

Standards in life sciences

Thomas Denecker
Institut Français de Bioinformatique



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



	A	B
1	Titre	
2	Auteur	
3	Date	
4	Résumé	
5	Mot-clés	
6	Identifiant	
7	Format	
8	Contexte de création	

Inspiré de <https://www.pasteur.fr/fr/file/20615/download>



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



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



Do you know any standard in life sciences ?

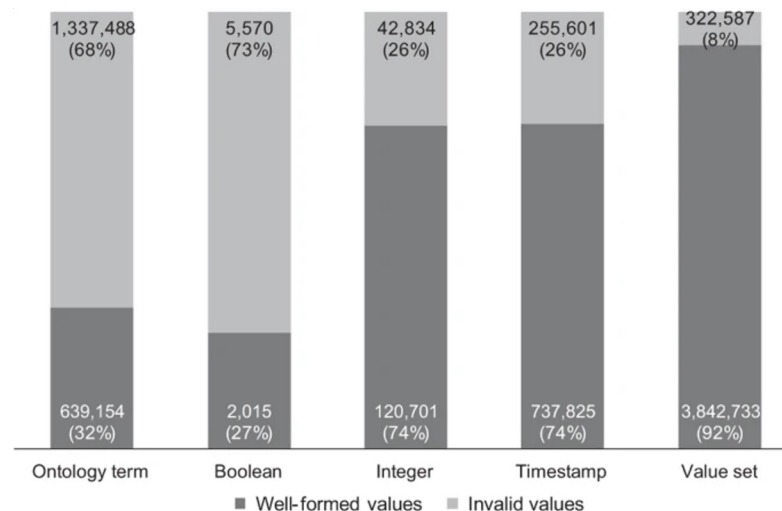




Submission in public resources is often a complex task

Submission procedures are heterogeneous

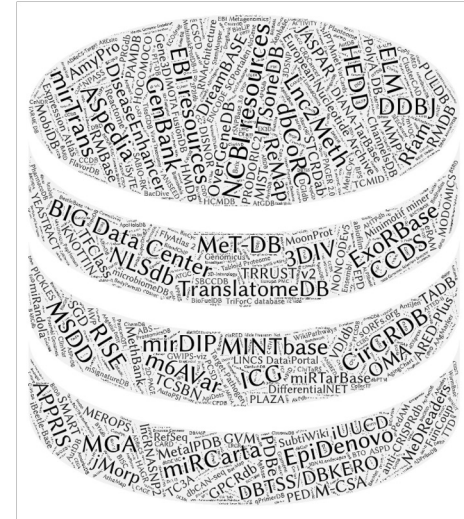
Metadata are often incomplete, inconsistent, redundant or not enough informative



Quality of dictionary attributes in NCBI BioSample according to their type, in [Gonçalves et al., 2019](#)

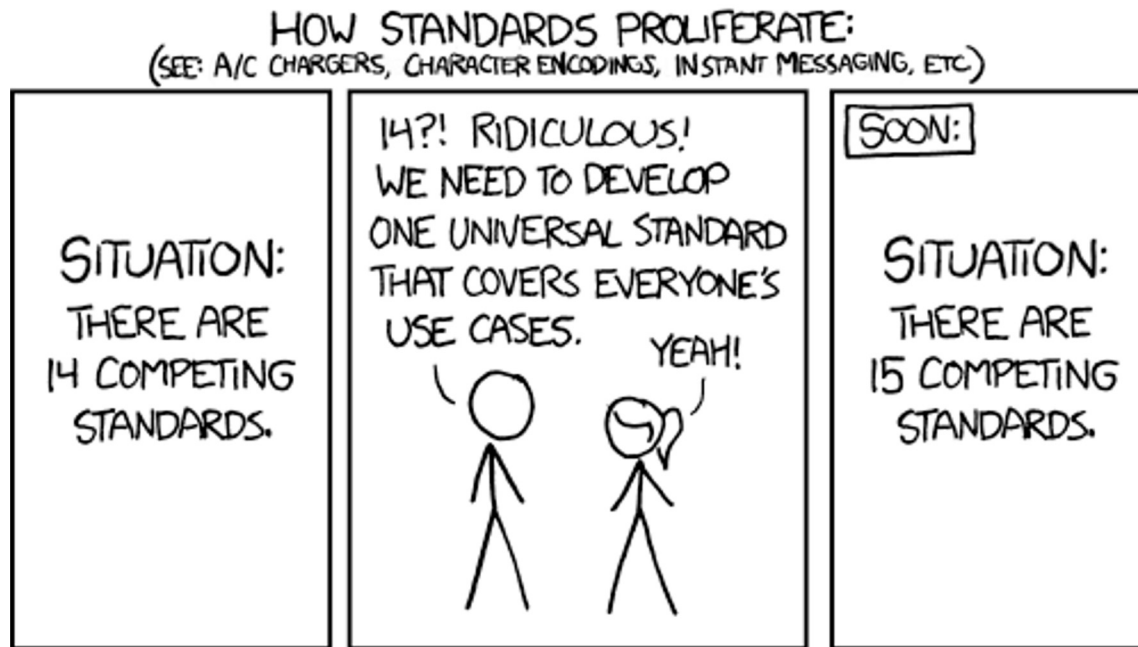
There are thousand of databases,
softwares and resources in biology with
unequal level of standard adoption

Is is not always easy for Life scientists
and bioinformaticians to identify and
use the most appropriate standards



1641 databases in NAR Database 2021

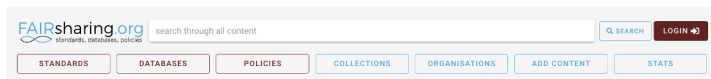
[Rigden et al, 2021](#)



Source: <https://xkcd.com/927/>

How do I find the standard I need ?

FAIR sharing & re3data



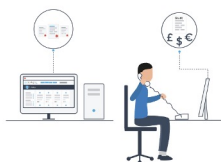
FAIRsharing.org
standards, databases, policies

search through all content

STANDARDS DATABASES POLICIES COLLECTIONS ORGANISATIONS ADD CONTENT STATS

A curated, informative and educational resource on data and metadata standards, inter-related to databases and data policies.

We guide consumers to discover, select and use these resources with confidence, and producers to make their resource more discoverable, more widely adopted and cited.

