# FAIR principles for 'omic sequence datasets



A.Ponsero - 24.06.25

QIB FAIR WORKSHOPS

# FAIR framework to enhance scientific data reusability



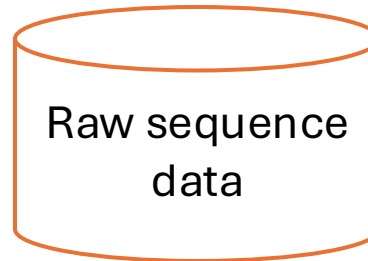**Since the publication of the FAIR data principles in 2016 :**
- Growing recognition of the importance of data management
- Push from institutions and funders to promote FAIR data
- Development of tools and methods for biological FAIR data

# FAIR framework to enhance scientific data reusability
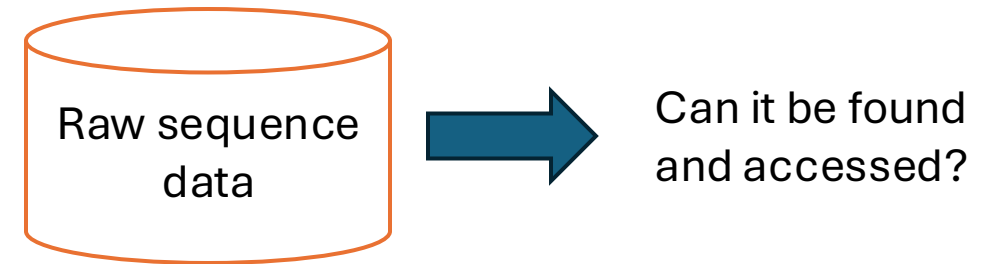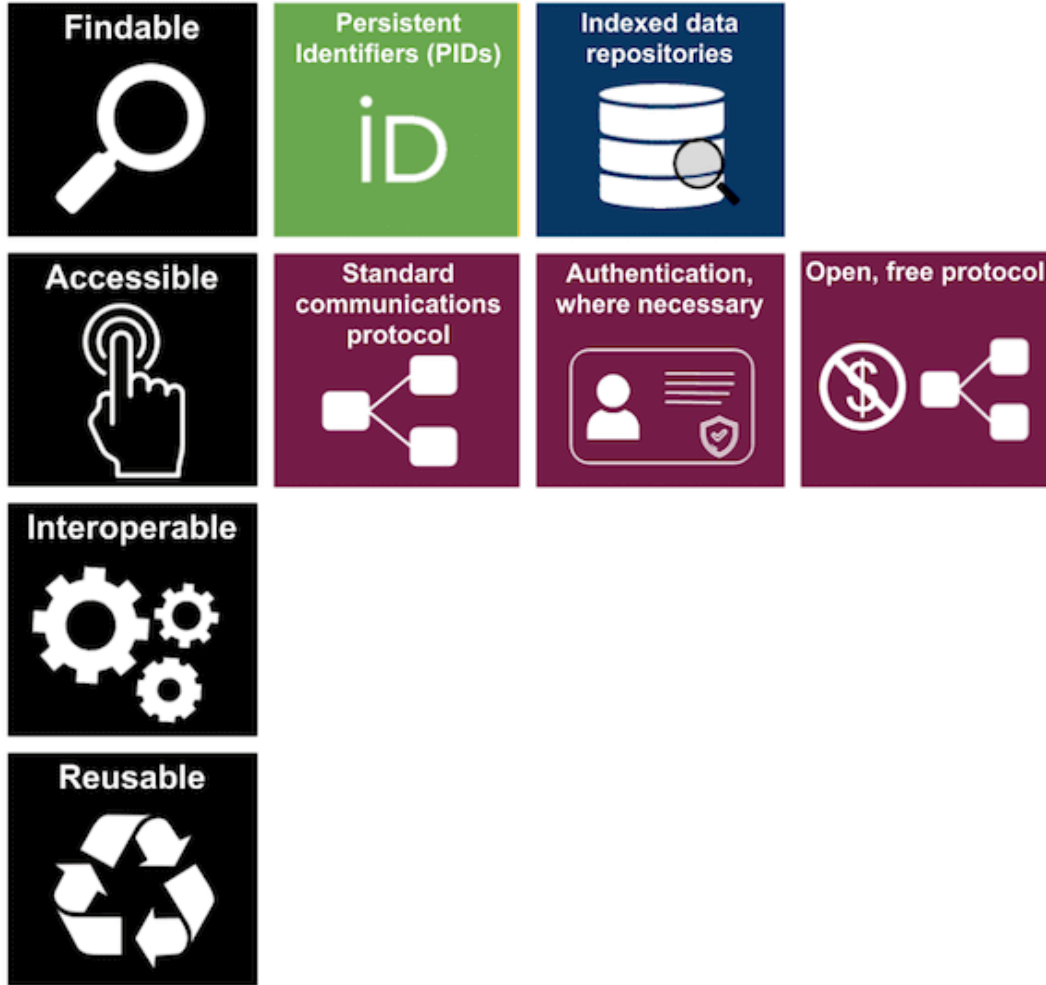


