# UniProtKB: manual and automated protein sequence annotation pipelines

# WELCOME !

**http://education.expasy.org/cours/SIB_UniProtKB_2018/**

**Basel, November 2018**

Elisabeth.Gasteiger@sib.swiss

Marie-Claude.Blatter@sib.swiss

**Swiss-Prot group
SIB Swiss Institute of Bioinformatics**

Elisabeth.Gasteiger@sib.swiss

Marie-Claude.Blatter@sib.swiss

The **goals of this course (theory + practicals)** are to give some basic theoretical and practical knowledge on protein sequence databases with a **focus on UniProtKB, on the different manual and automated annotation pipelines and on the optimum use of UniProt.**

We will discuss questions such as:

- Where do the protein sequences come from?

- What are the differences between the major protein sequence databases?

- What are the manual and automated gene / protein annotation pipelines?

- How to assess protein sequence accuracy and annotation quality?

- How to extract biological knowledge from a Blast result, a gene list or a multiple alignment?

- This afternoon: focus on automated annotation pipelines

UniProt

**Protein sequences**

**Fasta:**

• **SwissProt** is a high quality, curated protein database. Sequences are non-redundant, rather than non-identical, so you may get fewer matches for an MS/MS search than you would from a comprehensive database, such as NCBIprot. SwissProt is ideal for peptide mass fingerprint searches and MS/MS searches of well characterised organisms where it isn't essential to match every single spectrum.

• **NCBIprot** is a comprehensive, non-identical protein database maintained by NCBI for use with the search tools BLAST and Entrez. The entries have been compiled from GenBank CDS translations, PIR, SWISS-PROT, PRF, and PDB.

• **EMBL EST divisions** contain "single-pass" cDNA sequences, or Expressed Sequence Tags, from a number of organisms. During a Mascot search, the nucleic acid sequences are translated in all six reading frames. There are 10 divisions: Environmental_EST, Fungi_EST, Human_EST, Invertebrates_EST, Mammals_EST, Mus_EST, Plants_EST, Prokaryotes_EST, Rodents_EST, and Vertebrates_EST.

• **contaminants** is a database of common contaminants compiled by Max Planck Institute of Biochemistry, Martinsried

• **cRAP** is a database of common contaminants compiled by the Global Proteome Machine Organization

IPI
RefSeq
GenPept
…

Where do the protein sequences come from ?
What's about the sequence accuracy ?

UniProt

SIB
Swiss Institute of
Bioinformatics

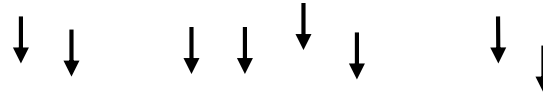**Table I.** *The identification of protein spots which were found in both MRC-5 and A549 cells.*

| Spot no. | Description | MW (Daltons) | pI | Phospho-staining | Phospho-peptides found | Phospho-proteins (SWISS-Prot) | Protein expression in A549 | Protein expression in MRC-5 | Protein function |
|---|---|---|---|---|---|---|---|---|---|
| D1 | 75 kDa Glucose-regulated protein | 73,734 | 5.9 | – | – | √ | 3.0 | –3.0 | Chaperone/stress response |
| D2 | 60 kDa Heat-shock protein, mitochondrial | 61,016 | 5.6 | – | √ | √ | 2.5 | –2.5 | Chaperone/stress response |
| D3 | Protein disulfide isomerase A3 | 56,747 | 5.9 | – | – | – | –3.9 | 3.9 | Metabolism |
| D4 | Lamin-A/C | 74,095 | 6.6 | – | – | √ | –1.1 | 1.1 | Cytoskeleton/mobility |
| D5 | Aldehyde dehydrogenase 1 | 54,696 | 6.3 | – | – | – | –1.9 | 1.9 | Metabolism |
| D6 | Beta-actin | 41,710 | 5.1 | √ | – | √ | –1.6 | 1.6 | Cytoskeleton/mobility |
| D7 | Keratin type II cytoskeletal 8 | 53,671 | 5.3 | – | – | √ | 1.7 | –1.7 | Cytoskeleton/mobility |
| D8 | Complement component 1 Q subcomponent-binding protein, mitochondrial | 31,342 | 4.5 | – | – | √ | 3.7 | –3.7 | Unknown |
| D9 | Tropomyosin alpha- | | 4.5 | – | – | √ | –3.1 | 3.1 | Cytoskeleton/mobility |
| D10 | Tropomyosin alpha- | | | – | – | √ | –2.0 | 2.0 | Cytoskeleton/mobility |
| D11 | Alpha-enolase | | | | √ | √ | 1.7 | –1.7 | Metabolism |
| D12 | Alpha-enolase | | | | | | 1.7 | –1.7 | Metabolism |
| D13 | Hemoglobin beta chain | 15,9.. | | | | | | –3.8 | Binding protein/folding |
| D14 | Non-POU domain-containing octamer-binding protein | 54,197 | 9.4 | | | | | | Binding protein/folding |
| D15 | Unnamed protein product | 65,980 | 7.6 | – | – | | | .2 | Unknown |
| D16 | Alpha-enolase | 47,008 | 7.2 | – | √ | √ | | 1.7 | Metabolism |
| D17 | Translationally-controlled tumor protein | 19,582 | 4.6 | – | – | √ | 1.1 | –1.1 | Binding protein/folding |
| D18 | DNA-binding protein | 35,801 | 8.9 | – | – | √ | 2.0 | –2.0 | Binding protein/folding |
| D19 | Unnamed protein product | 59,492 | 5.2 | – | – | – | –1.7 | 1.7 | Unknown |
| D20 | Heterogeneous nuclear ribonucleoprotein A1 | 38,822 | 9.3 | – | – | √ | 2.8 | –2.8 | Protein synthesis/degradation |
| D21 | Vimentin | 53,681 | 5.0 | – | √ | √ | –6.1 | 6.1 | Cytoskeleton/mobility |
| D22 | Peroxiredoxin 1 | 22,096 | 8.2 | – | √ | √ | 1.8 | –1.8 | Protection/detoxification |
| D23 | Peptidyl-prolyl cis-trans isomerase A, cyclophilin | 17,981 | 7.4 | – | √ | √ | 1.4 | –1.4 | Binding protein/folding |
| D24 | 78 kDa Glucose-regulated protein | 72,377 | 5.1 | – | √ | √ | –1.0 | 1.0 | Chaperone/stress response |
| D25 | Unnamed protein product | 65,980 | 7.6 | – | – | – | 1.5 | –1.5 | Unknown |
| D26 | Cofilin-1 | 18,490 | 8.2 | – | √ | √ | –1.9 | 1.9 | Cytoskeleton/mobility |

Where does the annotation come from ?
What's about the annotation accuracy ?

**Comparative proteomic analysis of lung cancer cell line and lung fibroblast cell line** (PMID: 19657000)

UniProt

SIB Swiss Institute of Bioinformatics

# Life of sequences …

## RNA, genes, genomes, …

**Nucleic acid sequence databases**    *Open access*

Annotated coding sequences (CDS) / gene prediction

**Protein sequence databases**    *Open access*

Protein sequences + biological knowledge (including GO terms)

**Nucleic acid sequence databases**
          **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

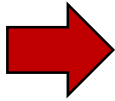          **UniProtKB/Swiss-Prot**
                    **Protein sequences**
                    **Biological knowledge**

          **UniProtKB/TrEMBL**
                    **Protein sequences**
                    **Biological knowledge**

          **GO annotation**

          **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

**Nucleic acid sequence databases**
        **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

    **UniProtKB/Swiss-Prot**
        **Protein sequences**
        **Biological knowledge**

    **UniProtKB/TrEMBL**
        **Protein sequences**
        **Biological knowledge**

    **GO annotation**

    **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# Protein sequence data:
# where does it come from ?

- > 190 billion 'different' proteins on earth ($\sum$ N species x M genes)

- ~ 130 million 'known and public' protein sequences in now
  - 50 % more by next year !

- More than 98% of the protein sequences are derived from the translation of nucleotide sequences (mRNA or DNA/genome)

- Less than 1 % come from direct protein sequencing (Edman, MS/MS…)

# Life of sequences …

## RNA, genes, genomes, …

| Nucleic acid sequence databases | Open access |

Annotated coding sequences (CDS) / gene prediction

mRNA / DNA - experimental / prediction

| Protein sequence databases | Open access |

Protein sequences + biological knowledge (including  GO terms)

# Life of sequences …

Data not submitted to public databases, delayed or canceled…

**RNA, genes, genomes, …**

**EMBL, GenBank, DDBJ**

*…if the submitters provide an annotated CoDing Sequence (CDS)*

INSDC

**Nucleic acid databases**

*no CDS*

*Gene prediction*
*RefSeq, Ensembl, other pipelines*

*RefSeq, Ensembl and other*

**Protein sequence databases**

UniProt

SIB
Swiss Institute of Bioinformatics

# Nucleic acid sequence databases

*open access*

*INSDC: ENA, GenBank, DDBJ*
*RefSeq (NCBI)*
*Ensembl (EBI)*

# EMBL-ENA/GenBank/DDBJ



**INSDC** International Nucleotide Sequence Database Collaboration

| ABOUT INSDC | POLICY | ADVISORS | DOCUMENTS |

### International Nucleotide Sequence Database Collaboration

- The International Nucleotide Sequence Database Collaboration (INSDC) is a long-standing foundational initiative that operates between DDBJ, EMBL-EBI and NCBI. INSDC covers the spectrum of data raw reads, though alignments and assemblies to functional annotation, enriched with contextual information relating to samples and experimental configurations.

| Data type | DDBJ | EMBL-EBI | NCBI |
|---|---|---|---|
| Next generation reads | Sequence Read Archive | European Nucleotide Archive (ENA) | Sequence Read Archive |
| Capillary reads | Trace Archive | | Trace Archive |
| Annotated sequences | DDBJ | | GenBank |
| Samples | BioSample | | BioSample |
| Studies | BioProject | | BioProject |

- The INSDC advisory board, the International Advisory Committee, is made up of members of each of the databases' advisory bodies. At their most recent meeting, members of this committee unanimously endorsed and reaffirmed the existing data-sharing policy of the three databases that make up the INSDC, which is stated below.
- Individuals submitting data to the international sequence databases should be aware of INSDC policy.

http://www.insdc.org/

# GenBank X02158

# GenBank X02158



**CDS annotation**
**CoDing Sequence**
**(proposed by submitters)**

```
CDS             join(615..627,1194..1339,1596..1682,2294..2473,2608..2763)
                /codon_start=1
                /product="erythropoietin"
                /protein_id="CAA26095.1"
                /db_xref="GDB:119110"
                /db_xref="GOA:P01588"
                /db_xref="HGNC:HGNC:3415"
                /db_xref="InterPro:IPR001323"
                /db_xref="InterPro:IPR003013"
                /db_xref="InterPro:IPR009079"
                /db_xref="InterPro:IPR012351"
                /db_xref="InterPro:IPR019767"
                /db_xref="PDB:1BUY"
                /db_xref="PDB:1CN4"
                /db_xref="PDB:1EER"
                /db_xref="UniProtKB/Swiss-Prot:P01588"
                /translation="MGVHECPAWLWLLLSLLSLPLGLPVLGAPPRLICDSRVLQRYLL
                EAKEAENITTGCAEHCSLNENITVPDTKVNFYAWKRMEVGQQAVEVWQGLALLSEAVL
                RGQALLVNSSQPWEPLQLHVDKAVSGLRSLTTLLRALGAQKEAISPPDAASAAPLRTI
                TADTFRKLFRVYSNFLRGKLKLYTGEACRTGDR"
sig_peptide     join(615..627,1194..1261)
mat_peptide     join(1262..1339,1596..1682,2294..2473,2608..2760)
                /product="erythropoietin"
intron          628..1193
                /number=1
exon            1194..1339
                /number=2
intron          1340..1595
                /number=2
exon            1596..1682
                /number=3
intron          1683..2293
                /number=3
exon            2294..2473
                /number=4
intron          2474..2607
                /number=4
exon            2608..3327
                /note="3' untranslated region"
                /number=5
ORIGIN
        1 agcttctggg cttccagacc cagctacttt gcggaactca gcaacccagg catctctgag
       61 tctccgccca agaccgggat gccccccagg aggtgtccgg gagcccagcc tttcccagat
      121 agcagctccg ccagtcccaa gggtgcgcaa ccggctgcac tcccctcccg cgacccaggg
      181 cccgggagca gcccccatga cccacacgca cgtctgcagc agccccgtca gcccccggagc
      241 ctcaacccag gcgtcctgcc cctgctctga ccccgggtgg cccctacccc tggcgacccc
      301 tcacgcacac agcctctccc ccacccccac ccgcgcacgc acacatgcag ataacagccc
      361 cgaccccccgg ccagagcccg agagtccctg ggccacccg gccgctcgct gcgctgcgcc
      421 gcaccgcgct gtcctccgg agccggaccg gggccaccgc gcccgctctg ctccgacacc
      481 gcgcccccctg gacagccgcc ctctcctcca ggcccgtggg gctggccctg caccgccgag
      541 cttcccggga tgagggcccc cggtgtggtc accgggcccc ccaggtcgct gagggacccc
      601 ggccaggcgc ggag atgggg gtgcacg gtg agtactcgcg ggctgggcgc tcccgcccgc
      661 ccgggtccct gtttgagcgg ggatttagcg ccccggctat tggccaggag gtggctgggt
      721 tcaaggaccg gcgacttgtc aaggaccccg gaagggggag gggggtgggg cagcctccac
      781 gtgccagcgg ggacttgggg gagtccttgg ggatggcaaa aacctgacct gtgaagggga
      841 cacagtttgg gggttgaggg gaagaaggtt tgggggggttc tgctgtgcca gtgagagga
      901 agctgataag ctgataacct gggcgctgga gccaccactt atctgccaga ggggaagcct
```