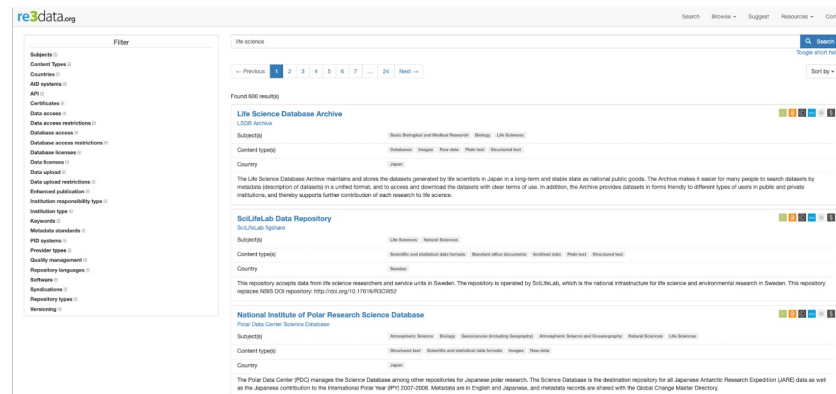


Funders & data policy makers

Recommend FAIRsharing to your awardees to inform the development of their data management plan, and select the appropriate resources to recommend in your data policy...

[read more](#)

1646 Standards	1980 Databases	161 Policies
Terminology Artifact 831	Repositories 1021	Journal 88
Model/Format 829	Knowledgebases 820	Funder 30



re3data.org

Search Browse Suggest Resource Contact

Filter

Subjects: 1
Content Types: 1
Countries: 1
AD systems: 1
API: 1
Certificates: 1
Data access: 1
Data access restrictions: 1
Database access: 1
Database access restrictions: 1
Data formats: 1
Data updates: 1
Data update restrictions: 1
Enhanced publications: 1
Institution responsibility type: 1
Institution type: 1
Keywords: 1
Metadata standards: 1
PID systems: 1
Provider type: 1
Quality management: 1
Repository languages: 1
Software: 1
Synchronisation: 1
Repository types: 1
Versioning: 1

life science

Found 600 result(s)

Life Science Database Archive
LSDA Archive
Subject(s): Basic Biological and Medical Research, Life Sciences
Content type(s): Database, Image, Plain text, Download text
Country: Japan
The Life Science Database Archive maintains and stores the datasets generated by the scientists in Japan in a long-term and stable state as national public goods. The Archive makes it easier for many people to search datasets by metadata (description of dataset) in a unified format, and to access and download the datasets with clear terms of use. In addition, the Archive provides datasets in forms friendly to different types of users in public and private institutions, and thereby supports further contribution of each research to life science.

SciLifeLab Data Repository
SciLifeLab Signet
Subject(s): Life Sciences, Natural Sciences
Content type(s): Scientific and statistical data formats, Document office document, Archived data, Plain text, Download text
Country: Sweden
This repository accepts data from life science researchers and service units in Sweden. The repository is operated by SciLifeLab, which is the national infrastructure for life science and environmental research in Sweden. This repository replaces NBS DCA repository <http://nbs.dca.se> <http://nbs.dca.se>

National Institute of Polar Research Science Database
Polar Data Center Science Database
Subject(s): Atmospheric Science, Biology, Oceanography (including Biological), Atmospheric Science and Oceanography, Natural Sciences, Life Sciences
Content type(s): Download text, Database and statistical data formats, Image, Plain text
Country: Japan
The Polar Data Center (PDC) manages the Science Database among other repositories for Japanese polar research. The Science Database is the destination repository for all Japanese Antarctic Research Expedition (JARE) data as well as the Japanese contribution to the International Polar Year (IPY) 2007-2008. Metadata are in English and Japanese, and metadata records are shared with the Global Change Master Directory.

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>

re3data.org - Registry of Research Data Repositories. <https://doi.org/10.17616/R3D> last accessed: 2023-01-25

European Nucleotide Archive (ENA)

10.25504/FAIRsharing.dj8n18

Type Repository

Registry Database

Description The European Nucleotide Archive (ENA) is a globally comprehensive data resource for nucleotide sequence, spanning raw data, alignments and assemblies, functional and taxonomic annotation and rich contextual data relating to sequenced samples and experimental design. Serving both as the database of record for the output of the world's sequencing activity and as a platform for the management, sharing and publication of sequence data, the ENA provides a portfolio of services for submission, data management, search and retrieval across web and programmatic interfaces. The ENA is part of the International Nucleotide Sequence Database Collaboration (INSDC), which comprises the DNA databank of Japan (DDBJ), the European Molecular Biology Laboratory (EMBL) and GenBank at the NCBI. These three organizations exchange data on a daily basis.

Homepage <http://www.ebi.ac.uk/ena>

Year of Creation 1980

Maintainers [cochrane](#)

Countries developing this resource Japan, United States, European Union

Subjects Functional Genomics, Metagenomics, Genomics, Bioinformatics, Data Management, Bioethics, Transcriptomics

Domains RNA Sequence Data, Annotation, Sequence Annotation, Genetic Assembly, Culture, Genomewide Anal., Microarray Anal., Nucleotide, Sequencing, Amino Acid Sequence, Data Storage

Taxonomic Range All

User Defined Tags Data Generation

[VIEW RELATION GRAPH](#)

How to cite this record

FAIRsharing.org/ENA, European Nucleotide Archive, DOI: 10.25504/FAIRsharing.dj8n18, Last Edited: Wednesday, May 4th 2022, 14:43, Last Editor: delphinedeaga, Last Accessed: Wednesday, January 25th 2023, 16:04

Publication for citation

European Nucleotide Archive in 2016. Toribio AL, Aizaki H, et al., (2016) [10.1093/nar/gkv1106](https://doi.org/10.1093/nar/gkv1106)

The following community curators have contributed to this record: [lccpaas@ghel.org](#)

Record ID: 1630 | [Record created at](#): Tuesday, November 4th 2014, 16:23 | [Record updated at](#): Wednesday, July 20th 2022, 11:30

COLLECTIONS & RECOMMENDATIONS

Search through in collections

IN COLLECTIONS (7)

IN POLICIES (48)

RELATED CONTENT

Search through related standards

RELATED STANDARDS (13)

RELATED DATABASES (32)

Data Citation Implementation

[Data Citation Implementation](#) [implements](#) [European Nucleotide Archive](#)

EDSCLife

[EDSCLife](#) [implements](#) [European Nucleotide Archive](#)

Short Read Archive eXtensible Markup Language

[European Nucleotide Archive](#) [implements](#) [Short Read Archive eXtensible Markup Language](#)

NCBI Taxonomy

[European Nucleotide Archive](#) [implements](#) [NCBI Taxonomy](#)

<https://fairsharing.org/FAIRsharing.e1byny>

Repository details

European Nucleotide Archive

[General](#) [Institutions](#) [Terms](#) [Standards](#)

