Alistair MacDougall[1], UniProt Consortium[1-4], Ensembl[1]

1 EMBL-EBI, Cambridge, UK; 2 Swiss Institute of Bioinformatics, Centre Medical Universitaire, Geneva, Switzerland
3 Protein Information Resource, Georgetown University Medical Center, Washington, USA; 4 Protein Information Resource, University of Delaware, Newark, USA
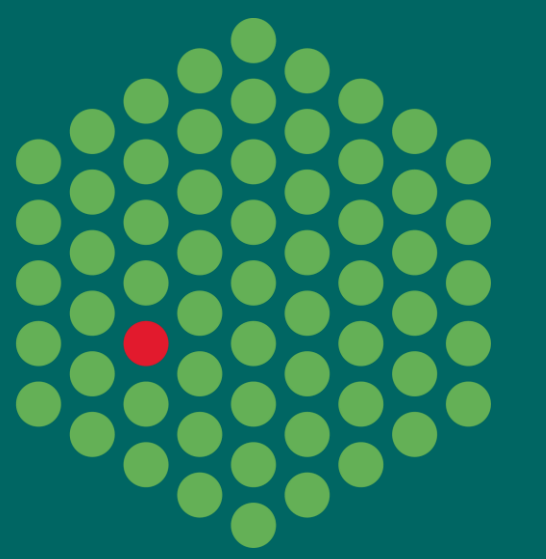
# Mapping Genes to Proteins in UniProtKB

**EMBL-EBI**

## UniProt Knowledgebase

The UniProt Knowledgebase (UniProtKB) endeavours to provide the scientific community with a comprehensive catalogue of protein sequence and functional information.

To make effective use of these data for genome studies, it is essential to have accurate mapping from gene to protein sequence, and in the reverse direction from protein to gene.

## Mapping UniProtKB:Ensembl

Ensembl and UniProtKB are closely aligned in the production of the annotated genome. UniProtKB data is used as one of the inputs for the Ensembl gene build process.

After Ensembl completes a gene build, cross-referencing in UniProtKB is carried out on the basis of exact matching between UniProtKB isoforms and the Ensembl transcripts.

## Current Status of Mapping

|  | Human | Mouse |
|---|---|---|
| **Ensembl (Release 93)** |  |  |
| Genes (protein coding) | 20,376 | 22,628 |
| Transcripts (protein coding) | 91,308 | 45,664 |
| Transcripts with no UniProt Xref | 922 | 28,639 |
| **UniProtKB (Release 2018_08)** |  |  |
| UniProtKB Entries in Proteome | 20,381 | 16,987 |
| Entries with Ensembl Transcript Xref | 19,374 | 15,065 |
| Entries with no Ensembl Xref | 1,007 | 1,922 |

Cross-referencing in UniProtKB to human and mouse transcripts is extensive, but not complete. Where entries, or isoforms within entries, appear to be missing in UniProtKB, this is almost always because the relevant protein sequences are held in the unreviewed section of UniProtKB, and have not yet been incorporated into the manually reviewed entries where they belong.

## Mapping Improvement

Discussion between UniProt and Ensembl on how to improve the extent and reliability of mappings has led to the creation of a common workspace in which the mapping between UniProtKB isoforms and Ensembl transcripts can be reviewed, agreed and updated.

## Planned Outcomes

An important outcome of this joint UniProt : Ensembl collaboration will be a publicly available mapping database, which will provide protein : transcript mappings that remain stable through the different release cycles of UniProtKB and Ensembl.

The database will also provide representative sets of proteins and transcripts that are reliably mapped for different genomes, starting with human and extending to mouse and other organisms.

The new database will allow UniProtKB to include Ensembl cross-referencing with the proteome datasets that are currently provided.

## Sequence Position Cross-referencing

Exact matching of sequence numbering plays a key role in making it possible to easily build datasets which combine genomic and protein data.