**Since the publication of the FAIR data principles in 2016 :**
- Growing recognition of the importance of data management
- Push from institutions and funders to promote FAIR data
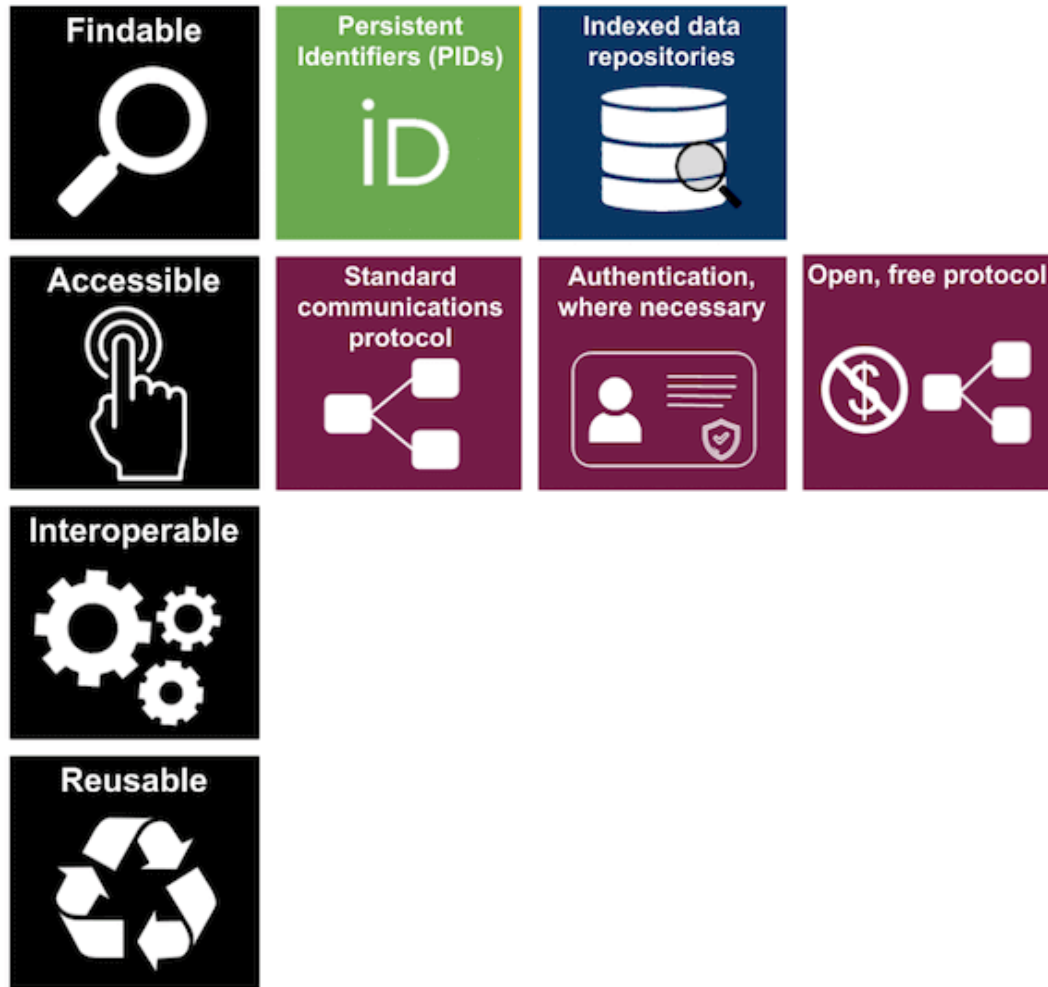- Development of tools and methods for biological FAIR data