**CDS translation**

Ca

**DNA sequence**

GenBank X02158

**CDS**
**CoDing Sequence**
**(proposed by submitters)**

```
join(615..627,1194..1339,1596..1682,2294..2473,2608..2763)
/codon_start=1
/product="erythropoietin"
/protein_id=" CAA26095.1 "
/db_xref="GDB:119110"
/db_xref="GOA: P01588 "
/db_xref="HGNC: HGNC:3415 "
/db_xref="InterPro: IPR001323 "
/db_xref="InterPro: IPR003013 "
/db_xref="InterPro: IPR009079 "
/db_xref="InterPro: IPR012351 "
/db_xref="InterPro: IPR019767 "
/db_xref="PDB: 1BUY "
/db_xref="PDB: 1CN4 "
/db_xref="PDB: 1EER "
/db_xref="UniProtKB/Swiss-Prot: P01588 "
/translation="MGVHECPAWLWLLLSLLSLPLGLPVLGAPPRLICDSRVLQRYLL
EAKEAENITTGCAEHCSLNENITVPDTKVNFYAWKRMEVGQQAVEVWQGLALLSEAVL
RGQALLVNSSQPWEPLQLHVDKAVSGLRSLTTLLRALGAQKEAISPPDAASAAPLRTI
TADTFRKLFRVYSNFLRGKLKLYTGEACRTGDR"
```

either generated by gene prediction programs or experimentally proven.

```
 541 cttcccggga tgagggcccc cggtgtggtc acccggcgcc ccaggtcgct gagggacccc
 601 ggccaggcgc ggagatgggg gtgcacggtg agtactcgcg ggctgggcgc tcccgcccgc
 661 ccgggtccct gtttgagcgg ggatttagcg ccccggctat tggccaggag gtggctgggt
 721 tcaaggaccg gcgacttgtc aaggaccccg gaaggggag ggggtggggg cagcctccac
 781 gtgccagcgg ggacttgggg gagtccttgg ggatggcaaa aacctgacct gtgaagggga
 841 cacagtttgg gggttgaggg gaagaaggtt tggggggttc tgctgtgcca gtggagagga
 901 agctgataag ctgataacct gggcgctgga gccaccactt atctgccaga ggggaagcct
 961 ctgtcacacc aggattgaag tttggccgga gaagtggatg ctggtagcct gggggtgggg
1021 tgtgcacacg gcagcaggat tgaatgaagg ccagggaggc agcacctgag tgcttgcatg
1081 gttgggggaca ggaaggacga gctggggcag agacgtgggg atgaaggaag ctgtccttcc
1141 acagccaccc ttctccctcc ccgcctgact ctcagcctgg ctatctgttc tagaatgtcc
1201 tgcctggctg tggcttctcc tgtccctgct gtcgctccct ctgggcctcc cagtcctggg
1261 cgccccacca cgcctcatct gtgacagccg agtcctgcag aggtacctct tggaggccaa
1321 ggaggccgag aatatcacgg tgagaccccct tccccagcac attccacaga actcacgctc
1381 agggcttcag ggaactcctc ccagatccag gaacctggca cttggtttgg ggtggagttg
1441 ggaagctaga cactgcccc ctacataaga ataagtctgg tggccccaaa ccatacctgg
1501 aaactaggca aggagcaaag ccagcagatc ctacgcctgt ggccagggcc agagccttca
1561 gggaccctttg actccccggg ctgtgtgcat ttcagacggg ctgtgctgaa cactgcagct
1621 tgaatgagaa tatcactgtc ccagacacca aagttaattt ctatgcctgg aagaggatgg
1681 aggtgagttc ctttttttttt ttttttcctt tctttttggag aatctcattt gcgagcctga
1741 ttttggatga aagggagaat gatcgaggga aaggtaaaat ggagcagcag agatgaggct
1801 gcctgggcgc agaggctcac gtctataatc ccaggctgag atggccgaga tgggagaatt
1861 gcttgagccc tggagtttca gaccaaccta ggcagcatag tgagatcccc catctctaca
1921 aacatttaaa aaaattagtc aggtgaagtg gtgcatggtg gtagtcccag atatttggaa
1981 ggctgaggcg ggaggatcgc ttgagcccag gaatttgagg ctgcagtgag ctgtgatcac
2041 accactgcac tccagcctca gtgacagagt gaggccctgt ctcaaaaaag aaaagaaaaa
2101 agaaaaataa tgagggctgt atggaatacg ttcattattc attcactcac tcactcactc
2161 attcattcat tcattcattc aacaagtctt attgcatacc ttctgtttgc tcagcttggt
2221 gcttgggget gctgaggggc aggagggaga gggtgacatc cctcagctga ctcccagagt
2281 ccactccctg taggtcgggc agcaggccgt agaagtctgg cagggcctgg ccctgctgtc
2341 ggaagctgtc ctgcggggcc aggccctgtt ggtcaactct tcccagccgt gggagcccct
2401 gcagctgcat gtggataaag ccgtcagtgg ccttcgcagc ctcaccactc tgcttcgggc
2461 tctggagcc caggtgagta ggagcggaca cttctgcttg ccctttctgt aagaaggggga
2521 gaagggtctt gctaaggagt acaggaactg tccgtattcc ttcccttttct gtggcactgc
2581 agcgacctcc tgttttctcc ttggcagaag gaagccatct ccctccaga tgcggcctca
2641 gctgctccac tccgaacaat cactgctgac actttccgca aactcttccg agtctactcc
2701 aatttcctcc gggggaaagct gaagctgtac acaggggagg cctgcaggac aggggacaga
2761 tgaccaggtg tgtccacctg ggcatatcca ccacctccct caccaacatt gcttgtgcca
2821 caccctcccc cgccactcct gaacccgtc gaggggctct cagctcagcg ccagcctgtc
2881 ccatggacac tccagtgcca ccaatgacat ctcaggggcc agaggaactg tccagagagc
2941 aactctgaga tctaaggatg tcacagggcc aacttgaggg cccagagcag gaagcattca
3001 gagagcagct ttaaactcag ggacagaccc atgctgggaa gacgcctgag ctcactcggc
3061 accctgcaaa attgatgcca ggacacgctt tggaggcgat ttacctgttt tcgcacctac
3121 catcagggac aggatgacct ggagaactta ggtggcaagc tgtgacttct ccaggtctca
3181 cgggcatggg cactcccttg gtggcaagag cccccttgac accggggtgg tgggaaccat
3241 gaagacagga tgggggctgg cctctggctc tcatggggtc caacttttgt gtattcttca
3301 acctcattga caagaactga aaccaccaat atgactcttg gcttttctgt tttctgggaa
3361 cctccaaatc ccctggctct gtcccactcc tggcagca
```

UniProt

SIB
Swiss Institute of
Bioinformatics

# Coding sequence (CDS) annotation

submit

**DNA sequence**  ACGCTCGTACGCATCGTCACTACTAGCTACGACGACGACACGCTACTACTCGACGATTCT

CDS
CoDing Sequence
(provided by submitters)

transcribe

**Derived mRNA sequence**  AUGCGUAGUGAUGAAUGCUGCUGUGCGAUGAGCUGC

translate

**Derived protein sequence**  MRSNECCCAMSC

Slide J. McDowall

UniProt

SIB Swiss Institute of Bioinformatics

# UCSC genome browser: human EPO

mRNAs and their corresponding CDS annotation
(from EMBL/GenBank/DDBJ)

# Coding sequence (CDS) annotation



Slide J. McDowall

# UCSC genome browser: another gene…



mRNAs and their corresponding CDS annotation (from EMBL/GenBank/DDBJ)

How to deal with these sequences: see later….

# EMBL/GenBank/DDBJ

- Archive: nothing goes out -> highly redundant !

- Most annotations (including CDS, gene and protein names, …) are done by the submitters: heterogeneity of the quality and of the completion

- Many errors: in sequences, in annotations (gene and protein names, …), in CDS attribution, no consistency of annotations

- Archive: all submitted information remains there; not updated

UniProt

SIB
Swiss Institute of
Bioinformatics

# EMBL/GenBank/DDBJ and annotation

"Beyond limited editorial control and some internal integrity checks (for example, proper use of INSD formats and translation of coding regions specified in CDS entries are verified), **the quality and accuracy of the record are the responsibility of the submitting author, not of the database.** The databases will work with submitters and users of the database to achieve the best quality resource possible."

http://www.insdc.org/policy

# EMBL/GenBank/DDBJ and annotation

- many scientists assume that GenBank annotation is kept up to date, and they are surprised to hear that it is not
- the **annotation has remained static**: a gene labeled 'hypothetical protein' a few years ago might now have a known function.
- **erroneous and inconsistent naming of genes**.
- a name is **transferred from one gene to another** on the basis of sequence similarity (usually from a BLAST search). As more genomes are annotated, and more BLAST searches are run, **the original source of the name quickly becomes lost**.
- scientists should fix errors that they find. But this would quickly destroy the **archival function of GenBank**, as original entries would be erased over time.

(PMID: 17274839)

# EMBL/GenBank/DDBJ and UniProtKB

What is transferred to UniProtKB:

- The protein sequence (translated CDS)

- Gene and protein names
- Origin of the sequence (tissues)
- Taxonomy
- Publication provided by the submitting author

# What is transferred to UniProtKB ('imported'):

DR EMBL; DQ339047; ABC68418.1; -; mRNA.

FT source 1..1397
FT /organism="Rattus norvegicus"
FT /strain="Sprague-Dawley"
FT /mol_type="mRNA"
FT /sex="female"
FT /tissue_type="ovary"
FT /db_xref="taxon:10116"
FT CDS 70..1329
FT /codon_start=1
FT /product="testis derived transcript"
FT /note="TES"
FT /db_xref="GOA:Q2LAP6"

```
RN   [1]
RP   NUCLEOTIDE SEQUENCE [MRNA].
RC   STRAIN=Sprague-Dawley; TISSUE=Ovary;
RA   Seo Y.M., Jang S.J., Chun S.Y.;
RL   Submitted (DEC-2005) to the EMBL/GenBank/DDBJ databases.
```

```
DE   RecName: Full=Testin;
GN   Name=Tes;
OS   Rattus norvegicus (Rat).
OC   Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC   Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia; Myomorpha;
OC   Muroidea; Muridae; Murinae; Rattus.
```

# Nucleic acid sequence databases

*open access*

*INSDC: ENA, GenBank, DDBJ*
*RefSeq (NCBI)*
*Ensembl (EBI)*

# RefSeq  @NCBI

- The Reference Sequence (RefSeq) database provides a non-redundant set of sequences, including <span style="color:red">genomic DNA, transcript (RNA), and protein products</span>, for major research organisms.

- Contains sequences constructed from INSDC sequences.

- Contains protein sequences derived from <span style="color:red">gene prediction</span>, not submitted to EMBL/GenBank/DDBJ

One record  for one transcript

**Announcements**

**September 17, 2018**
**RefSeq Release 90 is available for FTP**

This release includes:

| | |
|---|---|
| Proteins: | 121,138,769 |
| Transcripts: | 23,838,836 |
| Organisms: | 84,276 |
| Available at: | ftp://ftp.ncbi.nlm.nih.gov/refseq/release/ |
| Documentation: | Release Notes |

See previous announcements, follow NCBI on Twitter, or subscribe to NCBI's refseq-announce mail list to receive announcements.

UniProt

Nucleotide

Nucleotide ▼

Advanced

GenBank ▾

# Homo sapiens erythropoietin (EPO), mRNA

NCBI Reference Sequence: NM_000799.3

FASTA    Graphics

Go to: ⊘

```
LOCUS       NM_000799               1340 bp    mRNA    linear   PRI 08-MAR-2018
DEFINITION  Homo sapiens erythropoietin (EPO), mRNA.
ACCESSION   NM_000799
VERSION     NM_000799.3
KEYWORDS    RefSeq.
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Ve
            Mammalia; Eutheria; Euarchontoglires; Prim
            Catarrhini; Hominidae; Homo.
REFERENCE   1  (bases 1 to 1340)
  AUTHORS   Nishimura K, Matsumoto R, Yonezawa Y and N
  TITLE     Effect of quercetin on cell protection via
            injury of HepG2 cells
  JOURNAL   Arch. Biochem. Biophys. 636, 11-16 (2017)
   PUBMED   29080630
  REMARK    GeneRIF: these results suggested that quer
            effects in HepG2 cells are mediated via EP
REFERENCE   2  (bases 1 to 1340)
```

```
COMMENT     REVIEWED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence was derived from X02157.1, S65458.1 and
            AC009488.5.
            This sequence is a reference standard in the RefSeqGene project.
            On Aug 17, 2017 this sequence version replaced NM_000799.2.

            Summary: This gene encodes a secreted, glycosylated cytokine
            composed of four alpha helical bundles. The encoded protein is
            mainly synthesized in the kidney, secreted into the blood plasma,
            and binds to the erythropoietin receptor to promote red blood cell
            production, or erythropoiesis, in the bone marrow. Expression of
            this gene is upregulated under hypoxic conditions, in turn leading
            to increased erythropoiesis and enhanced oxygen-carrying capacity
            of the blood. Expression of this gene has also been observed in
            brain and in the eye, and elevated expression levels have been
            observed in diabetic retinopathy and ocular hypertension.
            Recombinant forms of the encoded protein exhibit neuroprotective
            activity against a variety of potential brain injuries, as well as
            antiapoptotic functions in several tissue types, and have been used
            in the treatment of anemia and to enhance the efficacy of cancer
            therapies. [provided by RefSeq, Aug 2017].

            Publication Note:  This RefSeq record includes a subset of the
            publications that are available for this gene. Please see the Gene
            record to access additional publications.

            ##Evidence-Data-START##
            Transcript exon combination :: X02157.1, BC093628.1 [ECO:0000332]
            RNAseq introns              :: single sample supports all introns
                                           SAMEA2158188, SAMEA2159368
                                           [ECO:0000348]
            ##Evidence-Data-END##
            COMPLETENESS: complete on the 3' end.
PRIMARY     REFSEQ_SPAN         PRIMARY_IDENTIFIER PRIMARY_SPAN        COMP
            1-597               X02157.1           1-597
            598-600             S65458.1           248-250
            601-1144            X02157.1           601-1144
            1145-1145           AC009488.5         64325-64325
            1146-1340           X02157.1           1146-1340
```

RefSeq: NP_000790.2. NM_000799.2.

```
CDS             182..763
                /gene="EPO"
                /gene_synonym="DBAL; ECYT5; EP; MVCD2"
                /note="epoetin"
                /codon_start=1
                /product="erythropoietin precursor"
                /protein_id="NP_000790.2"
                /db_xref="CCDS:CCDS5705.1"
                /db_xref="GeneID:2056"
                /db_xref="HGNC:HGNC:3415"
                /db_xref="MIM:133170"
                /translation="MGVHECPAWLWLLLSLLSLPLGLPVLGAPPRLICDSRVLERYLL
                EAKEAENITTGCAEHCSLNENITVPDTKVNFYAWKRMEVGQQAVEVWQGLALLSEAVL
                RGQALLVNSSQPWEPLQLHVDKAVSGLRSLTTLLRALGAQKEAISPPDAASAAPLRTI
                TADTFRKLFRVYSNFLRGKLKLYTGEACRTGDR"
```

One record  for one transcript

# Ensembl

- Creates, integrates and distributes reference datasets

- Joint project between EMBL-EBI and the Sanger Centre

Ensembl[i]    ENST00000252723; ENSP00000252723; ENSG00000130427.

