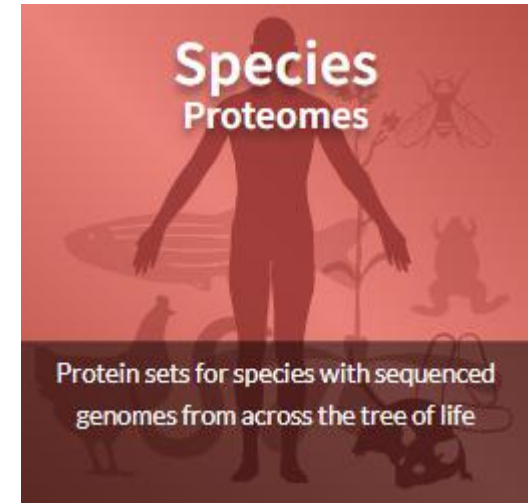# **Proteomes in UniProtKB**

Protein sequence databases and sequence annotation

Streamed from Geneva, 14 October 2022

# Proteomes

- A proteome is the set of proteins thought to be expressed by an organism (completely sequenced genomes).

- A proteome is formed by all UniProtKB entries mapping to a genome assembly from EMBL-ENA or Ensembl or RefSeq.

- Also from VectorBase, WormBase, Parasite

- Viral proteomes are manually checked and verified and periodically added to the database (SIB)



Species Proteomes

Protein sets for species with sequenced genomes from across the tree of life

# UniProtKB Proteomes

# Proteomes

## Proteomes · Amycolatopsis saalfeldensis

### Overview

**Status** 🔖 Reference proteome

**Protein count** 9,167

**Gene count** 9,167 Download one protein sequence per gene (FASTA)

**Proteome ID** UP000198582

**Taxonomy** Amycolatopsis saalfeldensis

**Strain** DSM 44993

**Genome assembly and annotation** GCA_9001105 5.1 from ENA/EMBL ⬏

**Genome representation** Full

**Completeness (CPD)** Close to standard (high value)

**BUSCO** ▫ Single ▫ Duplicated ▫ Fragmented ▫ Missing

n:356 · actinobacteria_phylum_odb10
C:100% (S:98.9% D:1.1%) F:0% M:0%

## Proteomes · Pan paniscus (Pygmy chimpanzee) (Bonobo)

### Overview

**Status** 🔖 Reference proteome

**Protein count** 42,996

**Gene count** 21,212 Download one protein sequence per gene (FASTA)

**Proteome ID** UP000240080

**Taxonomy** Pan paniscus (Pygmy chimpanzee) (Bonobo)

**Genome assembly and annotation** GCA_0001 58655.2 from Ensembl ⬏

**Genome representation** Full

**Completeness (CPD)** Outlier (high value)

**BUSCO** ▫ Single ▫ Duplicated ▫ Fragmented ▫ Missing

n:13780 · primates_odb10
C:95.2% (S:55.3% D:39.9%) F:1.2% M:3.6%

## Proteomes · Juglans regia (English walnut)

### Overview

**Status** 🔖 Reference proteome

**Protein count** 38,355

**Gene count** 32,573 Download one protein sequence per gene (FASTA)

**Proteome ID** UP000235220

**Taxonomy** Juglans regia (English walnut)

**Strain** cv. Chandler

**Genome assembly and annotation** GCF_001411 55.2 from Refseq ⬏

**Genome representation** Full

**Completeness (CPD)** Outlier (high value)

**BUSCO** ▫ Single ▫ Duplicated ▫ Fragmented ▫ Missing

n:2326 · eudicots_odb10
C:99.4% (S:60.8% D:38.6%) F:0.1% M:0.5%

# Proteomes

- The Proteomes webpage has been redesigned to enable users to view full details of their proteome(s) of interest in a single table view.

# CPD: Completeness in terms of number

We display the results of the 'Complete Proteome Detector' (CPD), an in-house algorithm which statistically evaluates the completeness and quality of each proteome by directly comparing it to those of a group of at least 3 closely taxonomically-related species. The CPD classifies each proteome as either 'standard', 'close to standard' or an 'outlier', according to protein count.

# BUSCO: Completeness in term of gene content

The BUSCO score includes percentages of complete single-copy genes, complete duplicated genes, fragmented and missing genes, as well as the total number of orthologous clusters used in the BUSCO assessment.

# BUSCO: Completeness in term of gene content

BUSCO v3 identifies complete, duplicated, fragmented, and potentially missing genes by comparison to a defined set of near-universal single copy orthologs. As a result of our adoption of the BUSCO algorithm, a number of incomplete and poor quality proteomes were identified and removed from UniProtKB. This, coupled with a concomitant clean-up of metagenomic data, resulted in a drop in the number of bacterial sequences in 2018.

