# Metadata standards and ontologies
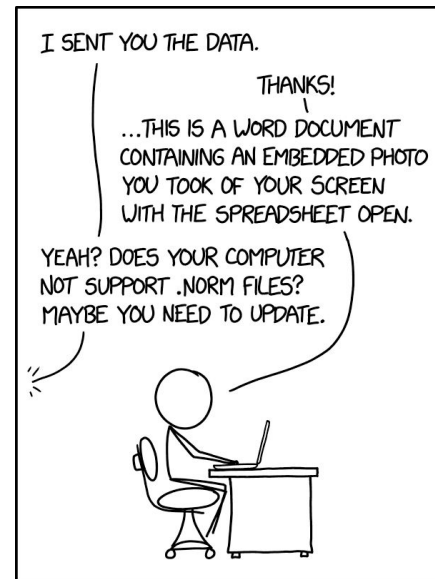
Claire Rioualen

IFB-core - Marseille

Biotic virtual lunch – February 24th, 2025

# Metadata issues and misuse: a never-ending struggle

- Ever tried downloading and analysing data from GEO or ArrayExpress?

  ○ Search for synonyms (ChIP-seq vs ChIP-sequencing)

  ○ Extract metadata with distinct field names (source, condition…)

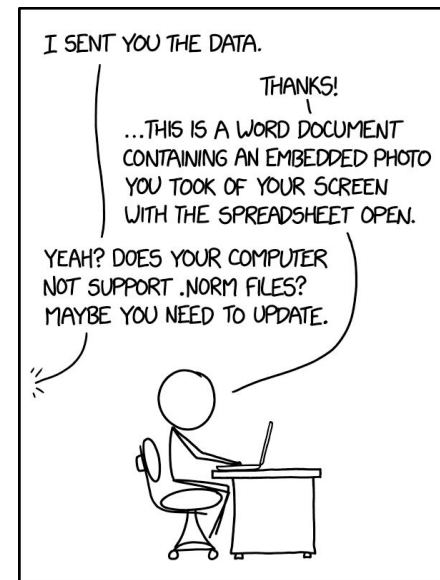  ○ Different terms for similar information (Escherichia coli K-12, MG1655…)

https://xkcd.com/2116/

# Metadata issues and misuse: a never-ending struggle

- Ever tried downloading and analysing data from GEO or ArrayExpress?

  - Search for synonyms (ChIP-seq vs ChIP-sequencing)

  - Extract metadata with distinct field names (source, condition…)

  - Different terms for similar information (Escherichia coli K-12, MG1655…)

- Ever tried integrating and re-analysing several public datasets?

  - Different gene or protein names

  - Metadata and actual data being inconsistent

  - Detailed methods or code lacking from publications

  - Authors not reachable, lab tech working elsewhere…

https://xkcd.com/2116/

*"Metadata, you see, is really a love note - it might be to yourself, but in fact it's a love note to the person after you, or the machine after you, where you've saved someone that amount of time to find something by telling them what this thing is"*

Jason Scott
http://ascii.textfiles.com/archives/3181

- Big data poses many challenges owing to its amount, complexity and heterogeneity

- Metadata is essential to annotate its **contents, origin and meaning**

- Metadata alone is not a guarantee that the data and associated results are **FAIR**

**Data**  **Metadata**

Michi
+777 236 7246

- Big data poses many challenges owing to its amount, complexity and heterogeneity

- Metadata is essential to annotate its **contents, origin and meaning**

- Metadata alone is not a guarantee that the data and associated results are **FAIR**

- **Standard metadata** improves the findability and reusability of the data for future users

- **Semantic metadata** improves data interoperability by using **machine-readable metadata schemes**

**Data**   **Metadata**

Michi
+777 236 7246

$F_{indable}$ $A_{ccessible}$ $I_{nteroperable}$ $R_{eusable}$

- **Controlled vocabulary**, thesaurus, subject headings
  - Preferred unique terminology

- **Taxonomy**
  - Terminology
  - Hierarchy of terms

- **Ontology**
  - Terminology
  - Alternative terms
  - Definition
  - Properties associated with terms
  - Hierarchical and semantic relations

https://xkcd.com/1179/

**Structured set of concepts** related to a given **field of knowledge**, their **definitions**, **relations**, unique and **permanent identifiers**, and other related **properties**
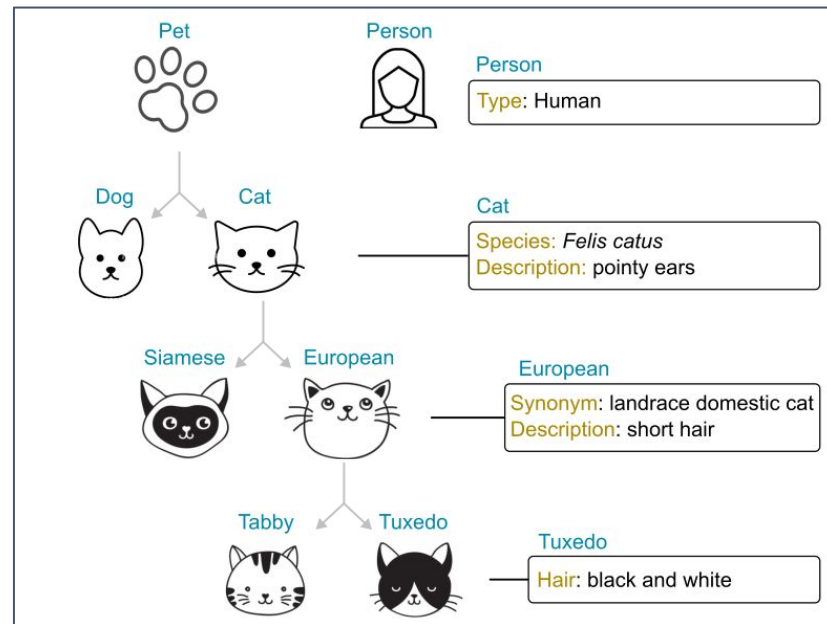
- **Class**: category of being

**Pet ontology**

**Structured set of concepts** related to a given **field of knowledge**, their definitions, **relations**, unique and **permanent identifiers**, and other related **properties**

- **Class**: category of being
- **Attributes**: properties associated with a class
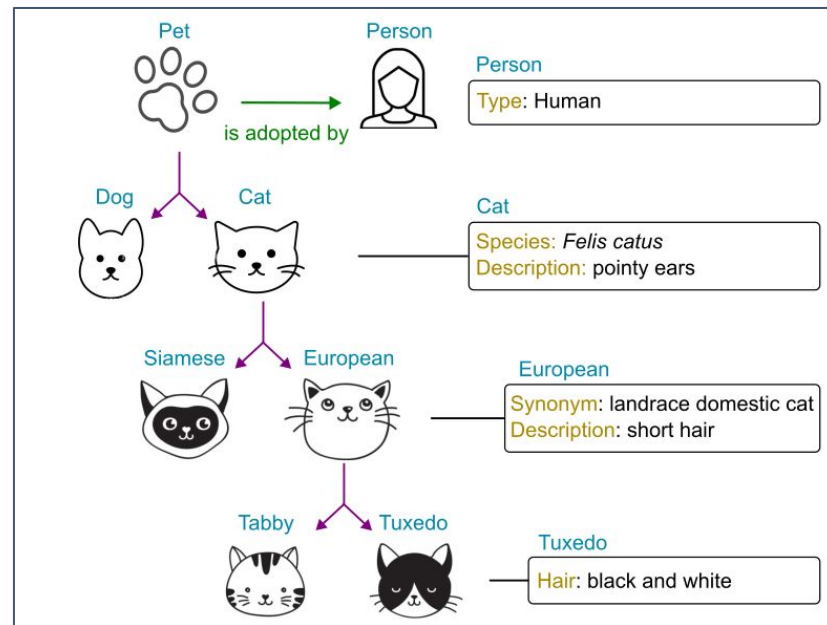
**Pet ontology**

