

Enhanced metadata standards supported by the National Microbiome Data Collaborative

Chris Mungall^{1*} (cjmungall@lbl.gov), Faiza Ahmed², Anubhav³, Jeffrey Baumes², Jonathan Beezley², Mark Borkum³, Lisa Bramer³, Shane Canon¹, Patrick Chain⁴, Danielle Christianson¹, Yuri Corilo³, Karen Davenport⁴, Brandon Davis², Meghan Drake⁵, William Duncan¹, Kjersten Fagnan¹, Mark Flynn⁴, David Hays¹, Bin Hu⁴, Marcel Huntemann¹, Julia Kelliher⁴, Sonya Lebedeva¹, Po-E Li⁴, Mary Lipton³, Chien-Chi Lo⁴, Douglas Mans³, Stanton Martin⁵, Lee Ann McCue³, David Millard³, Kayd Miller¹, Nigel Mouncey¹, Paul Piehowski³, Elais Player Jackson⁴, Anastasiya Prymolenna³, Samuel Purvine³, TBK Reddy¹, Rachel Richardson³, Migun Shakya⁴, Montana Smith³, Jagadish Chandrabose Sundaramurthi¹, Deepak Unni¹, Pajau Vangay¹, Bruce Wilson⁵, Donny Winston⁶, Elisha Wood-Charlson¹, Yan Xu⁴, **Emiley Eloë-Fadrosh¹**

¹ Lawrence Berkeley National Laboratory, Berkeley, CA; ² Kitware, Clifton Park, NY; ³ Pacific Northwest National Laboratory, Richland, WA; ⁴ Los Alamos National Laboratory, Los Alamos, NM; ⁵ Oak Ridge National Laboratory, Oak Ridge, TN; ⁶ Polyneme LLC, New York, NY

<https://microbiomedata.github.io/nmdc-metadata/>

Project Goals: Short statement of goals. (Limit to 1000 characters)

The National Microbiome Data Collaborative (NMDC) is a pilot initiative launched to support microbiome data exploration and discovery through a collaborative, integrative science gateway. With a community-centered design approach, the NMDC team is building an open-source, integrated data science ecosystem that leverages existing data standards, data resources, and infrastructure within the DOE complex.

Abstract

To understand microbiomes we need to integrate, analyze, and query large amounts of data, including multi-omics data (e.g., metagenome, metatranscriptome, metaproteome, and metabolome) and environmental data. This is challenging because these data are heterogeneous and complex, and existing standards and ontologies are lacking or incomplete.

For the NMDC project, we created a FAIR (findable, accessible, interoperable, and reusable) schema for handling data and metadata of multiple aspects of microbiome data, including environmental metadata about a sample and a study, metadata and provenance for all processing and workflows, and searchable information arising from annotation workflows (for example, functional annotations and results of binning).

The schema leverages and maps to existing standards where appropriate. For describing sample metadata and environmental characteristics, we leveraged the Genomics Standards Consortium (GSC) MIxS (Minimal Information about any Sequence) and use a combination of ENVO

(Environment Ontology) and GOLD used for classifying environments. This includes a mechanism for uniquely identifying source samples using identifier systems such as IGSN (International Geo Sample Number), allowing us to link together data from different omics processing pipeline connected to the same source sample. We extended the W3C PROV standard (<https://www.w3.org/TR/prov-overview/>) for metadata about computational workflows. For outputs of genomics/transcriptomics workflows, we built on standards such as GFF3, using standardized systems such as KEGG for functional annotation; and for metabolomics/metaproteomics we map to existing ontologies such as PSI-MS where possible.

Our schema weaves together these different standards into a coherent whole. It is rendered as JSON-Schema which allows for precise validation of data input streams using standard validators, as well . We also aim for FAIR compliance by also providing an RDF (Resource Description Framework) version of the schema, including mappings to existing standards.

We used this schema to integrate multiple diverse types of data into JSON-LD files, and to drive search in a web portal.

Funding statement.

This work is supported by the Genomic Science Program in the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research (BER) under contract numbers DE-AC02-05CH11231 (LBNL), 89233218CNA000001 (LANL), DE-AC05-00OR22725 (ORNL), and DE-AC05-76RL01830 (PNNL).