# BUSCO: Completeness in term of gene content

| Organism ▲ | Organism ID | Protein count ▲ | BUSCO |
|---|---|---|---|
| | | | ☐ Single ☐ Duplicated ☐ Fragmented ☐ Missing ▲ |
| **Escherichia coli (strain K12) (K12 / MG1655 / ATCC 47076)** | 83333 | 4,448 | n:440 · enterobacterales_odb10 C:100% (S:99.3% D:0.7%) F:0% M:0% |
| **Saccharomyces cerevisiae (strain ATCC 204508 / S288c) (Baker's yeast) (ATCC 204508 / S288c)** | 559292 | 6,060 | n:2137 · saccharomycetes_odb10 C:99.6% (S:97.4% D:2.2%) F:0.1% M:0.3% |
| **Homo sapiens (Human)** | 9606 | 80,581 | n:13780 · primates_odb10 C:99.5% (S:38.5% D:61.1%) F:0% M:0.4% |
| **Mus musculus (Mouse) (C57BL/6J)** | 10090 | 55,311 | n:13798 · glires_odb10 C:99.7% (S:51.2% D:48.5%) F:0% M:0.2% |
| **Drosophila melanogaster (Fruit fly) (Berkeley)** | 7227 | 22,071 | n:3285 · diptera_odb10 C:100% (S:41.8% D:58.2%) F:0% M:0% |
| **Bacillus subtilis (strain 168) (168)** | 224308 | 4,260 | n:450 · bacillales_odb10 C:99.3% (S:99.1% D:0.2%) F:0.2% M:0.4% |

# UniProtKB – scope

We additionally provide the assessment of the genome assembly status imported from the source of the genome assembly and annotation (e.g. Ensembl or RefSeq).

# Proteomes

- UniProt release 2022_04 contains over 230 million sequence records, with almost 470,000 proteomes originating from completely sequenced viral, bacterial, archaeal and eukaryotic genomes available through the UniProtKB Proteomes
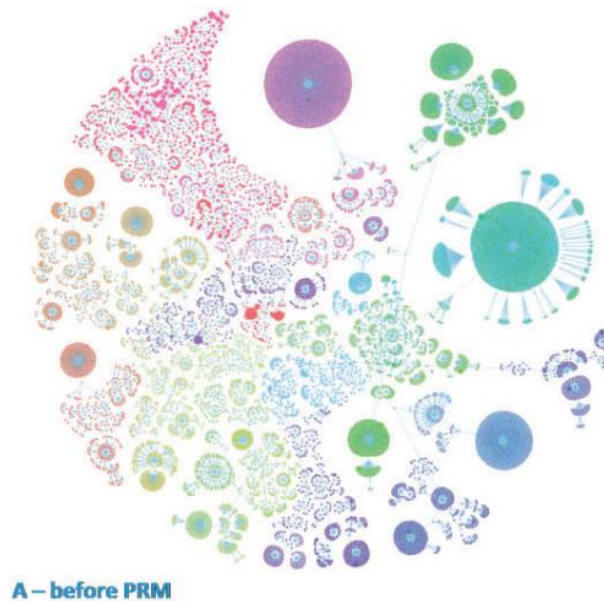
# Too many…