Raw sequence data

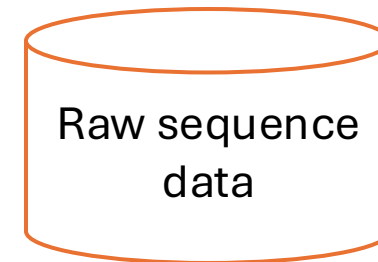What are the specificities of making 'omic data FAIR?

# FAIR framework to enhance 'omic data reusability
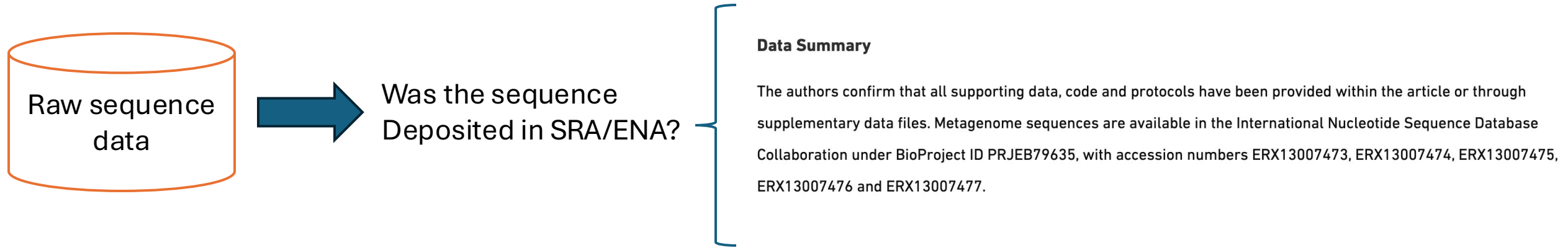
# FAIR framework to enhance 'omic data reusability



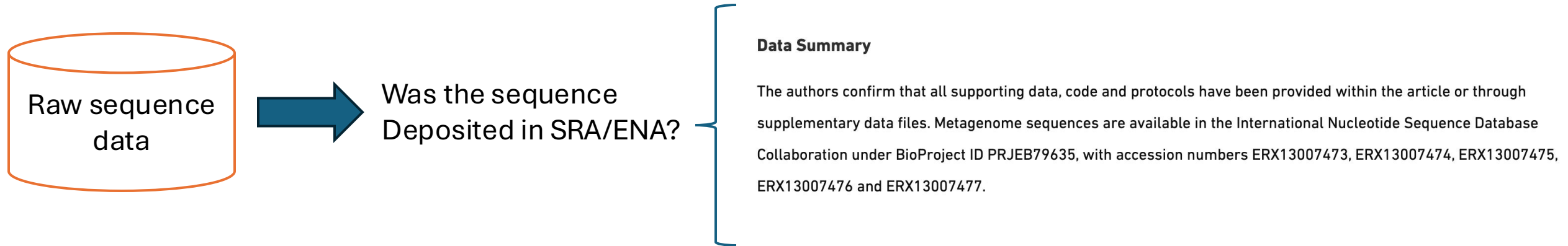Submission to Sequence repositories or generalist archives (Zenodo, Figshare...)

Raw sequence data → Was the sequence Deposited in SRA/ENA?

# The promise vs reality of 'Omic Data

Raw sequence data

Was the sequence Deposited in SRA/ENA?

**Data Summary**

The authors confirm that all supporting data, code and protocols have been provided within the article or through supplementary data files. Metagenome sequences are available in the International Nucleotide Sequence Database Collaboration under BioProject ID PRJEB79635, with accession numbers ERX13007473, ERX13007474, ERX13007475, ERX13007476 and ERX13007477.