- Contains protein sequences derived from gene prediction which are not submitted to EMBL/GenBank/DDBJ

  -> included in UniProtKB

# Ensembl and UniProtKB

- UniProtKB contain protein sequences derived from Ensembl gene prediction

https://www.ensembl.org/info/genome/genebuild/automatic_coding.html

- Proteome sets in UniProtKB are made of records with a cross-reference to Ensembl

- See later

**Nucleic acid sequence databases**
> **INSDC, RefSeq, Ensembl**

➡️ **Protein sequence databases**

**UniProtKB**

> **UniProtKB/Swiss-Prot**
>> **Protein sequences**
>> **Biological knowledge**

> **UniProtKB/TrEMBL**
>> **Protein sequences**
>> **Biological knowledge**

> **GO annotation**

> **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# The life of a sequence …

Data not submitted to public databases, delayed or canceled…

**RNA, genes, genomes, …**

**EMBL, GenBank, DDBJ**

*…if the submitters provide an annotated CoDing Sequence (CDS)*

**INSDC**

**Nucleic acid databases**

*no CDS*

*Gene prediction*
*RefSeq, Ensembl, other pipelines*

*RefSeq, Ensembl and other*

**Protein sequence databases**

UniProt

SIB
Swiss Institute of Bioinformatics

# Protein sequence databases

**UniProtKB/Swiss-Prot**: <u>manually</u> annotated protein sequences  (12'500 species)

**UniProtKB/TrEMBL**: submitted CDS (EMBL-ENA) + automated annotation; non redundant with Swiss-Prot (700'000 species)

GenPept: submitted CDS (GenBank); no annotation; redundant with Swiss-Prot (700'000 species);

PIR: Protein Information Ressource; archive since 2003; integrated into UniProtKB; still available in NCBInr

PDB: Protein Databank: 3D data and associated sequences

PRF: Journal scan of 'published' peptide sequences

**RefSeq**: Reference Sequence for DNA, RNA, protein + gene prediction  + <u>partial manual</u> annotation (68'200 species)

NextProt: human proteins from UniProtKB/Swiss-Prot

**NCBInr**: Swiss-Prot + GenPept + PIR + PDB + PRF  + RefSeq

# Major 'general' protein sequence database 'sources'

Ensembl
PIR     PRF
TPA     PDB

integrated databases
'cross-references'

**UniProtKB: Swiss-Prot + TrEMBL**

databases are kept
separated

**NCBI-nr**: **Swiss-Prot + TrEMBL** + **GenPept** + PIR + PDB + PRF + **RefSeq** + TPA

not complete !!!
(only entries created before 2007 ?)

Ensembl and Refseq: gene prediction

UniProt

SIB
Swiss Institute of
Bioinformatics

# Homo sapiens

## @ UniProt

## @ NCBInr

**Nucleic acid sequence databases**
    **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

➡️ **UniProtKB**

> **UniProtKB/Swiss-Prot**
> > **Protein sequences**
> > **Biological knowledge**
>
> **UniProtKB/TrEMBL**
> > **Protein sequences**
> > **Biological knowledge**
>
> **GO annotation**
>
> **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

www.uniprot.org

~130 millions of proteins/entries

derived from ~ 710'000 different species

7 million unique visitors/year

New release every month

UniProt consortium :    EMBL-EBI    PIR    SIB

EBI : European Bioinformatics Institute (UK)
SIB : Swiss Institute of Bioinformatics (CH)
PIR : Protein Information Resource (USA)

# www.uniprot.org

# UniProt databases

## UniProtKB

UniProt Knowledgebase

### Swiss-Prot (558,590)

Manually annotated and reviewed.

Records with information extracted from literature and curator-evaluated computational analysis.

### TrEMBL (126,780,198)

Automatically annotated and not reviewed.

Records that await full manual annotation.

## UniRef

The UniProt Reference Clusters (UniRef) provide clustered sets of sequences from the UniProt Knowledgebase (including isoforms) and selected UniParc records.

## UniParc

UniParc is a comprehensive and non-redundant database that contains most of the publicly available protein sequences in the world.

## Proteomes

A proteome is the set of proteins thought to be expressed by an organism. UniProt provides proteomes for species with completely sequenced genomes.

## Supporting data

Literature citations

Cross-ref. databases

Taxonomy

Diseases

XXX

Subcellular locations

Keywords

UniProt

Swiss Institute of Bioinformatics

# UniProtKB

Each record/entry contains:

**Protein sequence(s)**
&
**Biological knowledge**
i.e. function, catalytic activity, subcellular location,..
&
sequence annotation (features)
i.e. PTM, variants, domains, …

*Manually annotated <u>or</u> automatically annotated*

# UniProtKB is composed of 2 sections

## UniProtKB/Swiss-Prot

### Reviewed - Manually annotated

Records with information extracted from literature and <u>curator-evaluated</u> computational analysis.

*One gene / one record*

## UniProtKB/TrEMBL

### Unreviewed – Computationally analyzed

Records that await full manual annotation.

*One protein / one record*

some redundancy…

released every 4 weeks

# UniProtKB is composed of 2 sections



Major differences in the protein sequence and annotation accuracy !

**UniProtKB**

UniProt Knowledgebase

Swiss-Prot (558,590)

Manually annotated and reviewed.

TrEMBL (126,780,198)

Automatically annotated and not reviewed.

**0.5 %** of UniProtKB protein sequences

**99.5 %** of UniProtKB protein sequences

UniProt release 2018_09

# UniProtKB: source of protein sequences

- **INSDC (EMBL/ENA, GenBank, DDBJ) (> 95%)**
- **Ensembl (gene prediction)**
- RefSeq
- Sequences of PDB
- Sequences scanned from literature (PRF)
- Direct submission (including direct protein sequencing)

**No evidence nor quality statement for the protein sequence**

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

# UniProtKB: source of annotation/evidence statements

## UniProtKB/Swiss-Prot: Manual insertion, color in yellow

🏷 1 Publication ▾   🏷 By similarity ▾   🏷 UniRule annotation ▾   🏷 Imported ▾

🏷 Sequence analysis   🏷 Combined sources ▾   🏷 Curated

Computational analysis (curator-evaluated)

## UniProtKB/TrEMBL: Automated insertion, color in blue

🏷 UniRule annotation ▾   🏷 Imported ▾   🏷 SAAS annotation ▾

🏷 Sequence analysis ▾

Computational analysis (NOT curator-evaluated)

UniProt

SIB
Swiss Institute of
Bioinformatics

# UniProtKB/Swiss-Prot entry



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

# Function [i]

Receptor for glucocorticoids (GC) (PubMed:27120390). Has a dual mode of action: as a transcription factor that binds to glucocorticoid response elements (GRE), both for nuclear and mitochondrial DNA, and as a modulator of other transcription factors. Affects inflammatory responses, cellular proliferation and differentiation in target tissues. Involved in chromatin remodeling (PubMed:9590696). Plays a role in rapid mRNA degradation by binding to the 5' UTR of target mRNAs and interacting with PNRC2 in a ligand-dependent manner which recruits the RNA helicase UPF1 and the mRNA-decapping enzyme DCP1A, leading to RNA decay (PubMed:25775514). Could act as a coactivator for STAT5-dependent transcription upon growth hormone (GH) stimulation and could reveal an essential role of hepatic GR in the control of body growth (By similarity). ◆ By similarity ▾  ◆ 3 Publications ▾

Isoform Alpha: Has transcriptional activation and repression activity (PubMed:15866175, PubMed:19248771, PubMed:20484466, PubMed:23820903, PubMed:11435610, PubMed:15769988, PubMed:17635946, PubMed:19141540, PubMed:21664385). Mediates glucocorticoid-induced apoptosis (PubMed:23303127). Promotes accurate chromosome segregation during mitosis (PubMed:25847991). May act as a tumor suppressor (PubMed:25847991). May play a negative role in adipogenesis through the regulation of lipolytic and antilipogenic gene expression (By similarity). ◆ By similarity ▾  ◆ 11 Publications ▾

Isoform Beta: Acts as a dominant negative inhibitor of isoform Alpha (PubMed:7769088, PubMed:8621628, PubMed:20484466). Has intrinsic transcriptional activity independent of isoform Alpha when both isoforms are coexpressed (PubMed:19248771, PubMed:26711253). Loses this transcription modulator function on its own (PubMed:20484466). Has no hormone-binding activity (PubMed:8621628). May play a role in controlling glucose metabolism by maintaining insulin sensitivity (By similarity). Reduces hepatic gluconeogenesis through down-regulation of PEPCK in an isoform Alpha-dependent manner (PubMed:26711253). Directly regulates STAT1 expression in isoform Alpha-independent manner (PubMed:26711253). ◆ By similarity ▾  ◆ 5 Publications ▾

Isoform Alpha-2: Has lower transcriptional activation activity than isoform Alpha. Exerts a dominant negative effect on isoform Alpha trans-repression mechanism (PubMed:20484466).

Isoform GR-P: Increases activity of isoform Alpha. ◆ 1 Publication ▾

Isoform Alpha-B: More effective than isoform Alpha in transcriptional activation, but not repression activity. ◆ 2 Publications ▾

Isoform 10: Has transcriptional activation activity. ◆ 1 Publication ▾

Isoform Alpha-C1: Has transcriptional activation activity. ◆ 1 Publication ▾

Isoform Alpha-C2: Has transcriptional activation activity. ◆ 1 Publication ▾

Isoform Alpha-C3: Has highest transcriptional activation activity of all isoforms created by alternative initiation (PubMed:15866175, PubMed:23820903). Has transcriptional repression activity (PubMed:23303127). Mediates glucocorticoid-induced apoptosis (PubMed:23303127, PubMed:23820903). ◆ 3 Publications ▾

Isoform Alpha-D1: Has transcriptional activation activity. ◆ 1 Publication ▾

Isoform Alpha-D2: Has transcriptional activation activity. ◆ 1 Publication ▾

Isoform Alpha-D3: Has lowest transcriptional activation activity of all isoforms created by alternative initiation (PubMed:15866175, PubMed:23820903). Has transcriptional repression activity (PubMed:23303127). ◆ 3 Publications ▾

## Miscellaneous

Isoform Beta: High constitutive expression by neutrophils may provide a mechanism by which these cells escape glucocorticoid-induced cell death and up-regulation by proinflammatory cytokines such as IL8 further enhances their survival in the presence of glucocorticoids during inflammation. ◆ 1 Publication ▾

Can up- or down-modulate aggregation
localization of expanded polyglutamine

## Regions

| Feature key | Position(s) |
|---|---|
| DNA binding [i] | 418 – 493 |
| Zinc finger [i] | 421 – 441 |
| Zinc finger [i] | 457 – 476 |

# GO - Molecular function [i]

- core promoter binding ◆ Source: CAFA ▾
- DNA binding transcription factor activity ◆ Source: UniProtKB ▾
- glucocorticoid-activated RNA polymerase II transcription factor binding transcription factor activity ◆ Source: UniProtKB ▾
- glucocorticoid receptor activity ◆ Source: ProtInc ▾
- Hsp90 protein binding ◆ Source: UniProtKB ▾
- protein kinase binding ◆ Source: ARUK-UCL ▾
- RNA binding ◆ Source: UniProtKB-KW
- RNA polymerase II proximal promoter sequence-specific DNA binding ◆ Source: NTNU_SB ▾
- steroid binding ◆ Source: UniProtKB ▾
- steroid hormone binding ◆ Source: UniProtKB ▾
- SUMO binding ◆ Source: CAFA ▾
- transcriptional activator activity, RNA polymerase II proximal promoter sequence-specific DNA binding ◆ Source: UniProtKB ▾
- zinc ion binding ◆ Source: InterPro ▾

View the complete GO annotation on QuickGO ...

## Subcellular location



Graphics by Christian Stolte; Source: COMPARTMENTS

Manual annotation  Automatic computational assertion

UniProt annotation | GO - Cellular component

**Isoform Alpha :**

Mitochondrion

Mitochondrion ⓘ   ◈ 1 Publication ▾

**Nucleus**

Nucleus ⓘ   ◈ 5 Publications ▾

Cytoskeleton

spindle ⓘ   ◈ 1 Publication ▾
centrosome ⓘ   ◈ 1 Publication ▾

Other locations

Cytoplasm ⓘ   ◈ 5 Publications ▾
Note: After ligand activation, translocates from the cytoplasm to the nucleus.   ◈ 3 Publications ▾

**Isoform Beta :**

**Nucleus**

Nucleus ⓘ   ◈ 3 Publications ▾

Other locations

Cytoplasm ⓘ   ◈ 2 Publications ▾
Note: Expressed predominantly in the nucleus with some expression also detected in the cytoplasm.   ◈ 2 Publications ▾

Capt

UniProt

SIB
Swiss Institute of
Bioinformatics

# Sequence annotation (feature viewer)

P04150

# UniProtKB/TrEMBL entry



The same gene (NR3C1) in UniProtKB/TrEMBL

Due to some redundancy between Swiss-Prot and TrEMBL

# Function[i]

## Regions

| Feature key | Position(s) | Description | Actions | Graphical view | Length |
|---|---|---|---|---|---|
| DNA binding[i] | 418 – 493 | Nuclear receptor ⬗ PROSITE-ProRule annotation ▾ 🏠 Add 🔧 BLAST | | | 76 |

## GO - Molecular function[i]

- DNA binding transcription factor activity ⬗ Source: InterPro
- glucocorticoid receptor activity ⬗ Source: InterPro
- sequence-specific DNA binding ⬗ Source: InterPro
- steroid binding ⬗ Source: InterPro
- zinc ion binding ⬗ Source: InterPro

View the complete GO annotation on QuickGO ...

## GO - Biological process[i]

- transcription, DNA-templated ⬗ Source: UniProtKB-UniRule

UniProt

SIB
Swiss Institute of Bioinformatics

## Subcellular location



Graphics by Christian Stolte; Source: COMPARTMENTS

☐ Manual annotation  ☐ Automatic computational assertion

**UniProt annotation**    GO - Cellular component

**Nucleus**

Nucleus 🛈    ⬙ PROSITE-ProRule annotation ▾    ⬙ SAAS annotation ▾

**Keywords - Cellular component**

Nucleus   ⬙ PROSITE-ProRule annotation ▾    ⬙ SAAS annotation ▾

# B6ZGU6

# UniProtKB - B6ZGU6 (B6ZGU6_HUMAN)

**Display**

Entry

Publications

Feature viewer

Feature table

Domains & sites

ProtVista

**Nucleic acid sequence databases**
          **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**


        **UniProtKB/Swiss-Prot**
             **Protein sequences**
             **Biological knowledge**

        **UniProtKB/TrEMBL**
             **Protein sequences**
             **Biological knowledge**

        **GO annotation**

        **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# UniProtKB/Swiss-Prot

## Protein sequence

Manually annotated and reviewed.

