
- Select 'JGI' then select 'Other' and type something along the lines of 'test submission'
- Select Go to Next Step
- Name the study: "Ambassador Training: Validation activity"
- Put a fake email address or your own email address into the Principal Investigator Email field
- No need to fill out non-required information on this page
- Select Go to Next Step
- Select metagenome under JGI and type a fake study ID (123456)
- Select Go to Next Step
- Select the Soil Environment Package
- Select Go to Next Step
- **<u>Pause</u>** for discussion
- Upload the xlsx file, select the 'Validate' button and try to fix all of the invalid fields!

# Example Break Out Activity



**A more challenging example activity**
- Complete study and data information using a written study design
  - Ambassadors can write out a paragraph of study design (we have examples you can use)
  - Attendees read the paragraph and identify relevant metadata fields
  - Attendees complete the study and data information sections of the Submission Portal

# Discussion

- What do you currently use to track metadata and sample information?
- Which metadata standards or templates (if any) are your group currently using?
- What did you like about the NMDC Submission Portal?
- What could be improved about the Submission Portal?
- Did you encounter any bugs in the Submission Portal?

# **Submission Portal Experience**

- Do you have projects you would like to track in the Submission Portal?
  - Walk through your Submission Portal experience with our team!
  - Please add your name to shared notes to participate
- Live feedback from your workshops
  - Features people like/don't like
  - Things that aren't clear during the workshop
  - Bugs that come up while people are doing activities
- General reporting google form for issues: https://forms.gle/yxu9gkbufPigtbrB8

# Community-centered design process



Submission Portal Feedback
- Updates to how the templates are accessed and downloaded
- Feature and UI updates to improve experience with interface
- Help and tutorial guidance updates

Questions?

# NMDC Resources



**National Microbiome Data Collaborative**

**Website:** https://microbiomedata.org/
**Data Portal:** https://data.microbiomedata.org/
**Submission Portal:** https://data.microbiomedata.org/submission/home
**NMDC EDGE:** https://nmdc-edge.org/home
**Github:** https://github.com/microbiomedata
**Docker Hub:** https://hub.docker.com/u/microbiomedata
**Documentation:**
https://nmdc-documentation.readthedocs.io/en/latest/overview/nmdc_overview.html
**YouTube:** https://www.youtube.com/channel/UCyBqKc46NQZ_YgZlKGYeglw/featured

## Get involved!

**Sign up for our newsletter**
microbiomedata.org

**Become a NMDC Champion**
bit.ly/champions-program

**Find us on X/Twitter**
@microbiomedata

**Find us on LinkedIn**
https://bit.ly/NMDC_LinkedIn

**Find us on Instagram**
@microbiomedata

## Read more about the NMDC

Kelliher JM *et al.* Cohort-based learning for microbiome research community standards. *Nat Microbiol* (2023). doi.org/10.1038/s41564-023-01361-7.

Hu B, Canon S, Eloe-Fadrosh EA, et al.. Challenges in Bioinformatics Workflows for Processing Microbiome Omics Data at Scale. *Front Bioinform*. 1:826370. (2022) doi: 10.3389/fbinf.2021.826370.

Eloe-Fadrosh EA *et al*. The National Microbiome Data Collaborative Data Portal: an integrated multi-omics microbiome data resource. *Nucleic Acids Res.* 7;60(D1):D828–D836. (2022) doi: 10.1093/nar/gkab990.

Wood-Charlson, E.M., Anubhav, Auberry, D. *et al.* The National Microbiome Data Collaborative: enabling microbiome science. *Nat Rev Microbiol* **18,** 313–314 (2020). doi.org/10.1038/s41579-020-0377-0

Vangay, P *et al.* Microbiome metadata standards: Report of the National Microbiome Data Collaborative's workshop and follow-on activities. *mSystems* **6**, e01194-20 (2021). doi.org/10.1128/mSystems.01194-20

# Data Stewardship and NMDC Data Portal
2024 Ambassador Cohort

# Purpose of this training



- Provide an overview of data stewardship and relevant guiding principles
- Discuss the importance of proper data stewardship for the future of microbiome data
- Provide template slides for events
- Introduce audience to the NMDC Data Portal
  - The hands-on scavenger hunt can be used (or modified and used) for use in your events too!

