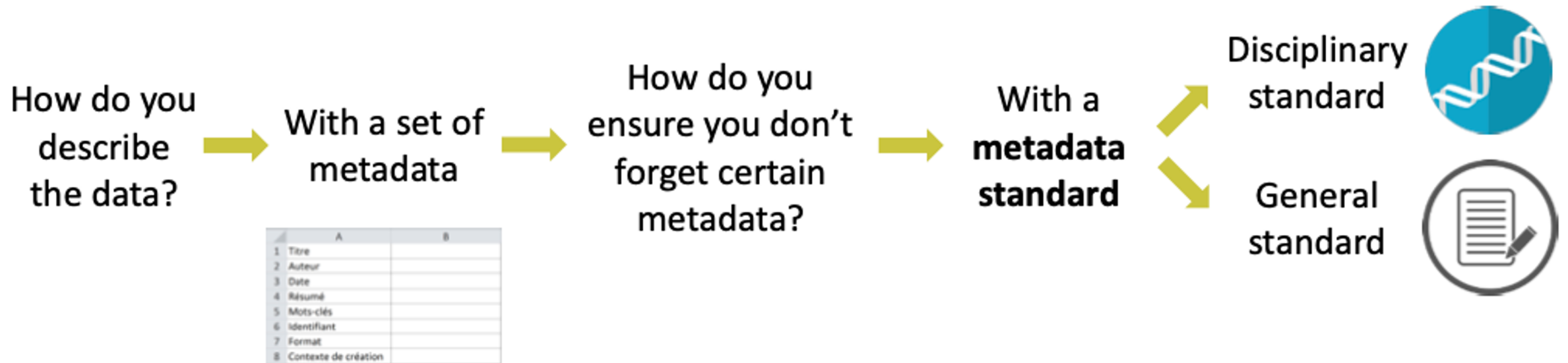# Life science standards and ENA submission

helene.chiapello@inrae.fr  &  thomas.denecker@france-bioinformatique.fr

Hélène Chiapello   : https://orcid.org/0000-0001-5102-0632
Thomas Denecker  : https://orcid.org/0000-0003-1421-7641

# Metadata & standards in life sciences

# Metadata standards help describing data



How do you describe the data? → With a set of metadata → How do you ensure you don't forget certain metadata? → With a **metadata standard** → Disciplinary standard / General standard

3

# Definition of a standard

In essence, a standard is an **agreed way of doing something**.

A standard provides the **requirements**, **specifications**, **guidelines** or **characteristics** that can be used for the **description**, **interoperability**, **citation**, **sharing**, **publication**, or **preservation** of all kinds of **digital objects** such as data, code, algorithms, workflows, software, or papers.

*source: https://fairsharing.org/educational/*

**Example of standard in biology :** Gene Ontology

# The standards concern both data and metadata

Why do I have to use a **data standard**?

- to analyse, compare and exchange data
- to publish datasets in international resources

And a **metadata standard**?

- To describe data richly and accurately, with the same vocabulary as the rest of your scientific community
- To make your metadata interoperable and to allow other systems to exploit them

The Gene Ontology is a **metadata** standard

Question: What do you know as standard in life sciences ?

*10 minutes to find an example (one for data and one for metadata) and write a note in*

*https://scrumblr.ethibox.fr/standard*

# Metadata exhibit questionable quality in biology

Submission in public resources is often a complex task

Submission procedures are heterogeneous

**Metadata are often incomplete, inconsistent, redundant or not enough informative**



**Quality of dictionary attributes in NCBI BioSample according to their type, in Gonçalves et al., 2019**

# Standard adoption and perenity

- There are thousand of databases, softwares and resources in biology with **unequal level of standard adoption**
- Is is not easy for Life scientists and bioinformaticians to identify and use the most appropriate standards



**1641** databases in NAR Database 2021

Rigden *et al*, 2021

# How do I find the standard I need?

# The FAIRsharing portal

Sansone, *et al.* FAIRsharing as a community approach to standards, repositories and policies.