# UniProtKB/Swiss-Prot

**gene-centric** /protein-centric

- all protein products encoded by one gene are described in a single entry

- One or several protein sequences (isoforms) per entry (canonical / isoform)

- curated sequences in UniProtKB/Swiss-Prot are based on the reference genome (when available)

- collaboration with other databases

The protein sequences

found in UniProtKB/Swiss-Prot

are often derived from different different nucleotide sequences from EMBL/GenBank/DDBJ entries

+

mapped to the corresponding genomic sequence

## Cross-references[i]

### Sequence databases

| | |
|---|---|
| Select the link destinations: ○EMBL[i] ○GenBank[i] ○DDBJ[i] | U50743 mRNA. Translation: AAB09425.1. Frameshift.<br>AF316896 Genomic DNA. Translation: AAG37906.1.<br>AF316896 Genomic DNA. Translation: AAG37907.1.<br>AF241235 Genomic DNA. Translation: AAG34359.1.<br>AF241235 Genomic DNA. Translation: AAG34360.1.<br>AF241236 mRNA. Translation: AAG34361.1.<br>X86400 mRNA. Translation: CAA60152.1. Sequence problems.<br>BT006721 mRNA. Translation: AAP35367.1.<br>BC013289 mRNA. Translation: AAH13289.1.<br>BC005302 mRNA. Translation: AAH05302.1. |
| CCDS[i] | CCDS8385.1. [P54710-2]<br>CCDS8386.1. [P54710-1] |
| PIR[i] | S54159. |
| RefSeq[i] | NP_001671.2. NM_001680.4. [P54710-1]<br>NP_067614.1. NM_021603.3. [P54710-2] |

### Genome annotation databases

| | |
|---|---|
| Ensembl[i] | ENST00000260287; ENSP00000260287; ENSG00000137731. [P54710-2]<br>ENST00000292079; ENSP00000292079; ENSG00000137731. [P54710-1]<br>ENST00000528014; ENSP00000432430; ENSG00000137731. [P54710-2]<br>ENST00000532119; ENSP00000436414; ENSG00000137731. [P54710-2] |

**gene-centric** / protein-centric

Used to construct the UniProtKB canonical sequence

Automatically mapped to the UniProtKB record

http://www.uniprot.org/uniprot/P54710#cross_references

- **At least 20% of UniProtKB/Swiss-Prot entries required curation effort to "correct" the sequences.**

- **Typical problems**
  - **unsolved conflicts;**
  - **uncorrected initiation sites;**
  - **frameshifts;**
  - **other 'problems'**

UniProt

SIB
Swiss Institute of
Bioinformatics

# UCSC genome browser: another gene



mRNAs and their corresponding
CDS annotation
(from EMBL/GenBank/DDBJ)

# P54710 - ATNG_HUMAN

| | |
|---|---|
| Protein | **Sodium/potassium-transporting ATPase subunit gamma** |
| Gene | **FXYD2**, ATP1C, ATP1G1 |
| Organism | *Homo sapiens (Human)* |
| Status | Reviewed - ⊙⊙⊙⊙⊙ - Experimental evidence at protein level[i] |

## Manually checked (comparison with orthologs, etc.)

```
          MPPGGHIESGASFSPSAGSDGDKGLKPREEGEAPEAGPCLASQLAPVQWVAPPRLLLRLI
X86400
AF241236  -------MTGLSMDGGGSPKGDVDPFYYDYETVRNGG--------------------LI
BC005302  --------MDRWYLGGSPKGDVDPFYYDYETVRNGG--------------------LI
                .....**  .      :     . :.*                      **

X86400    FIDLPALLIAPTAESSAEEDEEPHDEGQSSEDQAPIANGLIVIVERVHVPLGAAATVHRQ
AF241236  FAGLAFIV----------------------------GLLILLSR-RFRCGGNKKRRQI
BC005302  FAGLAFIV----------------------------GLLILLSR-RFRCGGNKKRRQI
          *  .*. ::                            **:::::.*  ..   *.  . ..

X86400    PSHFPR
AF241236  NEDEP-
BC005302  NEDEP-
            .  *
```

incorrect → AF241236

correct → BC005302

### Cross-references[1]

**Sequence databases**

Select the link destinations:
- ● EMBL[i]
- ○ GenBank[i]
- ○ DDBJ[i]

U50743 mRNA. Translation: AAB09425.1. Frameshift.
AF316896 Genomic DNA. Translation: AAG37906.1.
AF316896 Genomic DNA. Translation: AAG37907.1.
AF241235 Genomic DNA. Translation: AAG34359.1.
AF241235 Genomic DNA. Translation: AAG34360.1.
AF241236 mRNA. Translation: AAG34361.1.
X86400 mRNA. Translation: CAA60152.1. Sequence problems.
BT006721 mRNA. Translation: AAP35367.1.
BC013289 mRNA. Translation: AAH13289.1.
BC005302 mRNA. Translation: AAH05302.1.

UniProt

# P54710 - ATNG_HUMAN

| | |
|---|---|
| Protein | **Sodium/potassium-transporting ATPase subunit gamma** |
| Gene | **FXYD2**, ATP1C, ATP1G1 |
| Organism | *Homo sapiens (Human)* |
| Status | Reviewed - ⊚⊚⊚⊚⊚ - Experimental evidence at protein level[i] |

## Sequences (2)[i]

Sequence status[i]: Complete.

This entry describes **2** isoforms[i] produced by **alternative splicing**. ☰ Align

---

### Isoform 1 (identifier: P54710-1) [UniParc] ⬇ FASTA

*Also known as: A*

*This isoform has been chosen as the 'canonical' sequence. All positional information in this entry refers to it. This is also the sequence that appears in the downloadable versions of the entry.*

« Hide

```
MTGLSMDGGG SPKGDVDPFY YDYETVRNGG LIFAGLAFIV GLLILLSRRF   50
RCGGNKKRRQ INEDEP                                        66
```

---

### Isoform 2 (identifier: P54710-2) [UniParc] ⬇ FASTA

*Also known as: B*

*The sequence of this isoform differs from the canonical sequence as follows:*

    1-8: MTGLSMDG → MDRWYL

« Hide

```
MDRWYLGGSP KGDVDPFYYD YETVRNGGLI FAGLAFIVGL LILLSRRFRC   50
GGNKKRRQIN EDEP                                          64
```

UniProt

Swiss Institute of Bioinformatics

# UniProtKB/Swiss-Prot

## Entry vs Protein sequence(s)

One entry – one gene – one species

One or several protein sequences (isoforms) per entry

canonical & isoform

gene-centric / protein-centric

UniProt

SIB
Swiss Institute of
Bioinformatics

UniProtKB/Swiss-Prot protein knowledgebase release 2018_02 statistics

1. INTRODUCTION

Release 2018_02 of 28-Feb-18 of UniProtKB/Swiss-Prot contains 556825 sequence entries,

comprising 199652254 amino acids abstracted from 258419 references.

270 sequences have been added since release 2018_01, the sequence data of

41 existing entries has been updated and the annotations of

274847 entries have been revised.

Number of fragments: 9129

Number of additional sequences produced by alternative splicing, initiation or promoter usage, or ribosomal frameshifting: 39825

556'825 'canonical'

+ 39'825 'isoforms'
(7 %)

http://web.expasy.org/docs/relnotes/relstat.html

**Beware**

The isoform sequences are <u>not</u> included in all datasets,

Examples:

- Complete proteome -> download Fasta (canonical & isoform)
- Blast@ NCBI (NCBInr)

# UniProtKB/Swiss-Prot
## Biological knowledge / annotation

Manually annotated and reviewed.

Knowledge:
- comprehensive summary (free text) that provides a complete overview of the information available
- standardized vocabularies to facilitate subsequent retrieval whenever possible

# Source of annotation/Evidence statements

- Selected Publication (experimental)  `1 Publication ▾`  `Curated`

- Another UniProtKB entry (orthologs):  `By similarity ▾`

- An entry from another database:  `Imported ▾`

- Curator-evaluated computational analysis*  `UniRule annotation ▾`
  `Sequence analysis`

- Combined sources  `Combined sources ▾`

*more details later…*

# Source of annotation/Evidence statements

- Selected Publication (experimental)  🏷 1 Publication ▾   🏷 Curated

- Another UniProtKB entry (orthologs): 🏷 By similarity ▾

- An entry from another database: 🏷 Imported ▾

- Curator-evaluated computational analysis* 🏷 UniRule annotation ▾
  🏷 Sequence analysis

- Combined sources  🏷 Combined sources ▾

*more details later…*

# comprehensive and computer friendly representation of biological knowledge

PubMed=16595657

Selected Publication



FIGURE 3. **Subcellular localization of Alb3-GFP and Alb4-GFP fusion proteins.** *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3- or Alb4-GFP. Maximum intensity signals from confocal images are shown for chlorophyll autofluorescence, GFP fluorescence, and an overlay of both. OEP7-GFP is included as a

corresponds to the predicted length of the Alb4 mRNA. Even after prolonged exposure of the blots treated with the Alb4 probe, no signal could be found at ~3.5 kb, the predicted size of the Artemis transcript.

*Alb4 Is a Thylakoid Membrane Protein*—Alb4 is predicted to have a chloroplast transit peptide with a processing site after amino acid resi-

due 45 based on the ChloroP prediction program (33). To test this prediction, *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3-GFP or Alb4-GFP. Merging of the GFP and autofluorescence images indicated a thylakoid localization of Alb4. The GFP distribution for Alb4 is similar to that of Alb3 and not to that of outer envelope protein AtOEP7 (Fig. 3). To test this assumption, *in vitro* translated radiolabeled Alb4 was imported into isolated pea

# comprehensive and computer friendly representation of biological knowledge

## PubMed=16595657

## UniProtKB  Q9FYL3



FIGURE 3. **Subcellular localization of Alb3-GFP and Alb4-GFP fusion proteins.** *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3- or Alb4-GFP. Maximum intensity signals from confocal images are shown for chlorophyll autofluorescence, GFP fluorescence, and an overlay of both. OEP7-GFP is included as a

### Subcellular location[1]

- Plastid > chloroplast thylakoid membrane  ✏ 1 Publication ▾  ; Multi-pass membrane protein  ✏ 1 Publication ▾

**Topology**

| Feature key | Position(s) | Length | Description | Graphical view |
|---|---|---|---|---|
| Transmembrane[i] | 115 – 135 | 21 | Helical ✏ Sequence analysis | |
| Transmembrane[i] | 184 – 204 | 21 | Helical ✏ Sequence analysis | |
| Transmembrane[i] | 263 – 283 | 21 | Helical ✏ Sequence analysis | |
| Transmembrane[i] | 302 – 322 | 21 | Helical ✏ Sequence analysis | |

**GO - Cellular component[i]**

- chloroplast ✏ Source: TAIR ▾
- chloroplast thylakoid membrane ✏ Source: TAIR ▾
- integral component of membrane ✏ Source: UniProtKB-KW
- thylakoid ✏ Source: TAIR ▾

### Controlled vocabulary

2. **"A second thylakoid membrane-localized Alb3/OxaI/YidC homologue is involved in proper chloroplast biogenesis in Arabidopsis thaliana."**
   Gerdes L., Bals T., Klostermann E., Karl M., Philippar K., Huenken M., Soll J., Schuenemann D.
   J. Biol. Chem. 281:16632-16642(2006) [PubMed] [Europe PMC] [Abstract]
   **Cited for**: SEQUENCE REVISION, TISSUE SPECIFICITY, SUBCELLULAR LOCATION.

corresponds to the predicted length of the Alb4 mRNA. Even after prolonged exposure of the blots treated with the Alb4 probe, no signal could be found at ~3.5 kb, the predicted size of the Artemis transcript.

*Alb4 Is a Thylakoid Membrane Protein*—Alb4 is predicted to have a chloroplast transit peptide with a processing site after amino acid resi-

due 45 based on the ChloroP prediction program (33). To test this prediction, *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3-GFP or Alb4-GFP. Merging of the GFP and autofluorescence images indicated a thylakoid localization of Alb4. The GFP distribution for Alb4 is similar to that of Alb3 and not to that of outer envelope protein AtOEP7 (Fig. 3). To test this assumption, *in vitro* translated radiolabeled Alb4 was imported into isolated pea

# Source of annotation/Evidence statements

- Selected Publication (experimental)   [1 Publication ▾]

- Another UniProtKB entry (orthologs):   [By similarity ▾]

- An entry from another database:   [Imported ▾]

- Curator-evaluated computational analysis   [UniRule annotation ▾] *
  [Sequence analysis]

- Combined sources   [Combined sources ▾]

*more details later…*

UniProt

SIB
Swiss Institute of
Bioinformatics

# comprehensive and computer friendly representation of biological knowledge

Selected Publication

PubMed=16595657



FIGURE 3. **Subcellular localization of Alb3-GFP and Alb4-GFP fusion proteins.** *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3- or Alb4-GFP. Maximum intensity signals from confocal images are shown for chlorophyll autofluorescence, GFP fluorescence, and an overlay of both. OEP7-GFP is included as a

corresponds to the predicted length of the Alb4 mRNA. Even after prolonged exposure of the blots treated with the Alb4 probe, no signal could be found at ~3.5 kb, the predicted size of the Artemis transcript.

*Alb4 Is a Thylakoid Membrane Protein*—Alb4 is predicted to have a chloroplast transit peptide with a processing site after amino acid resi-

UniProtKB  Q9FYL3

**Molecule processing**

| Feature key | Position(s) | Length | Description |
|---|---|---|---|
| Transit peptide[i] | 1 – 45 | 45 | Chloroplast  Sequence analysis |
| Chain[i] | 46 – 499 | 454 | ALBINO3-like protein 1, chloroplastic |

due 45 based on the ChloroP prediction program (33). To test this prediction, *Arabidopsis* mesophyll protoplasts were transiently transformed with constructs for Alb3-GFP or Alb4-GFP. Merging of the GFP and autofluorescence images indicated a thylakoid localization of Alb4. The GFP distribution for Alb4 is similar to that of Alb3 and not to that of outer envelope protein AtOEP7 (Fig. 3). To test this assumption, *in vitro* translated radiolabeled Alb4 was imported into isolated pea

# Protein sequence analysis: in-house resource

*Table 1. Sequence analysis tools used during the UniProtKB manual curation process*

| Program | Version | Prediction |
|---------|---------|------------|
| **Topology** | | |
| Signal P (6) | 3.0 | Presence and location of signal peptides |
| TargetP (6) | 1.1 | Presence and location of transit peptides |
| Predotar (7) | 1.03 | Mitochondrial, plastid or ER targeting sequences |
| ESKW* (8) | UniProt-modified version 1.0 | Transmembrane domains |
| MEMSAT (9) | UniProt-modified version 1.8a | Transmembrane domains |
| TMHMM (10) | 2.0 | Transmembrane domains |
| Phobius (11) | Unknown | Discriminates transmembrane and signal regions |
| **Post-translational modifications** | | |
| GPI-predictor (12) | 1.0 | GPI lipid anchor sites |
| NetNGlyc (13) | 1.0 | N-glycosylation sites |
| NetOGlyc (14) | 3.1 | O-glycosylation sites |
| NMT Predictor (15) | 1.0 | N-terminal myristoylation sites |
| Sulfinator (16) | 1.0 | Tyrosine sulfation sites |
| **Domains** | | |
| ps_scan | 1.0 | Internal PROSITE profile, pattern and rule scanning program |
| InterPro (17) | Uses latest versions of InterPro and InterProScan | Retrieves non-PROSITE motif matches using InterPro database or InterProScan |
| Coils (18) | 2.2 | Coiled-coil regions |
| polyAA | 1.0 | Internal program which identifies homopolymeric stretches of amino acids |
| REPEAT (19) | 1.1 | Identifies the following repeats: Ankyrin, Armadillo, HAT, HEAT, Kelch, Leucine-rich, PFTA, PFTB, RCC1, TPR, WD40 |

*ESKW = transmembrane prediction algorithm by Eisenberg, Schwarz, Komaromy and Wall

**Nucleic acid sequence databases**
  **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

  **UniProtKB/Swiss-Prot**
    **Protein sequences**
    **Biological knowledge**

➡ **UniProtKB/TrEMBL**
    **Protein sequences**
    **Biological knowledge**

  **GO annotation**

  **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# UniProtKB/TrEMBL



One protein sequence per entry

gene-centric /**protein-centric**

**99,5 % of UniProtKB protein sequences**

# P73722 - P73722_SYNY3

| | |
|---|---|
| Protein | Submitted name: **SOS function regulatory protein** |
| Gene | **lexA** |
| Organism | *Synec... sp. (strain PCC 6803 / Kazusa)* |
| Status | Unreviewed - ◉◉○○○ - Protein inferred from homology[i] |

## Function[i]

Represses a number of genes involved in the response to DNA damage (SOS response), including recA and lexA. In the presence of single-stranded DNA, RecA interacts with LexA causing an autocatalytic cleavage which disrupts the DNA-binding part of LexA, leading to derepression of the SOS regulon and eventually DNA repair 🏷 By similarity .