Name of repository	European Nucleotide Archive
Additional name(s)	ENA
Repository URL	https://www.ebi.ac.uk/ena/browser/home
Subject(s)	Biology Medicine General Genetics Bioinformatics and Theoretical Biology Microbiology, Virology and Immunology Life Sciences Basic Biological and Medical Research
Description	<p>The European Nucleotide Archive (ENA) captures and presents information relating to experimental workflows that are based around nucleotide sequencing. A typical workflow includes the isolation and preparation of material for sequencing, a run of a sequencing machine in which sequencing data are produced and a subsequent bioinformatic analysis pipeline. ENA records this information in a data model that covers input information (sample, experimental setup, machine configuration), output machine data (sequence traces, reads and quality scores) and interpreted information (assembly, mapping, functional annotation). Data arrive at ENA from a variety of sources. These include submissions of raw data, assembled sequences and annotation from small-scale sequencing efforts, data provision from the major European sequencing centres and routine and comprehensive exchange with our partners in the International Nucleotide Sequence Database Collaboration (INSDC). Provision of nucleotide sequence data to ENA or its INSDC partners has become a central and mandatory step in the dissemination of research findings to the scientific community. ENA works with publishers of scientific literature and funding bodies to ensure compliance with these principles and to provide optimal submission systems and data access tools that work seamlessly with the published literature.</p>
Contact	https://www.ebi.ac.uk/ena/browser/support
Content type(s)	Scientific and statistical data formats Network-based data Structured graphics Plain text Software applications Raw data Structured text
Keyword(s)	GRAM DNA, RNA spigenomics genome human genetics metagenomics protein coding samples taxonomy data transcript/seq
Persistent identifier(s) of the repository	RRID:rrf-0000-32981 RRID:SCR_006515 OMICS_D1029 MIR-00000372 FAIRsharing_doi:10.25504/FAIRsharing.dj8n18
Repository type(s)	disciplinary
Mission statement for designated community	https://www.ebi.ac.uk/ena/browser/about
Research data repository language(s)	English
Data and/or service provider	data provider service provider

[Back to search](#) [Submit a change request](#) [Get a badge](#)

Cite this re3data.org record:

re3data.org: European Nucleotide Archive; editing status 2021-12-13; re3data.org - Registry of Research Data Repositories. <https://doi.org/10.17616/R3dHW3J> last accessed: 2023-01-25

<https://www.re3data.org/repository/r3d100010527>

PRIDE
Proteomics IDentifications database (PRIDE)

10.25561/FAIRsharing.e1byny

GENERAL INFORMATION

Type: Repository

Registry: Database

Description: The PRIDE Proteomics IDentifications (PRIDE) Archive database is a centralized, standards compliant, public data repository for mass spectrometry proteomics data, including protein and peptide identifications and the corresponding expression values, post-translational modifications and supporting mass spectra evidence (both as raw data and peak list files). PRIDE is a core member in the ProteomeXchange (PX) consortium, which provides a standardised way for submitting mass spectrometry based proteomics data to public-domain repositories.

Homepage: <http://www.ebi.ac.uk/pride/>

Year of Creation: 2004

Maintainers: pridebi

Countries developing this resource: Austria, Belgium, Croatia, Czech Republic, Denmark, Finland, France, Germany, Greece, Hungary, Iceland, Ireland, Israel, Italy, Lithuania, Luxembourg, Malta, Montenegro, Netherlands, Norway, Portugal, Slovakia, Spain, Sweden, Switzerland, United Kingdom

Subjects: Proteomics

Domains: Mass Spectrometry, Peptide Identification, Protein Identification, Peptide, Post-translational Protein Modification, Mass Spectrometry Assay, Curated Information, Protein

Taxonomic Range: All

User Defined Tags: None

[VIEW RELATION GRAPH](#)

How to cite this record

Faloutsos et al. PRIDE: Proteomics IDentifications database, DOI: 10.25561/FAIRsharing.e1byny, Last Edited: Monday, October 3rd 2022, 9:07, Last Editor: aliyositor, Last Accessed: Wednesday, January 25th 2023, 15:58

Publication for citation

The PRIDE database and related tools and resources in 2019: Improving support for identification data. Perez-Ruiz M, Coudane A, Bai J, Berni-Luque M, Vespignani S, Kundi D, Hegerl A, Gira J, Mayer U, Mescher M, Perez E, Jochims J, Pfeiffer J, Sachdev R, Vlasov S, Tsvetkov S, Cox J, Audin E, Walter M, Janczuk A, Tennent T, Buzza A, Visciano JA (2019) 10.1093/nar/gky1108

RECORD INFO

Record ID: 1863 | [in Record created at](#) Tuesday, November 4th 2014, 16:23 | [in Record updated at](#) Monday, October 3rd 2022, 9:07

COLLECTIONS & RECOMMENDATIONS

Search through in collections

IN COLLECTIONS (4)

IN POLICIES (34)

IN COLLECTIONS (4)

IN POLICIES (34)

IN COLLECTIONS (4)

IN POLICIES (34)

RELATED CONTENT

Search through related standards

RELATED STANDARDS (15)

RELATED DATABASES (14)

RELATED STANDARDS (15)

RELATED DATABASES (14)

<https://fairsharing.org/FAIRsharing.e1byny>

Repository details

Proteomics IDentifications Database

General	Institutions	Terms	Standards
Name of repository	PRoteomics IDentifications Database		
Additional name(s)	PRIDE		
Repository URL	https://wwwdev.ebi.ac.uk/pride/		
Subject(s)	Biology Biochemistry Bioinformatics and Theoretical Biology Analytical Chemistry, Method Development (Chemistry) Life Sciences Basic Biological and Medical Research Analytical Chemistry, Method Development (Chemistry) Chemistry Natural Sciences		
Description	The PRIDE Proteomics IDentifications database is a centralized, standards compliant, public data repository for proteomics data, including protein and peptide identifications, post-translational modifications and supporting spectral evidence. PRIDE encourages and welcomes direct user submissions of mass spectrometry data to be published in peer-reviewed publications.		
Content type(s)	Structured text Standard office documents Images Structured graphics Plain text Scientific and statistical data formats Software applications Archived data		
Keyword(s)	amino acids mass spectra mass spectrometry peptide identifications protein identifications proteomics sequencing		
Persistent identifier(s) of the repository	RRID:rrf-0000-03336 RRID:SCR_003411 OMICS_02456 MIR:00100094 FAIRsharing_doi:10.25504/FAIRsharing.e1byny		
Repository size	13,226 datasets		
Repository type(s)	disciplinary		
Mission statement for designated community	https://wwwdev.ebi.ac.uk/pride/markdownpage/citationpage		
Research data repository language(s)	English		
Data and/or service provider	data provider		
Back to search Submit a change request Get a badge			



Cite this re3data.org record:

re3data.org: PRoteomics IDentifications Database; editing status 2021-11-09; re3data.org - Registry of Research Data Repositories. <http://doi.org/10.17716/R3JG6V> last accessed: 2023-01-25

<https://www.re3data.org/repository/r3d100010137>

BiolImage Archive (BiolImage Archive)

10.25504/FAIRsharing.x38D2k

Type: Repository

Registry: Database

Description: The BiolImage Archive stores and distributes biological images that are useful to life-science researchers. Its development will provide data archiving services to the broader biomedicine database community. This includes added-value biomedicine data resources such as EMPAR, Cell-IR and Tissue-IR.