Findable 🔍
Accessible 🔒
Interoperable ⚙️
Reusable ♻️

# Data Stewardship & FAIR Data

# Microbiome Data

PubMed Results for "Microbiome" over Time

# The immense scale of omics data



Advances in sequencing and omics technologies have **far outpaced** data infrastructure

Source A

Target C

Source B

Target D

# The immense scale of omics data



Advances in sequencing and omics technologies have **far outpaced** data infrastructure

# Data Stewardship
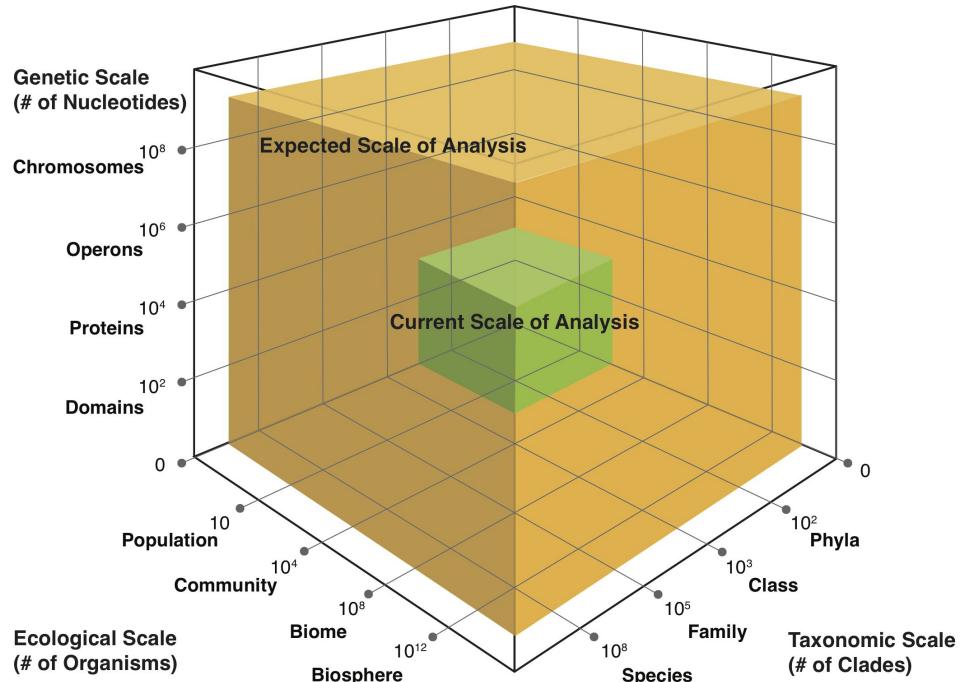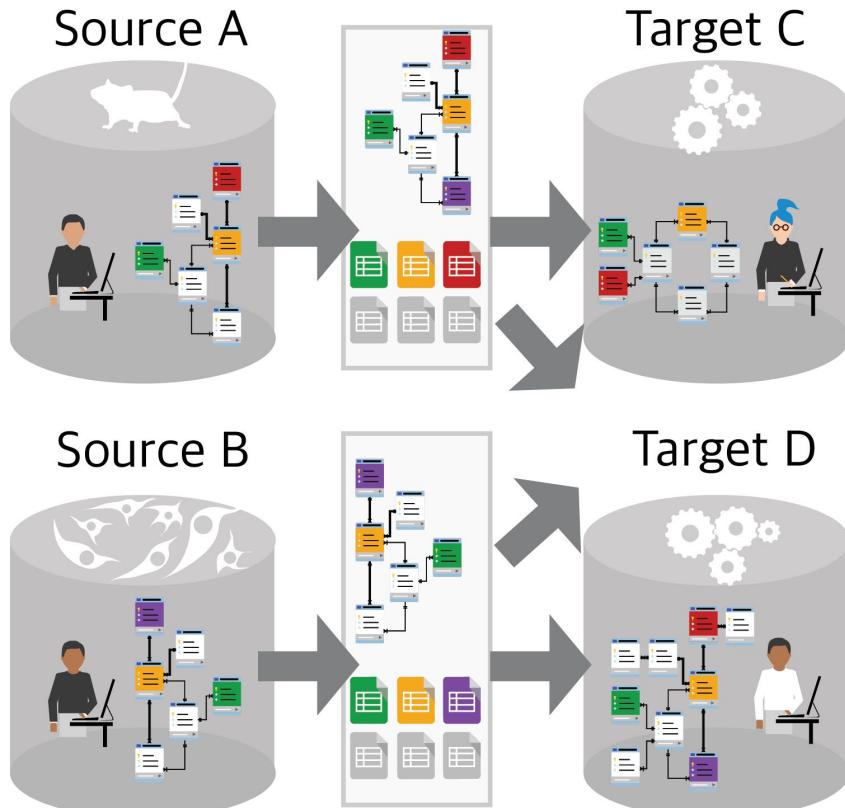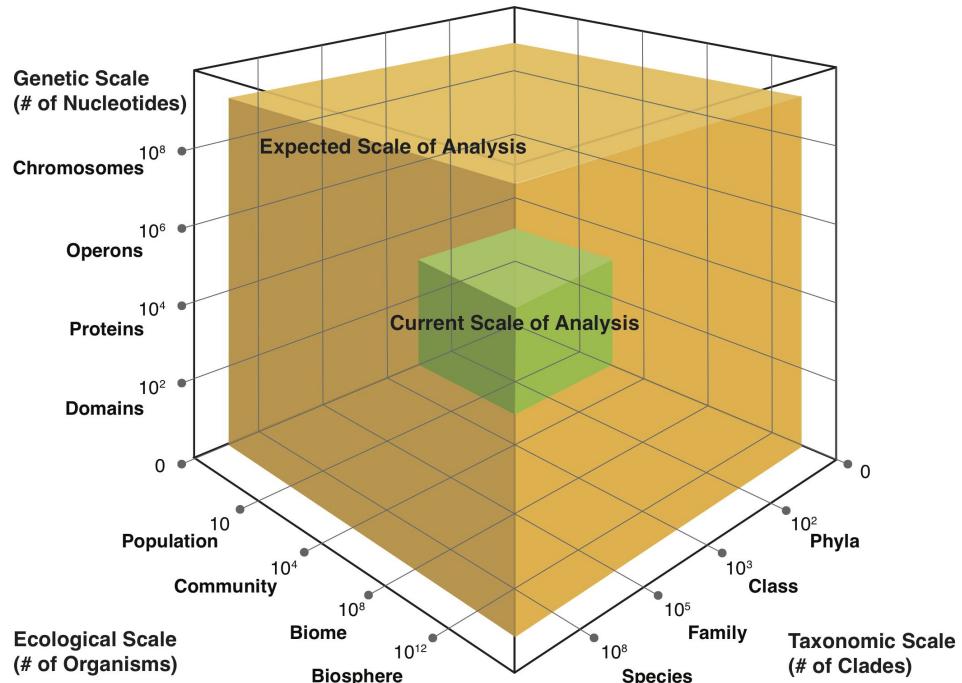
Aspects of data stewardship:

- Data generation and acquisition
- Data storage
- Data processing
- Data publication, release
- Data reuse

Everyone's responsibility to ensure data is properly managed, trustworthy, and FAIR

# Accountability in data stewardship



- All researchers
- Funders
- User facilities
- Publishers
- Societies
- Institutions
- Data storage facilities

Important to implement data stewardship best practices as early in the research process as possible

# Why should you care about data management and data stewardship?

Streamlines your research process
- Less data loss and waste, easier to find information
- Can answer new scientific questions, easily share data
- Publication of data
- Data preservation in the future

Makes data accessible for others
- Collecting and providing data and metadata allows other researchers to understand your full study context for data reuse
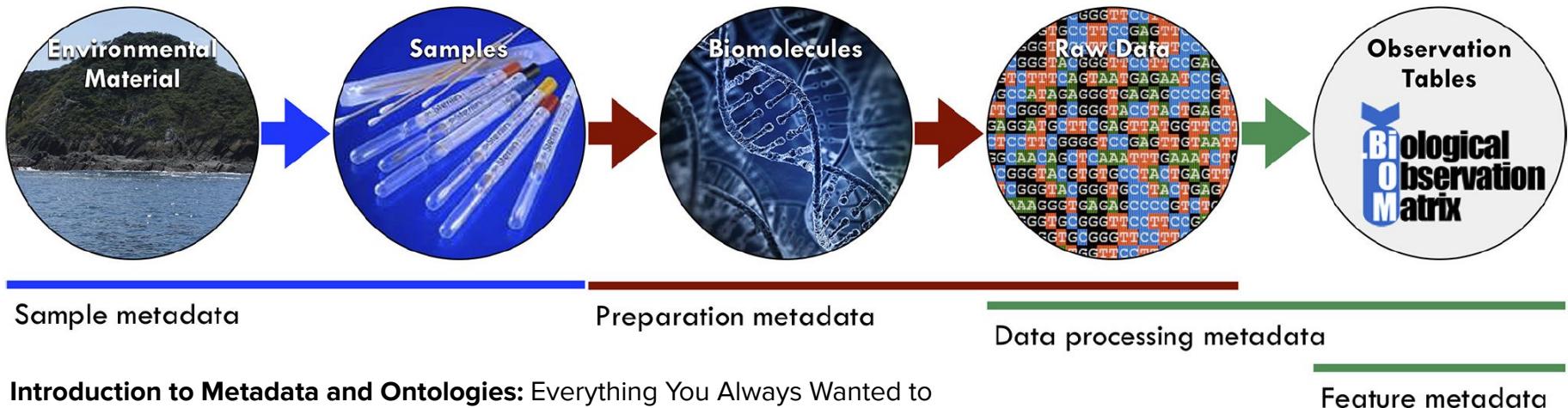- Increased scientific exposure

Ensures you receive credit for your work
- Datasets can get DOIs, included in publications
- Poorly managed data may not retain provenance or may not be able to be reused, limiting your contributions to your field
- Impact is quantifiable, measurable

# Metadata standards
## *A pillar for data stewardship best practices*

**nmdc**
National Microbiome
Data Collaborative

Environmental Material → Samples → Biomolecules → Raw Data → Observation Tables

Sample metadata

Preparation metadata

Data processing metadata

Feature metadata

**Introduction to Metadata and Ontologies:** Everything You Always Wanted to Know About Metadata and Ontologies (But Were Afraid to Ask)

DOI: 10.25979/1607365

# FAIR Data

- FAIR is about:
  - Data and metadata
    - Metadata and metadata standards should be articulated and made publicly available to the greatest extent possible
  - Machine-actionability
    - Relevant on all levels of data aggregation
    - Human and machine readable considerations
  - Controlled data access
    - Explicit, well-defined and readily available terms and conditions under which data will be shared or made accessible

FAIR guiding principles: https://doi.org/10.1038/sdata.2016.18

0101 **Findable**
Ensure all data registered within NMDC are human and machine readable

**Accessible**
Identify data sets that are available, including any authentication and authorization requirements

**Interoperable**
Provide provenance, metadata, and uniformly processed data, we are lowering the barriers to making data interoperable

**Reusable**
Enable download of data, data products, and workflows for external reprocessing

# Findable

The first step in (re)using data is to find datasets. Metadata and data should be machine-readable and easy to find by the community.

- (Meta)data are assigned a globally unique and persistent identifier
- Data are described with rich metadata
- Metadata clearly and explicitly include the identifier of the data they describe
- (Meta)data are registered or indexed in a searchable resource

# Accessible



Once the user finds the required data, they need to know how the data can be accessed, possibly including authentication and authorization.

- (Meta)data are retrievable by their identifier using a standardized communications protocol
- The protocol is open, free, and universally implementable
- The protocol allows for an authentication and authorization procedure, where necessary
- Metadata are accessible, even when the data are no longer available

# Interoperable

Data usually needs to be integrated with other data. In addition, data needs to interoperate with applications or workflows for analysis, storage, and processing.

