# Topic 9. Data integration

**Why, how, promise and pitfalls. Standards, ontologies, ID mapping and metadata**

# The Genotype to Phenotype challenge

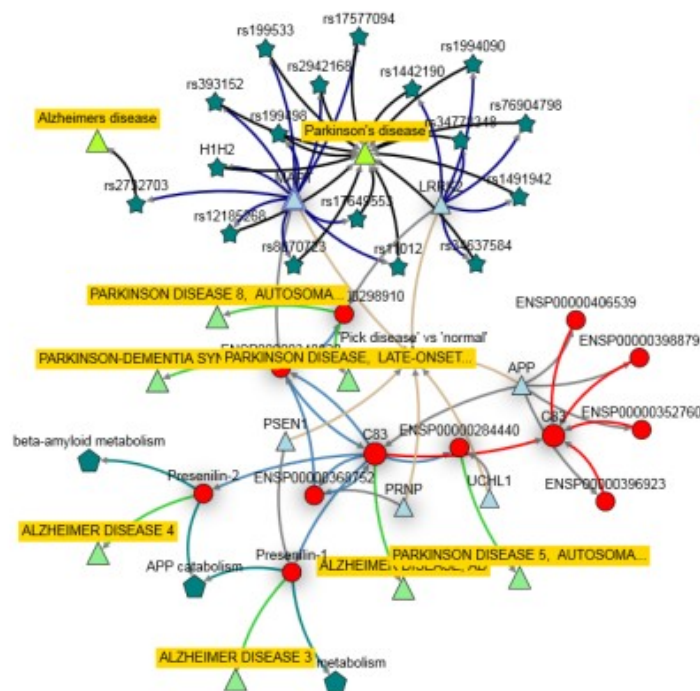1. **You need to understand the limitations/context of the data you start with, to be able to get something useful from a larger integrated dataset**

2. Even when you have the same data to work with, integrating it in different ways can give completely different views of the bigger biological picture
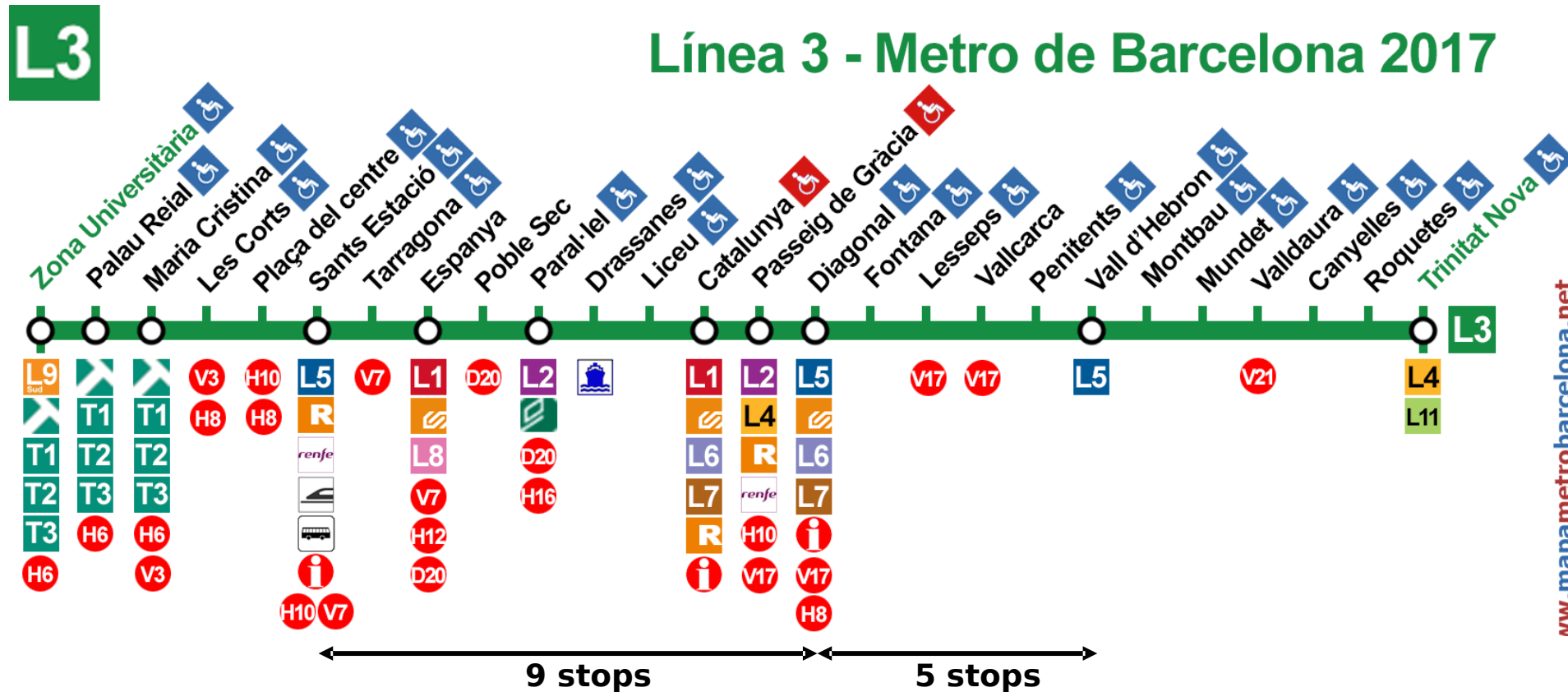
3.  Even if you understand the data and have a good knowledge of the biology, integrating data without the right skills can give a poor result or false impression of reality