🏷 SAAS annotations ▼

## Catalytic activity[i]

Hydrolysis of Ala-|-Gly bond in repressor LexA. 🏷 SAAS annotations ▼

**Keywords - Molecular function[i]**

Hydrolase 🏷 SAAS annotations ▼ , Repressor 🏷 SAAS annotations ▼

**Keywords - Biological process[i]**

DNA damage, DNA repair, DNA replication 🏷 SAAS annotations ▼ , SOS response 🏷 SAAS annotations ▼ , Transcription, Transcription regulation 🏷 SAAS annotations ▼

**Keywords - Ligand[i]**

DNA-binding 🏷 SAAS annotations ▼

Capture                    Ctrl+Ins

UniProt

SIB
Swiss Institute of Bioinformatics

# UniProtKB/TrEMBL

## Protein sequence

Automatically annotated and not reviewed.

## Protein sequence

- The quality of the protein sequences is dependent on the information provided by the submitter of the original nucleotide entry (EMBL-ENA CDS) or of the gene prediction pipeline (i.e. Ensembl).

- 100% identical sequences (same length, same organism are merged automatically).

gene-centric **/protein-centric**

One protein sequence per entry

UniProt

SIB
Swiss Institute of
Bioinformatics

# UniProtKB/TrEMBL

## Biological knowledge / annotation

publication

reviewed entries (Swiss-Prot)

unreviewed entries (TrEMBL)

INFORMATION CONTENT

Automated annotation

UniProt

SIB
Swiss Institute of Bioinformatics

# Automated annotation

**Automated generated rules (SAAS)**

Generates a set of decision trees using data mining (new set every UniProtKB release)

SAAS annotation

**Manually generated rules (UniRule)**
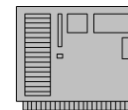
Maintains a set of manual annotation rules

UniRule annotation

UniRule = PIR + HAMAP + Rulebase

**Sequence analysis methods (SAM)**

Signal, transmembrane, coils prediction

Sequence analysis

**InterPro**

Domains & GO terms

InterPro
Protein sequence analysis & classification

InterPro annotation

UniProt

SIB
Swiss Institute of
Bioinformatics

# SAAS



SAAS learns on the properties present in the reviewed UniProtKB (Swiss-Prot) entries and uses the following attribute types to define the learning entries: InterPro protein family, taxonomy and sequence length. This combination allows SAAS to generate rules to annotate protein properties such as **function, catalytic activity, pathway membership, subcellular location, protein names and feature predictions.**

## SAAS: SAAS00002149

👁 View all proteins annotated by this rule   ⟳ Remove highlights

### If a protein meets these conditions... [i]

#### Common conditions

taxon = Bacteria
Matches InterPro signature IPR002146
Matches InterPro signature IPR028987
Matches Pfam signature PF00430
Matches SCOP Superfamily signature SSF81573

### ... then these annotations are applied [i]

#### Function [i]

$F_1F_0$ ATP synthase produces ATP from ADP in the presence of a proton or sodium gradient. F-type ATPases consist of two structural domains, $F_1$ containing the extramembraneous catalytic core and $F_0$ containing the membrane proton channel, linked together by a central stalk and a peripheral stalk. During catalysis, ATP synthesis in the catalytic domain of $F_1$ is coupled via a rotary mechanism of the central stalk subunits to proton translocation.

http://insideuniprot.blogspot.ch/2016/10/automatic-learning-based-annotation-in.html

# SAAS



Query: source: SAAS00002149, November 2018

Beware: the SAAS number may change, if the rules changes…

# Automated annotation

**Automated generated rules (SAAS)**

Generates a set of decision trees using data mining (new set every UniProtKB release)

SAAS annotation

**Manually generated rules (UniRule)**

Maintains a set of manual annotation rules

UniRule = PIR + HAMAP + Rulebase

UniRule annotation

**Sequence analysis methods (SAM)**

Signal, transmembrane, coils prediction

Sequence analysis

**InterPro**

Domains & GO terms

InterPro

Protein sequence analysis & classification

InterPro annotation

UniProt

SIB
Swiss Institute of Bioinformatics

UniProt

UniRule▾                    Advanced▾  🔍 Search

BLAST   Align   Retrieve/ID mapping                          Help   Contact

## UniRule: UR000068985

**Source ID:** RU361160

👁 View all proteins annotated by this rule   ↻ Remove highlights

### If a protein meets these conditions... ⁱ

#### Common conditions

Matches PRINTS signature PR00078
Matches PROSITE signature PS00071
sequence length = 0 - 500
taxon = Eukaryota, Bacteria
Matches TIGRFAM signature TIGR01534
Does not match InterPro signature IPR006422

#### Special conditions

taxon ≠ Viridiplantae, Cryptophyta, Rhodophyta, Bacteria

taxon = Viridiplantae, Cryptophyta, Rhodophyta, Bacteria

### ... then these annotations are applied ⁱ

**Protein names**ⁱ

*Recommended name:*
   **Glyceraldehyde-3-phosphate dehydrogenase**

(EC:1.2.1.12)

(EC:1.2.1.-)

**Sequence similarities**ⁱ

Belongs to the glyceraldehyde-3-phosphate dehydrogenase family.

**Catalytic activity**ⁱ

D-glyceraldehyde 3-phosphate + phosphate + NAD$^+$ = 3-phospho-D-glyceroyl phosphate + NADH.

**Subunit structure**ⁱ

Homotetramer.

UniProt

SIB Swiss Institute of Bioinformatics

Query: source:RU361160, Nov 2018

## Profiles and annotation rules are manually curated:  HAMAP

### General rule information [?]

| | |
|---|---|
| Accession | MF_01578 |
| Dates | 23-FEB-2007 (Created)<br>19-OCT-2016 (Last updated, Version 15) |
| Name | Shikimate_DH_YdiB |
| Scope | Bacteria; Enterobacterales |
| Templates | P0A6D5 (YDIB_ECOLI); Q8ZPR4 (YDIB_SALTY): [Recover all] |
| Triggered by | HAMAP; MF_01578 (Get profile general information and statistics) |

### Propagated annotation [?]

#### Identifier, protein and gene names [?]

| | | |
|---|---|---|
| Identifier | YDIB | |
| Protein name | RecName: | Full=Quinate/shikimate dehydrogenase; |
| | | EC=1.1.1.282; |
| | AltName: | Full=NAD-dependent shikimate 5-dehydrogenase; |
| Gene name | ydiB | |

#### Comments [?]

| | |
|---|---|
| Function | The actual biological function of YdiB remains unclear, nor is it known whether 3-dehydroshikimate or quinate represents the natural substrate. Catalyzes the reversible NAD-dependent reduction of both 3-dehydroshikimate (DHSA) and 3-dehydroquinate to yield shikimate (SA) and quinate, respectively. It can use both NAD or NADP for catalysis, however it has higher catalytic efficiency with NAD. |
| Catalytic activity | L-quinate + NAD(P)(+) = 3-dehydroquinate + NAD(P)H.<br>Shikimate + NAD(P)(+) = 3-dehydroshikimate + NAD(P)H. |
| Pathway | Metabolic intermediate biosynthesis; chorismate biosynthesis; chorismate from D-erythrose 4-phosphate and phosphoenolpyruvate: step 4/7. |
| Subunit | Homodimer. |
| Similarity | Belongs to the shikimate dehydrogenase family. |

UniProt

SIB
Swiss Institute of
Bioinformatics

## Profiles and annotation rules are manually curated:  **HAMAP**



Query: source:"mf 01578"

## Profiles and annotation rules are manually curated: **HAMAP**

# Automated annotation

**Automated generated rules (SAAS)**

Generates a set of decision trees using data mining (new set every UniProtKB release)

 SAAS annotation

**Manually generated rules (UniRule)**

Maintains a set of manual annotation rules

UniRule = PIR + HAMAP + Rulebase

 UniRule annotation

**Sequence analysis methods (SAM)**

Signal, transmembrane, coils prediction

 Sequence analysis

**InterPro**

Domains & GO terms

InterPro — Protein sequence analysis & classification

InterPro annotation

UniProt

SIB — Swiss Institute of Bioinformatics

## SAM - Sequence Analysis Methods

**Last modified October 16, 2015**

UniProt's Automatic Annotation pipeline enhances the unreviewed records in UniProtKB by enriching them with automatic classification and annotation. In this context, we use a suite of Sequence Analysis Methods (SAM) to enrich the unreviewed TrEMBL records in the UniProt Knowledgebase with extra sequence-specific information.

## Methods

Predictions of sequence features such as Signal, Transmembrane and Coil regions are generated using the following software from external providers:

- TMHMM
- SignalP
- Phobius
- Coils

These methods are applied to UniProtKB sequences by InterPro to predict sequence features. More annotations (mainly keywords) are then added automatically to enrich the generated predictions. The new predictions are propagated to all the UniProtKB/TrEMBL records that do not already contain such feature predictions from the UniRule automatic annotation system.

No GO (cellular component) annotation in this case…
http://www.uniprot.org/help/sam

SAM: Transmembrane

annotation:(type:transmem) AND reviewed:no

SAM: Signal peptide

# Automated annotation

**Automated generated rules (SAAS)**

Generates a set of decision trees using data mining (new set every UniProtKB release)

SAAS annotation

**Manually generated rules (UniRule)**

Maintains a set of manual annotation rules

UniRule annotation

UniRule = PIR* + HAMAP + Rulebase

**Sequence analysis methods (SAM)**

Signal, transmembrane, coils prediction

Sequence analysis

**InterPro**

Domains & GO terms

InterPro
Protein sequence analysis & classification

InterPro annotation

# UniProtKB - F1MSM3 (F1MSM3_BOVIN)

## Display

- Entry
- Publications
- Feature viewer
- Feature table

None

**BLAST** | Align | **Format** | **Add to basket** | **History**

| | |
|---|---|
| Protein | Submitted name: **Uncharacterized protein** |
| Gene | **NOTCH1** |
| Organism | *Bos taurus (Bovine)* |
| Status | Unreviewed · Annotation score: ●●●●● · Protein predicted[i] |

```
>tr|F1MSM3|F1MSM3_BOVIN Uncharacterized protein OS=Bos taurus GN=NOTCH1 PE=4 SV=2
MPPLLAPLLCLALLPALAARGLRCSQPGETCLNGGKCEVFPNGTEACICGGAFAGQQCQA
PNPCLSAPCKNGGTCHTTEREGLVDYVCGCRLGFSGPLCLTPRDHACLASPCLNGGTCDL
LTLTEYKCLCTPGWSGKTCQQADPCASNPCANGGQCLPFEASYICHCPPGFHGPTCRQDV
NECSQSPGLCHHGGTCLNEVGSYRCVCRPTHTGPHCELPYVPCSPSPCQNGGTCRPTGDT
THECACLPGFTGQNCEENIDDCPGNSCKNGGACVDGVNTYNCRCPPEWTGQYCTEDVDEC
QLMPNACQNGGTCHNTHGGYNCVCVNGWTGEDCSENIDDCASASCFQGATCHDRVASFYC
ECPHGRTGLLCHLNDACISNPCNEGSNCDTNPVNGKAICTCPSGYTGPACSQDVDECSLG
ANPCEHAGKCINTLGSFECQCLQGYTGPRCEIDVNECVSNPCQNDATCLDQIGEFQCICM
PGYEGLHCEVNTDECASSPCLQNGRCLDKINEFVCECPTGFTGHLCQYDVDECASTPCKN
GAKCLDGPNTYTCVCTEGYTGPHCEVDIDECDPDPCHYGSCKDGVATFTCLCQPGYTGHH
CESNINECHSQPCRHGGTCQDRDNAYLCFCLKGTTGPNCEINLDDCASNPCDSGTCLDKI
DGYECACEPGYTGSMCNINIDECADSPCHNGGTCEDGINGFTCRCPEGYHDPTCLSEVNE
CSSNPCIHGACRDSLNGYKCDCDPGWSGANCDVNNDECESNPCINGGTCKDMTSGYVCAC
REGFSGPNCQTNINECASNPCLNQGTCIDDVAGYKCNCLLPYTGATCEVVLAPCAPGPCR
NGGECRESEDYESFSCACPAGWQGQTCEIDINECVKSPCRAGASCQNTNGSYRCHCQAGY
TGRNCETDIDDCRPNPCHNGGSCTDGINTAFCDCLPGFQGAFCEEDINECASSPCRNGAN
CTDCVDSYTCTCPTGFSGIHCENNTPDCTESSCFNGGTCVDGINSFTCLCPPGFTGSYCQ
HDVNECDSRPCLHGGTCHDSYGTYTCTCPQGYTGLNCQTLVRWCDSSPCKNDGRCWQTNA
LYRCECHSGWTGLYCDVPSVSCEVAARQQGVNVTHLCRNGGLCMNAGNTHRCHCQAGYTG
SYCEEQVDECSPSPCQNGATCTDYPGGYSCECVAGYHGVNCSEEVNECLSQPCRNGGTCI
DLTNTYKCSCPRGTQGVHCEINVDDCNPPIDPVSRGPKCFNNGTCVDQVGGYSCSCPPGF
VGERCEGDVNECLSNPCDARGTQNCVQHVNAFHCECRAGHTGRRCESVINGCKDRPCKNG
GSCAVASNTARGFICKCPAGFEGATCENDARSCGSLRCLNGGTCIAGPRSPTCLCLGPFT
GPECQFPASSPCVGGNPCYNQGVCEPTAESPFYRCRCPAKFNGLLCHILDYSFGGGVGLD
IPPPQIEETCELPGCREEAGNKVCSLQCNSHACGWDGGDCSLDFDDPWQNCTQSLQCWKY
FSNGRCDSQCNSAGCLFDGFDCQRAEGQCNPLYDQYCKDHFRDGHCDQGCNSAECEWDGL
DCAEHVPERLAAGTLVLVVLMPPEQLRNRSLHFLRELSRLLHTNVVFKRDASGQQMIFPY
YGQEPHCRQGSAPRSVGVSTTHALLVLDKASPQGHCAPPGLFLSIVYLEIDNRQCVQSSS
QCFQSATDVAAFLGALASLGSLNIPYKIEAVQSETVEPPPPPPLHFMYVAVVAFVLLFFV
GCGVLLSRKRRRQHGQLWFPEGFKVSEASKKKRREPLGEDSVGLKPLKNSSDGALMDDNQ
NEWGDEGLEAKKFRFEEPVVLPDLDDQTDHRQWTQQHLDAADLRVSAMAPTPPQGEADAD
CMDVNVRGPDGFTPLMIASCSGGGLETGNSEEEEDAPAVISDFIYQGASLHNQTDRTGET
ALHLAARYSRSDAAKRLLEASADANIQDNMGRTPLHAAVSADAQGVFQILIRNRATDLDA
RMHDGTTPLILAARLAVEGMLEDLINSHADVNAVDDLGKSALHWAAAVNNVEAAVVLLKN
GANKDMQNNKEETPLFLAAREGSYETAKVLLDHFANRDITDHMDRLPRDIAQERMHHDIV
RLLDEYSLVRSPPLHGATLGGTPTLSPPLCSPNGYLGNLKPPMQGKKARKPSTKGLACGG
KEPKDLKARRKKSQDGKGCLLDSGSVMSPVDSLESPHGYLSDVASPPLLPSPFQPSPSVP
LNHLPGMPETHLGVSHLSVAAKPEMAVLSGGSRLAFEAGPPRLSHLPVASSTSTILGSGG
SGGSGAVNFTVGGAAGLNGQCEWLSRLQNGLVPNQYNPLRGGVTPGTLSTQAAGLQHGTV
GPLHAPALSQVMTYQALPSTRLASQPHLVQPQQNLQMQPPSMPPQPNLQPHLGVSSAASG
HLGRSFLGGELSQADMQPLGPGNLAAHTVLPQDGQVLPTSLPSTLAPPTMAPPMTTAQFL
TPPSQHSYSSSPVDNTPSHQLQVPEHPFLTPSPESPDQWSSSSPHSNISDWSEGISSPPT
SVPSQIAHVPEAFK
```

