

Open Science

Publicly Available Resources

Alberto Santos & Yesid Cuesta-Astroz



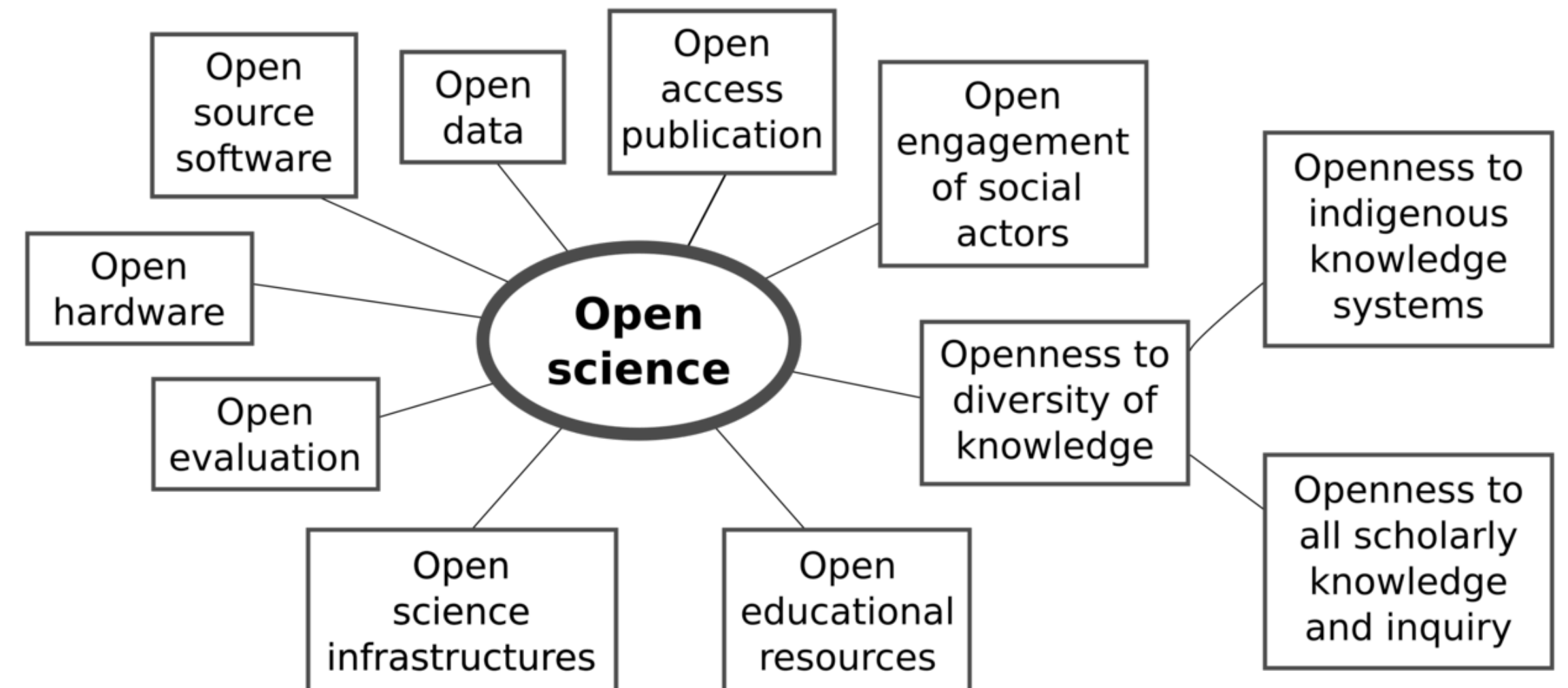
Outline

- What is Open Science?
- FAIR Data and Software
- Challenges sharing and reusing data
- Standardisation and Ontologies
- Publicly available resources