- (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- (Meta)data use vocabularies that follow FAIR principles
- (Meta)data include qualified references to other (meta)data

# Reusable

The ultimate goal of FAIR is to optimize the reuse of data. To achieve this, metadata and data should be well-described so that they can be replicated and/or combined in different settings.

- (Meta)data are richly described with a plurality of accurate and relevant attributes
- (Meta)data are released with a clear and accessible data usage license
- (Meta)data are associated with detailed provenance
- (Meta)data meet domain-relevant community standards

# Open Science & Open Data

Why is openly accessible data important?

- Part of **FAIR** data practices
- Scientific advances can be made by synthesizing new and old studies
- Open data can create new research avenues
- Open data increases accessibility to all researchers

# CARE principles

For inclusive development and innovation

C1

Collective Benefit

C2 For improved governance and citizen engagement

C3 For equitable outcomes

Recognizing rights and interests

A1

Authority to Control

A2 Data for governance

A3 Governance of data

For positive relationships

R1

Responsibility

R2 For expanding capability and capacity

R3 For Indigenous languages and worldviews

For minimizing harm and maximizing benefit

E1

Ethics

E2 For justice

E3 For future use

- Indigenous data: "Data generated by Indigenous Peoples or by other governments and institutions on or about Indigenous Peoples and territories"

- CARE principles aim for data stewardship practices that align with Indigenous interests and governance needs
    - Making data FAIR while acknowledging power differentials and historical contexts
    - People- and person-oriented, reflecting the crucial role of data in advancing Indigenous innovation and self-determination

CARE principles for Indigenous Data Governance

# CARE Principles



The CARE Principles are available in Spanish, Vietnamese, Māori, German, and Khmer

CARE principles for Indigenous Data Governance

# IDEA and Data Stewardship

*Inclusion* means creating environments, large and small, that foster welcoming and belonging.

*Diversity* refers to the variety of backgrounds, cultures, disciplines, approaches, perspectives, and ways in which we solve problems.

*Equity* means achieving the aspirational state of 'opportunity parity' for all. To achieve equity, we must increase access to and remove barriers to opportunity, taking into consideration individual needs wherever feasible.

*Accountability* means taking individual and collective responsibility for our actions, behavior, and impact on others.

How can **data stewardship** incorporate IDEA principles?

- ○ Openly accessible data makes resources more equitable
- ○ Training and educational resources make research more accessible
- ○ From CARE principles - how does your data impact groups involved, who benefits from your data, is anyone excluded who should not be?

https://diversity.lbl.gov/ideaberkeleylab/

**Data Management**

# Data Lifecycle and Metadata



Best Practices in Data Management support the entire Data Lifecycle

# Data Management

For research: the *intentional* process of collecting, storing, processing, and protecting data

For data preservation: data are **FAIR** now and into the future

Benefits of good data management:
- Decreased data loss
- Streamlined data deposition and publication
- Work is appropriately credited
- Impact is measurable / quantifiable
- Increased scientific exposure



From: Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, Cain KE, et al. (2014) Troubleshooting Public Data Archiving: Suggestions to Increase Participation. PLoS Biol 12(1): e1001779. https://doi.org/10.1371/journal.pbio.1001779

# What is a Data Management Plan?

**Required part of any federally funded grant proposal**

- How you will collect, categorize, store, and share any data produced during the duration of a grant
- How that data will be preserved and made accessible after the completion of a project

**Impact:**

- Well-managed data can be published as a product
  - DOIs assigned to datasets
- Data management plans extend beyond the researchers
- Leads to greater re-use of data (internally and externally)

# What to include in your DMP?

**SAMPLE AND DATA TYPES AND SOURCES**

Outlines what kinds of data will be produced throughout the project.

**DATA STANDARDS AND FORMATS**

Defines all variables of interest and communicates that you are aware of and will abide by community best practices whenever possible.

## What goes into a data management plan?

**DATA DISSEMINATION & ARCHIVING**

Describes what the final data products will be and how you will protect data, if applicable.

**POLICIES FOR DATA SHARING, PUBLIC ACCESS, AND RE-USE**

Communicates that you understand your funders data sharing policies and that you have a plan to ensure public availability.