Prosite, Smart, PFAM

# UniProtKB - F1MSM3 (F1MSM3_BOVIN)

**Display**

Entry

Publications

Feature viewer

Feature table

None

BLAST  Align  Format  Add to basket  History

Protein | Submitted name: **Uncharacterized protein**

Gene | **NOTCH1**

Organism | *Bos taurus (Bovine)*

Status | Unreviewed - Annotation score: ●●●●● - Protein predicted[i]

## GO - Molecular function[i]

- calcium ion binding  Source: InterPro
- chromatin DNA binding  Source: Ensembl
- core promoter binding  Source: Ensembl
- enzyme inhibitor activity  Source: Ensembl
- receptor activity  Source: InterPro
- sequence-specific DNA binding  Source: Ensembl
- transcriptional activator activity, RNA polymerase II transcription factor binding  Source: Ensembl
- transcription factor activity, sequence-specific DNA binding  Source: Ensembl

## Family & Domains[i]

**Domains and Repeats**

| Feature key | Position(s) | Description |
|---|---|---|
| Domain[i] | 20 – 59 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 60 – 100 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 103 – 140 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 141 – 177 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 179 – 217 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 219 – 256 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 258 – 294 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 296 – 334 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 336 – 372 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 373 – 411 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 413 – 451 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 453 – 489 | EGF-like  InterPro annotation ▼ |
| Domain[i] | 491 – 527 | EGF-like  InterPro annotation ▼ |

UniProt

SIB Swiss Institute of Bioinformatics

# UniProtKB/TrEMBL
# &
# Automated annotation

Important remarks

# UniRule

UniProtKB ▾ | source:pir* OR source:hamap OR source:rulebase OR source:saas OR source:sam

Total number of records in UniProtKB

**UniProtKB**

UniProt Knowledgebase

**Swiss-Prot (558,590)**
Manually annotated and reviewed.
Records with information extracted from literature and curator-evaluated computational analysis.

**TrEMBL (126,780,198)**
Automatically annotated and not reviewed.
Records that await full manual annotation.

Records with house-made automated annotation

⭐ Reviewed (313,777)
Swiss-Prot

📄 Unreviewed (69,577,937)
TrEMBL

Only ~ 55 % of TrEMBL records contain automated annotation

# Differences between TrEMBL and Swiss-Prot

|  | TrEMBL | Swiss-Prot |
|---|---|---|
| annotation | automatic | manual |
| Annotation = complete ? | Partial annotation (~55 % of the entries) | As complete and systematic as possible |
| Set of sequences = complete ? | As complete as possible; does not contain Swiss-Prot sequences ! | Complete sets only for a few organisms |
| Number of entries | 127 000 000 | 550 000 |
| Number of species | 700 000 | 13 000 |

**When you compare biological information of given datasets of proteins beware the ratio of TrEMBL vs Swiss-Prot entries in your dataset: the results might not be only 'biological'!**

UniProt

# Automated annotation

# Set of mouse proteins with N-glycosylation – Nov 2018

Reviewed (17,001)
Swiss-Prot

Unreviewed (67,833)
TrEMBL

Popular organisms

Mouse (84,834)

Proteomes

UP000000589 (53,857)

organism:"Mus musculus (Mouse) [10090]"

Reviewed (3,706)
Swiss-Prot

Unreviewed (28)
TrEMBL

Popular organisms

Mouse (3,734)

Proteomes

UP000000589 (3,720)

annotation:(type:carbohyd "n linked glcnac ellipsis") AND organism:"Mus musculus (Mouse) [10090]"

21.7 %

0.04 %

6.0 %

UniProt

SIB
Swiss Institute of Bioinformatics

# UniProtKB sequence annotation

- Feature viewer

- Genome browser

# UniProtKB annotation & Feature viewer



https://www.uniprot.org/uniprot/P15056/protvista

# UniProtKB annotation & Feature viewer

# UniProtKB annotation & Genome browsers

http://www.uniprot.org/uniprot/Q92667

The UniProt release 2018_11 will be publicly available on **December 5th**.

# Rhea in UniProt



Improved annotation precision and coverage compared to reactions provided by the IUBMB Enzyme Classification

Improved usability, interoperability and consistency of UniProt enzyme data

Metabolic and metabolomic data mining and integration

**A link between chemistry and biology**

# What will change?



UniProt 2018_11

# Display on the UniProt website



**Catalytic activity** [i]

agmatine [i] + H2O [i] = N-carbamoylputrescine [i] + NH4(+) [i]   ⬢ 2 Publications ▾
EC:3.5.3.12   ⬢ 2 Publications ▾   « Hide Reaction
Source: RHEA:18037

agmatine                    H2O                  N-carbamoylputrescine              NH4(+)

# Display on the UniProt website

- **Under development…**

# A complete integration of data



A key to understand metabolic phenotype

# What kind of searches are now possible ?

*Searches on the UniProt website or SPARQL queries*

- What are the human proteins that catalyze reactions using D-glucose ? *(a specific compound)*

- What are the enzymes in species X involved in lipid metabolism ? *(a class of compounds)*

- What is a reaction network for a specified organism of interest ?

- Retrieve the links between genes, transcripts and proteins to relevant metabolites, in a specified organism

- Identify putative enzymes acting on a specific metabolite

**Nucleic acid sequence databases**
　　　　**INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

　　　　**UniProtKB/Swiss-Prot**
　　　　　　　**Protein sequences**
　　　　　　　**Biological knowledge**

　　　　**UniProtKB/TrEMBL**
　　　　　　　**Protein sequences**
　　　　　　　**Biological knowledge**

　　**GO annotation**

　　　　**UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# Gene Ontology

- The Gene Ontology is a controlled vocabulary, a set of standard terms—words and phrases—used for indexing and retrieving information.

- Created by the GO consortium

- Contains ~50'000 terms.

- GO also defines the relationships between the terms (hierarchy), making it a structured vocabulary.

- `http://www.geneontology.org`

UniProt

# 3 categories of GO terms

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

- Cell division
- Mitosis
- Organelle fission

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

- Protein kinase activity
- Insulin binding
- Insulin receptor activity

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

- Mitochondrion
- Mitochondrial matrix
- Mitochondrial membrane

# Each GO term has two definitions

**Gene Ontology Browser**
Term Detail

GO term: **cell differentiation**
GO id: **GO:0030154**
Definition: **The process whereby relatively unspecialized cells, e.g. embryonic or regenerative cells, acquire specialized structural and/or functional features that characterize the cells, tissues, or organs of the mature organism or some other relatively stable phase of the organism's life history.**

A textual definition written by a biologist

Gene_Ontology
  ⓅBiological_process
    ⓘcellular process
      ⓘcell communication +
      ⓘcell differentiation [GO:0030154] *(493 genes, 649 annotations)*
        ⓘadipocyte differentiation +
        ⓘantipodal cell differentiation +
        ⓘcardiac cell differentiation +
Gene_Ontology
  Ⓟbiological_process
    ⓘdevelopment
      ⓘabscission +
      ⓘaging +
      ⓘblastocyst development +
      ⓘblastocyst hatching
      ⓘcell development +
      Ⓟcell differentiation [GO:0030154] *(493 genes, 649 annotations)*
        ⓘadipocyte differentiation +

Graph structure (~hierarchy): **formal and computable**

31

UniProt

SIB
Swiss Institute of Bioinformatics

# GO databases



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

http://www.geneontology.org

# Gene Ontology (GO) annotation

# GO Annotation

1. **Experimental annotations (EXP)**

   Biocurators read papers and summarize the content with a list of appropriated GO terms

2. **Curated non experimental annotation**

   Evidence such as sequence similarity, database entries

3. **Automatically assigned annotation (IEA)**

   No curator oversight

   Inferred by Electronic Annotation: **> 98 %**

# Annotations are supported by Evidence codes



**Figure from**: Gaudet, Skunka, Hu, Dessimoz, The Gene Ontology Handbook 2017

# GO annotation: experimental (manual)

PUBMED 1: Protein kinase R (PKR) recognizes double stranded
RNA in the cytoplasm

PUBMED 2: PKR acts as a kinase in the nucleus

GO annotation:

Molecular Function:    -double stranded RNA binding
-kinase activity

Cellular Component:    -nucleus
-cytoplasm

# GO annotation: automated

## Examples of IEA Usage    >98 %

## External Mappings

→ • InterPro2GO

• EC2GO

• SwissProt Keywords

• UniProtKB Subcellular Localization

## Automated Annotation by Orthology

→ • Ensembl Compara

Experimental annotations by assigner

98 % of the GO annotation is 'non-experimental'

(1) You can have different set of GO terms for a same gene depending on the database

# UniProtKB/Swiss-Prot: P0DMC2 (ELA_DANRE)

**GO - Biological process** [i]

- apelin receptor signaling pathway — Source: UniProtKB ▾
- cell migration involved in mesendoderm migration — Source: UniProtKB ▾
- endoderm development — Source: UniProtKB ▾
- heart development — Source: UniProtKB ▾
- mesendoderm migration — Source: UniProtKB ▾

# UniProtKB/TrEMBL: A0A0R4IR01_DANRE

**GO - Biological process** [i]

- angioblast cell migration from lateral mesoderm to midline — Source: ZFIN ▾
- apelin receptor signaling pathway — Source: ZFIN ▾
- cell migration involved in gastrulation — Source: ZFIN ▾
- cell migration involved in mesendoderm migration — Source: ZFIN ▾
- chordate embryonic development — Source: ZFIN ▾
- endodermal cell differentiation — Source: ZFIN ▾

# Panther: P0DMC2 (ELA_DANRE)

## GENE ONTOLOGY DATABASE ANNOTATIONS ② Close

GO MF Complete: -

GO BP Complete: chordate embryonic development, apelin receptor signaling pathway, cell migration involved in gastrulation, cell migration involved in mesendoderm migration, endodermal cell differentiation, angioblast cell migration from lateral mesoderm to midline

GO CC Complete: extracellular space

UniProt

SIB
Swiss Institute of
Bioinformatics

# (2) The 'quantity' and 'quality' of GO annotation varie from one species to another



http://amigo.geneontology.org/amigo/base_statistics

# (2) The 'quantity' and 'quality' of GO annotation varie from one species to another



| | | | | | |
|---|---|---|---|---|---|
| ☐ | P33709 | EPO_SHEEP | ⭐ | **Erythropoietin** | Ovis aries (Sheep) | apoptotic process; cell proliferation; cellular hyperosmotic response; embryo implantation; erythrocyte differentiation; erythrocyte maturation; erythropoietin-mediated signaling pathway; hemoglobin biosynthetic process; negative regulation of calcium ion transport into cytosol; negative regulation of cation channel activity; negative regulation of erythrocyte apoptotic process; negative regulation of intrinsic apoptotic signaling pathway in response to osmotic stress; negative regulation of transcription from RNA polymerase II promoter; peptidyl-serine phosphorylation; positive regulation of cell proliferation; positive regulation of DNA replication; positive regulation of Ras protein signal transduction; positive regulation of transcription, DNA-templated; positive regulation of tyrosine phosphorylation of STAT protein; response to hypoxia |
| ☐ | Q2XNF5 | EPO_DANRE | ⭐ | **Erythropoietin** | Danio rerio (Zebrafish) (Brachydanio rerio) | erythrocyte maturation; hemopoiesis; nucleate erythrocyte development; response to activity |
| ☐ | G9JKG7 | G9JKG7_HUMAN | | **Erythropoietin** | Homo sapiens (Human) | acute-phase response; aging; apoptotic process; embryo implantation; erythrocyte maturation; hemoglobin biosynthetic process; negative regulation of myeloid cell apoptotic process; negative regulation of neuron death; peptidyl-serine phosphorylation; positive regulation of activated T cell proliferation; positive regulation of ERK1 and ERK2 cascade; positive regulation of neuron projection development; regulation of transcription from RNA polymerase II promoter; response to axon injury; response to dexamethasone; response to electrical stimulus; response to estrogen; response to hyperoxia; response to hypoxia; response to interleukin-1; response to lipopolysaccharide; response to salt stress; response to testosterone; response to vitamin A |

## Where do they come from ?

# (3) timeliness of GO annotation…





Impact of outdated gene annotations on pathway enrichment analysis
(August 2016)
Analysis of 75 mutated glioblastoma (GBM) genes using annual annotations from 2009–2016
Pathway analysis assesses the statistical enrichment of biological processes and pathways in a given gene list on the basis of information in Gene Ontology (GO) and pathway databases such as Reactome and PathwayCommons. GO is updated daily and Reactome versions are released quarterly, but many software tools interpret gene lists using functional information that has not been updated for years.
http://www.nature.com/nmeth/journal/v13/n9/full/nmeth.3963.html

**Nucleic acid sequence databases**
        **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

        **UniProtKB/Swiss-Prot**
                **Protein sequences**
                **Biological knowledge**

        **UniProtKB/TrEMBL**
                **Protein sequences**
                **Biological knowledge**

        **GO annotation**

➡️ **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# UniProtKB Proteomes

## UniProtKB
UniProt Knowledgebase

### Swiss-Prot (558,590)
⭐ Manually annotated and reviewed.

Records with information extracted from literature and curator-evaluated computational analysis.