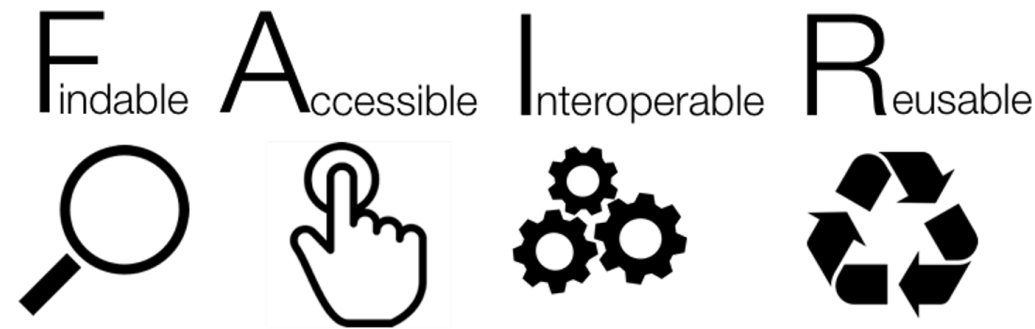
What is Open Science

Impact, Contribution, Trust

- Make scientific research **accessible** to all levels of society:
 - Publications
 - Samples
 - Methods
 - **Software**
 - **Data**
- Advantages:
 - **Reproducibility** and **replicability**
 - Societal **responsibility** — publicly funded, publicly available
 - **Multi-purpose** of research outputs
- Disadvantages: concerns of data **misuse**



FAIR Data and Software



- **F**indable and **A**ccessible

- Add enough **metadata** — data about your data
- Deposit your data in **public repositories** or make them available in **databases**

[Minimum Information for Biological and Biomedical Investigations](#)

[Zenodo](#)
[Figshare](#)
[Pride](#)
[Metabolights](#)
[GEO](#)
[GitHub](#)

- **I**nteroperable:

- Use **standard** and **open formats**
- Provide **all data needed** to reproduce your analysis

- **R**eusable:

- **Describe** your data well, e.g., good metadata but also
- Attach a **license**

Provide README files describing the data
Use descriptive column headers for the data tables

Challenges Sharing and Reusing

The marshmallow test — delayed gratification

- Open does not mean **FAIR**
- Requires an **effort**
- **Metadata** becomes the most important data
- In many cases there are **no standards or multiple ones**
- **Most of the data** out there **not FAIR**



<https://imgflip.com/memegenerator>

https://en.wikipedia.org/wiki/Stanford_marshmallow_experiment

Standardisation and Ontologies

- Data **standardisation** requires defining **terminologies** and **vocabularies** that:
 - Assign **unique identifiers** to entities/concepts such as proteins, genes, diseases
 - **Describe** those entities/concepts and **provide meaning**
 - **Relate** those concepts to other terms
 - Classify those entities/concepts into **categories**