Homepage: <https://www.ebi.ac.uk/biolimage-archive/>

Year of Creation: 2019

Maintainers: [ugr](#), [matthew](#)

Countries developing this resource: Austria, Belgium, Croatia, Czech Republic, Denmark, Finland, France, Germany, Greece, Hungary, Iceland, Ireland, Israel, Italy, Lithuania, Luxembourg, Malta, Montenegro, Netherlands, Norway, Portugal, Slovakia, Spain, Sweden, Switzerland, United Kingdom

Subjects: Molecular Biology, Clinical Studies, Life Sciences, Cell Biology, Biomedical Sciences, Preclinical Studies

Domains: Biomedicine, Cell, Microscopy, Light Microscopy, Electron Microscopy, Imaging, Super-resolution Microscopy, Image, Histology, White Mount Tissue, Tissue, High-content Screen

Taxonomic Range: All

User Defined Tags: Pre-Clinical Imaging

[VIEW RELATION GRAPH](#)

How to cite this record

FAIRsharing.org/BiolImage Archive, DOI: 10.25504/FAIRsharing.x38D2k, Last Edited: Thursday, November 24th 2022, 14:48, Last Editor: ramon granell, Last Accessed: Tuesday, January 31st 2023, 13:49

Publication for citation

A call for public archives for biological image data: Elberberg J, Svedberg J, Rutar M, Cook CE, Sarkans U, Pothardshorn A, Biazma A, Bimay E (2018) 10.1038/s41592-018-0195-8

Record ID: 2690 | Record created at: Thursday, October 17th 2019, 18:23 | Record updated at: Wednesday, December 7th 2022, 16:57

COLLECTIONS & RECOMMENDATIONS

Search through in collections

IN COLLECTIONS (5) | IN POLICIES (1)

EOSC Life

EOSC Life collects BiolImage Archive

RELATED CONTENT

Search through related standards

RELATED STANDARDS (1) | RELATED DATABASES (3)

Recommended Metadata for Biological Images

BiolImage Archive implements Recommended Metadata for Biological Images

<https://fairsharing.org/FAIRsharing.x38D2k>

Repository details

BiolImage Archive

General	Institutions	Terms	Standards
Name of repository	BiolImage Archive		
Additional name(s)	BIA		
Repository URL	https://www.ebi.ac.uk/biolimage-archive		
Subject(s)	Basic Biological and Medical Research Life Sciences Bioinformatics and Theoretical Biology Biology		
Description	The BiolImage Archive stores and distributes life sciences imaging datasets. It supports deposition of biological imaging data associated with publications for the whole research community, as well as reference imaging datasets. All data deposited to the BiolImage Archive is made openly accessible to the scientific community.		
Contact	biolimage-archive@ebi.ac.uk https://www.ebi.ac.uk/biolimage-archive/contact-us/		
Content type(s)	Images Raw data Scientific and statistical data formats Archived data		
Keyword(s)	biomedicine high-content screening microscopy		
Persistent identifier(s) of the repository	FAIRsharing_doi:10.25504/FAIRsharing.x38D2k		
Repository size	694 datasets		
Repository type(s)	disciplinary		
Mission statement for designated community	https://www.ebi.ac.uk/biolimage-archive/scope/		
Research data repository language(s)	English		
Data and/or service provider	data provider service provider		

[Back to search](#)

[Submit a change request](#)

[Get a badge](#)



Cite this re3data.org record:

re3data.org: BiolImage Archive; editing status 2022-08-13; re3data.org - Registry of Research Data Repositories.
<http://doi.org/10.17616/R31NJN99> last accessed: 2023-01-31