### TrEMBL (126,780,198)
Automatically annotated and not reviewed.

Records that await full manual annotation.

## UniRef
The UniProt Reference Clusters (UniRef) provide clustered sets of sequences from the UniProt Knowledgebase (including isoforms) and selected UniParc records.

## UniParc
UniParc is a comprehensive and non-redundant database that contains most of the publicly available protein sequences in the world.

## Proteomes
A proteome is the set of proteins thought to be expressed by an organism. UniProt provides proteomes for species with completely sequenced genomes.

## Supporting data

| Literature citations | Taxonomy | Subcellular locations |
|---|---|---|
| Cross-ref. databases | Diseases | Keywords |
| | XXX | |

# Proteomes / Reference proteomes

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

- **A proteome is formed from all UniProtKB/Swiss-Prot entries + those UniProtKB/TrEMBL entries mapping to Ensembl (genomes).**

- Pseudogenes and other dubious uncharacterized ORFs are removed

**R** 16,547 Reference proteomes

169,798 Other proteomes

## Superkingdom

95,840 Bacteria

87,840 Viruses

933 Archaea

1,731 Eukaryota

UniProt

SIB
Swiss Institute of
Bioinformatics

# Human proteins

(records)

# Human proteome



Reviewed (20,316)
Swiss-Prot

Unreviewed (141,875)
TrEMBL

Popular organisms

Human (162,191)  ✖

Proteomes

UP000005640 (71,772)

Reviewed (20,303)
Swiss-Prot

Unreviewed (51,469)
TrEMBL

Popular organisms

Human (71,772)  ✖

Proteomes

UP000005640 (71,772)  ✖

Query: organism:"Homo sapiens (Human) [9606]"

Query: organism:"Homo sapiens (Human) [9606]" AND proteome:up000005640

⚠ **+ Swiss-Prot isoforms**: ~22 000 (link)

| | Entry | Entry name | | Protein names |
|---|---|---|---|---|
| ☐ | P0DOX6 | IGM_HUMAN | ⭐ | **Immunoglobulin mu heavy chain** |
| ☐ | P0DOX5 | IGG1_HUMAN | ⭐ | **Immunoglobulin gamma-1 heavy chain** |
| ☐ | P0DOX4 | IGE_HUMAN | ⭐ | **Immunoglobulin epsilon heavy chain** |
| ☐ | P0DOX3 | IGD_HUMAN | ⭐ | **Immunoglobulin delta heavy chain** |
| ☐ | P0DOX8 | IGL1_HUMAN | ⭐ | **Immunoglobulin lambda-1 light chain** |
| ☐ | P0DOX2 | IGA2_HUMAN | ⭐ | **Immunoglobulin alpha-2 heavy chain** |
| ☐ | P0DOX7 | IGK_HUMAN | ⭐ | **Immunoglobulin kappa light chain** |
| ☐ | P69208 | MORN_HUMAN | ⭐ | **Morphogenetic neuropeptide** |
| ☐ | P22103 | PNEU_HUMAN | ⭐ | **Pneumadin** |
| ☐ | P02728 | GLEM_HUMAN | ⭐ | **Erythrocyte membrane glycopeptide** |
| ☐ | P02729 | GLUR_HUMAN | ⭐ | **Urine glycopeptide** |
| ☐ | P01358 | GAJU_HUMAN | ⭐ | **Gastric juice peptide 1** |
| ☐ | P01858 | TUFT_HUMAN | ⭐ | **Phagocytosis-stimulating peptide** |

**NOT proteome:up000005640** AND reviewed:yes AND organism:"Homo sapiens (Human) [9606]"

UniProt

# Proteomes / Reference proteomes

- Some proteomes have been (manually and algorithmically) selected as **reference proteomes** (useful for biomedical research and phylogeny).

- Regularly updated

http://www.uniprot.org/help/proteome
http://www.uniprot.org/help/reference_proteome

# Proteomes / Reference proteomes



The 'redundant' proteomes are excluded from UniProtKB (since 2015).
The corresponding protein sequences are only available in UniParc …

Proteomes 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
Andrea.Auchincloss@sib.swiss or Ivo.Pedruzzi@sib.swiss

# Downloads

# Downloads

## UniProtKB

Parent directory

Reviewed (Swiss-Prot) / FAQ — Flat file (.dat)

Unreviewed (TrEMBL) / FAQ — Flat file

Isoform sequences / FAQ — Flat file

Taxonomic divisions / README — Flat file

Reference proteomes / README — Fasta format (.fasta)

Pan proteomes / README — Fasta format

ID mapping / README

Proteomics mapping / README

Variants / README

Genome annotation tracks / README

Documents

XML schema

Does not include «Swiss-Prot isoforms»

http://www.uniprot.org/downloads

# Query + Download: homo sapiens

updated

| HUMAN (TaxID: 9606) **WEB** | Swiss-Prot (canonical) | Swiss-Prot (isoforms) | TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (without Swiss-Prot isoforms) |
|---|---|---|---|---|---|
| **Total** | 20 316* | 22 014 | 141 875[#] | 184 205 | 162 191 |
| **Proteome** | 20 303 | 22 014 | 51 469 | **93 786** | 71 772 |

*uniprot_sprot_human.dat.gz          [#]uniprot_trembl_human.dat.gz

| HUMAN **FTP** (README) **Reference proteome** | Main fasta Swiss-Prot & TrEMBL (« canonical ») | Additional fasta Swiss-Prot & TrEMBL (« isoforms ») | Swiss-Prot & TrEMBL (canonical + isoforms) |
|---|---|---|---|
| **Total** | 20 998* | 72 788[#] | **93 786** |

*UP000005640_9606.fasta.gz

[#]UP000005640_9606_additional.fasta.gz

UniProt release 2018_02

UniProt

* # ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/README

SIB Swiss Institute of Bioinformatics

# Query + Download: mus musculus

| MOUSE (TaxID: 10090) WEB | Swiss-Prot (canonical) | Swiss-Prot (isoforms) | TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (without Swiss-Prot isoforms) |
|---|---|---|---|---|---|
| Total | 16 877 | 8 140 | 64 732 | 89 749 | 81 609 |
| Proteome | 16 873 | 8 138 | 34 070 | **59 081** | 50 943 |

| MOUSE FTP (README) Reference proteome | Main fasta Swiss-Prot & TrEMBL (« canonical ») | Additional fasta Swiss-Prot & TrEMBL (« isoforms ») | Swiss-Prot & TrEMBL (canonical + isoforms) |
|---|---|---|---|
| Total | 22 281* | 36 800# | **59 081** |

*UP000000589_10090.fasta.gz

#UP000000589_10090_additional.fasta.gz

UniProt release 2017_05

UniProt

* # ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/README

SIB Swiss Institute of Bioinformatics

# Query + Download: ARATH

| ARATH (TaxID: 3702) WEB | Swiss-Prot (canonical) | Swiss-Prot (isoforms) | TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (canonical + isoforms) | Swiss-Prot & TrEMBL (without Swiss-Prot isoforms) |
|---|---|---|---|---|---|
| Total | 15 333 | 2 122 | 74 101 | 91 546 | 89 434 |
| Proteome | 15 260 | 2 122 | 23 969 | **41 341** | 39 229 |

| ARATH FTP (README) Reference proteome | Main fasta Swiss-Prot & TrEMBL « canonical » | Additional fasta Swiss-Prot & TrEMBL « isoforms » | Swiss-Prot & TrEMBL (canonical + isoforms) |
|---|---|---|---|
| Total | 27 510* | 13 831# | **41 341** |

*UP000006548_3702.fasta.gz

#UP000006548_3702_additional.fasta.gz

UniProt release 2017_05

UniProt

* # ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/README

Swiss Institute of Bioinformatics

# UniProt website and tools

# The UniProt web site – www.uniprot.org

- Powerful search engine, google-like and easy-to-use, but also supports very directed field searches

- Entry views, search result views and downloads are customizable

- The URL of a result page reflects the query; all pages and queries are **bookmarkable**, supporting programmatic access

- Tools: Blast, Align, Retrieve/Idmapping, Peptide search

# Result pages: highly customizable (also available for Blast)

# Result page: download



**Different formats (fasta, txt, excell, RDF, etc.)**

# Highlight sequence annotation in alignment (BLAST or multiple alignment)

# Peptide search



http://www.uniprot.org/peptidesearch/

# help@uniprot.org

# How can you increase the impact of your research papers and contribute to UniProt?

**Nucleic acid sequence databases**
        **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

        **UniProtKB/Swiss-Prot**
                **Protein sequences**
                **Biological knowledge**

        **UniProtKB/TrEMBL**
                **Protein sequences**
                **Biological knowledge**

        **GO annotation**

        **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**

# Major 'general' protein sequence database 'sources'

Ensembl

PIR          PRF

TPA    PDB

integrated databases 'cross-references'

**UniProtKB: Swiss-Prot + TrEMBL**

databases are kept separated

**NCBI-nr: Swiss-Prot + TrEMBL + GenPept + PIR + PDB + PRF + RefSeq + TPA**

not complete !!!
(only entries created before 2007 ?)

Ensembl and Refseq: gene prediction

# NCBI nr - Entrez 'protein'

# NCBI-nr

- GenPept (source: GenBank; translated CDS)

- RefSeq

- TPA (third part annotation)


- Swiss-Prot (does not include isoform sequences)

- PIR (not updated since 2003)

- PRF (journal scan of 'published' peptide)

- PDB  (Protein Data Bank, 3D structure)

- TrEMBL (some entries….)

NCBI-nr: Swiss-Prot + GenPept + PIR + PDB + PRF + RefSeq + TPA

# GenPept

**Translation from annotated CDS in GenBank**
Contains all translated CDS annotated in
GenBank/EMBL/DDBJ sequences

- equivalent to UniProtKB/TrEMBL,
except that it is
redundant with other databases
(Swiss-Prot, RefSeq, PIR….)

**GenPept**: 'translations from all annotated coding regions (CDS) in GenBank

```
LOCUS       AF312033_10              192 aa            linear   ROD 10-DEC-2009
DEFINITION  EPO [Mus musculus].
ACCESSION   AAK28825 AAK28053
VERSION     AAK28825.1  GI:13517500
DBSOURCE    accession AF312033.1
KEYWORDS    .
SOURCE      Mus musculus (house mouse)
  ORGANISM  Mus musculus
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Glires; Rodentia;
            Sciurognathi; Muroidea; Muridae; Murinae; Mus; Mus.
REFERENCE   1  (residues 1 to 192)
  AUTHORS   Wilson,M.D., Riemer,C., Martindale,D.W., Schnupf,P., Boright,A.P.,
            Cheung,T.L., Hardy,D.M., Schwartz,S., Scherer,S.W., Tsui,L.-C.,
            Miller,W. and Koop,B.F.
  TITLE     Comparative analysis of the gene-dense ACHE/TFR2 region on human
            chromosome 7q22 with the orthologous region on mouse chromosome 5
  JOURNAL   Nucleic Acids Res. 29 (6), 1352-1365 (2001)
  PUBMED    11239002
REFERENCE   2  (residues 1 to 192)
  AUTHORS   Wilson,M.D. and Koop,B.F.
  TITLE     Direct Submission
  JOURNAL   Submitted (06-OCT-2000) Biology, Centre for Environmental Health,
            University of Victoria, P.O. Box 3020, Victoria, B.C. V8W 3N5,
            Canada
REFERENCE   3  (residues 1 to 192)
  AUTHORS   Wilson,M.D., Martindale,D.W., Schnupf,P. and Koop,B.F.
  TITLE     Direct Submission
  JOURNAL   Submitted (06-OCT-2000) Biology, Centre for Environmental Health,
            University of Victoria, P.O. Box 3020, Victoria, B.C. V8W 3N5,
            Canada
COMMENT     On Dec 10, 2009 this sequence version replaced gi:13492038.
            Method: conceptual translation supplied by author.
FEATURES             Location/Qualifiers
     source          1..192
                     /organism="Mus musculus"
                     /strain="129/Sv"
                     /db_xref="taxon:10090"
                     /chromosome="5"
     Protein         1..192
                     /product="EPO"
     Region          29..192
                     /region_name="EPO_TPO"
                     /note="Erythropoietin/thrombopoietin; pfam00758"
                     /db_xref="CDD:189705"
     CDS             1..192
                     /gene="Epo"
                     /coded_by="complement(join(AF312033.1:202181..202336,
                     AF312033.1:203077..203256,AF312033.1:203583..203669,
                     AF312033.1:204132..204274,AF312033.1:204830..204842))"
ORIGIN
        1 mgvperptll lllslllipl glpvlcappr licdsrvler yileakeaen vtmgcaegpr
       61 lsenitvpdt kvnfyawkrm eveeqaievw qglsllseai lqaqallans sqppetlqlh
      121 idkaisglrs ltsllrvlga qkelmsppdt tppaplrtlt vdtfcklfrv yanflrgklk
      181 lytgevcrrg dr
//
```

Annotation according to the submitter

No GO term !

# RefSeq

Produced by NCBI and NLM

http://www.ncbi.nlm.nih.gov/RefSeq/

http://www.ncbi.nlm.nih.gov/books/bookres.fcgi/handbook/ch18.pdf

FAQ: http://www.ncbi.nlm.nih.gov/books/NBK50679/
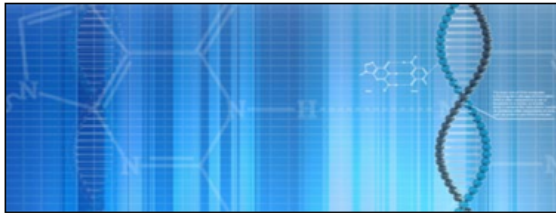
# RefSeq: NCBI Reference Sequence Database

A comprehensive, integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and protein.

## Using RefSeq

About RefSeq

Human Reference Genome

Prokaryotic RefSeq Genomes

FAQ

NCBI Handbook

Factsheet

## RefSeq Access

Human Genome Resources and Download

RefSeq FTP

RefSeq genomes FTP

New RefSeq genomic (last 30 days)

New RefSeq transcripts (last 30 days)

New RefSeq proteins (last 30 days)

Searching for RefSeq records (Queries)

## RefSeq projects

Consensus CDS (CCDS)

RefSeq Functional Elements

RefSeqGene

Targeted Loci

Virus Variation

## Announcements

**January 12, 2018**
**RefSeq Release 86 is available for FTP**

This release includes:

Proteins:     102,133,844
Transcripts:  21,370,778
Organisms:    75,218
Available at:  ftp://ftp.ncbi.nlm.nih.gov/refseq/release/
Documentation: Release Notes

See previous announcements, follow NCBI on Twitter, or subscribe to NCBI's refseq-announce mail list to receive announcements.