# Data integration: some important considerations

4. How (and which) data is represented has a significant effect on our understanding - aiding us in some areas, but can mis-represent the truth in others

4. How (and which) data is represented has a significant effect on our understanding - aiding us in some areas, but can mis-represent the truth in others

4. How (and which) data is represented has a significant effect on our understanding - aiding us in some areas, but can mis-represent the truth in others
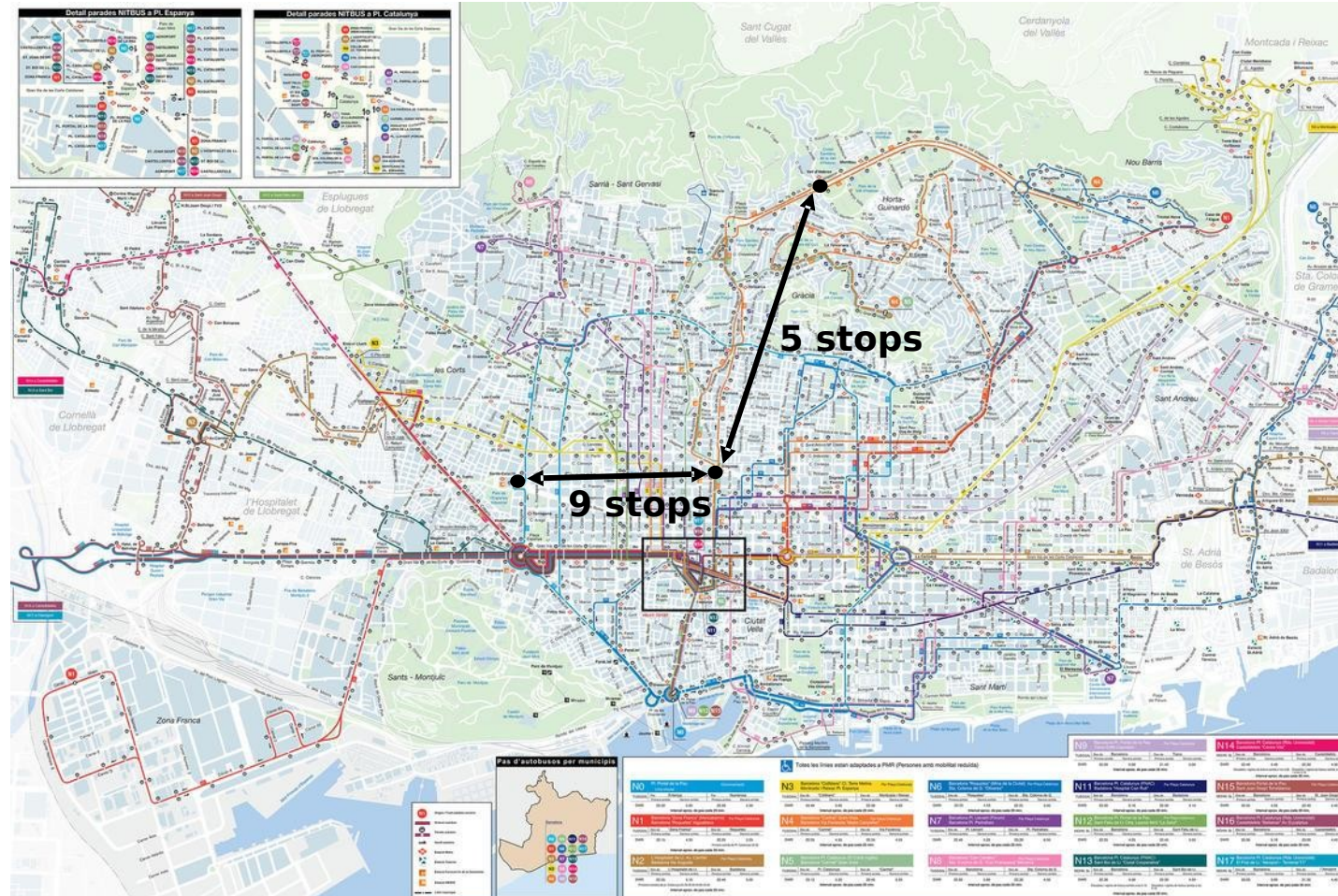
5. We do not (usually) have all the data - unless we recognize this we can be led to false conclusions.



Debunking Lunar Landing Conspiracies with Maxwell and VXGI
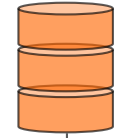https://www.youtube.com/watch?v=O9y_AVYMEUs

Big data, bioinformatics, and particularly data integration have the potential to help us understand biology better.