<https://www.re3data.org/repository/r3d100013949>

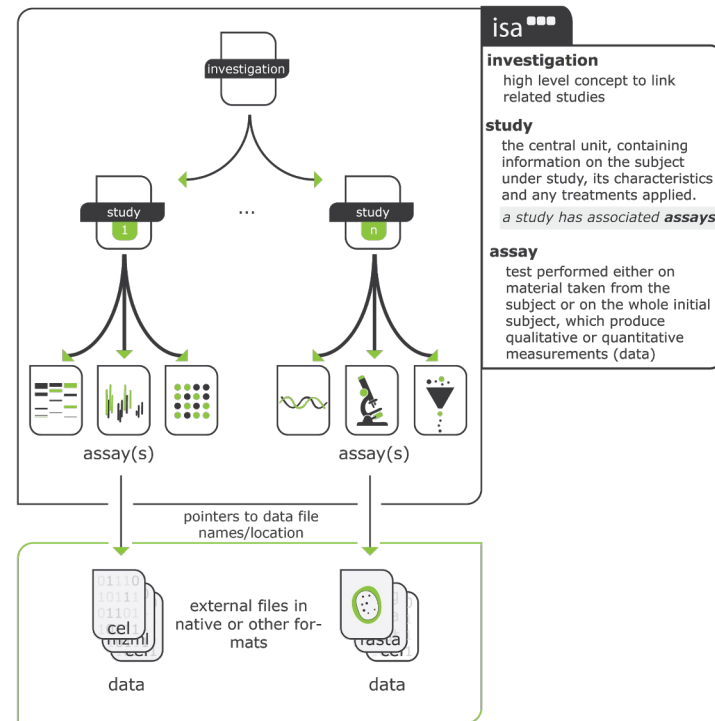
Focus

A standard for Life Science Data

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

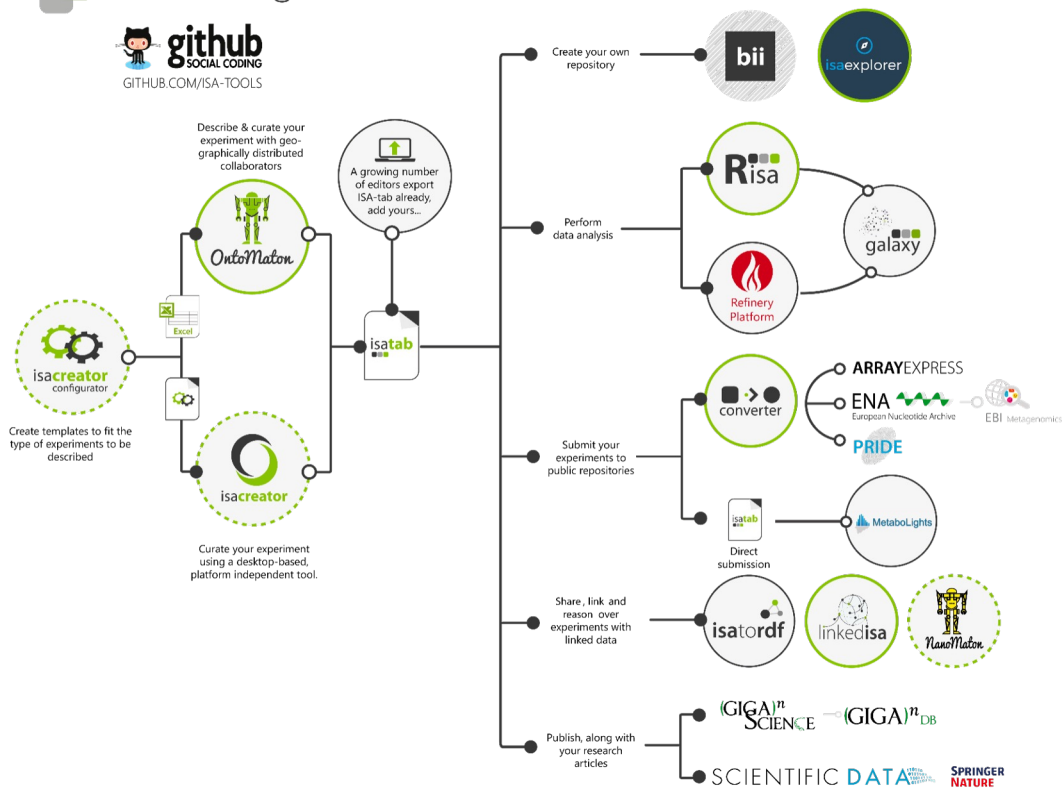
- **Investigation**: the project context
- **Study**: an experimentation in one location
- **Assay**: a specific measurement that targets a trait with a method and a scale

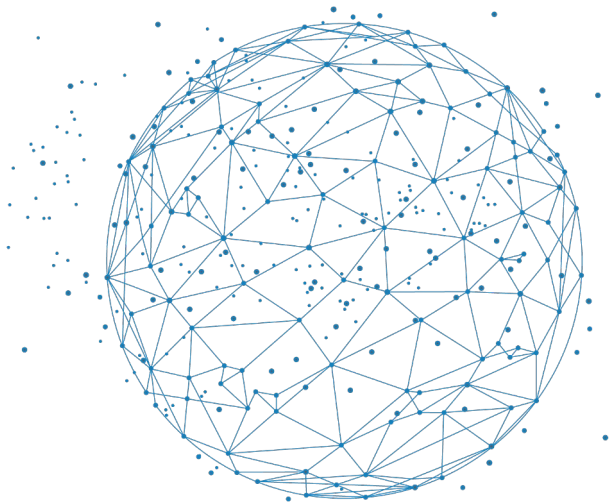
ISA software suite: supporting standards-compliant experimental annotation and enabling curation at the community level. Rocca-Serra P et al. **Bioinformatics** 2010. <https://doi.org/10.1093/bioinformatics/btq415>



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

isa-tools.org Curate, store, analyse, share and publish your bioscience experiment