Nat Biotech. 2019

https://doi.org/10.1038/s41587-019-0080-8



https://fairsharing.org

# The FAIRsharing portal

Citable *DOI* for all records

Accessible via *API* or *web interface*

*Curation*

# Standard maintenance is a key point

**Standard records that have maintainers**



40.7%

59.3%

■ Yes ■ No

**Standards that have a publication**



40.6%

59.4%

■ Yes ■ No

59.3 % of standards have no maintainer

59.4% of standard has no publication

https://fairsharing.org/summary-statistics/?collection=standards

# Types of data standards

**Conceptual model, schema, exchange formats,etc…**
e.g. FASTA

**Minimum information reporting requirements, checklists…**
e.g. MIAME guidelines



| Formats | Terminologies | Guidelines | Identifiers |

**Controlled vocabularies, taxonomies, ontologies…**
e.g. Gene Ontology

**Formal systems for resources and digital objects that allow their identification**
e.g. DOI

# The landscape of standards in life sciences

FASTQ

FASTA

SBML

GFF

MIAME

Newick

EC number

BAM

VCF

MINSEQE

n=432  n=539  n=180  n=18

Formats → Terminologies → Guidelines → Identifiers

ID

**COMMUNITY STANDARDS**
for metadata and identifiers

Source: https://fairsharing.org/search/?q=Life+science

genomic
gsc STANDARDS *consortium*

isatools

GENEONTOLOGY
Unifying Biology

Crop Ontology
for agricultural data

DISEASE
ONTOLOGY

# Collections in the FAIRsharing portal

A *collection* include standards and/or databases *grouped by domain, species or organization*

*Graph view* to visualize relationship links between resources

https://fairsharing.org/collections/

# Collections in Life Sciences

63 collections related to Life Science standards in FAIRsharing

Example 1: the *FAIRdom community Standards collection* (System biology)

https://fairsharing.org/collection/FAIRDOM

# Some collections are recent

Example 2: The *Covid-19* collection



https://fairsharing.org/collection/COVID19Resources



https://fairsharing.org/graph/#/collection/bsg-c000070

# What about the minimum required metadata in biology?

Example 3: the *Minimum Information for Biological and Biomedical Investigations* collection
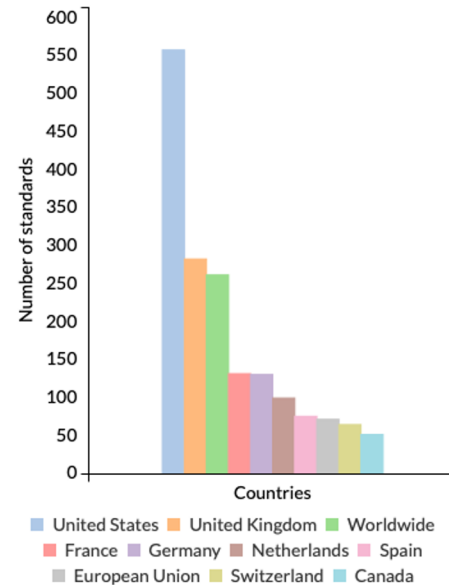
https://fairsharing.org/collection/MIBBI