However, it never provides a definitive answer - it can be used to guide experimentation; and ideas developed from *in silico* approaches should be validated by laboratory (and preferably orthogonal) methods.
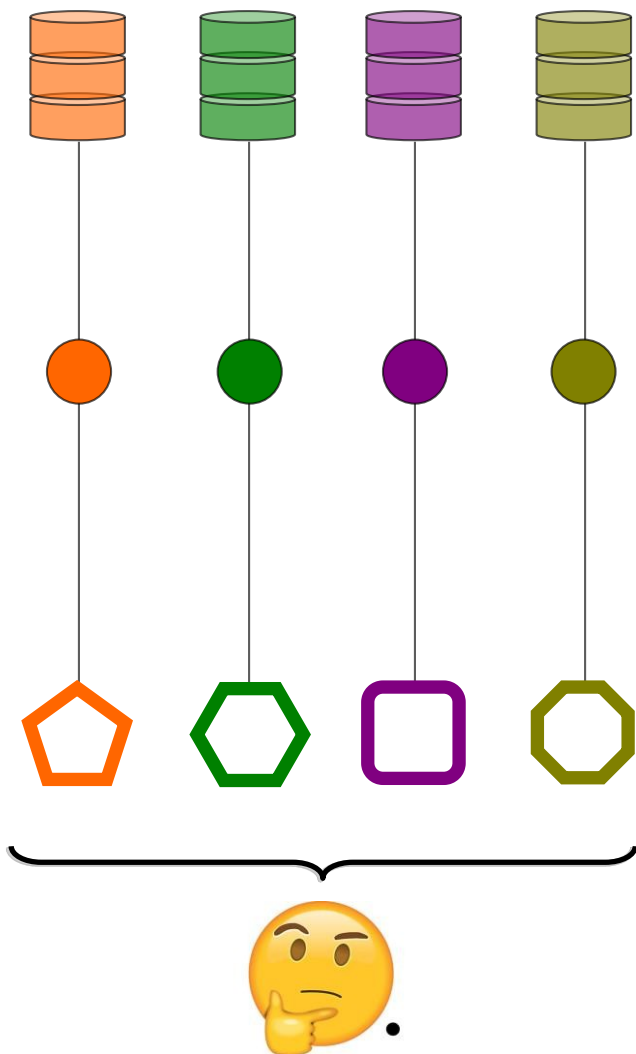
**Ideally**

## Reality



## Many data resources
- **Many to maintain**
- **New appearing**
- **Few have a sustained future**
- **Not easy to find them**

## Different query interfaces

## Variable results
- **Formats**
- **Schemas**
- **Data content**

## Data integration
- **Redundancy**
- **Inconsistency**

FAIRsharing.org
standards, databases, policies
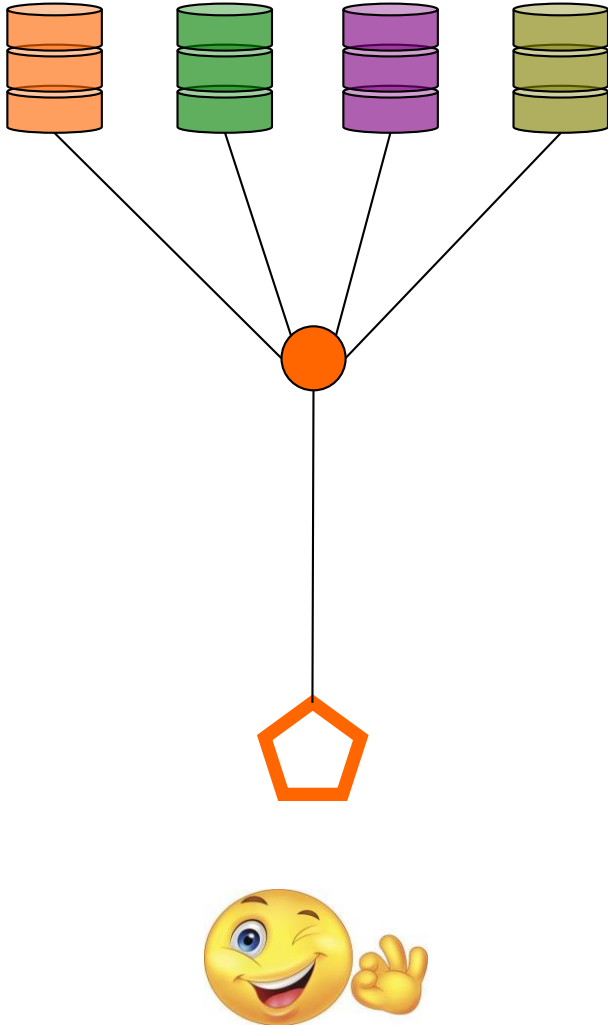
A catalogue of databases, described according to the BioDBcore guidelines, along with the standards used within them; partly compiled with the support of Oxford University Press (NAR Database Issue and DATABASE Journal).