The UniProt Knowledgebase (UniProtKB) has witnessed exponential growth in the last few years with a two-fold increase in the number of entries in 2014. This follows the vastly increased submission of multiple genomes for the same or closely related organisms. This increase was accompanied by a high level of redundancy in unreviewed UniProtKB (TrEMBL), and many sequences were over-represented in the database.
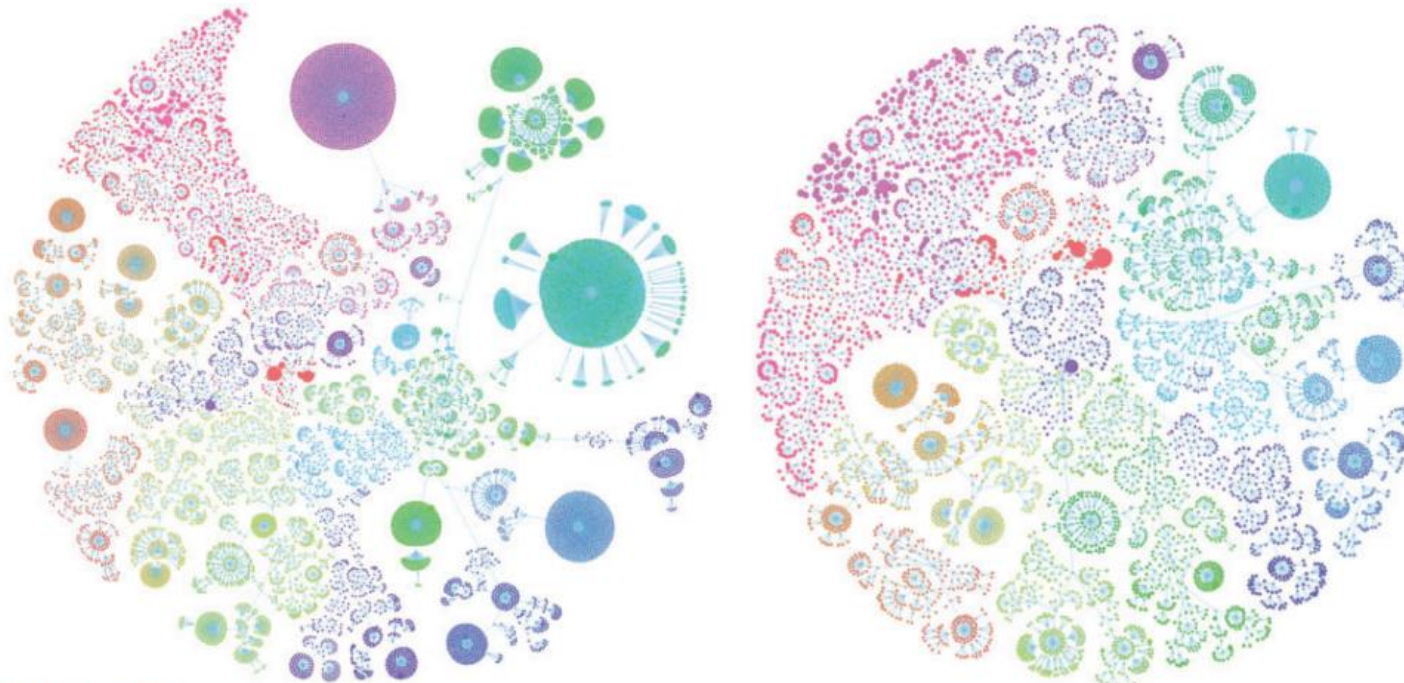
# How to make useful proteomes

This was especially true for bacterial species where different strains of the same species have been sequenced and submitted (e.g. 1,692 strains of Mycobacterium tuberculosis, corresponding to 5.97 million entries). High redundancy led to an increase in the size of UniProtKB, and thus to the amount of data to be processed internally and by our users, but also to repetitive results in BLAST searches for over-represented sequences.



A — before PRM

# Redundant proteomes

A redundancy removal process was first introduced in 2015. This process identifies and removes almost identical proteomes of the same species before their inclusion in UniProtKB and places their sequences in UniParc. Currently this process has removed ~38% of all complete proteomes (~241 million proteins) from UniProtKB.



A – before PRM

B – after PRM

# Redundant proteomes

The redundant proteome sequences are available through UniParc to researchers and stable proteome identifiers (of the form UPXXXXXXXXX, where Xs are integers) are maintained for each redundant proteome to ensure findability.
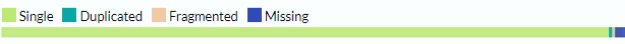
labelled with a specific icon



https://www.uniprot.org/proteomes/UP000008521

# Redundant proteomes

# Redundant proteomes

# Redundant proteomes

# Redundant proteomes

- Proteome's redundancy is only applied to bacteria and fungi for the moment.

- See FAQ:

  https://www.uniprot.org/help/redundancy

- https://www.uniprot.org/help/proteome_redundancy

- http://insideuniprot.blogspot.ch/2015_05_01_archive.html

- If you need 'to protect' a proteome for a good reason, please contact us !

- help@uniprot.org

# Still too many to be useful…

- 'Redundant' proteomes are excluded from UniProtKB (since 2015) (bacteria & fungi).

- Corresponding protein sequences are only available in UniParc …



Number of entries in UniProtKB/TrEMBL

https://www.uniprot.org/help/proteome_redundancy

# Reference proteomes

For the remaining proteomes we provide a Reference Proteome set (~9% of total proteomes) selected by the research community and supplemented with selected proteomes from a computational clustering to provide the best annotated proteome in their cluster.