**DATA AND SAMPLE PRESERVATION**

Communicates the sustainability plan for your data, showing your funder that the data products will last after the completion of the project.

**ROLES AND RESPONSIBILITIES**

Shows how your data management plan will be executed and ensures that your team's data management responsibilities are clearly defined.

# NMDC DMP Resources

## DMPTool Template



## DMP Consultancies

**One-on-one guidance on how to make a FAIR Data Management Plan**

NMDC team has the expertise to offer guidance on the creation of DMPs for the **Department of Energy Office of Science** proposals



https://microbiomedata.org/data-management/

# Tools for Data Management

Raise your hand if you use the following for managing your data:

- Lab notebooks
- Field notes, pieces of paper
- Remembering the information
- Google sheets
- Excel
- Tablet, iPad
- Laboratory Information Management Systems (LIMS) systems
- Online tools, data management plan software
- Post-it notes, paper towels, the back of a glove in the lab

# Data Management Reflection

- How do you implement data management best practices in your organization?
  - Is there anything you would add/change to make your data more FAIR?
  - What tools work for you for data management? Are there tools that could improve this process for you?
  - How can we implement the principles we discussed and ensure accountability?
- Are the FAIR principles enough?

Questions?

# NMDC Data Portal

# NMDC Data Portal



https://data.microbiomedata.org

# 'Omics Type

# Geographic Location

# Collection Date

# Combination of 'omics types

# Consortia and Studies



National Microbiome Data Collaborative

# Samples

# KEGG Term Search

# Other searchable metadata fields

# Environmental Metadata

# Environment tab

# Study Pages

# Data Download

Need to log in with ORCiD to download data

Can directly download all files from the processed multi-omics data

# Instructional Content

Data Portal User Guide:

https://nmdc-documentation.readthedocs.io/en/latest/howto_guides/portal_guide.html

Data Portal Tutorial Video:

https://nmdc-documentation.readthedocs.io/en/latest/tutorials/nav_data_portal.html

Data Portal Documentation:

https://nmdc-documentation.readthedocs.io/en/latest/reference/data_portal.html

More information found by clicking the "?"s

# NMDC API

The NMDC API
- Auto-generated documentation at and UI at
  https://api.microbiomedata.org/docs#
- User friendly detailed documentation



sites

A site corresponds to a physical place that may participate in job execution.

A site may register data objects and capabilties with NMDC. It may claim jobs to execute, and it may update job operations with execution info.

A site must be able to service requests for any data objects it has registered.

A site may expose a "put object" custom method for authorized users. This method facilitates an operation to upload an object to the site and have the site register that object with the runtime system.

| GET | /sites | List Sites |
| POST | /sites | Create Site |
| GET | /sites/{site_id} | Get Site |
| GET | /sites/{site_id}/capabilities | List Site Capabilities |
| PUT | /sites/{site_id}/capabilities | Replace Site Capabilities |

Endpoints are color coded based on function.

# The NMDC API



find **Find NMDC metadata entities.**

| GET | /studies | Find Studies |
| GET | /studies/{study_id} | Find Study By Id |
| GET | /biosamples | Find Biosamples |
| GET | /biosamples/{sample_id} | Find Biosample By Id |
| GET | /data_objects | Find Data Objects |
| GET | /data_objects/{data_object_id} | Find Data Object By Id |
| GET | /activities | Find Activities |
| GET | /activities/{activity_id} | Find Activity By Id |
| GET | /search | Search Page |
| GET | /pipeline_search | Pipeline Search |
| POST | /pipeline_search | Pipeline Search |
| POST | /pipeline_search_form | Pipeline Search |