# What is an ontology?

**Structured set of concepts** related to a given **field of knowledge**, their definitions, **relations**, unique and **permanent identifiers**, and other related **properties**
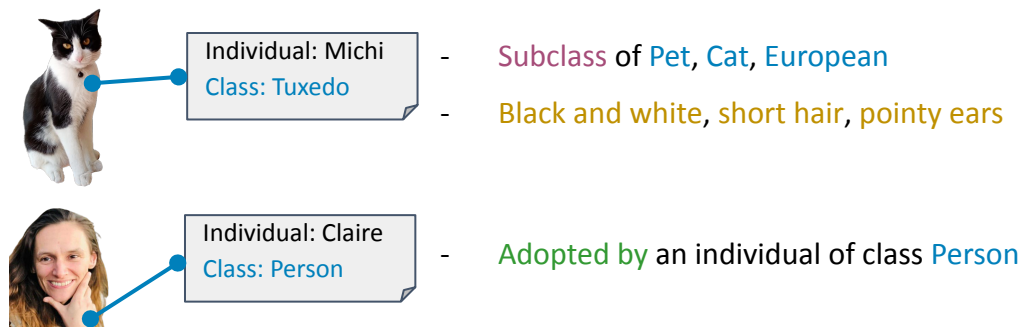
- **Class**: category of being
- **Attributes**: properties associated with a class
- Relations: **hierarchical** or **semantic**

**Pet ontology**
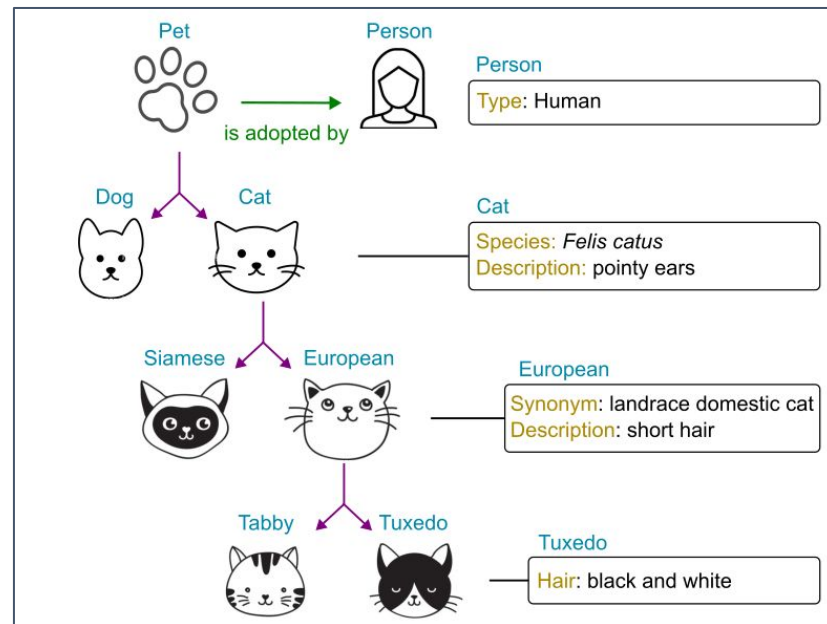
# What is an ontology?

**Structured set of concepts** related to a given **field of knowledge**, their definitions, **relations**, unique and **permanent identifiers**, and other related **properties**

- **Class**: category of being
- **Attributes**: properties associated with a class
- Relations: **hierarchical** or **semantic**
- **Individual**: instance of a class

Individual: Michi
Class: Tuxedo

- Subclass of Pet, Cat, European
- Black and white, short hair, pointy ears

Individual: Claire
Class: Person

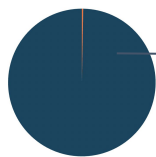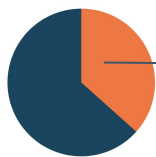- Adopted by an individual of class Person

**Pet ontology**

**schema.org:** a terminology and **metadata scheme** for *Things*

schema.org

- Used for web page annotations, it improves indexation by search engines

0.3% of websites use schema.org annotations in their metadata

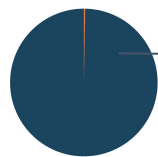36.6% of Google search results include websites with schema.org metadata

Source: Beard et al. 2016

# Metadata schemes: the example of schema.org

**schema.org:** a terminology and **metadata scheme** for *Things*

schema.org

- Used for web page annotations, it improves indexation by search engines

0.3% of websites use schema.org annotations in their metadata

36.6% of Google search results include websites with schema.org metadata
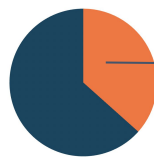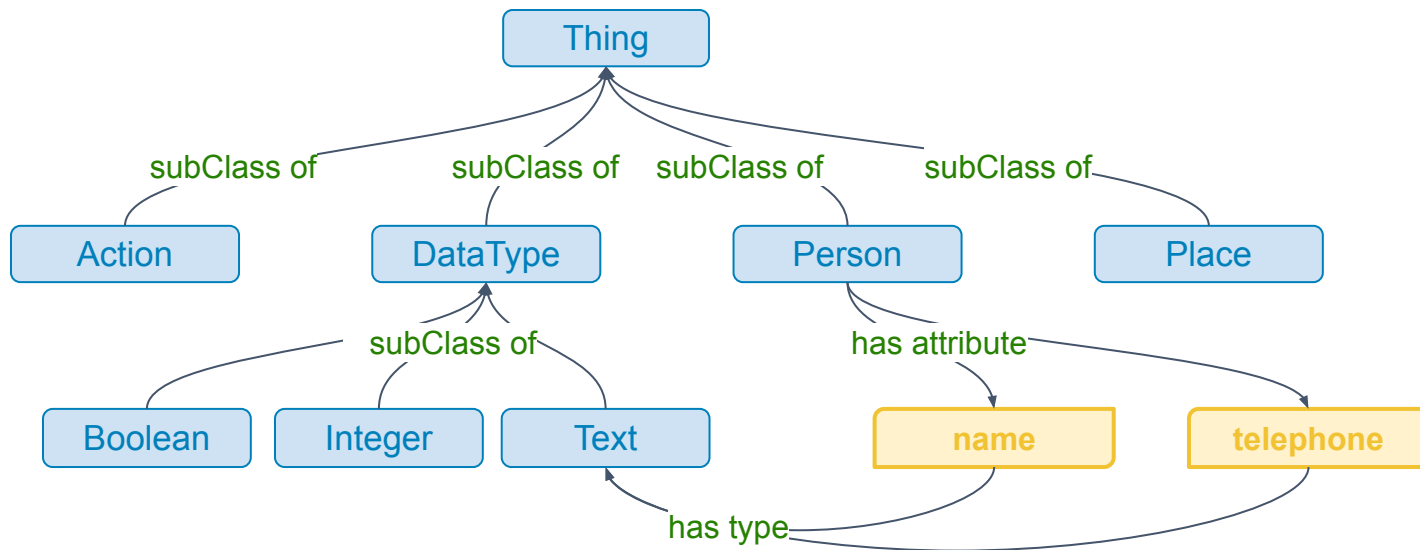
Source: Beard et al. 2016