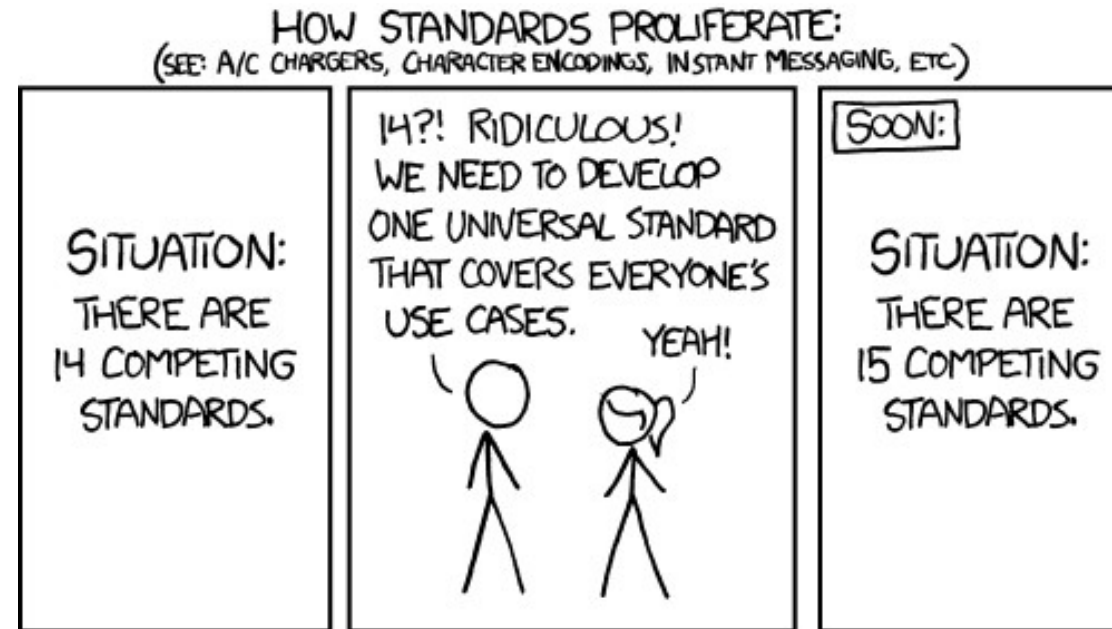**Compiles 2062 databases, 1710 standards & 168 policies***
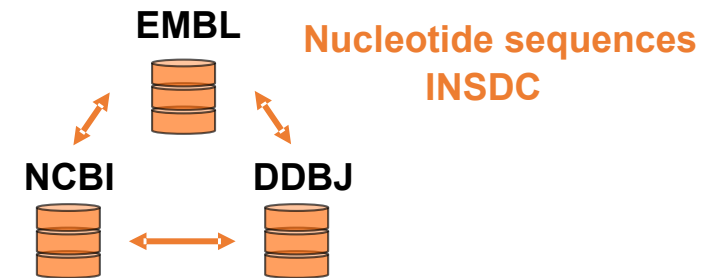
*20/11/2023

## Compromise



**Standards & Ontologies**

# Standards



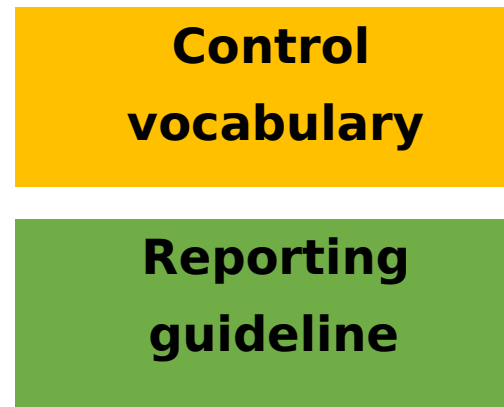**Collaboration among data providers**

- **More data coverage**
- **Less redundancy**
- **Less inconsistency**
- **Better data management: access, exchange, sharing, portability, interoperability, annotation, comparison, verification, representation, integration, reusability**
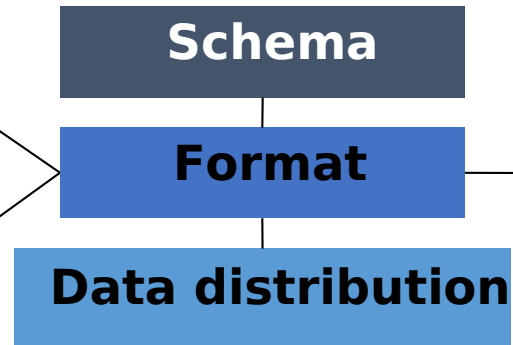
EMBL

Nucleotide sequences
INSDC

NCBI          DDBJ

# Standards

**How to annotate it?**
**Controlled vocabularies (ontologies)**

**How/where to store it?**
**Database, common schemas**

**Control vocabulary**

**Reporting guideline**

**Schema**

**Format**

**Data distribution**

**Identifiers**

**How to identify items?**
**Common identifiers**

**What to store? How to put the data in context?**
**Minimum information guidelines, Metadata**

**How to make it available?**
**Common formats, common query interfaces**

# Standards

**How to annotate it?**
**Controlled vocabularies (ontologies)**

How/where to store it?
Database, common schemas

**Control vocabulary**

Schema

Format

Identifiers

**Reporting guideline**

Data distribution

How to identify items?
Common identifiers

What to store? How to put the data in context?
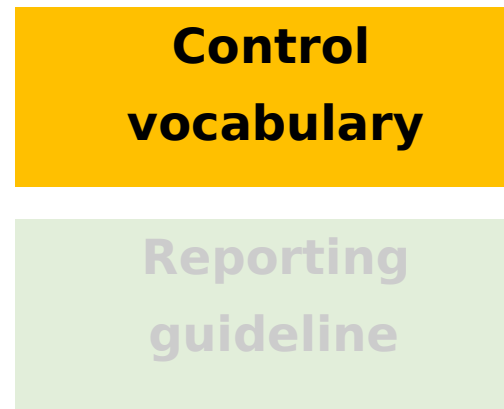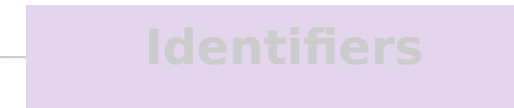Minimum information guidelines, Metadata