# The promise vs reality of 'Omic Data

Raw sequence data

Was the sequence Deposited in SRA/ENA?

**Data Summary**

The authors confirm that all supporting data, code and protocols have been provided within the article or through supplementary data files. Metagenome sequences are available in the International Nucleotide Sequence Database Collaboration under BioProject ID PRJEB79635, with accession numbers ERX13007473, ERX13007474, ERX13007475, ERX13007476 and ERX13007477.

## Every fifth published metagenome is not available to science

Ester M. Eckert [co], Andrea Di Cesare [co], Diego Fontaneto [co], Thomas U. Berendonk, Helmut Bürgmann, Eddie Cytryn, Despo Fatta-Kassinos, Andrea Franzetti, D. G. Joakim Larsson, Célia M. Manaia, Amy Pruden, Andrew C. Singer, Nikolina Udikovic-Kolic, Gianluca Corno [✉]
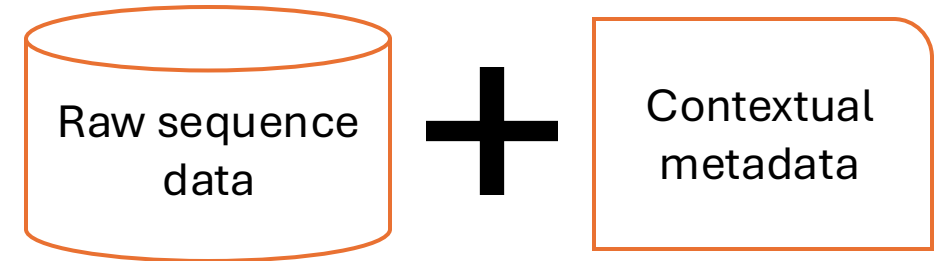
13% of papers don't mention availability
8%   have broken/nonexistent links

# FAIR framework to enhance 'omic data reusability



Raw sequence data

Where do the sample/isolate come from?
How was it collected and when?

# FAIR framework to enhance 'omic data reusability



Findable | Persistent Identifiers (PIDs) | Indexed data repositories

Accessible | Standard communications protocol | Authentication, where necessary | Open, free protocol

Interoperable | Vocabularies | Vocabularies are FAIR | Linked metadata

Reusable | Provenance | Community standards

Raw sequence data + Contextual metadata

Can it be found and accessed?

Additional efforts necessary to ensure interoperability and reusability of research outputs

# A metadata desert

## COVID-19 pandemic reveals the peril of ignoring metadata standards

Lynn M. Schriml ✉, Maria Chuvochina, Neil Davies, Emiley A. Eloe-Fadrosh, Robert D. Finn, Philip Hugenholtz, Christopher I. Hunter, Bonnie L. Hurwitz, Nikos C. Kyrpides, Folker Meyer, Ilene Karsch Mizrachi, Susanna-Assunta Sansone, Granger Sutton, Scott Tighe & Ramona Walls

less than 33% of 2.1 million COVID genomes had basic contextual metadata

**When** the next global outbreak crisis occurs, we need a predefined, widely adopted multidimensional approach to organize critical genomic data. Our strategy to broadly inform how to clearly describe genomic metadata and the tools to prepare genomic metadata datasets needs to be expanded now. Our community needs the organizational ability and coordination to respond to the imminent need well in advance. Opportunities for coordination of reported data types are critical for data interoperability as contact tracing efforts and outbreak resources, such as Nextstrain[19] and GISAID[20] are being developed.

In the words of Benjamin Franklin: "By failing to prepare, you are preparing to fail."