ALLIANCE of GENOME RESOURCES

QfO Quest for Orthologs



**Compute pair-wise co-membership value (X) in UniRef50 for all proteomes**

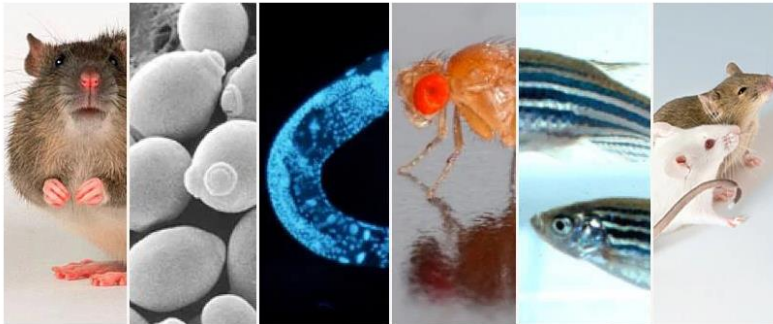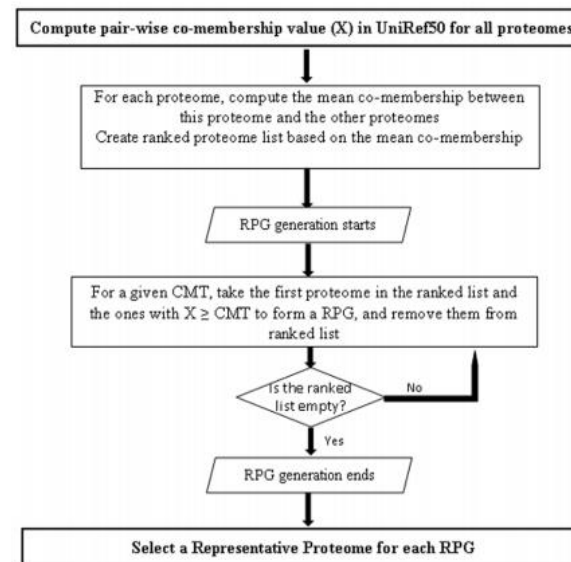For each proteome, compute the mean co-membership between this proteome and the other proteomes
Create ranked proteome list based on the mean co-membership

RPG generation starts

For a given CMT, take the first proteome in the ranked list and the ones with X ≥ CMT to form a RPG, and remove them from ranked list

Is the ranked list empty? — No

Yes

RPG generation ends

**Select a Representative Proteome for each RPG**

**Figure 1. Flow chart of the method used to select Representative Proteomes.** For details please see [materials and methods] section.
doi:10.1371/journal.pone.0018910.g001

**Citation:** Chen C, Natale DA, Finn RD, Huang H, Zhang J, et al. (2011) Representative Proteomes: A Stable, Scalable and Unbiased Proteome Set for Sequence Analysis and Functional Annotation. PLoS ONE 6(4): e18910. doi:10.1371/journal.pone.0018910

# Reference proteomes

labelled with a specific icon 

## Proteomes · Bacillus subtilis (strain 168)

### Overview

| | |
|---|---|
| Status 🔖 Reference proteome | Genome representation Full |
| Protein count[i] 4,260 | Pan proteome This proteome is part of the Bacillus subtilis (strain 168) pan proteome (FASTA) |
| Gene count 4,260 Download one protein sequence per gene (FASTA) | Completeness (CPD)[i] Standard |
| Proteome ID[i] UP000001570 | BUSCO[i] ▇ Single ▇ Duplicated ▇ Fragmented ▇ Missing |
| Taxonomy Bacillus subtilis (strain 168) | |
| Strain 168 | n:450 · bacillales_odb10 |
| Genome assembly and annotation[i] GCA_000009045.1 from ENA/EMBL ⧉ | C:99.3% (S:99.1% D:0.2%) F:0.2% M:0.4% |

Gram-positive harmless bacterium found in soil. Its cell envelope consists of a thick peptidoglycan wall and and cell membrane. It is capable of producing endospores resistant to adverse environmental conditions such as heat and desiccation and is widely used for the production of enzymes and specialty chemicals.

### Components

⬇ Download   View proteins

| ☐ Component name | Genome accession(s) | Protein count |
|---|---|---|
| ☐ Chromosome | AL009126 ⧉ | 4,260 |

### Publications

The complete genome sequence of the Gram-positive bacterium Bacillus subtilis.

Kunst F., Ogasawara N., Moszer I., Albertini A.M., Alloni G., Azevedo V., Bertero M.G., Bessieres P., Bolotin A. [...] , Danchin A.

PubMed ⧉
Europe PMC ⧉
Nature 390:249-256 (1997) ⧉



UniProtKB contains over 22,000 reference

proteomes from all branches of the ToL at

https://www.uniprot.org/proteomes/

# UniProtKB – scope



**E.coli (3 reference proteomes; 2 chromosomes, 1 plasmid) & E.coli phages (3 reference proteomes)**

# Exercise



- @UniProtKB: look for the Candida albicans proteomes

- How many proteomes ?

- What is the Proteome ID of the Candida albicans Reference Proteome ?


- To which strain does it correspond ?

- How many Swiss-Prot records ?

- How many TrEMBL records ?

# Exercise

# Exercise