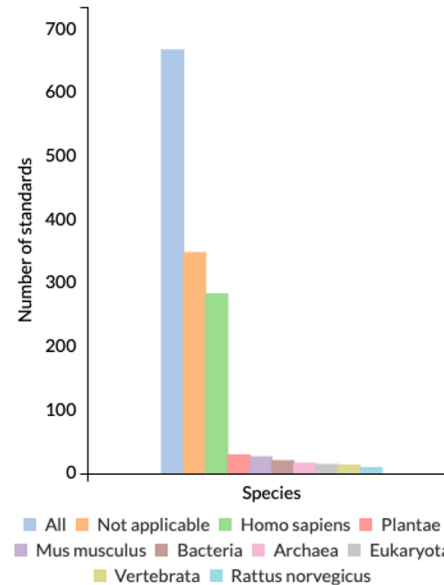# Summary statistics about standards



Top 10 disciplines covered by standards

Top 10 standard producing countries

Top 10 species covered by standards

**Life Science is one of the best covered discipline**

US and UK are the main standards producers

Human species is the best covered species

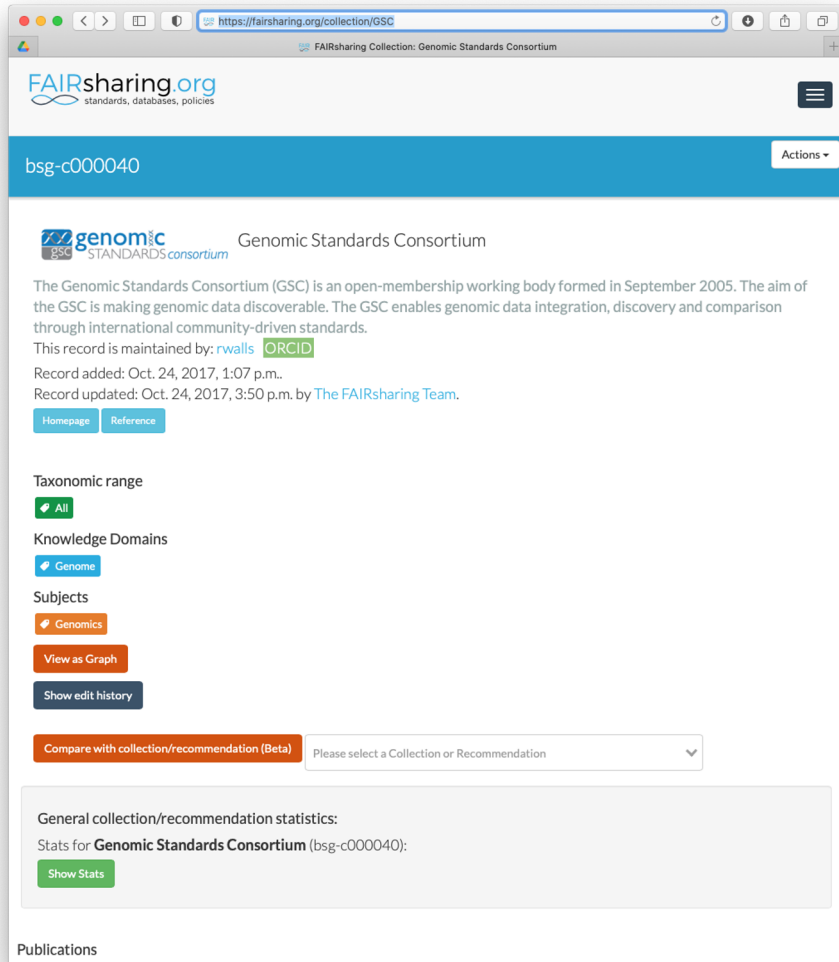https://fairsharing.org/summary-statistics/?collection=standards

# Practice

Find the *Genomics Standard Consortium (GSC) used by both ENA and SRA databases* in the FAIRsharing resource

Use both the record summary and the Graph visualization to interpret and answer the questions in zoom:

1. How many records (*i.e.* standards) are associated to the *GSC* ?

2. What type of standard is *Minimum Information about any (x) Sequence (MiXS)* ?

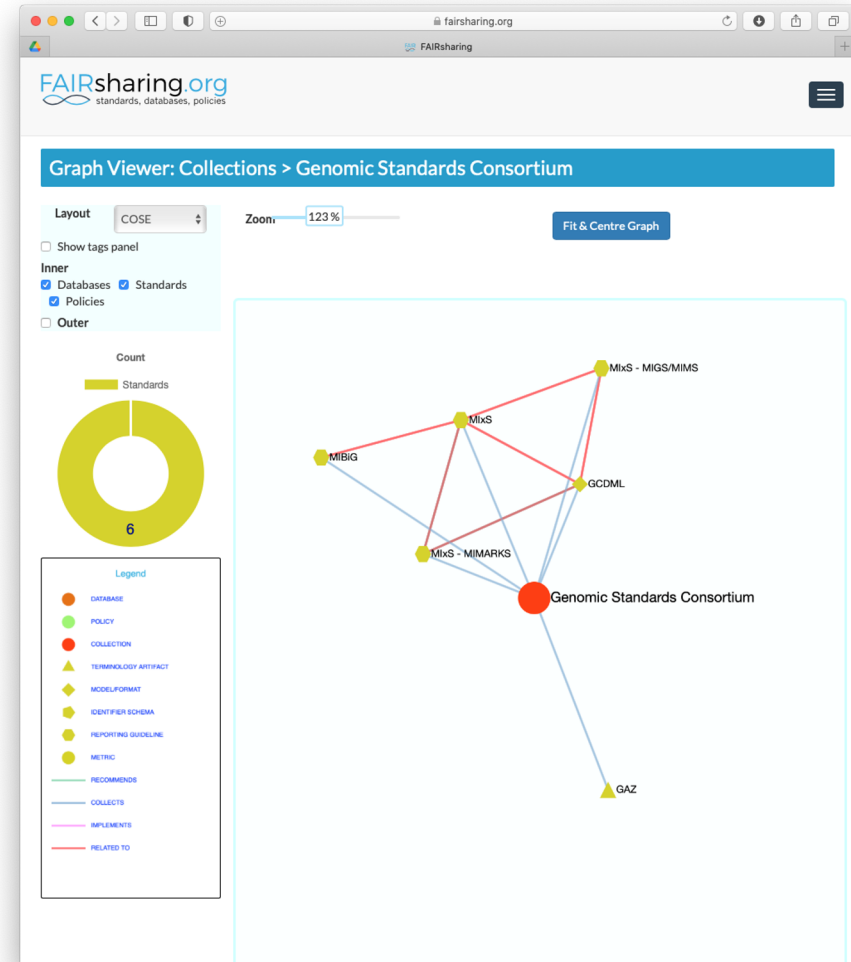3. What is the record status of the *GAZ* record ?