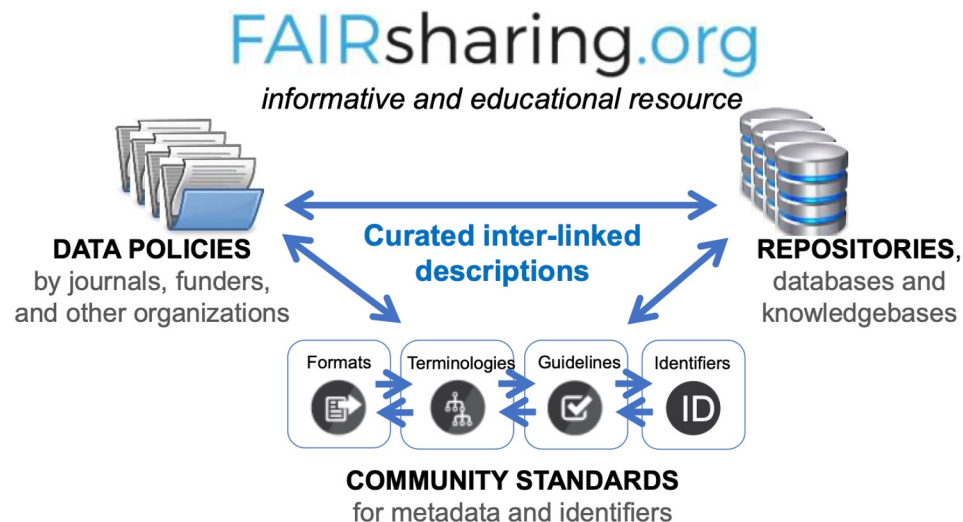


Supplementary slides

Citable *DOI* for all records

Accessible via *API* or *web interface*

Curation

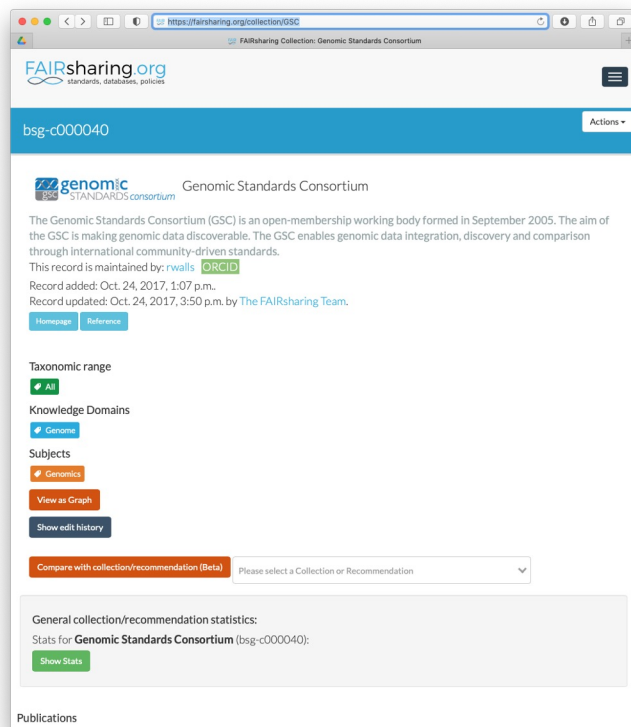


**RECORD
STATUS**

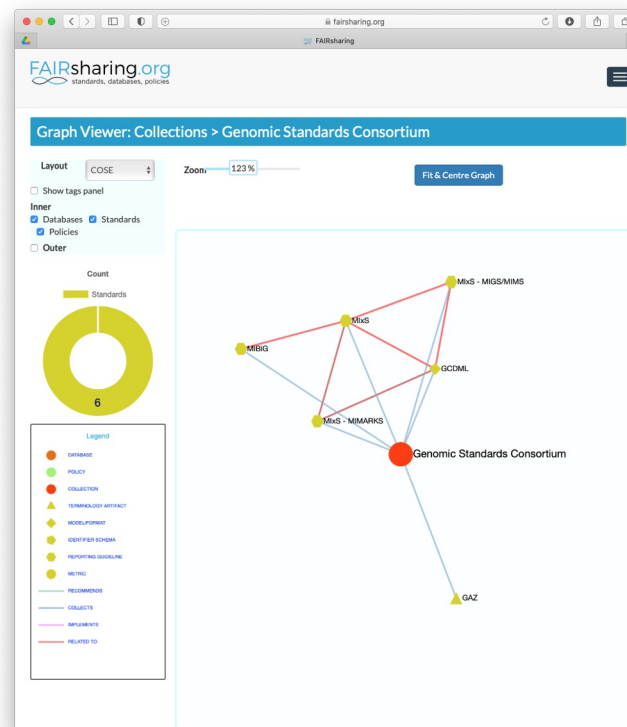
- R** Ready for use, implementation, or recommendation
- Dev** In development
- U** Status uncertain
- D** Deprecated as subsumed or superseded

All records are manually **curated**
in-house, **verified** and **claimed** by the
community behind each resource

The Genomic Standards Consortium (GSC)



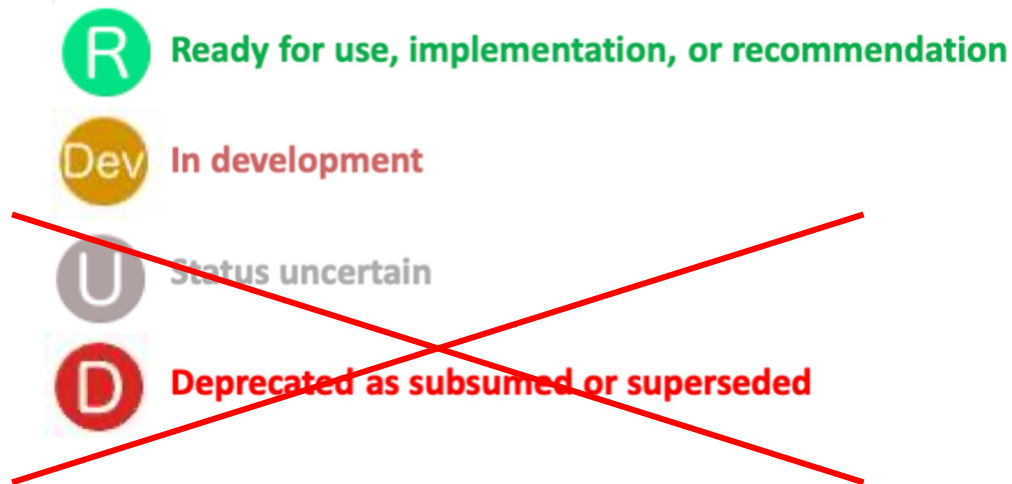
The screenshot shows the FAIRsharing.org interface for the collection 'bsg-c000040' (Genomic Standards Consortium). The page includes a description of the GSC, its mission, and its founding date (September 2005). It also features a 'Taxonomic range' section with a dropdown menu set to 'All', a 'Knowledge Domains' section with a dropdown menu set to 'Genome', and a 'Subjects' section with a dropdown menu set to 'Genomics'. A 'View as Graph' button is visible. At the bottom, there is a 'General collection/recommendation statistics' section with a 'Show Stats' button.



<https://fairsharing.org/collection/GSC>

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

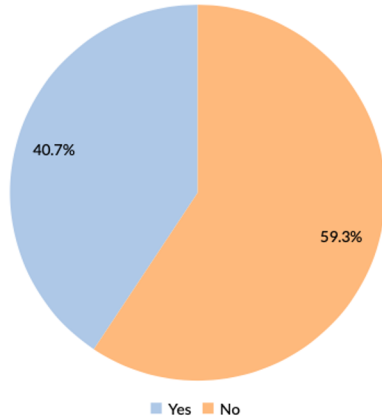
Please don't use “Uncertain” or “Deprecated” standards



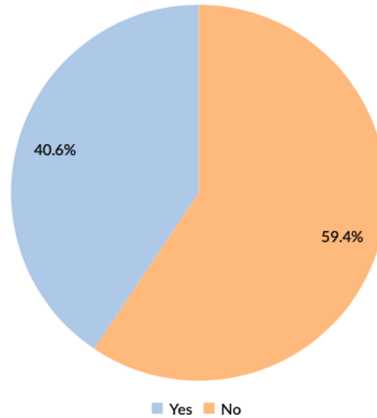
Standard maintenance is a key point



Standard records that have maintainers



Standards that have a publication

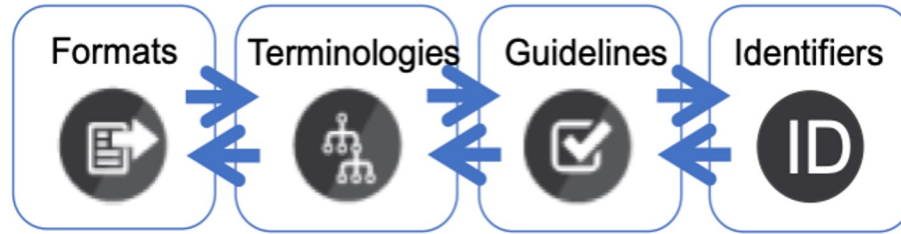


59.3 % of standards have no maintainer

59.4% of standard has no publication

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

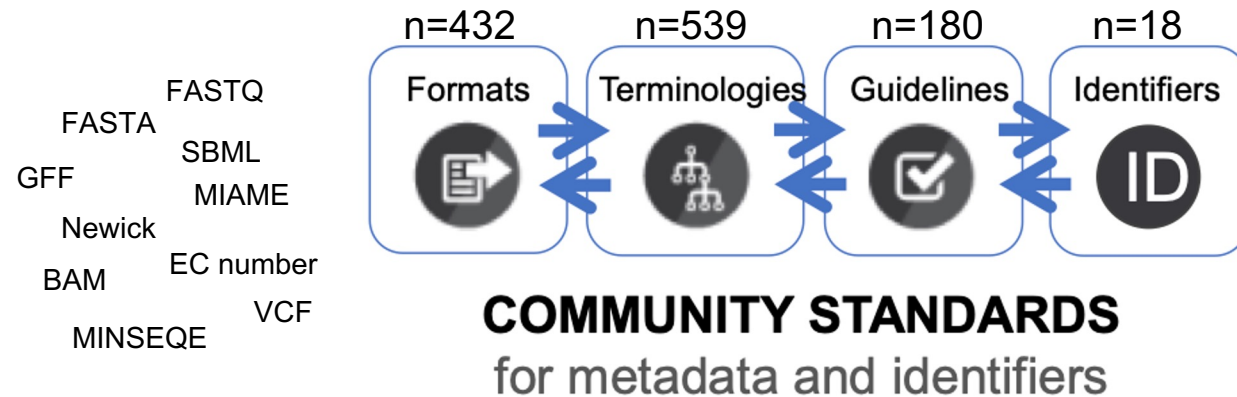
**Conceptual
model, schema,
exchange
formats, etc...**
e.g. SBML,
FASTA