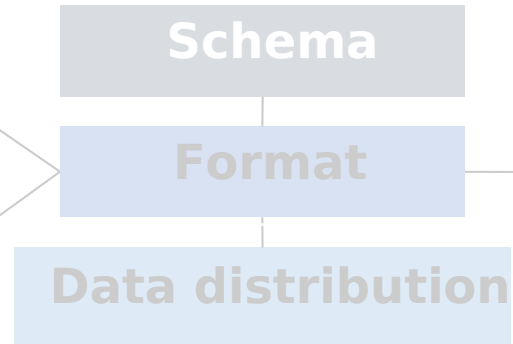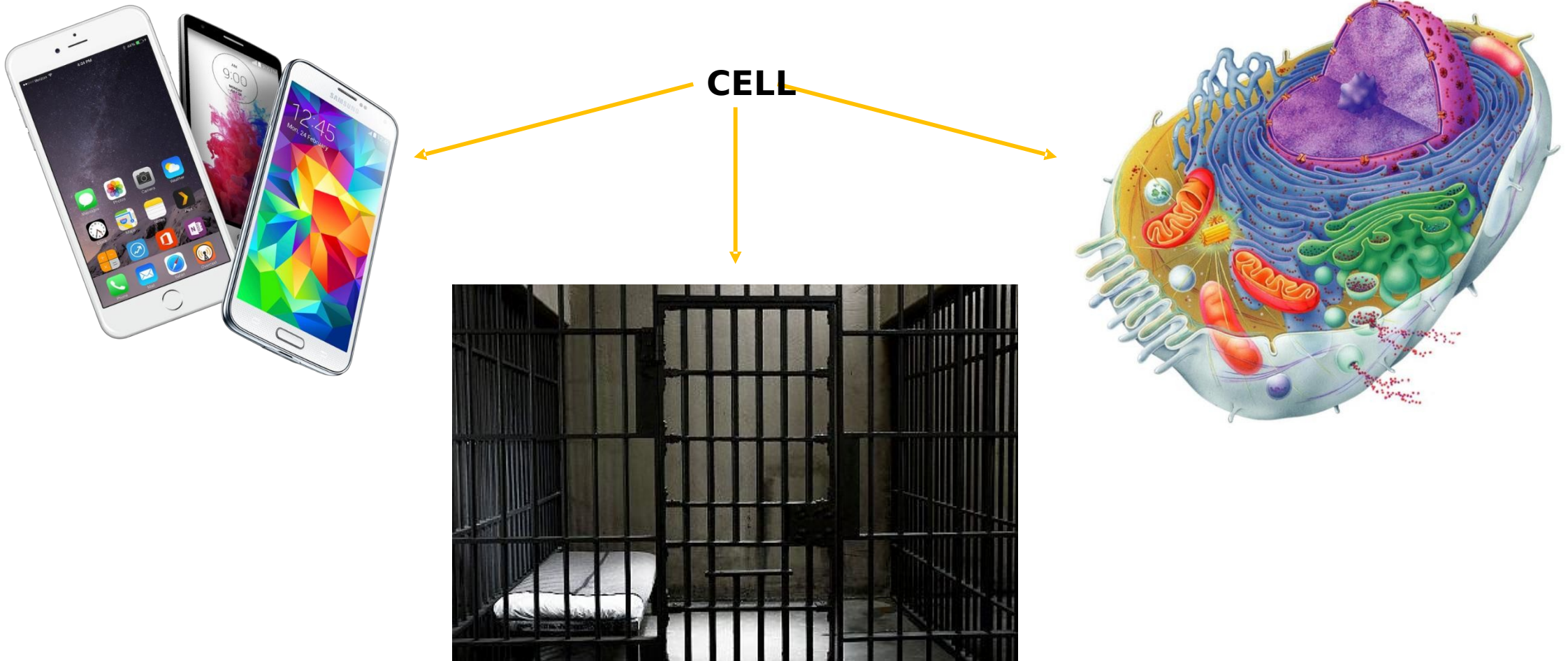How to make it available?
Common formats, common query interfaces

# Controlled vocabularies (ontologies)

- **Inconsistency in natural language: same name, different concepts**

**CELL**

# Controlled vocabularies (ontologies)

- **Controlled vocabularies** provide a way to organize knowledge for subsequent retrieval

- Mandate the use of defined preselected terms

- The Gene Ontology (GO)

- A way to capture biological knowledge for individual gene products in a written and computable form