# The Genomics Standard Consortium (GSC)



https://fairsharing.org/collection/GSC



https://fairsharing.org/graph/#/collection/bsg-c000040

# The Genomics Standard Consortium (GSC)

Genomic data integration, and comparison through international community-driven standards

**Producer of the *Minimum Information Standards* (Checklists) used by ENA (EBI) & SRA (NCBI)**

**Ex: MIxS** : Minimum Information about any (x) Sequence



Yilmaz et *al*, 2011

# The ISA model

**A standard for Life ScienceData**

A **model** to capture **experimental metadata** through **3 core entities**:

- **Investigation**: the project context
- **Study**: an experimentation in one location
- **Assay:** a specific measurement that targets a trait with a method and a scale

ISA software suite: supporting standards-compliant experimental annotation and enabling curation at the community level. Rocca-Serra P et al. **Bioinformatics 2010**. https://doi.org/10.1093/bioinformatics/btq415



Sources: https://isa-tools.org and : https://isa-specs.readthedocs.io/en/latest/isamodel.html

# European Nucleotide Archive (ENA) submission

# Why do I need to submit my data and metadata to ENA ?

- Open Science and reproducibility of experiments
- 3$^{rd}$ party access
- Archival
- Publication
- Analyses, example: MGinfy

# The ENA metadata model

ISA compliant !

All **samples** submitted to ENA must conform to a **Checklist**

# THE ENA Checklists

- A **checklist** defines the **minimum and optional metadata** expected to describe biological samples

- ENA are based on the **Genomic Standards Consortium (GSC)** recommandations
- The **most suitable checklist** depends on the type of the sample: https://www.ebi.ac.uk/ena/browser/checklists
- All ENA checklist are defined by an **access number** like ERCxxx (Ena R Checklist xxx)
  - example: GSC MIxS plant associated https://www.ebi.ac.uk/ena/browser/view/ERC000020

# Data brokering at IFB

# Why developing data brokering at IFB?

**Observations:**

- Submissions are often complex and difficult to perform by individual teams
- Metadata are often poorly understood resulting in incomplete, redundant and inconsistent submissions
- ENA asks that IFB becomes the French national broker

**Main idea:** offer a national service of **data brokering at IFB** to simplify and rationalize data exchange between international resources and the french Elixir node IFB.

**Brokering include 3 types of activities:** tools development, training and support to users

# Data Brokering service developed by IFB

*IFB services to manage and centralize data and metadata of a project*

*IFB services to submit data and metadata of a project to international resources*

# The omicsBroker tool

- **omicsBroker** is a tool to easily annotate and submit **omics data** to **international repositories**
- For now, only available as a **PROTOTYPE**
  - based on **R Shiny** technology
  - allowing to test submission of genomic and transcriptomic samples and reads to **ENA test instance**
- The final tool will be developed using Django technology and will **manage data and metadata from different sources** to make submission to international resources easier

https://github.com/IFB-ElixirFr/omicsBroker

# Practice

Use omicsBroker prototype (134.158.247.213:443) to test submission of samples to ENA

Use information of the corresponding DMP to associate relevant metadata to data https://dmp.opidor.fr

3 groups

- bacterial genome (IFB_Training_salivarius)
- plant transcriptome (IFB Training : Sars-CoV-2)
- SARS-Cov2 genome (IFB_Training_plant)

https://ifb-elixirfr.github.io/IFB-FAIR-data-training/sequences/module3_sequence3_tp.html

# To conclude: sources & useful links

| Description | Name | URL |
| --- | --- | --- |
| A curated, informative and educational resource on data and metadata standards, inter-related to databases and data policies. | FAIRsharing portal | https://fairsharing.org |
| Investigation, Study, Assay (ISA) ressource: A standard model an a set of tools to capture experimental data in life sciences | ISAtools | https://isa-tools.org |
| Genomics Standard Consortium (GSC): An international consortium developing standards and checklists in genomics | GSC | https://gensc.org |
| European National Archive Checklists | ENA Checklists | https://www.ebi.ac.uk/ena/browser/checklists |
| European National Archive submission documentation | ENA submission guide | https://ena-docs.readthedocs.io/en/latest/submit/general-guide.html |
| A prototype to test submission of samples and DNAseq or RNAseq reads to ENA | omicsBroker | https://github.com/IFB-ElixirFr/omicsBroker |

# Thanks



Paulette Lieby



Jean-François Dufayard



Frédéric de Lamotte