- **Solution —> Ontologies**

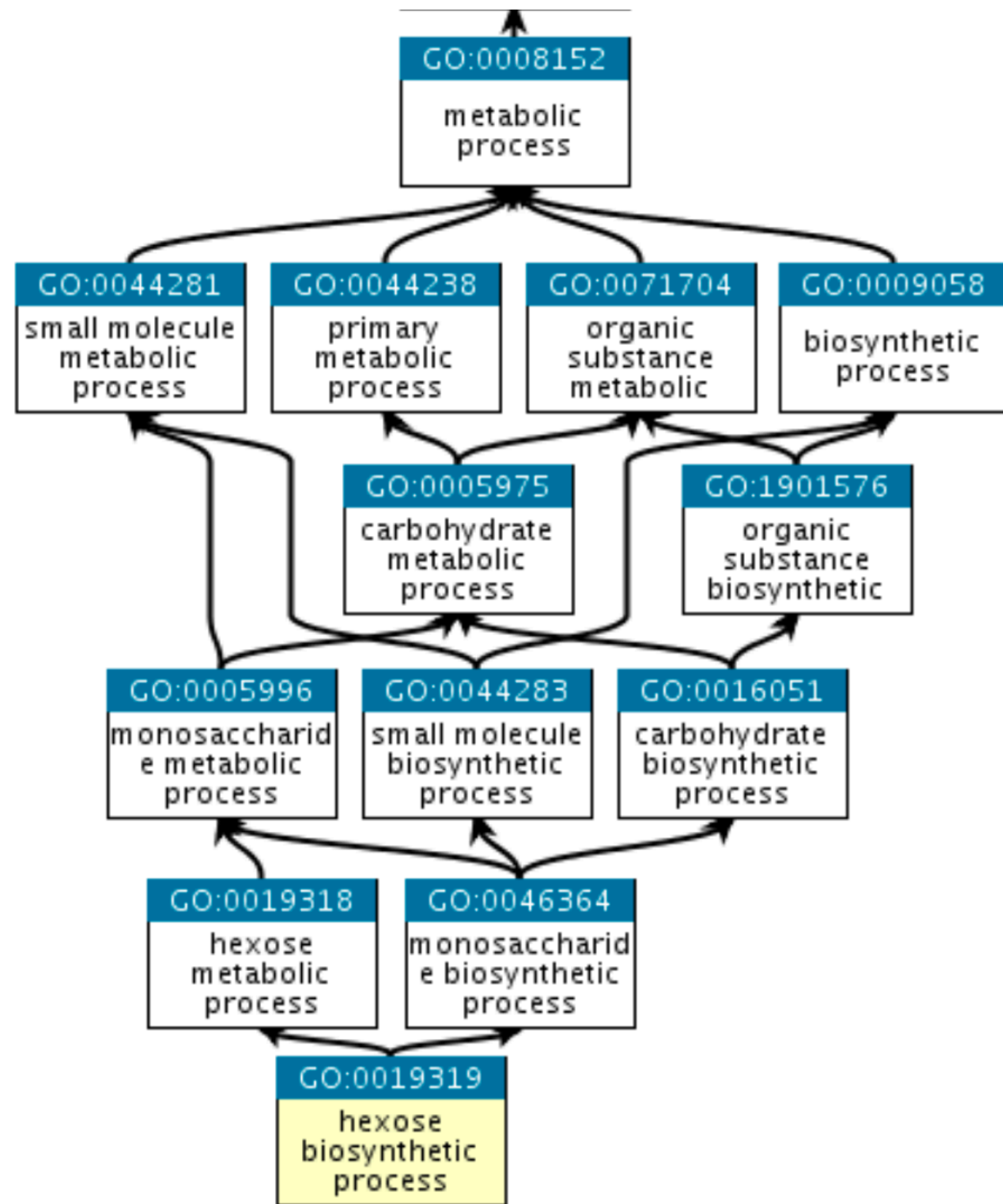
- **Ontology:**

formal way of representing knowledge in which concepts are described both by their meaning and their relationship to each other