# When Standard Frameworks Don't Fit Complex Science



## Human gut MixS
- Limited contextual metadata data check
- Difficult to capture the complexity of study variables

# When Standard Frameworks Don't Fit Complex Science
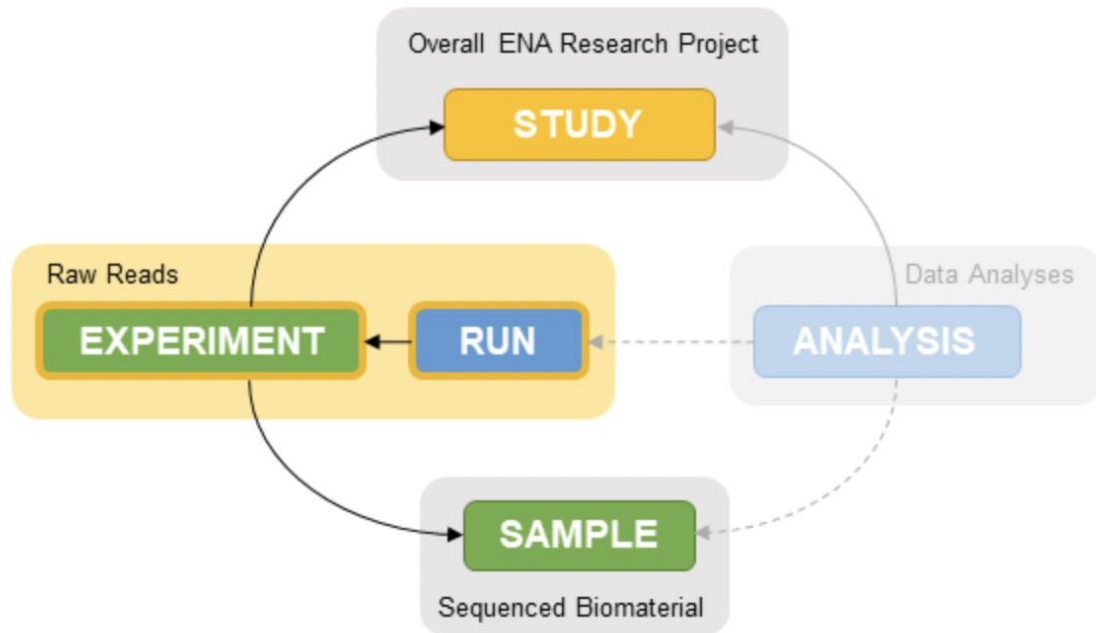


Data model design to accommodate most study design
→ Difficult for longitudinal or complex study design
→ How to link to non-sequence measurements (multi-omic)?
→ Non-sample sequences submission (blanks)?

# When Standard Frameworks Don't Fit Complex Science

## Biosample: SAMEA2579905 ↗

Stool sample from danish

| | |
|---|---|
| **Organism:** | human gut metagenome |
| **Scientific Name:** | human gut metagenome |
| **Sample Accession:** | SAMEA2579905 |
| **Location:** | 55.676097 N 12.568337 E |
| **Center Name:** | BGI |
| **Sample Alias:** | 10_12M |
| **Checklist:** | ERC000015 |
| **Sample Title:** | Stool sample from danish |
| **ENA-CHECKLIST:** | ERC000015 |
| **Environment (Material):** | faeces |
| **Collection Date:** | 2008/2010 |
| **Geographic Location (Longitude):** | 12.568337 |
| **Geographic Location (Latitude):** | 55.676097 |
| **Geographic Location (Country And/or Sea,region):** | Denmark |

## Biosample: SAMEA2579915 ↗

Stool sample from danish

| | |
|---|---|
| **Organism:** | human gut metagenome |
| **Scientific Name:** | human gut metagenome |
| **Sample Accession:** | SAMEA2579915 |
| **Location:** | 55.676097 N 12.568337 E |
| **Center Name:** | BGI |
| **Sample Alias:** | 10_B |
| **Checklist:** | ERC000015 |
| **Sample Title:** | Stool sample from danish |
| **ENA-CHECKLIST:** | ERC000015 |
| **Environment (Material):** | faeces |
| **Collection Date:** | 2008/2010 |
| **Geographic Location (Longitude):** | 12.568337 |
| **Geographic Location (Latitude):** | 55.676097 |
| **Geographic Location (Country And/or Sea,region):** | Denmark |

Additional information from supplemental materials necessary to link these two samples from the same infant

# FAIR framework to enhance 'omic data reusability



**Goal: not replacing standards, but enriching them to fit the institute/field needs**

→ how could we improve FAIR principles implementation at QIB?

