# Xenopus and Zebrafish Annotation in the UniProt Knowledgebase (UniProtKB)

EMBL-EBI

European Bioinformatics Institute is an Outstation of the European Molecular Biology Laboratory.

UniProt

European Bioinformatics Institute

**Rebecca E. Foulger[1] and UniProt Consortium[1,2,3]**

[1]European Bioinformatics Institute, Wellcome Trust Genome Campus, Cambridge, UK
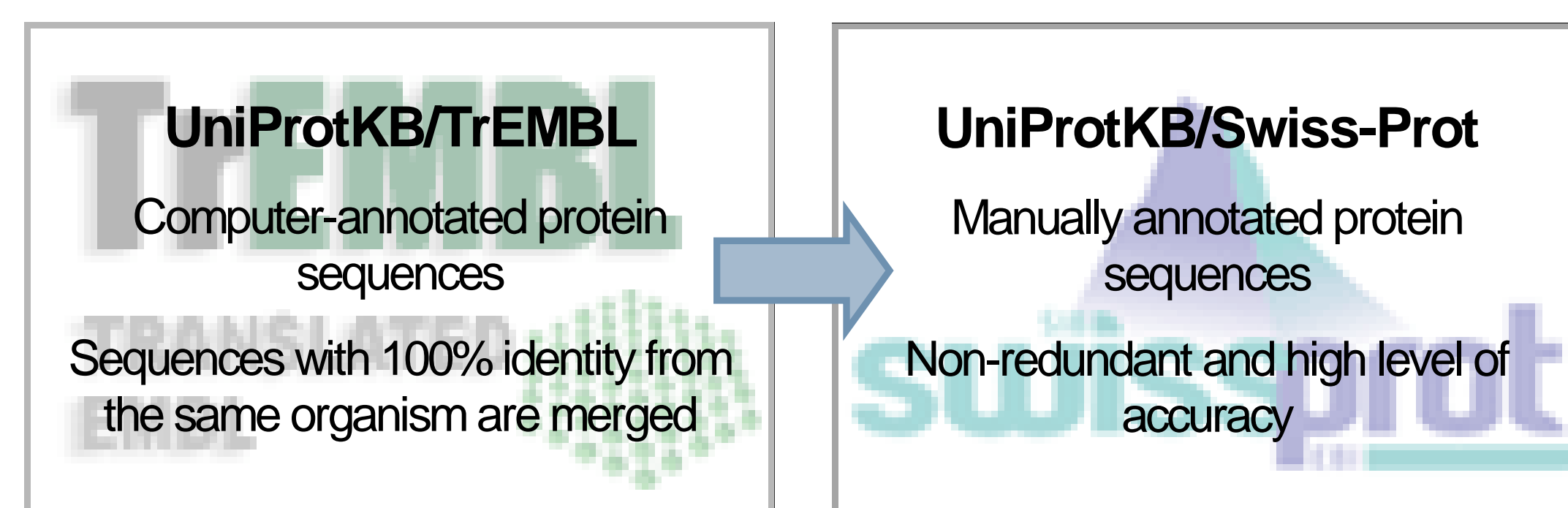[2]Swiss Institute of Bioinformatics, Geneva, Switzerland
[3]Protein Information Resource, Georgetown University, Washington DC, USA

## Introduction to UniProt

• UniProt (Universal Protein Resource: http://www.uniprot.org) provides a central resource of protein sequences with functional annotation. The UniProt Consortium is a collaboration between the Swiss Institute of Bioinformatics (SIB), the European Bioinformatics Institute (EBI) and the Protein Information Resource (PIR).

• UniProt Knowledgebase (UniProtKB) contains the manually annotated UniProtKB/Swiss-Prot section and the automatically annotated UniProtKB/TrEMBL section.

**UniProtKB/TrEMBL**
Computer-annotated protein sequences
Sequences with 100% identity from the same organism are merged

**UniProtKB/Swiss-Prot**
Manually annotated protein sequences
Non-redundant and high level of accuracy

## Why frogs and fish?

• *Xenopus laevis* (African clawed frog) and *Danio rerio* (zebrafish) are both model organisms in the laboratory, as they are easy to accommodate and can produce large numbers of externally-developing embryos that can be easily viewed and manipulated.

• With a generation time of 2 to 3 months, *D. rerio* is amenable to genetic analysis, and the Sanger Institute began sequencing the zebrafish genome in 2001.The extensive similarity between the zebrafish and human genomes means many human developmental and disease genes have counterparts in the zebrafish.