The **World Wide Web Consortium (W3C)** endorses multiple standards for metadata

- **Resource Description Framework** (RDF)
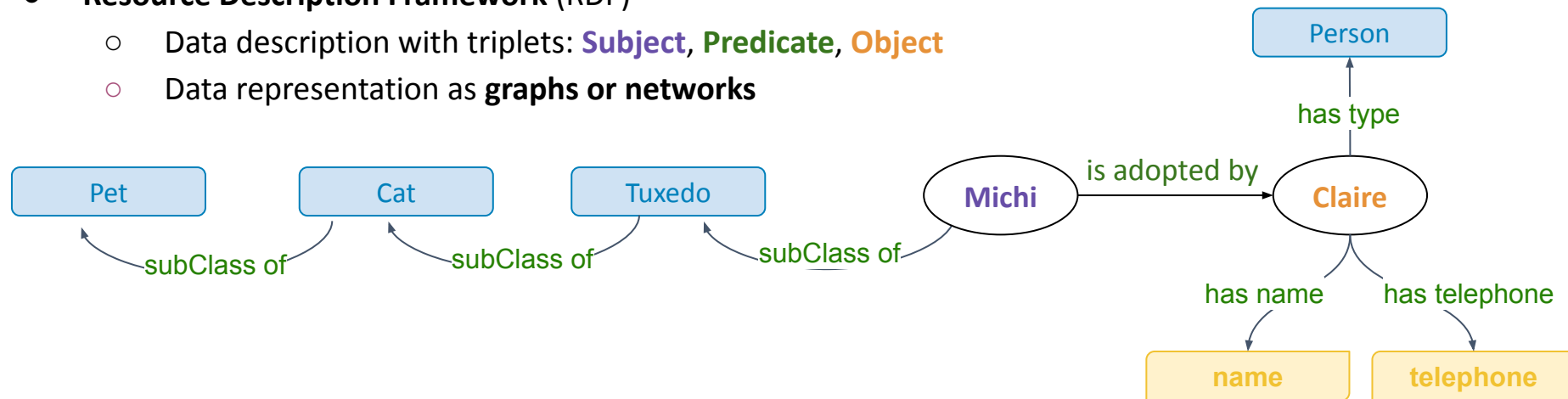  - Data description with triplets: **Subject**, **Predicate**, **Object**

The **World Wide Web Consortium (W3C)** endorses multiple standards for metadata

- **Resource Description Framework** (RDF)
    - Data description with triplets: **Subject**, **Predicate**, **Object**
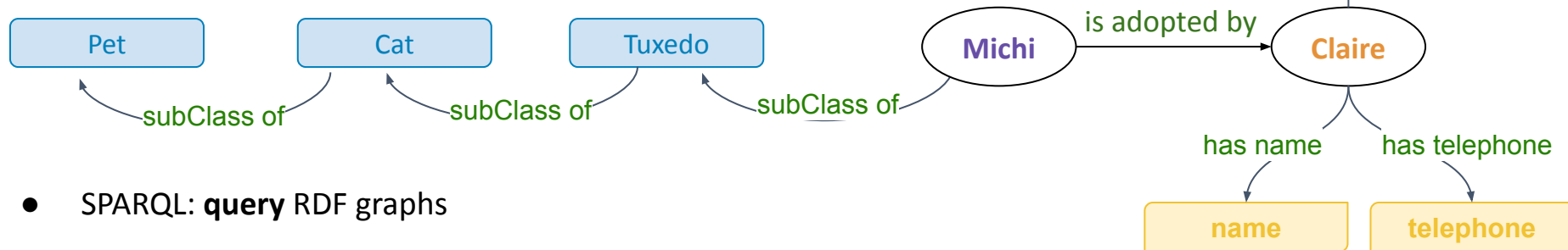    - Data representation as **graphs or networks**

The **World Wide Web Consortium (W3C)** endorses multiple standards for metadata

- **Resource Description Framework** (RDF)
  - Data description with triplets: **Subject**, **Predicate**, **Object**
  - Data representation as **graphs or networks**



- SPARQL: **query** RDF graphs

```
PREFIX po: <http://pet_ontolgy.org>
PREFIX sch: <http://schema.org>

SELECT   ?pet ?owner
WHERE {  ?pet rdf:subClass_of po:Cat
         ?pet po:is_adopted_by ?owner
         ?owner sch:has_name Claire }
```

The **World Wide Web Consortium (W3C)** endorses multiple standards for metadata

W3C®

- **Resource Description Framework** (RDF)
  - Data description with triplets: **Subject**, **Predicate**, **Object**
  - Data representation as **graphs or networks**