GET function retrieves records

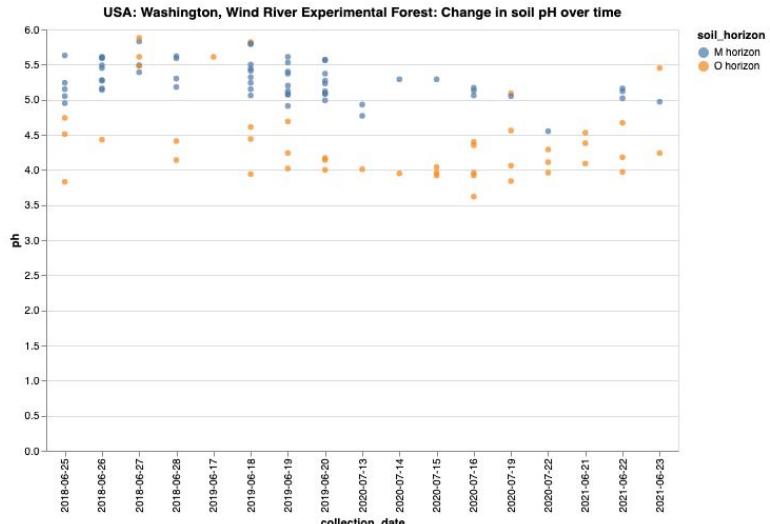POST function creates records

Find study or biosample metadata by ID

# Available Jupyter notebooks

The NMDC has several publicly available Jupyter notebooks that researchers can explore, use, and modify

https://github.com/microbiomedata/notebook_hackathons/tree/main

Questions?

# Adding Data

- Do you have data you'd like to add to the data portal?
- We are accepting environmental microbiome omics data with sufficient metadata
- NMDC team can process for you and upload to the Portal
- Why add your data?
  - Increased citations, increased visibility, inclusion in meta-analyses, comparisons with other datasets in the portal, ethical benefits
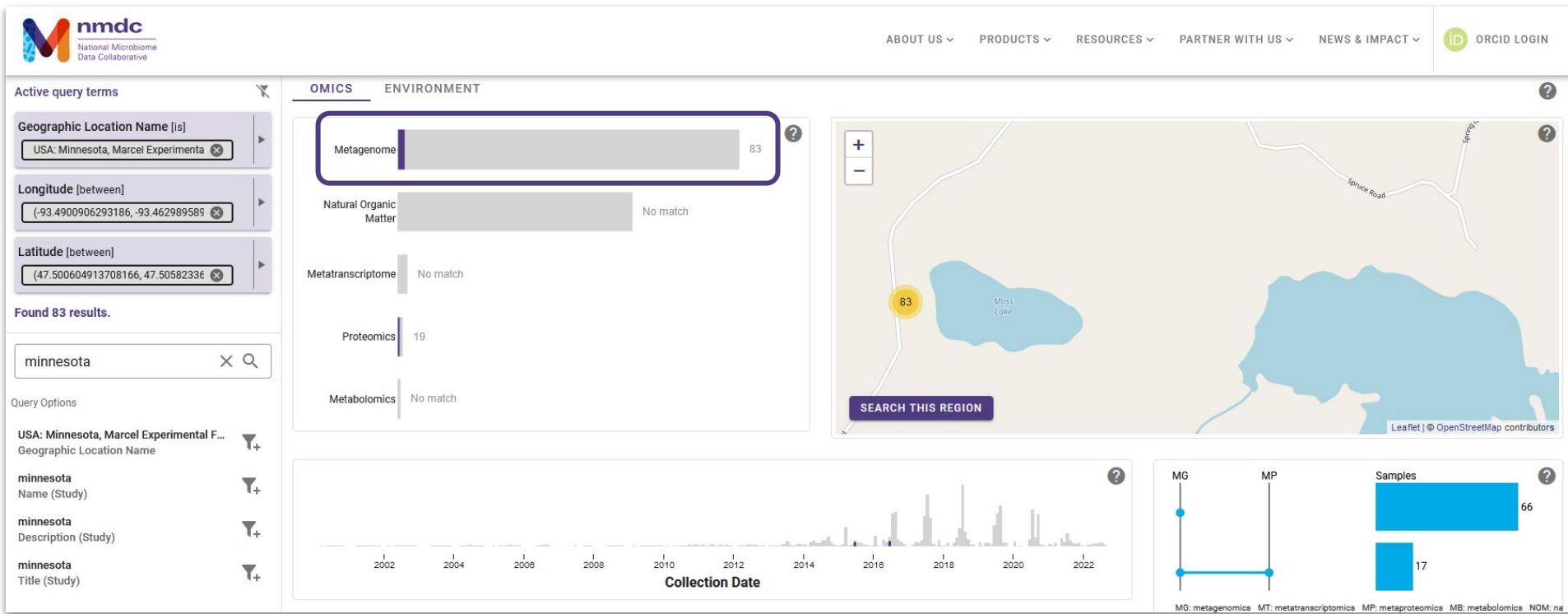
# Scavenger Hunt

How many metagenome samples were collected next to Moss Lake, Minnesota?