→ What are the current needs of the institute in FAIR data?

Hurwitz lab


Salonen lab

THE UNIVERSITY OF ARIZONA

HELSINGIN YLIOPISTO
HELSINGFORS UNIVERSITET
UNIVERSITY OF HELSINKI

## Marine environments

- Cover >70% of Earth's surface
- Largest continuous ecosystem
- Physical forces create many niches

## Marine microbes

- Critical to food webs
- Drive elemental cycles
- Impact atmosphere & climate

## Ecological context required

- Need well-integrated data

# Great Datasets Exists



>30 years of sample collections

TBs of sequence data

Most oceans and sea sampled

# Great Datasets Exists But Can't Be Used Easily



**Sequence and physiochemical contextual data deposited in different repositories**

**Lack of common vocabularies**

**Different Units**

**Disparate data types**

# Semantic harmonization

| SampleID | NITRATE | Depth |
|----------|---------|-------|
| TARA_1 | 3.082 | 25.6 |
| TARA_2 | 0.193 | 125.2 |
| TARA_3 | 1.967 | 45.7 |

| SampleID | NO3 | Depth |
|----------|-----|-------|
| BATS_1 | 40.02 | 2 |
| BATS_2 | 19.2 | 50.2 |
| BATS_3 | 36.67 | 75.7 |

# Semantic & unit harmonization



| SampleID | NITRATE | Depth |
|----------|---------|-------|
| TARA_1 | 3.082 | 25.6 |
| TARA_2 | 0.193 | 125.2 |
| TARA_3 | 1.967 | 45.7 |

Contextual metadata
associated to the sequence
(tsv)

Semantic layer (json)

Blumberg *et al.* 2021

# Catering for a complex Data Model



Contextual water column measure from CTD & Niskin

Large scale oceanographic information from cruise tracks

Ponsero et al. 2021

# Frictionless data packages

# Frictionless data packages



Relational database reconnecting all contextual information to the sample

# Frictionless data packages

# Frictionless data packages



Frictionless datapackage

Semantic layer and unit description

**Name:** Temp
**Type:** numeric
**Rdftype:** ENVO:Temperature of water
**Unit:** UO:degree Celsius

FP

Github & Zenodo

**Frictionless Data**

**Packaging Data**

Package data with its metadata and schema for increased usability and clarity.

**Transforming Data**

Data often requires some transformations, like cleaning or conversions from one format to another.

**Pushing and Storing Data**

Frictionless has several plugins for accessing and storing data, for example in a SQL database.

# Frictionless data packages



Semantic layer and unit description

Name: Temp
Type: numeric
Rdftype: ENVO:Temperature of water
Unit: UO:degree Celsius

Frictionless datapackage

FP

Data processing, validation and unit conversion

Experiments and Runs data retrieval

SRA

Planet Microbe

Planet Microbe web-application

Ponsero *et al.* 2021

# Increase reusability of datasets for meta-analysis



**4D search**

**map search**

**semantic search**

Ponsero *et al.* 2021

# Increased reusability of datasets for meta-analysis

**In a nutshell:**
- FAIR metagenomics is about capturing the complexity of the data

- Adding a semantic layer to your contextual data will help with interoperability and reusability

- Taking the effort to create curated data packages allows to make valuable 'omics datasets interoperable and reusable

**In a nutshell:**
- FAIR metagenomics is about capturing the complexity of the data

- Adding a semantic layer to your contextual data will help with interoperability and reusability

- Taking the effort to create curated data packages allows to make valuable 'omics datasets interoperable and reusable

What if your data is sensitive or is not yet publicly available, why making your data FAIR matters?