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

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

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



Source:

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



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

Graph view to visualize relationship links between resources

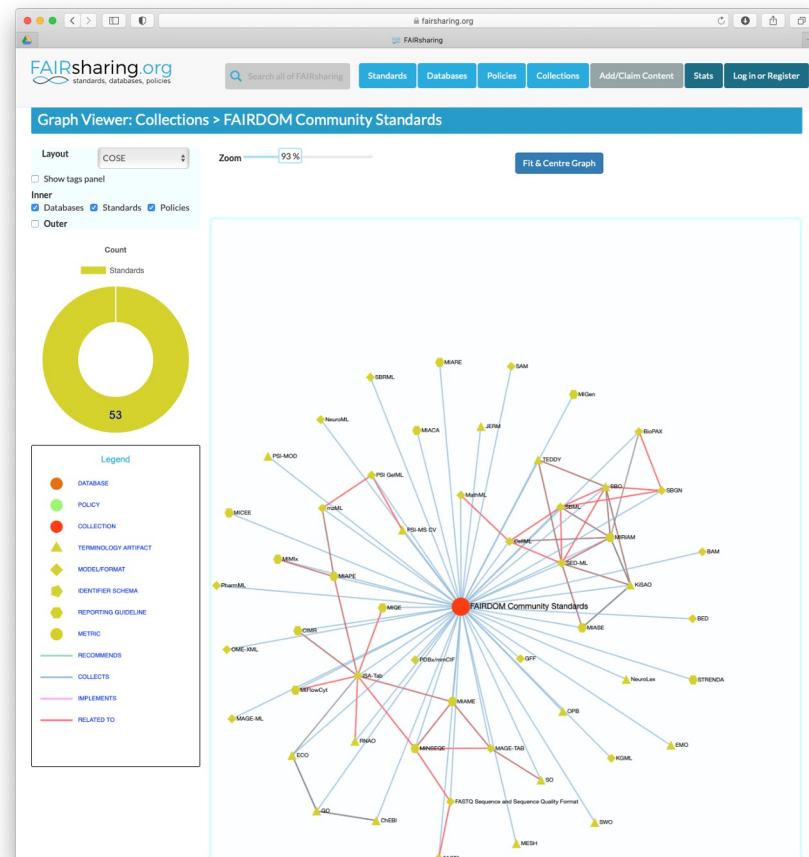
<https://fairsharing.org/collections/>

The screenshot shows the FAIRsharing.org interface for the 'COVID-19 Resources' collection. The top navigation bar includes links for Standards, Databases, Policies, Collections, Add Claim Content, Stats, and Log in or Register. The main content area displays a graph view of the collection, showing records 1 to 50 of 80. The graph visualizes the relationships between various resources, including databases, standards, and policies. A sidebar on the left provides filters for Recommended Records, Associated Publications, Claimed?, Record Status, Standard Type, Registry, and Domains. The main table lists records with columns for Registry, Name, Abbreviation, Type, Subject, Domain, Taxonomy, Related Database, Related Standard, and In Collection.

Registry	Name	Abbreviation	Type	Subject	Domain	Taxonomy	Related Database	Related Standard	In Collection
American Type Culture Collection Database	ATCC	Database	ATCC	ATCC	ATCC	ATCC	None	None	Yes
Australian New Zealand Clinical Trials Registry	ANZCTR	Database	ANZCTR	ANZCTR	ANZCTR	ANZCTR	None	None	Yes
EBM-DBS	EBM-DBS	Database	EBM-DBS	EBM-DBS	EBM-DBS	EBM-DBS	None	None	Yes

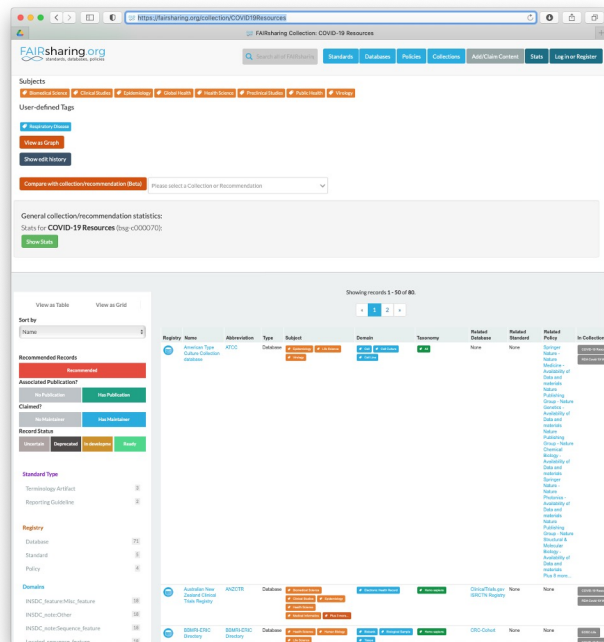
Example 1: the *FAIRdom community Standards collection* (System biology)