- SPARQL: **query** RDF graphs

```
PREFIX po: <http://pet_ontolgy.org>
PREFIX sch: <http://schema.org>

SELECT   ?pet ?owner
WHERE { ?pet rdf:subClass_of po:Cat
        ?pet po:is_adopted_by ?owner
        ?owner sch:has_name Claire }
```
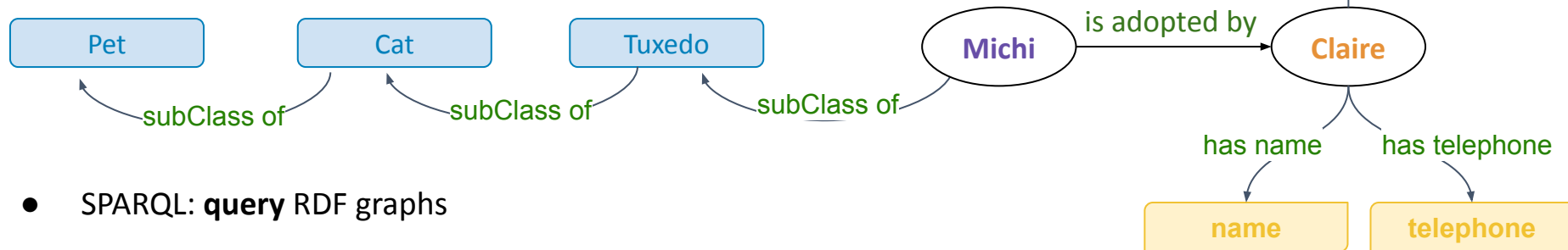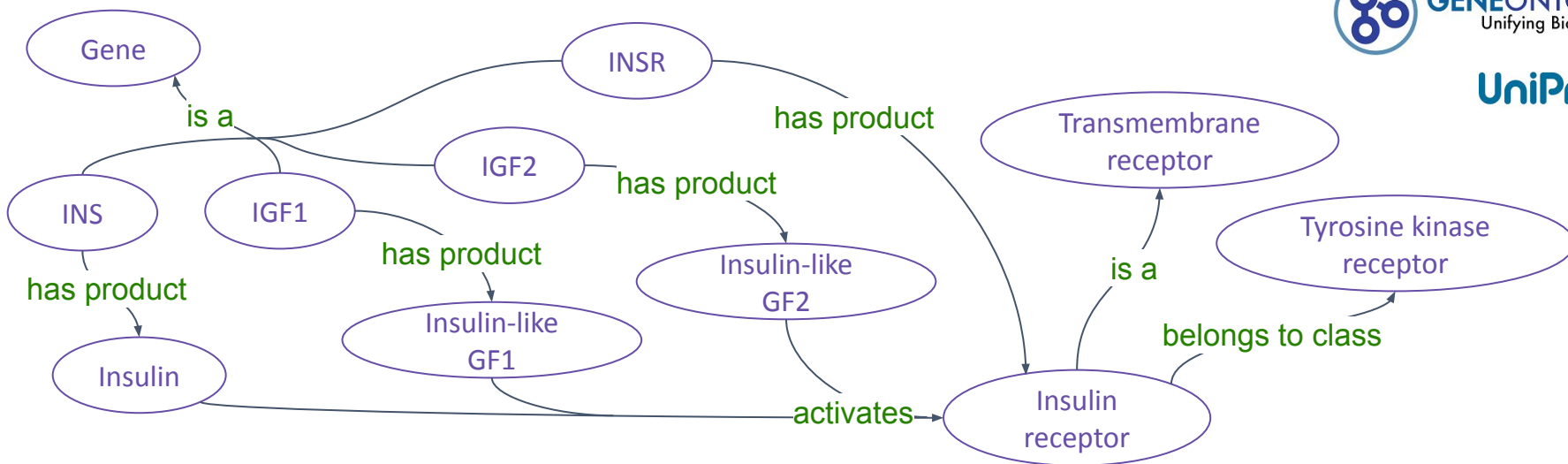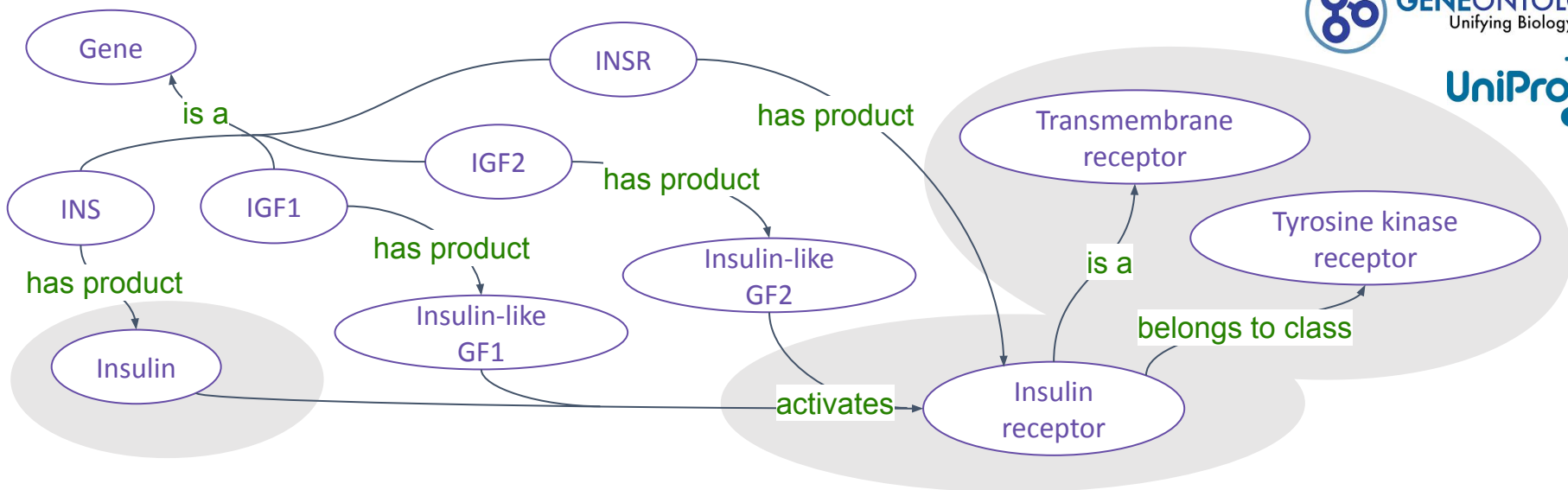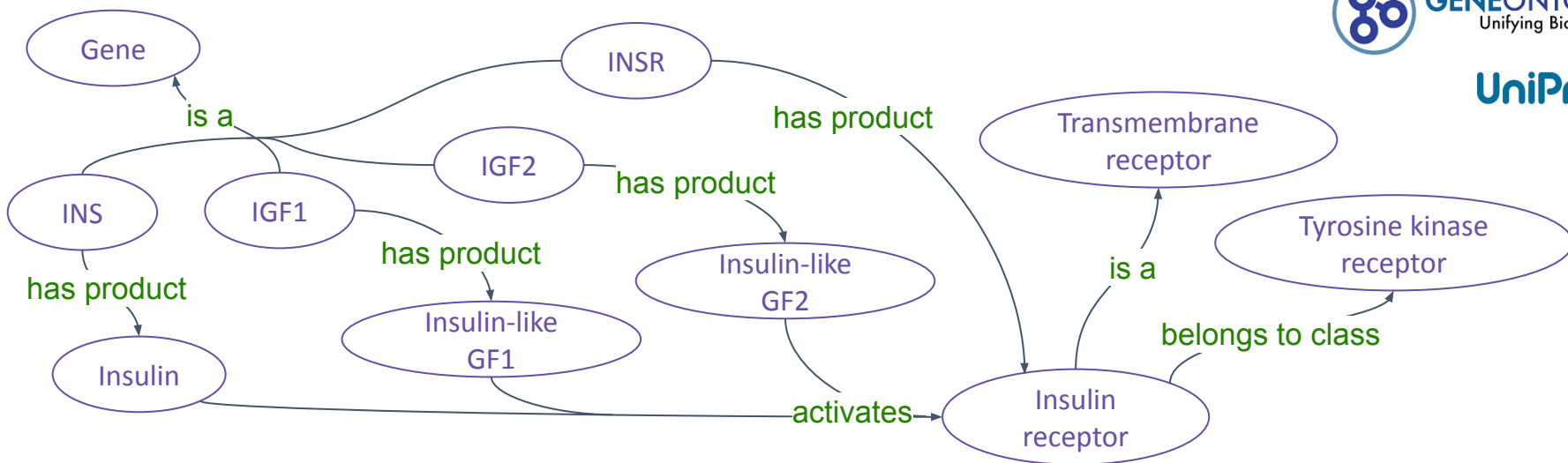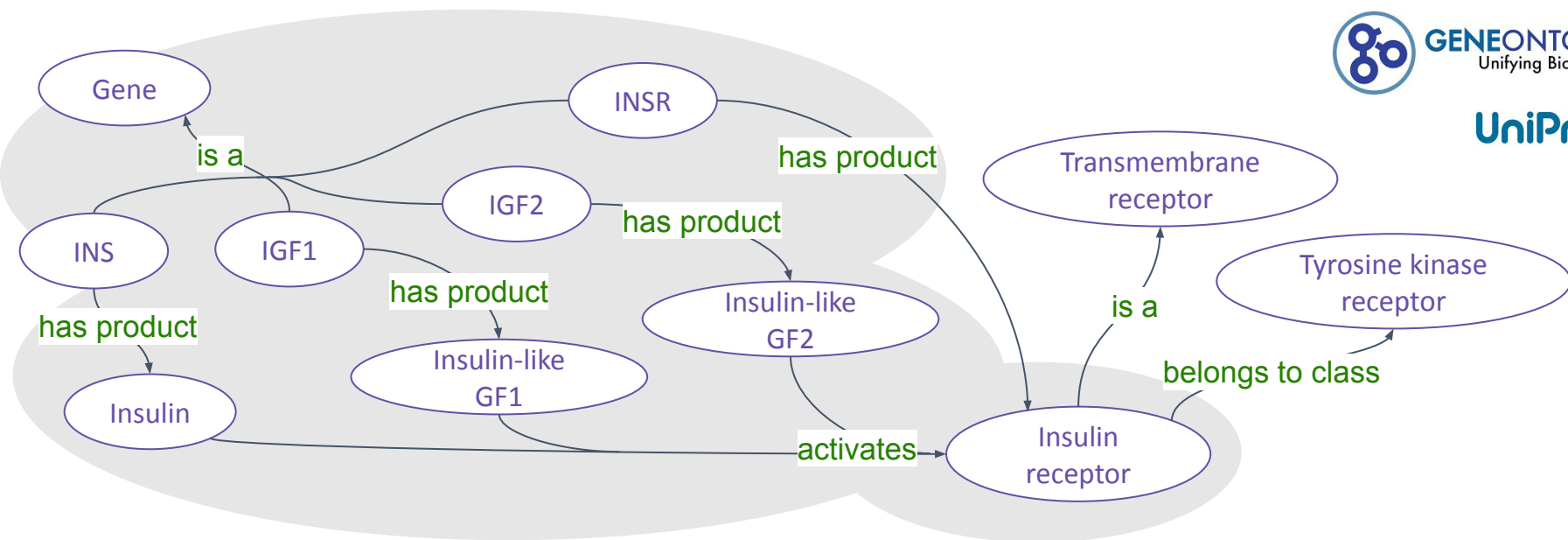
16