- A set of concepts and their relationships to each other arranged as a hierarchy
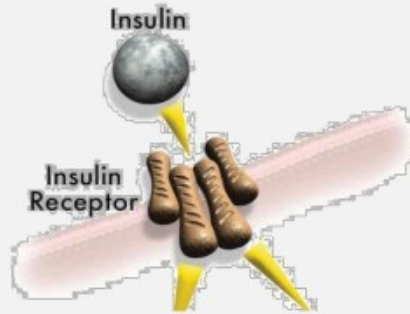


http://www.geneontology.org/

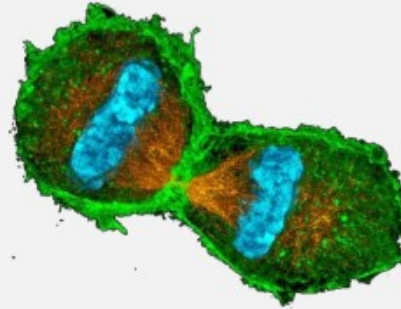# Controlled vocabularies: The Gene Ontology (GO)

## 1. Molecular Function
An elemental activity or task or job

Insulin

Insulin Receptor

- protein kinase activity
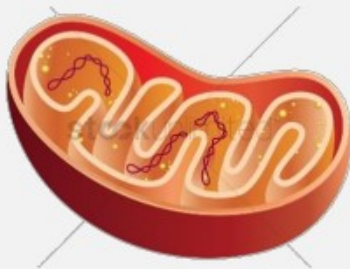- insulin receptor activity

## 2. Biological Process
A commonly recognized series of events

- cell division

## 3. Cellular Component
Where a gene product is located

- mitochondrion
- mitochondrial matrix
- mitochondrial inner membrane

# Controlled vocabularies: The Gene Ontology (GO)



https://www.ebi.ac.uk/QuickGO/

# Controlled vocabularies: The Gene Ontology (GO)



GO:0051961

**Unique identifier**

negative regulation of nervous system development

**Biological Process**

**Term name**

Definition (GO:0051961 GONUTS page)

Any process that stops, prevents, or reduces the frequency, rate or extent of nervous system development, the origin and formation of nervous tissue.

**Definition**

12,371 annotations

## Synonyms

Synonyms are alternative words or phrases closely related in meaning to the term name, with indication of the relationship between the name and synonym given by the synonym scope.

**Synonyms**

| Synonym | Type |
|---|---|
| inhibition of nervous system development | narrow |
| downregulation of nervous system development | exact |
| down-regulation of nervous system development | exact |
| down regulation of nervous system development | exact |

https://www.ebi.ac.uk/QuickGO/term/GO:0051961

# Controlled vocabularies: The Ontology Lookup Servi



https://www.ebi.ac.uk/ols/index

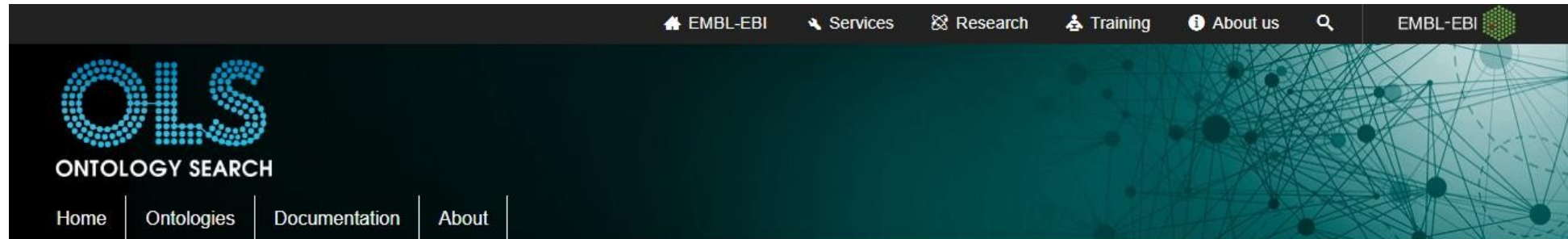# Controlled vocabularies: The Ontology Lookup Service (OLS)

# Standards

How to annotate it?
Controlled vocabularies (ontologies)

How/where to store it?
Database, common schemas

| Control vocabulary |
|---|

| **Reporting guideline** |
|---|

| Schema |
|---|

| Format |
|---|

| Data distribution |
|---|

| Identifiers |
|---|

How to identify items?
Common identifiers

**What to store? How to put the data in context?**
Minimum information guidelines, Metadata

How to make it available?
Common formats, common query interfaces

# Minimum information guidelines

- **The reader can understand, and potentially reproduce, the authors experiments**
- **The data can be moved into a database by a manual curator**
- **The data can be captured, in some form, by text-mining**

*Nature Genetics* **29**, 365 - 371 (2001)
doi:10.1038/ng1201-365

**Minimum information about a microarray experiment (MIAME)—toward standards for microarray data**

nature
biotechnology

The minimum information about a proteomics experiment (MIAPE)

Gigascience. 2: 13.
Published online . doi: 10.1186/2047-217X-2-13                    PMCID: PMC3853013

The role of reporting standards for metabolite annotation and identification in metabolomic studies

# Metadata

- **Metadata is data that describes other data**

- **In the case of experimental data, this is a description of experimental conditions; e.g., patient/environmental details, cell lines, instrumentation**