## Issues and needs for data stewardship

- Curation and harmonization (and translations) of the questionnaires data

- Tracking metagenomes, metabolomes and isolates obtained from the samples

- Coordinating multi-location collaborative projects on the dataset

- Tracking issues, contaminations and labelling errors

- Ensuring long-term reusability of the dataset even after students/postdocs have left

Korpela et al. 2021

# Database to improve data management of complex datasets

# Database to improve data management of complex datasets

# Database to improve data management of complex datasets

# Database to improve data management of complex datasets

# Database to improve data management of complex datasets



Sequence/sample quality and known issues

Family

Sample collection and handling

Sequencing run and sequence file

Controlled vocabulary
-> Finnish original wording + translations

HELMi questionnaires

Questionnaires

Diaries

Registry

Publicly available along with MixS compliant metadata

→ Modular database that adds more and more information as the project develops

# Harmonization with other datasets

- Systematic search to identify relevant infant cohorts

- Controlled vocabulary for variables of interests

- Curation of relevant metadata across 14 cohorts

Bargheet et al. 2025

**In a nutshell:**
- FAIR standards needs to account for sensitive data protection

- Implementing extensive data stewardship allows for long terms project management resilience
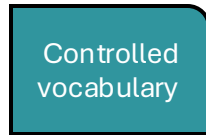
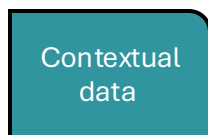- Help supporting collaborative works and meta-analysis

## Before publication



Sequence Data **+** Contextual data **+** Controlled vocabulary

- Internal data management that document the dataset
- Internal controlled vocabulary
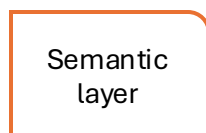
- Modular and comprehensive
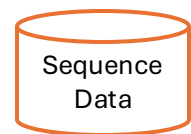- Limit data loss risks during the project
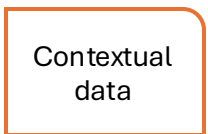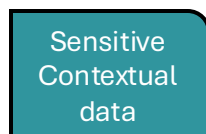
**Before publication**

Sequence Data + Contextual data + Controlled vocabulary
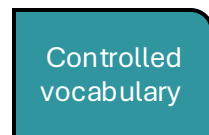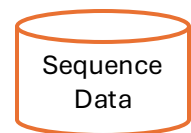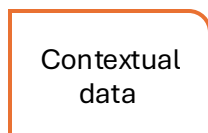
- Internal data management that document the dataset
- Internal controlled vocabulary
- Modular and comprehensive
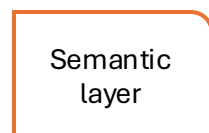- Limit data loss risks during the project

**At publication**

Sequence Data + Contextual data + Semantic layer

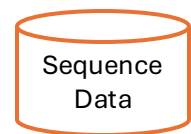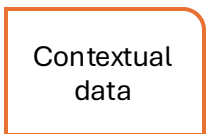Sequence Data + Contextual data + Sensitive Contextual data + Controlled vocabulary
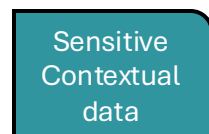
- Fixed data package that contained all relevant contextual data
- Sensitive data may be accessible upon request
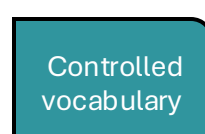- Long term storage and safety of the dataset needs to be addressed
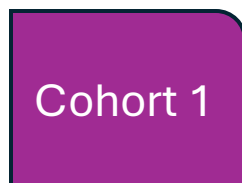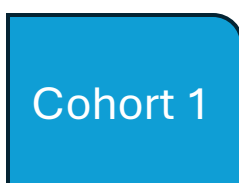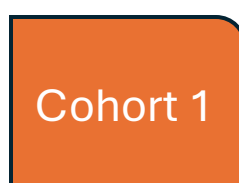
**After publication**

Cohort 1 + Cohort 1 + Cohort 1 — Harmonization

- Curation and harmonization efforts to make high-value datasets re-usable for meta-analysis or cross-study analysis

- **What are the current pain points in data management in your lab?**
    - Drawing the data management plan → a DS meeting in October will focus on this
    - Data management during the project?
    - Making data accessible? Including secondary outputs (protein catalogues,... etc)
    - Reusing old data with current datasets?

- **What use cases would you be interested in applying FAIR principles ?**
    - What high-value datasets would benefit in being curated in a data package?
    - Highly collaborative projects that would benefit in having a data stewardship in place?