# Scavenger Hunt

How many metagenome samples were collected next to Moss Lake, Minnesota?

# Scavenger Hunt

How many metagenome samples were collected next to Moss Lake, Minnesota?

- Find where to download the metagenome Reads QC statistics file from the first resulting sample

# Scavenger Hunt



## How many metagenome samples were collected from peatland soil in 2016?

- ○ Find the link to the data in IMG

How many Natural Organic Matter samples are from Freshwater biomes?

# Scavenger Hunt



## How many Natural Organic Matter samples are from Freshwater Biomes?

# Scavenger Hunt

How many SPRUCE biosamples have been characterized by both metagenomics _and_ metaproteomics?

# Scavenger Hunt

How many SPRUCE biosamples have been characterized by both metagenomics _and_ metaproteomics?

# Scavenger Hunt

How many samples are from Puerto Rico?

# Scavenger Hunt



nmdc
National Microbiome
Data Collaborative

## How many samples are from Puerto Rico?

# Data Portal Feedback

- Were the metadata search fields comprehensive? Are there any other metadata terms you would want to search or sort by?
- Favorite feature(s)? Least favorite?
- Was the Data Portal intuitive to use?
- What did you look for? What were you hoping to find? Could you always find what you were looking for?
- Any sticking points?
- Other datasets or data types you'd like to see?
- Other questions?

# Community-centered design process



Data Portal Feedback
- Bulk data download updates
- Visualization updates - map, upset plot
- Implementing taxonomy-based search
- Improved search and filtering

# NMDC Resources

**Website:** https://microbiomedata.org/
**Data Portal:** https://data.microbiomedata.org/
**Submission Portal:** https://data.microbiomedata.org/submission/home
**NMDC EDGE:** https://nmdc-edge.org/home
**Github:** https://github.com/microbiomedata
**Docker Hub:** https://hub.docker.com/u/microbiomedata
**Documentation:**
https://nmdc-documentation.readthedocs.io/en/latest/overview/nmdc_overview.html
**YouTube:** https://www.youtube.com/channel/UCyBqKc46NQZ_YgZlKGYeglw/featured

## Get involved!

**Sign up for our newsletter**
microbiomedata.org

**Become a NMDC Champion**
bit.ly/champions-program

**Find us on X/Twitter**
@microbiomedata

**Find us on LinkedIn**
https://bit.ly/NMDC_LinkedIn

**Find us on Instagram**
@microbiomedata

## Read more about the NMDC

Kelliher JM *et al*. Cohort-based learning for microbiome research community standards. *Nat Microbiol* (2023). doi.org/10.1038/s41564-023-01361-7.

Hu B, Canon S, Eloe-Fadrosh EA, et al.. Challenges in Bioinformatics Workflows for Processing Microbiome Omics Data at Scale. *Front Bioinform*. 1:826370. (2022) doi: 10.3389/fbinf.2021.826370.

Eloe-Fadrosh EA *et al*. The National Microbiome Data Collaborative Data Portal: an integrated multi-omics microbiome data resource. *Nucleic Acids Res.* 7;60(D1):D828–D836. (2022) doi: 10.1093/nar/gkab990.

Wood-Charlson, E.M., Anubhav, Auberry, D. *et al.* The National Microbiome Data Collaborative: enabling microbiome science. *Nat Rev Microbiol* **18,** 313–314 (2020). doi.org/10.1038/s41579-020-0377-0

Vangay, P *et al.* Microbiome metadata standards: Report of the National Microbiome Data Collaborative's workshop and follow-on activities. *mSystems* **6**, e01194-20 (2021). doi.org/10.1128/mSystems.01194-20