- Common problems: incomplete and inconsistent metadata

| | | | |
|---|---|---|---|
| 37 year old male | initial phase male | male fetus | six males mixed |
| 600 yr. old male | m | male plant | stallion |
| adult male | make | male, 8 weeks old | steer |
| bull | makle | male, castrated | sterile male |
| castrated male | mal e | male, pooled | strictly male |
| cm | male | males | tetraploide male |
| dioecious male | male (7-2872) | man | type i males |
| diploid male | male (7-3074) | men | type ii males |
| drone | male (m-a) | normale male | virgin male |
| engorged male | male (m-o) | ram | winged and wingless males |
| fertile male | male caucasian | rooster | young male |
| four males mixed | male child | s1 male sterile | |
| individual male | male fertile | sex: male | female* |

# Standards

How to annotate it?
Controlled vocabularies (ontologies)

How/where to store it?
Database, common schemas

| Control vocabulary |
| Reporting guideline |

| Schema |
| **Format** |
| Data distribution |

| Identifiers |

How to identify items?
Common identifiers

What to store? How to put the data in context?
Minimum information guidelines, Metadata

**How to make it available?**
**Common formats, common query interfaces**

# Data formats

- **Merge data from different data sources with minimal effort**
- **Parse the data**
- **Build tools to use the data**

- **PSI-MI XML format (molecular interactions)**

## Overview of Sequence Data Formats.

Zhang H[1].

⊕ **Author information**

**Abstract**

Next-generation sequencing experiment can generate billions of short reads for each sample and processing of the raw reads will add more information. Various file formats have been introduced/developed in order to store and manipulate this information. This chapter presents an overview of the file formats including FASTQ, FASTA, SAM/BAM, GFF/GTF, BED, and VCF that are commonly used in analysis of next-generation sequencing data.

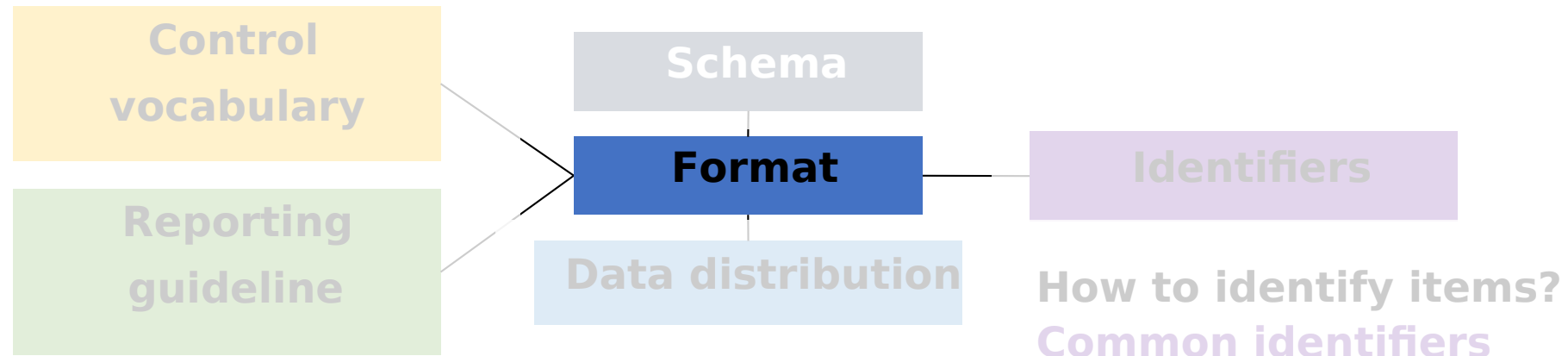**KEYWORDS:** BED; FASTA; FASTQ; GFF/GTF; Next-generation sequencing; SAM/BAM; Sequencing data; Sequencing data file format; VCF

# Standards

How to annotate it?
Controlled vocabularies (ontologies)

How/where to store it?
Database, common schemas

**Control vocabulary**

**Schema**

**Format**

**Identifiers**

**Reporting guideline**

**Data distribution**

**How to identify items?**
**Common identifiers**

What to store? How to put the data in context?
Minimum information guidelines, Metadata

How to make it available?
Common formats, common query interfaces

# Identifiers & ID mapping

- **Identifiers: The use of identifiers allows for unambiguous identifications of molecules or conceptual entities, and their representation in databases**

- **ID mapping: There is a large number of identifiers that aim to represent the same entity**

ENSP00000269305     **P04637**     NP_001119584.1

**P53_HUMAN**

```
MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGP
DEAPRMPEAAPPVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK
SVTCTYSPALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHE
RCSDSDGLAPPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNYMCNS
SCMGGMNRRPILTIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEEENLRKKGEPHHELP
PGSTKRALPNNTSSSPQPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPG
GSRAHSSHLKSKKGQSTSRHKKLMFKTEGPDSD
```