• *X. laevis* has a much longer generation time of 1-2 years, and its tetraploid genome creates difficulties for genetic analysis. The related species *Xenopus tropicalis* (Western clawed frog) is diploid with a relatively small genome

|  | *X. laevis* | *X. tropicalis* |
|---|---|---|
| Ploidy | Allotetraploid | Diploid |
| Genome Size | $3.1 \times 10^9$ bp | $1.7 \times 10^9$ bp |
| Adult size | 10 cm | 4-5 cm |
| Generation time | 1-2 years | 4 months |
| Egg size | 1-1.3 mm | 0.7 – 0.8 mm |

Xenopus photos by Enrique Amaya

and a shorter generation time, so is much better suited for genetic studies. Sequencing of the *X. tropicalis* genome has begun and is supported by cDNA and EST sequencing projects, making it an ideal time to focus on Xenopus curation.

## The Challenge of Duplicated Genomes

• Both *Danio rerio* and *Xenopus laevis* have undergone whole genome duplications, resulting in duplicated genes for multiple loci.

• Protein sequences encoded by duplicated *X. laevis* and *D. rerio* genes are annotated as separate UniProtKB/Swiss-Prot entries.

| Accession | Entry name | Status | Protein names | Gene names | Organism |
|---|---|---|---|---|---|
| Q9W707 | FXF1B_XENLA | ★ | Forkhead box protein F1-B (FoxF1-B) (FoxF1b) (Fork head domain-related protein 13) (xFD-13') | foxf1-B | Xenopus laevis (African clawed frog) |
| Q9W706 | FXF1A_XENLA | ★ | Forkhead box protein F1-A (FoxF1-A) (FoxF1a) (Fork head domain-related protein 13) (xFD-13) | foxf1-A | Xenopus laevis (African clawed frog) |

A and B nomenclature is used to distinguish the duplicated sequences

Synonyms from the literature are recorded

| Accession | Entry name | Status | Protein names | Gene names | Organism |
|---|---|---|---|---|---|
| Q6PC54 | SN25B_DANRE | ★ | Synaptosomal-associated protein 25-B (SNAP-B) (Synaptosome-associated protein 25.2) (SNAP-25.2) | snap25b (Snap) (snap25.2) | Danio rerio (Zebrafish) (Brachydanio rerio) |
| Q5TZ66 | SN25A_DANRE | ★ | Synaptosomal-associated protein 25-A (SNAP-25A) (Synaptosome-associated protein 25.1) (SNAP-25.1) | snap25a (Snap) (snap25.1) (si:dkeyp-8f4.6) | Danio rerio (Zebrafish) (Brachydanio rerio) |

• For *D. rerio*, the UniProtKB A/B nomenclature matches ZFIN. For *X. laevis*, the A/B naming is taken from the literature where available, or assigned by a curator when the literature doesn't specify.

## Curation of a *X. tropicalis* TrEMBL Entry

For both Xenopus and zebrafish, protein and gene nomenclature is generally propagated from the human ortholog, with species-specific names added as synonyms

Identifiers from large-scale projects are recorded

• We use a protein-by-protein approach to curation.

• PubMed searches, requests from users, cross-reference updates and sequence revisions are all used to identify which proteins are of priority to curate.

• Entries are curated using our flat-file editor, CRiSP.

Functional annotation is attributed to individual papers

References include both large scale projects (including XGC/ZGC and Sanger cDNA projects), and sequences and papers from individual labs. A PubMed search is performed to find additional characterization papers. All available papers are manually curated and added to a TrEMBL entry

Evidence tags are used throughout the entry to show the source of the annotation

A wide range of functional data is taken from papers. Data is also propagated from similar UniProtKB/Swiss-Prot entries, and predicted by sequence analysis tools

Isoforms are annotated based on all available sequences and literature

GO annotation is added manually using the Protein2GO editor. IEA annotations are generated from automatic methods, including KW:GOterm mapping

Most cross-references are added automatically. Some (E.g. PROSITE, EMBL/DDBJ/GenBank, and MOD X-ref lines) can be added or modified by a curator

Reciprocal links to Xenbase gene pages were introduced at the end of 2008

Xenbase

Computer-generated keywords (KWs) are verified by a curator, and further keywords are added manually

Key residues, domains and motifs are annotated based on sequence analysis tools (such as Anabelle), similarity to existing UniProtKB/Swiss-Prot entries, and information in the literature

Evidence tags show the source of annotation

Internal comments are used to explain annotation and point to related annotated entries

The longest isoform is usually displayed. Where conflicts exist, for frogs and fish the consensus sequence is normally displayed

European Bioinformatics Institute (EMBL-EBI)
Swiss Institute of Bioinformatics (SIB)
Protein Information Resource (PIR)

Email: help@uniprot.org

URL: www.uniprot.org

Nature Precedings : doi:10.1038/npre.2009.3953.1 : Posted 23 Apr 2009