# UniprotKB and the Gene Ontology



"The insulin receptor is a transmembrane receptor activated by insulin, IGF-I, IGF-II, and belongs to the class of tyrosine kinase receptors."

# UniprotKB and the Gene Ontology



"The insulin receptor is a transmembrane receptor activated by insulin, IGF-I, IGF-II, and belongs to the class of tyrosine kinase receptors."

What are the genes coding for the activators of insulin receptors?

# UniprotKB and the Gene Ontology



"The insulin receptor is a transmembrane receptor activated by insulin, IGF-I, IGF-II, and belongs to the class of tyrosine kinase receptors."

What are the genes coding for the activators of insulin receptors?

- INS, IGF1, IGF2

```
SELECT ?gene
WHERE {
  ?gene go:has_product ?protein
  ?protein upr:activates upr:insulin_receptor
}
```

# Ontologies for life sciences

Numerous ontologies, from very general to very specific domains

- **MeSH - Medical Subject Headings**
  - Pubmed indexing and search

- **HPO - Human Phenotype Ontology**
  - Multi-lingual support
  - Layperson synonyms

  🇬🇧Macroencephaly
  🇯🇵大頭
  🙌Big head

- **SO - Sequence Ontology**
  - Sequence attributes, sequence features…

- **EFO - Experimental Factor Ontology**
  - Cell types, biological processes, protocols..

- Experimental biology: Cellosaurus, Microbial Conditions Ontology…

- Species specific: fish ontology, potato ontology, banana ontology…

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…

Organisation in **4 main classes**:

- **Topic**: field of study or technology          → *Genomics, Sequencing*

# The EDAM ontology for bioinformatics analyses

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…

Organisation in **4 main classes**:

- **Topic**: field of study or technology     → *Genomics, Sequencing*

- **Operation**: process or function     → *DNA mapping, Peak calling*

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…

Organisation in **4 main classes**:

- **Topic**: field of study or technology  → *Genomics, Sequencing*

- **Operation**: process or function  → *DNA mapping, Peak calling*

- **Data**: data type  → *Gene ID, DNA sequence*

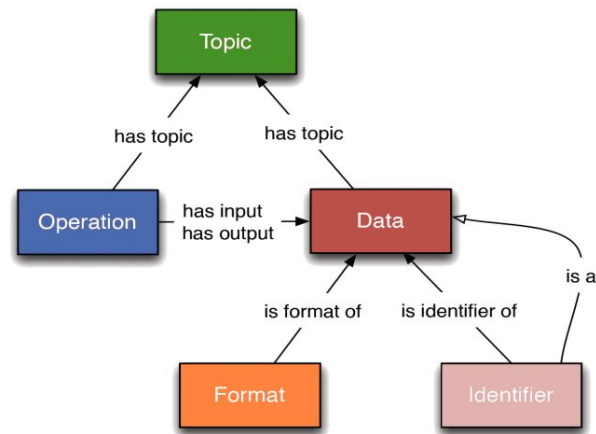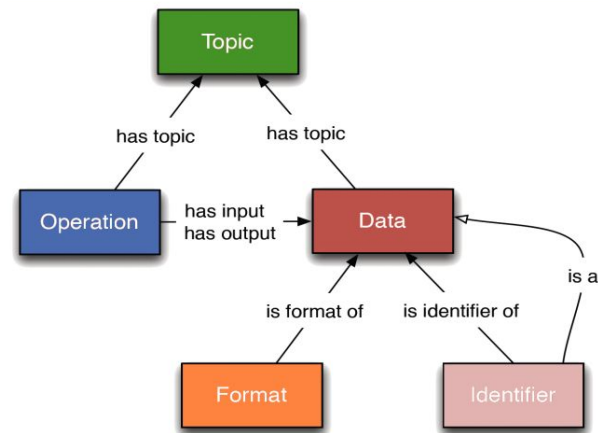# The EDAM ontology for bioinformatics analyses

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…
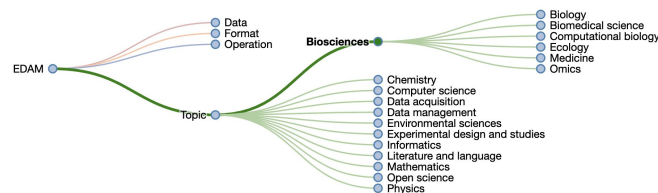
Organisation in **4 main classes**:

- **Topic**: field of study or technology      → *Genomics, Sequencing*

- **Operation**: process or function      → *DNA mapping, Peak calling*

- **Data**: data type      → *Gene ID, DNA sequence*

- **Format**: file format      → *FASTQ, BAM, JSON*

# The EDAM ontology for bioinformatics analyses

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…

Organisation in **4 main classes**:

- **Topic**: field of study or technology    → *Genomics, Sequencing*

- **Operation**: process or function    → *DNA mapping, Peak calling*

- **Data**: data type    → *Gene ID, DNA sequence*

- **Format**: file format    → *FASTQ, BAM, JSON*

# The EDAM ontology for bioinformatics analyses

**EDAM**

Concepts about **data management and analysis** for life sciences, their **relations** and **attributes**
→ definition, permanent identifier, synonyms, references…

Organisation in **4 main classes**:

- **Topic**: field of study or technology  → *Genomics, Sequencing*
- **Operation**: process or function  → *DNA mapping, Peak calling*
- **Data**: data type  → *Gene ID, DNA sequence*
- **Format**: file format  → *FASTQ, BAM, JSON*



Development

- **Open and collaborative** via Github or the GUI Protégé/WebProtégé
- **Visualisation** and navigation with the EDAM browser
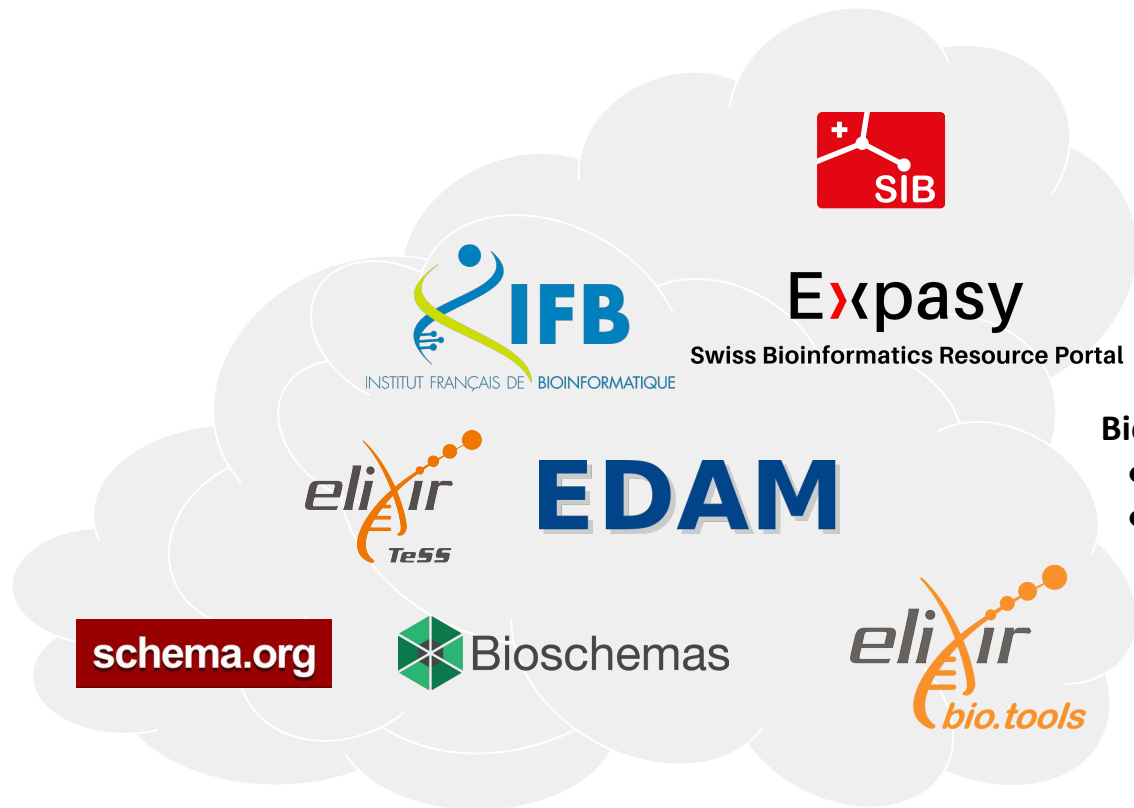- **Uninterrupted**: knowledge never stops growing!