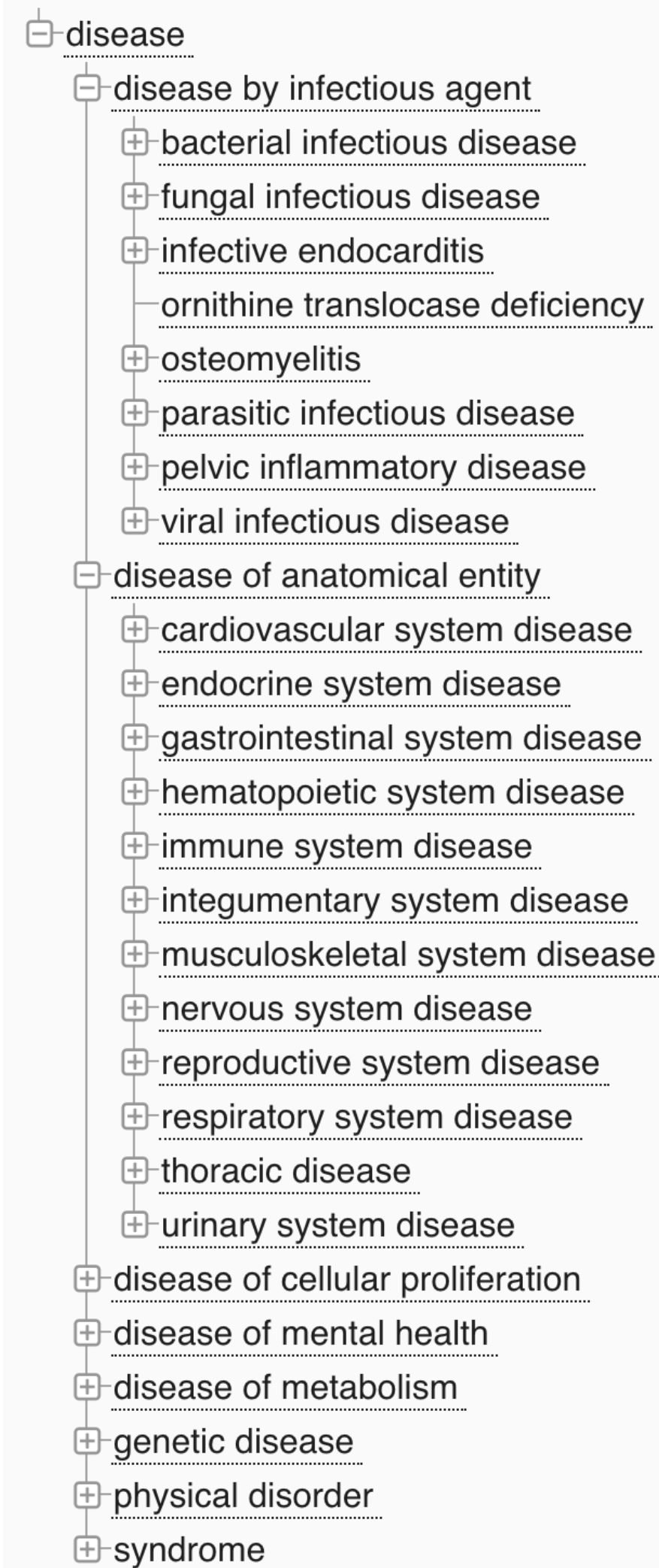
A collection of terms and their definitions for a specific domain