**HGNC:11998**

**NX_P04637**

**Antigen NY-CO-13**     **Cellular tumor antigen p53**     **TP53**
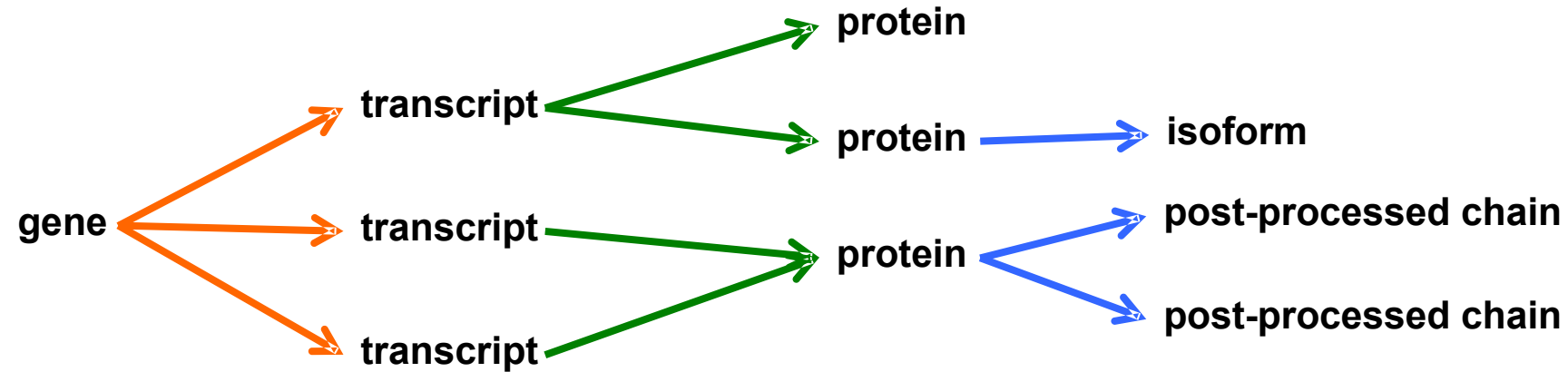
**P04637-1**     uc002gii.2     CHEMBL2221344

# Identifiers & ID mapping

- Most commonly used accessions:

  - **Entrez GeneIDs:** Gene-centered identifier: DNA consensus sequence, no isoform or variants.

  - **UniProt:** Represents proteins, taking into account isoforms. Additional identifiers for variants and post-processed chains.

  - **RefSeq:** Represents sequences of DNA, RNA and proteins.

  - **Ensembl:** Identifiers that represent genes and their different products: gene, gene tree, protein, regulatory feature, transcript, exon and protein family.

  - **International Protein Index:** Proteomics reference database (protein sequences). Now obsoleted, but still used in proteomics.

  - **HUGO gene symbols:** Unique symbols and names for human loci (protein-coding genes, RNA genes and pseudogenes).

  - **Organism centered databases:** FlyBase, TAIR, WormBase, SGD, …

# Identifiers & ID mapping

- Identifiers most of the times do not have a 1:1 relationship: **gene ≠ transcript ≠ protein**



- Identifiers (and sequences!) **disappear** and **get updated**

- Identifiers are **misused**; *e.g.,* gene identifiers used to represent proteins

- ID mapping tools / mapping tables

**UniProt ID mapping**            https://www.uniprot.org/uploadlists/

**Ensembl BioMart**               https://www.ensembl.org/biomart/martview

**DAVID GeneID Conversion Tool**  https://david.ncifcrf.gov/conversion.jsp

**HGNC project**                  http://www.genenames.org/

- Common problems:

  - You are not able to map all the IDs

  - You get 1:n mappings

# Identifiers & ID mapping



http://www.uniprot.org/uploadlists/

# Identifiers & ID mapping

Frameworks for finding and mapping equivalent database identifiers for genes, proteins, and metabolites

- **BridgeDb**

  - Cytoscape app: http://apps.cytoscape.org/apps/bridgedb
  - bridgeDbR Bioconductor package: http://bioconductor.org/packages/devel/bioc/html/BridgeDbR.html
    Web service API: http://webservice.bridgedb.org/

- **BioMart**

  - Cytoscape app: http://apps.cytoscape.org/apps/biomartwebserviceclient
  - biomaRt Bioconductor package: https://bioconductor.org/packages/release/bioc/html/biomaRt.html
  - Web service API: http://www.biomart.org/martservice.html

# Points to remember

- **In order to merge/integrate data it needs to have as many similarities as possible** – the same entity/concept described in different ways will not merge

- **Data transfer/conversion almost always means data loss** – if you want detail, go to the source

- **More resources does not necessarily mean more data** – many databases import from the few that curate

- **If using a resource which imports data** – check the date of the last import

- **Web-services run directly from the most recent release** – may be a better source of data than a data-warehouse, which imports infrequently

- **One advantage of database submission** – data will then be available in standard formats
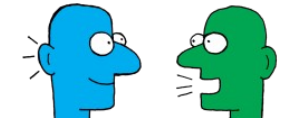
# Some final recommendations

- **Curation, curation, curation**
  - **Contribute to databases that accept submissions**
  - **Use credible data sources**

- **Collaboration, constant feedback, agree on standards**
  - **Use standard data formats**
  - **Use ontology terms to annotate metadata**

- **Openness, reproducibility**
  - **Use open-source software**
  - **Manage and share code**
  - **Transparent data sources and analysis procedures**
  - **Document as much as you can**