**Bioschemas**
- Metadata schemes using schema.org
- Definition of biology-oriented profiles using EDAM terms
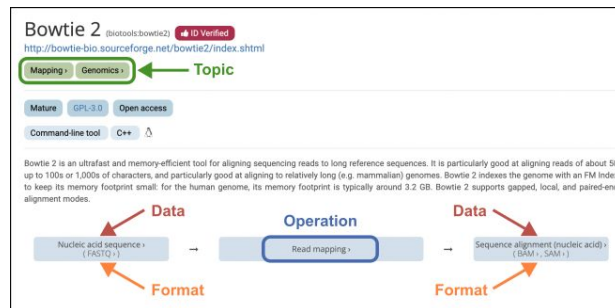
**TeSS portal**
- Automated scraping from selected resources using Bioschemas annotations
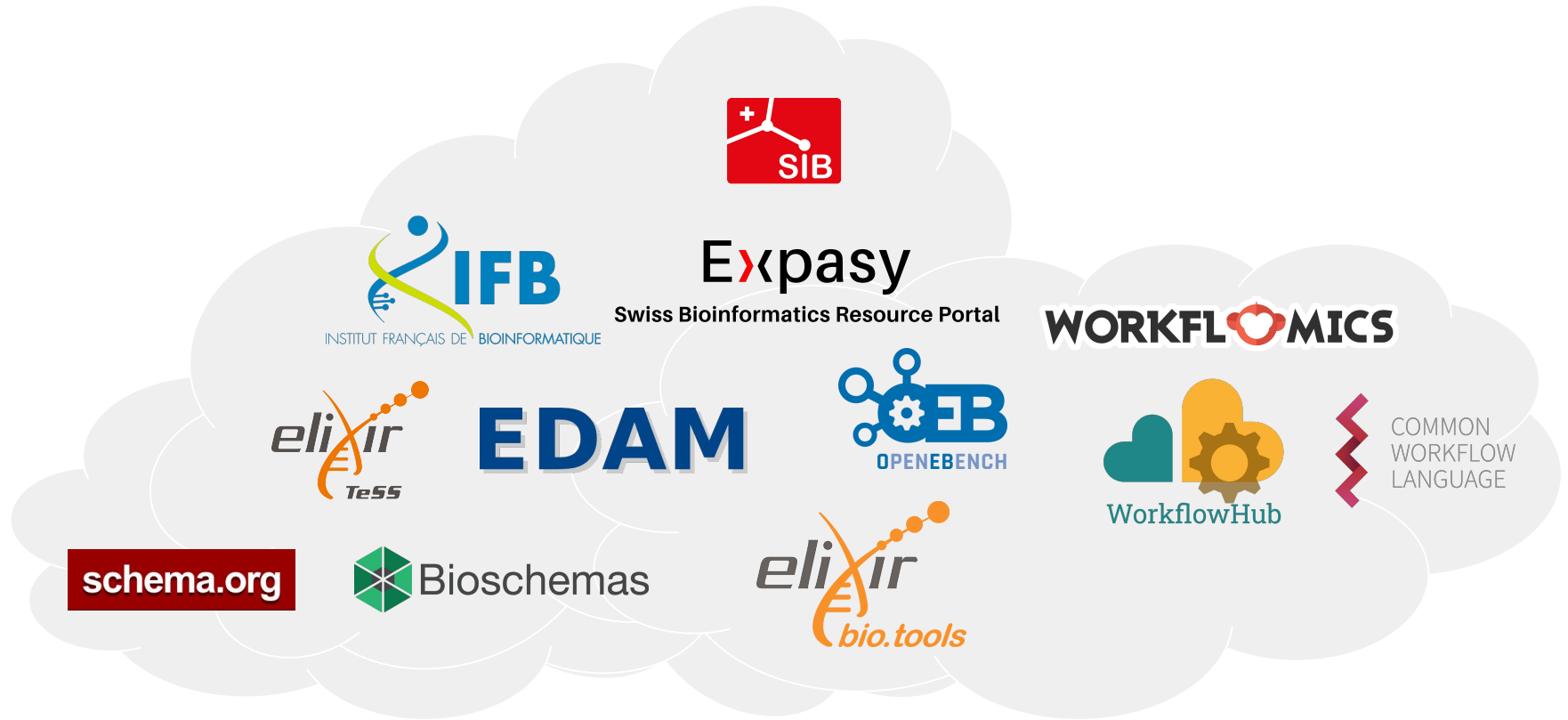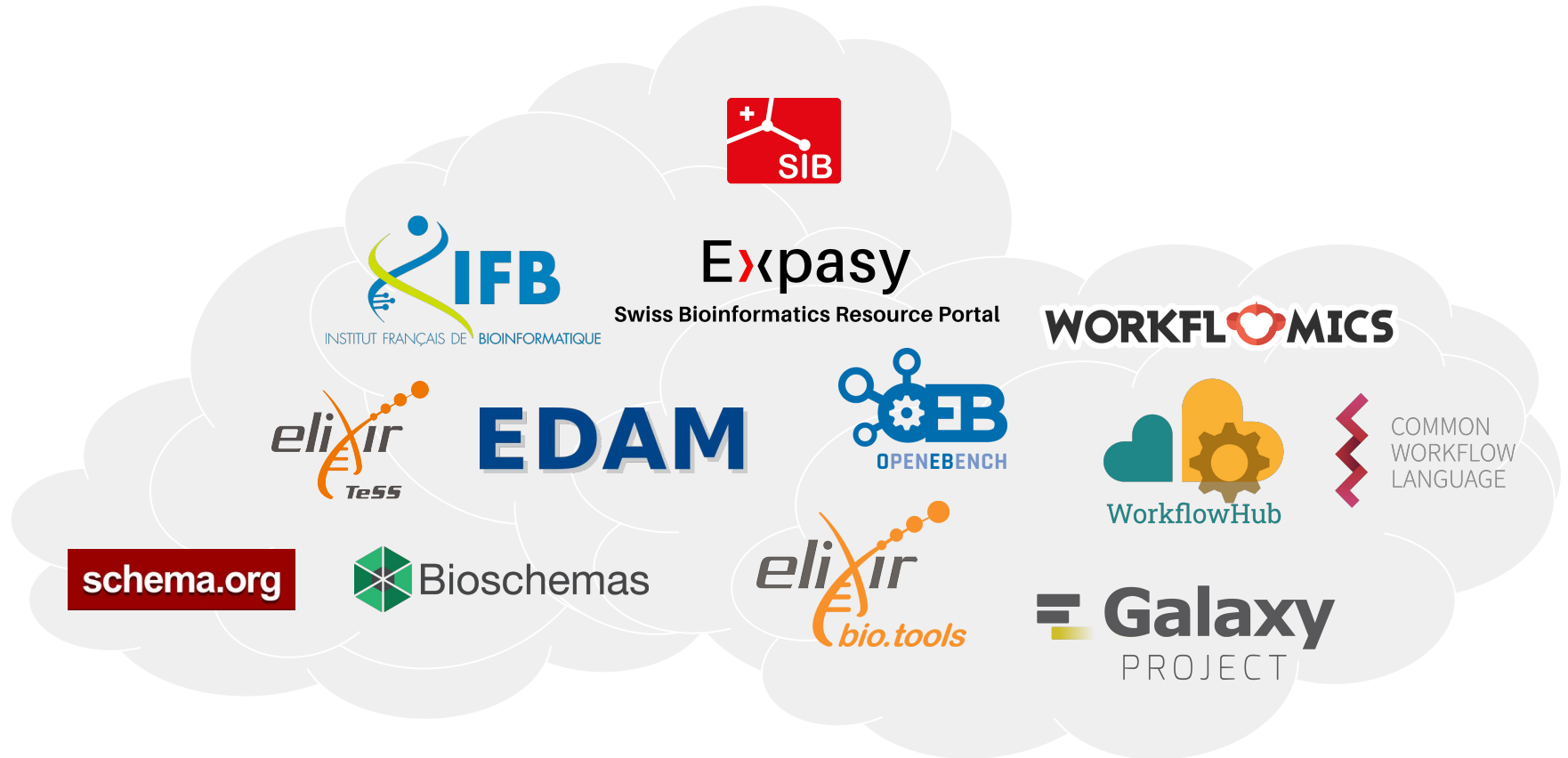- Aggregation, cataloguing, filtering training materials using EDAM keywords
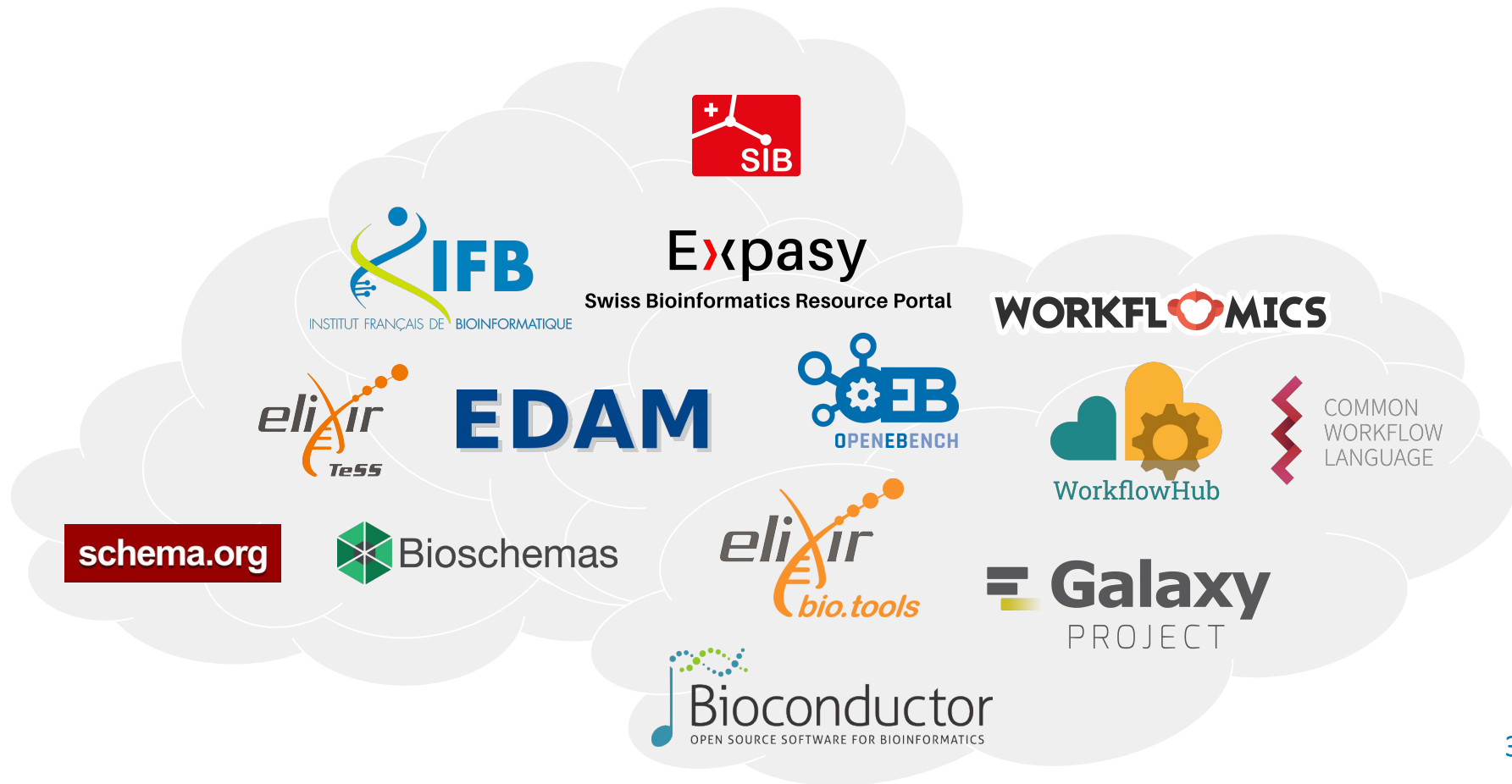
**Bio.tools catalogue**
- 30,000 tools annotated
- Cataloguing, filtering, extracting via an API or a web portal