Ontologies

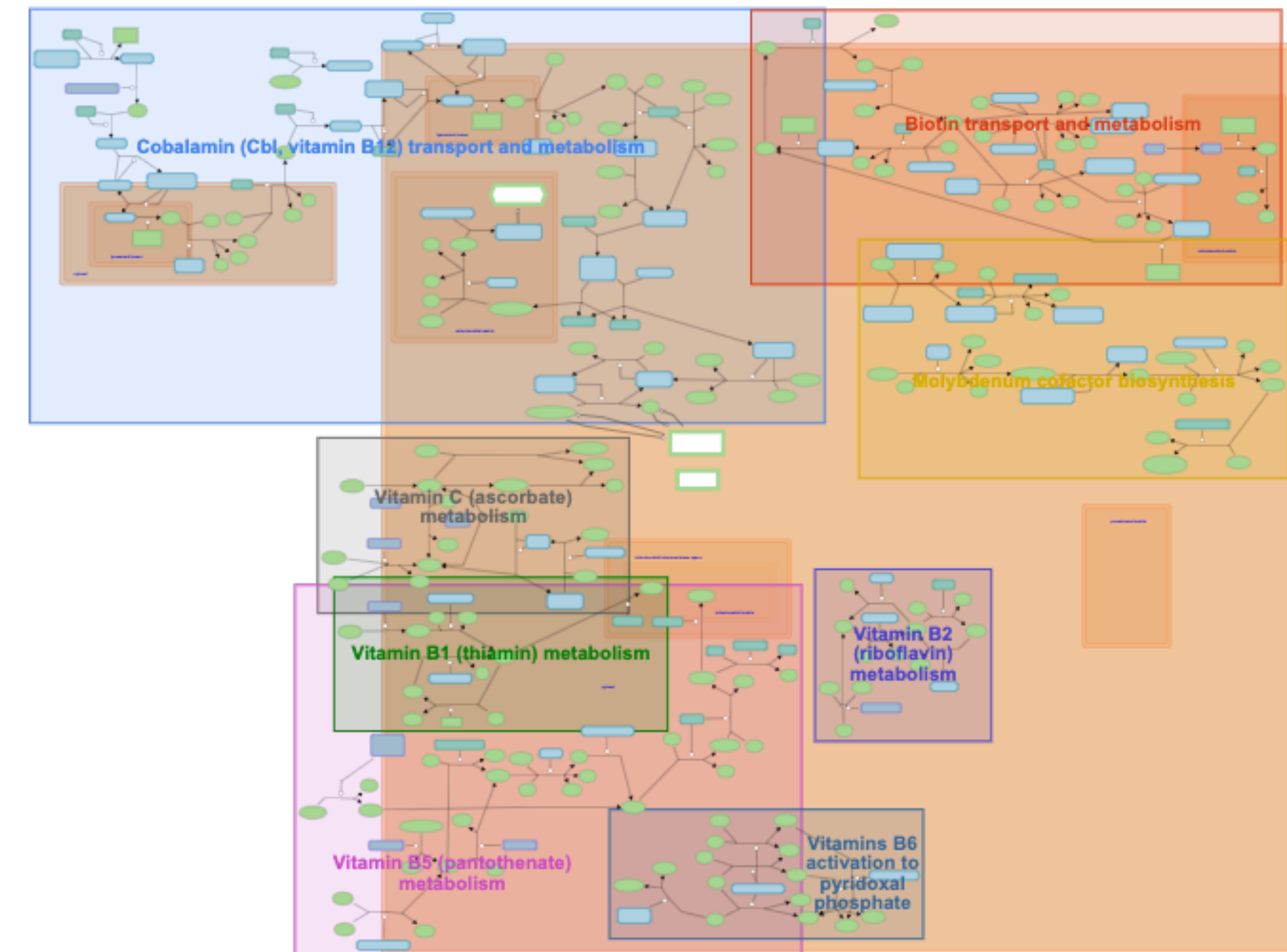
Gene Ontology



Disease Ontology



REACTOME Pathways



<https://www.ebi.ac.uk/ols/ontologies>

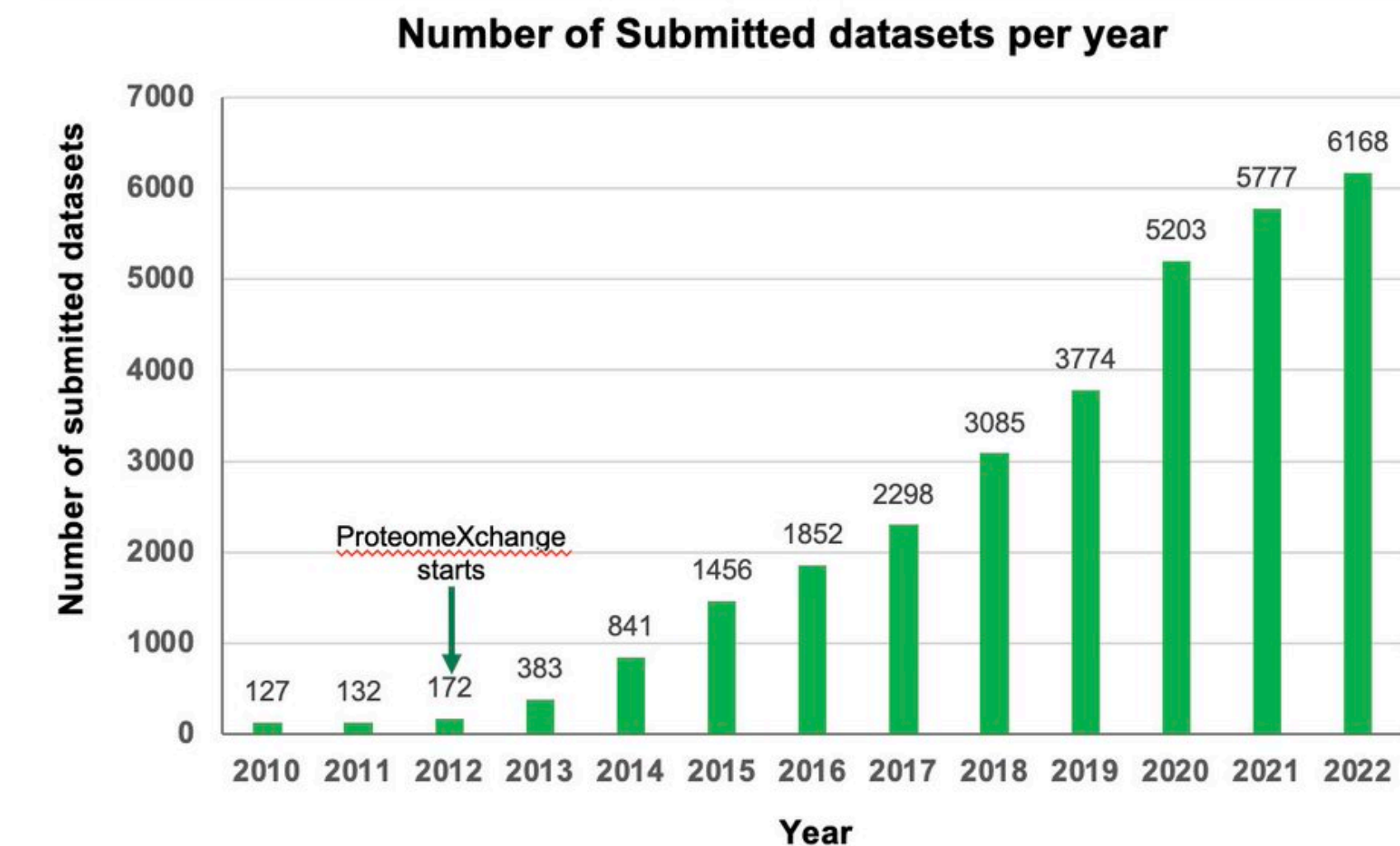
<https://reactome.org/>

<http://geneontology.org/>

Publicly Available Resources

Be a Data Parasite

- Do not reinvent the wheel
- **Extend the life and purpose** of publicly available **data**
- Build **in-silico hypotheses** before jumping into experiments (cheaper, higher success rate)
- **Download — Use — Test — Transform — Upload**
- **Growing number of resources and datasets** available



Examples of Microbes Resources

ALEdb 1.0: a database of mutations from adaptive laboratory evolution experimentation <https://aledb.org/>

MiMeDB: the Human Microbial Metabolome Database <https://mimedb.org/>

Web of microbes (WoM): a curated microbial exometabolomics database for linking chemistry and microbes <https://metatlas.nersc.gov/wom/project-begin.view>

MicroPhenoDB Associates Metagenomic Data with Pathogenic Microbes, Microbial Core Genes, and Human Disease Phenotypes <http://www.liwzlab.cn/microphenodb>

BacDive in 2022: the knowledge base for standardized bacterial and archaeal data <https://bacdive.dsmz.de/>

MASI: microbiota—active substance interactions database <http://www.aiddlab.com/MASI/>

iModulonDB: a knowledgebase of microbial transcriptional regulation derived from machine learning <https://imodulondb.org/index.html>

MIBiG 3.0: a community-driven effort to annotate experimentally validated biosynthetic gene clusters <https://mibig.secondarymetabolites.org/>

The Microbial Metabolites Database

MiMeDB

- The **human microbiome** is believed to produce or process **>55,000 different compounds** – many of which **affect human health, behavior and disease**
- **Microbes synthesise primary metabolites** required for their own survival, but they also **produce other compounds** arising **from substrates or host-derived food sources**

E.g., microbes transform xenobiotics from food constituents, food additives, phytochemicals, drugs, cosmetics and other exogenous or man-made chemicals

- **MiMeDB is a database of small molecule metabolites found in the human microbiome**
- Provides links between **metabolites, microbes, hosts, health and exposure** data

Microbe and Disease Phenotype Association Database

MicroPhenoDB

- Manually **curated database** integrating **microbe-disease associations**
- Provides **5677** non-redundant **associations** between **1781 microbes** and **542 human diseases** across more than **22 human tissues**
- Disease phenotypes are classified using **Experimental Factor Ontology (EFO)** (<https://www.ebi.ac.uk/efo/>)
- Aims to accelerate **metagenomic data analysis**