<https://fairsharing.org/collection/FAIRD0M>

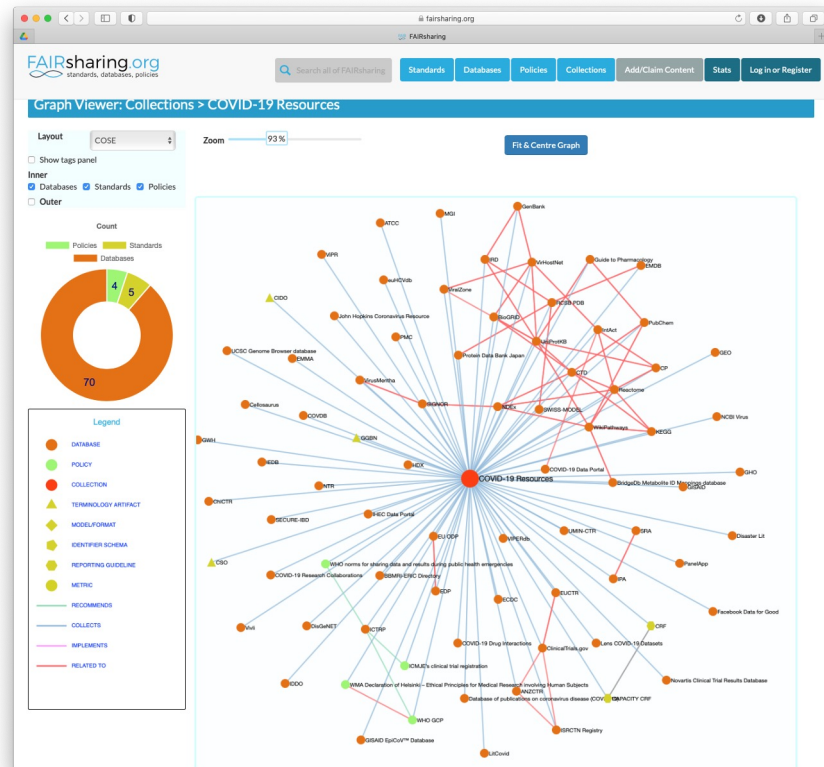


Some collections are recent

Example 2: The Covid-19 collection

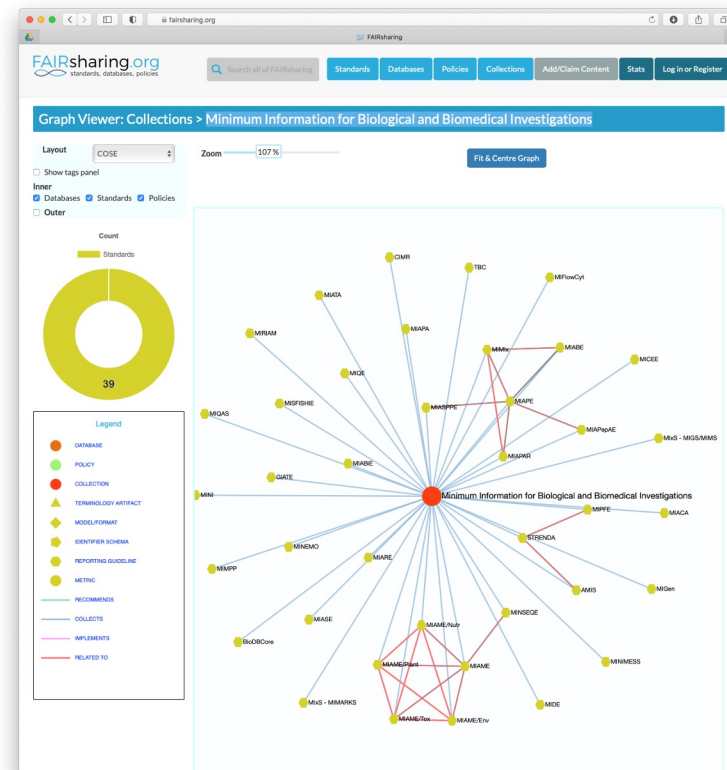


<https://fairsharing.org/collection/COVID19Resources>



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

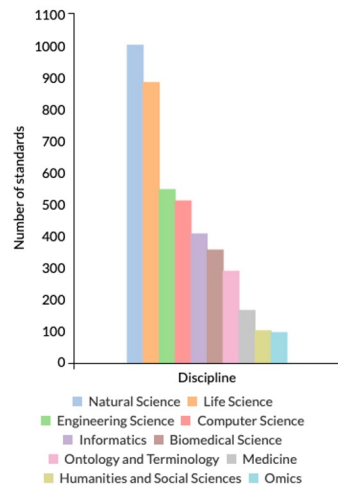
Example 3: the *Minimum Information for Biological and Biomedical Investigations* collection



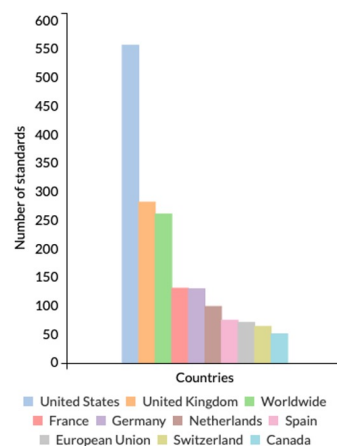
<https://fairsharing.org/collection/MIBBI>



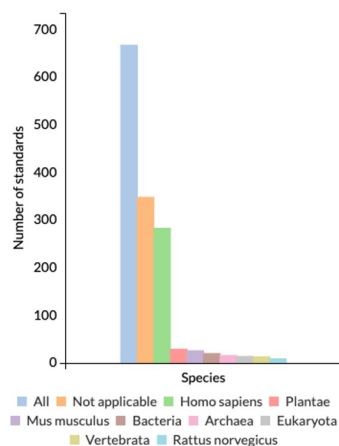
Top 10 disciplines covered by standards



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>

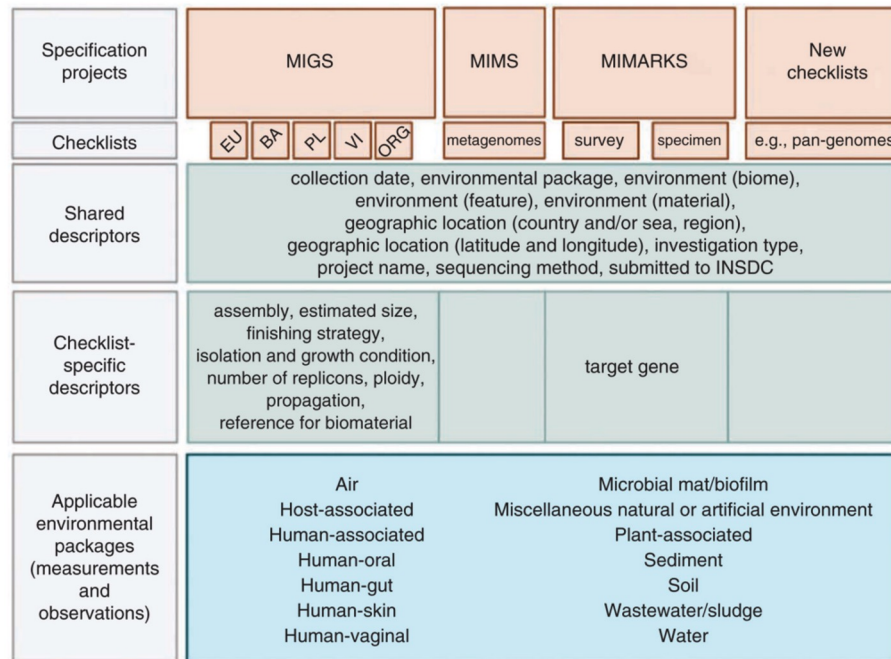


Find the *Genomic Standards Consortium (GSC)* used by both *ENA* and *SRA* databases in the **FAIRsharing collections**

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 ?

- An international community-driven standard in **Genomics** producer of the ***MlxS: Minimum Information Standards about any(X) Sequence***
- MlxS includes **technology-specific checklists** (MIGS, MIMS, MIMARKS,...) and also allows **annotation of sample data** using environmental packages



[Yilmaz et al, 2011](#)

Source: <https://gensc.org>