https://xkcd.com/1406

Using ontology-based metadata has many more applications

- Data **integration** from heterogeneous sources
- Knowledge representation as **graphs or networks**
- Knowledge **discovery**, **predictions**, hypotheses
- **Inferring**, **querying**, **reasoning**
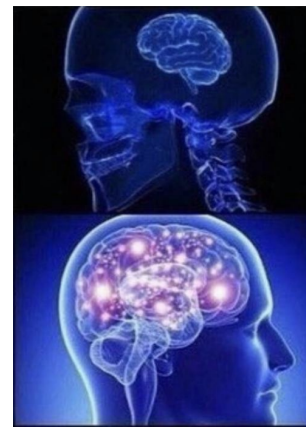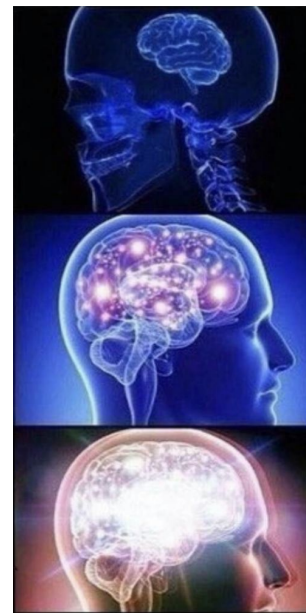
# Ontologies: to infinity and beyond

Using ontology-based metadata has many more applications

- Data **integration** from heterogeneous sources
- Knowledge representation as **graphs or networks**
- Knowledge **discovery**, **predictions**, hypotheses
- **Inferring**, **querying**, **reasoning**

Ontologies provide a basis for **semantic web** and technologies

- Defined by Tim Berners-Lee as a "**web of data**" ≠ "web of documents"
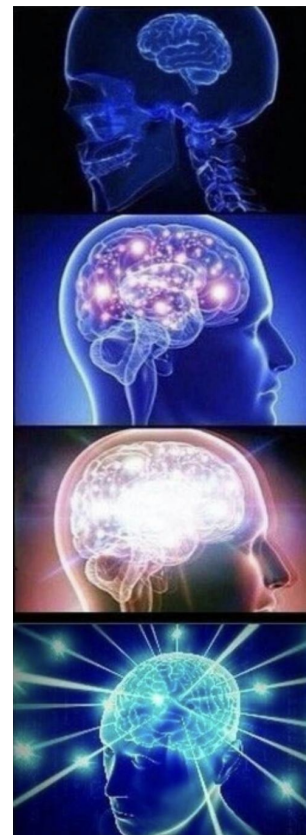- Allows for the *meaning* of data to be **machine-readable**

Using ontology-based metadata has many more applications

- Data **integration** from heterogeneous sources
- Knowledge representation as **graphs or networks**
- Knowledge **discovery**, **predictions**, hypotheses
- **Inferring**, **querying**, **reasoning**

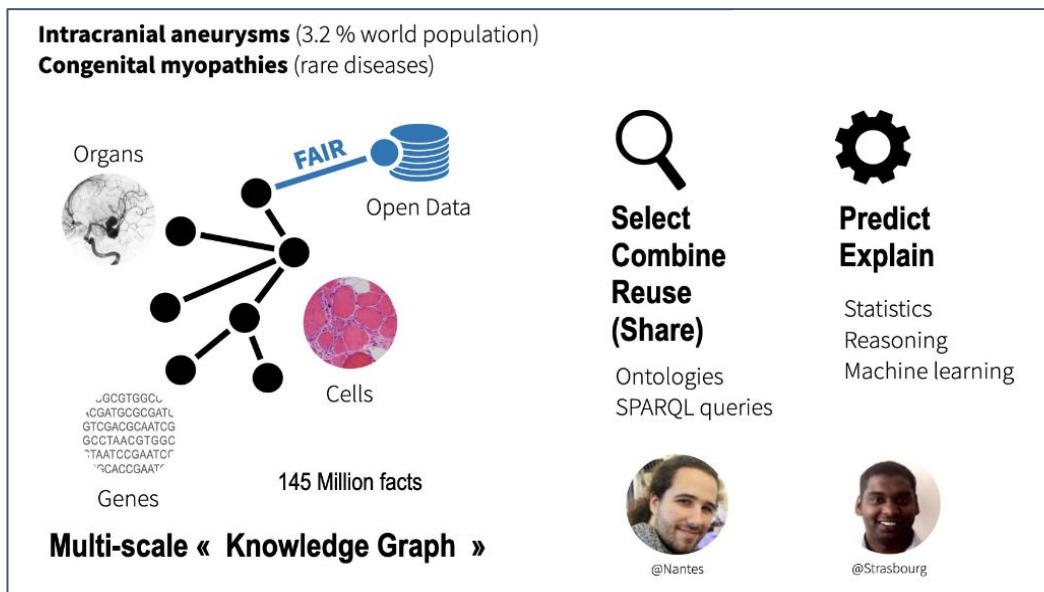Ontologies provide a basis for **semantic web** and technologies

- Defined by Tim Berners-Lee as a "**web of data**" ≠ "web of documents"
- Allows for the *meaning* of data to be **machine-readable**
- Semantic technologies:

  - **Natural Language Processing** (NLP): speech recognition, biocuration
  - **Large Language Models** (LLMs): generative AI tools such as ChatGPT, Gemini
  - **Image analysis**: mass spec, spatial omics, microscopy…

INEX-MED: bridging imaging-omics-clinical data for the study of intracranial aneurysms

- ICAN cohort: 3,400 subjects, 3,000 MRIs, 800 whole genomes
- Association between imaging phenotypes and omics signatures?
- Patients with higher/lower risks of aneurysm rupturing?

Ontologies and semantic metadata

- Annotation of diverse resources: data, articles, software tools, training materials and events…
- Improve data interoperability and accessibility
- Integration and analysis of heterogeneous data

## Ontologies and semantic metadata

- Annotation of diverse resources: data, articles, software tools, training materials and events…
- Improve data interoperability and accessibility
- Integration and analysis of heterogeneous data

## Useful resources

- **OBO foundry**: community development of interoperable ontologies for biological sciences
- **Ontology Lookup Service (OLS)** and **BioPortal**: access and search biomedical ontologies
- **Ontotext** learning resources: what are ontologies and semantic web?
- FAIR-Checker: improve the FAIRness of your web resources (Gaignard et al., 2023)