## Related Links

Assembly

Gene

Genome

Genome Data Viewer

Annotated Eukaryotic Genomes

## Feedback & Credits

Publications and Citing RefSeq

Contact RefSeq Help Desk

Contact CCDS Help Desk

Submit a GeneRIF

Collaborators

http://www.ncbi.nlm.nih.gov/refseq/

# RefSeq

**RefSeq**: The Reference Sequence (RefSeq) collection aims to provide a comprehensive, integrated, non-redondant set of sequences, including <span style="color:red">genomic DNA, transcript (RNA), and protein products</span>, for major research organisms.

one mRNA sequence -> one entry

*Different entries for identical protein sequences*

## Announcements

### January 12, 2018
### RefSeq Release 86 is available for FTP

This release includes:

Proteins:        102,133,844
Transcripts:     21,370,778
Organisms:       75,218
Available at:    ftp://ftp.ncbi.nlm.nih.gov/refseq/release/
Documentation:Release Notes

See previous announcements, follow NCBI on Twitter, or subscribe to NCBI's refseq-announce mail list to receive announcements.

GenPept ▾

# erythropoietin precursor [Homo sapiens]

NCBI Reference Sequence: NP_000790.2

Identical Proteins    FASTA    Graphics

Go to: ☑

```
LOCUS       NP_000790                193 aa            linear   PRI 08-MAR-2018
DEFINITION  erythropoietin precursor [Homo sapiens].
ACCESSION   NP_000790
VERSION     NP_000790.2
DBSOURCE    REFSEQ: accession NM_000799.3
KEYWORDS    RefSeq.
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   1  (residues 1 to 193)
  AUTHORS   Nishimura K, Matsumoto R, Yonezawa Y and Nakagawa H.
  TITLE     Effect of quercetin on cell protection via erythropoietin and cell
            injury of HepG2 cells
  JOURNAL   Arch. Biochem. Biophys. 636, 11-16 (2017)
   PUBMED   29080630
  REMARK    GeneRIF: these results suggested that quercetin's cytoprotective
            effects in HepG2 cells are mediated via EPO production.
REFERENCE   2  (residues 1 to 193)
  AUTHORS   Flamme I, Ellinghaus P, Urrego D and Kruger T.
  TITLE     FGF23 expression in rodents is directly induced via erythropoietin
            after inhibition of hypoxia inducible factor proline hydroxylase
  JOURNAL   PLoS ONE 12 (10), e0186979 (2017)
   PUBMED   29073196
  REMARK    GeneRIF: EPO dependent regulation pathway of FGF23 gene expression
            Publication Status: Online-Only
```

AC number

Protein: NP_
mRNA: NM_
DNA: NC_

references

# Accession number

| Accession prefix | Molecule type | Comment |
| --- | --- | --- |
| AC_ | Genomic | Complete genomic molecule, usually alternate assembly |
| NC_ | Genomic | Complete genomic molecule, usually reference assembly |
| NG_ | Genomic | Incomplete genomic region |
| NT_ | Genomic | Contig or scaffold, clone-based or WGS[a] |
| NW_ | Genomic | Contig or scaffold, primarily WGS[a] |
| NS_ | Genomic | Environmental sequence |
| NZ_[b] | Genomic | Unfinished WGS |
| NM_ | mRNA | |
| NR_ | RNA | |
| XM_[c] | mRNA | Predicted model |
| XR_[c] | RNA | Predicted model |
| AP_ | Protein | Annotated on AC_ alternate assembly |
| NP_ | Protein | Associated with an NM_ or NC_ accession |
| YP_[c] | Protein | |
| XP_[c] | Protein | Predicted model, associated with an XM_ accession |
| ZP_[c] | Protein | Predicted model, annotated on NZ_ genomic records |

[a] Whole Genome Shotgun sequence data.

[b] An ordered collection of WGS sequence for a genome.

[c] Computed.

UniProt

COMMENT    REVIEWED REFSEQ: This record has been curated by NCBI staff. The
reference sequence was derived from X02157.1, S65458.1 and
AC009488.5.
This sequence is a reference standard in the RefSeqGene project.
On Apr 6, 2005 this sequence version replaced NP_000790.1.

Summary: This gene encodes a secreted, glycosylated cytokine
composed of four alpha helical bundles. The encoded protein is
mainly synthesized in the kidney, secreted into the blood plasma,
and binds to the erythropoietin receptor to promote red blood cell
production, or erythropoiesis, in the bone marrow. Expression of
this gene is upregulated under hypoxic conditions, in turn leading
to increased erythropoiesis and enhanced oxygen-carrying capacity
of the blood. Expression of this gene has also been observed in
brain and in the eye, and elevated expression levels have been
observed in diabetic retinopathy and ocular hypertension.
Recombinant forms of the encoded protein exhibit neuroprotective
activity against a variety of potential brain injuries, as well as
antiapoptotic functions in several tissue types, and have been used
in the treatment of anemia and to enhance the efficacy of cancer
therapies. [provided by RefSeq, Aug 2017].

Publication Note:  This RefSeq record includes a subset of the
publications that are available for this gene. Please see the Gene
record to access additional publications.

##Evidence-Data-START##
Transcript exon combination :: X02157.1, BC093628.1 [ECO:0000332]
RNAseq introns              :: single sample supports all introns
                               SAMEA2158188, SAMEA2159368
                               [ECO:0000348]

##Evidence-Data-END##

**Full support of intron position by RNA-seq alignment evidence used in automatic assertion** ECO:0000348

http://purl.obolibrary.org/obo/ECO_0000348

A type of full support of intron position by RNA-seq alignment evidence that is used in an automatic assertion.

**Ontology:** Evidence ontology ECO

Annotation
No GO term !

Sequence origin
+ ECO code

# RefSeq

| | manual annotation |
|---|---|
| GENOME ANNOTATION | No |
| INFERRED | No |
| MODEL | No |
| PREDICTED | No |
| PROVISIONAL | No |
| REVIEWED | **Yes** (sequence + functional information and features) |
| VALIDATED | **Yes** (initial sequence) |
| Whole Genome Sequencing (WGS) | No |

http://www.ncbi.nlm.nih.gov/RefSeq/

UniProt

SIB
Swiss Institute of
Bioinformatics

# Automated annotation: CDD

```
                         propagated from UniProtKB/Swiss-Prot (P04150.1)"
     Region              98..115
                         /region_name="Required for high transcriptional activity
                         of isoform Alpha-C3. {ECO:0000269|PubMed:23820903}"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="propagated from UniProtKB/Swiss-Prot (P04150.1)"
     Site                113
                         /site_type="other"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="Phosphoserine. {ECO:0000250|UniProtKB:P06537};
                         propagated from UniProtKB/Swiss-Prot (P04150.1)"
     Site                134
                         /site_type="other"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="Phosphoserine. {ECO:0000244|PubMed:18669648};
                         propagated from UniProtKB/Swiss-Prot (P04150.1)"
     Site                141
                         /site_type="other"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="Phosphoserine. {ECO:0000250|UniProtKB:P06537};
                         propagated from UniProtKB/Swiss-Prot (P04150.1)"
     Site                203
                         /site_type="other"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="Phosphoserine. {ECO:0000244|PubMed:24275569,
                         ECO:0000269|PubMed:12000743, ECO:0000269|PubMed:18483179,
                         ECO:0000269|PubMed:25847991}; propagated from
                         UniProtKB/Swiss-Prot (P04150.1)"
     Site                211
                         /site_type="other"
                         /experiment="experimental evidence, no additional details
                         recorded"
                         /note="Phosphoserine. {ECO:0000269|PubMed:12000743,
                         ECO:0000269|PubMed:18483179, ECO:0000269|PubMed:25847991};
                         propagated from UniProtKB/Swiss-Prot (P04150.1)"
```

Annotation
- automated (CDD)
- derived from Swiss-Prot
- in-house (ab initio)

UniProtKB/Swiss-Prot:     One gene -> one entry (9 isoforms)
RefSeq:                   One mRNA sequence -> one entry

## Results

1 out of 1 identifier from UniProtKB AC/ID was successfully mapped to 15 RefSeq Protein IDs.

Format

| From | To |
|------|-----|
| P04637 | NP_000537.3 |
| P04637 | NP_001119584.1 |
| P04637 | NP_001119585.1 |
| P04637 | NP_001119586.1 |
| P04637 | NP_001119587.1 |
| P04637 | NP_001119588.1 |
| P04637 | NP_001119589.1 |
| P04637 | NP_001119590.1 |
| P04637 | NP_001263624.1 |
| P04637 | NP_001263625.1 |
| P04637 | NP_001263626.1 |
| P04637 | NP_001263627.1 |
| P04637 | NP_001263628.1 |
| P04637 | NP_001263689.1 |
| P04637 | NP_001263690.1 |

ID/AC mapping: see later

# Query: organism:"Homo sapiens (Human) [9606]"



Reviewed (20,171)
Swiss-Prot

Number of canonical and isoform protein sequences: 42,147 (download data in FASTA format)

Unreviewed (136,539)
TrEMBL

## Popular organisms

Human (156,710) ✖

## Proteomes

UP000005640 (70,952)

## RefSeq

| Nucleotide | | |
|---|---|---|
| | Nucleotide ▾ | "Homo sapiens"[Organism] |
| | | Create alert   Advanced |

| | |
|---|---|
| **Species** | Summary ▾  20 per page ▾  Sort by Default order ▾      Send: ▾ |
| Animals (199,675) | |
| Customize ... | **Items: 1 to 20 of 199675** |
| **Molecule types** | << First  < Prev  Page 1  of 9984  Next >  Last >> |
| genomic DNA/RNA (21,643) | |
| mRNA (111,384) | ⓘ Filters activated: RefSeq. Clear all |
| rRNA (27) | |
| Customize ... | ☐  **Homo sapiens** kinase D interacting substrate 220 (KIDINS220), transcript variant 16, non-coding |
| **Source databases**    clear | 1.  RNA |
| INSDC (GenBank) (0) | 8,692 bp linear transcribed-RNA |
| ✓ RefSeq (199,675) | Accession: NR_145965.1  GI: 1149123051 |
| Customize ... | GenBank   FASTA   Graphics |

UniProt

SIB
Swiss Institute of Bioinformatics

# NCBI nr

## query & BLAST

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein

# Look for human FOXP2 @ NCBI nr

# BLAST @ NCBI

# BLAST @ NCBI

**Sequences producing significant alignments:**

| Accession | Description | Max score | Total score | Query coverage | E value | Max ident | Links |
|---|---|---|---|---|---|---|---|
| NP_000790.2 | erythropoietin precursor [Homo sapiens] >sp\|P01588.1\|EPO_HUMAN F | 392 | 392 | 100% | 2e-137 | 100% | U G M |
| CAA26095.1 | erythropoietin [Homo sapiens] | 391 | 391 | 100% | 7e-137 | 99% | G |
| ACI13657.1 | SP(EPO)-EPO(Nat)-ELP [synthetic construct] >gb\|ACI13658.1\| SP(EI | 396 | 396 | 100% | 2e-136 | 99% | |
| XP_519268.2 | PREDICTED: erythropoietin [Pan troglodytes] | 389 | 389 | 100% | 4e-136 | 99% | U G M |
| CAA26094.1 | unnamed protein product [Homo sapiens] >gb\|ACI13661.1\| secEPO [ | 389 | 389 | 100% | 4e-136 | 99% | G M |
| ACI13653.1 | SP(EPO)-EPO(Nat) [synthetic construct] >gb\|ACI13654.1\| SP(EPO)- | 390 | 390 | 100% | 9e-136 | 99% | |
| ACI13659.1 | SP(EPO)-EPO(R103)-ELP [synthetic construct] | 394 | 394 | 100% | 1e-135 | 99% | |
| ACI13655.1 | SP(EPO)-EPO(R103) [synthetic construct] | 387 | 387 | 100% | 6e-135 | 99% | |
| ACI13663.1 | secEPO-GFP [synthetic construct] | 394 | 394 | 100% | 3e-134 | 99% | |
| ACI13664.1 | SP(EPO)-EPO(Nat)-GFP [synthetic construct] | 394 | 394 | 100% | 3e-134 | 99% | |
| EAW76494.1 | erythropoietin, isoform CRA_b [Homo sapiens] | 382 | 382 | 97% | 1e-131 | 100% | G |
| AAI43226.1 | EPO protein [Homo sapiens] | 345 | 345 | 100% | 8e-119 | 99% | G M |
| ACJ06770.1 | erythropoietin [synthetic construct] | 340 | 340 | 86% | 6e-117 | 100% | |
| XP_003278152.1 | PREDICTED: erythropoietin-like [Nomascus leucogenys] | 340 | 340 | 100% | 8e-117 | 97% | G M |
| ACI13646.1 | SP(Tob)-EPO(Opt) [synthetic construct] | 341 | 341 | 97% | 8e-117 | 89% | |
| ACI13645.1 | SP(Tob)-EPO(Nat) [synthetic construct] | 341 | 341 | 97% | 1e-116 | 89% | |
| ACI13656.1 | SP(EPO)-EPO(Agly) [synthetic construct] | 340 | 340 | 100% | 3e-116 | 97% | |
| ACI13662.1 | chEPO [synthetic construct] | 339 | 339 | 89% | 8e-116 | 97% | |
| ACI13647.1 | SP(Tob)-EPO(R103) [synthetic construct] | 338 | 338 | 97% | 2e-115 | 89% | |
| ACI13650.1 | SP(Tob)-EPO(Opt)-ELP [synthetic construct] | 340 | 340 | 97% | 2e-114 | 89% | |
| ACI13665.1 | SP(Tob)-EPO(Nat)-GFP [synthetic construct] | 343 | 343 | 97% | 3e-114 | 89% | |
| 1CN4_C | Chain C, Erythropoietin Complexed With Extracellular Domains Of Eryt | 333 | 333 | 86% | 3e-114 | 98% | S |
| ACI13660.1 | SP(EPO)-EPO(Agly)-ELP [synthetic construct] | 340 | 340 | 100% | 3e-114 | 97% | |
| ACI13649.1 | SP(Tob)-EPO(Nat)-ELP [synthetic construct] | 340 | 340 | 97% | 4e-114 | 89% | |
| ACI13651.1 | SP(Tob)-EPO(R103)-ELP [synthetic construct] | 338 | 338 | 97% | 2e-113 | 89% | |
| ACI13648.1 | SP(Tob)-EPO(Agly) [synthetic construct] | 331 | 331 | 97% | 1e-112 | 87% | |
| 1EER_A | Chain A, Crystal Structure Of Human Erythropoietin Complexed To Its | 328 | 328 | 86% | 4e-112 | 97% | S |

## UniProtKB/Swiss-Prot alternative isoform sequences are not included !

UniProt

SIB Swiss Institute of Bioinformatics

# NCBI nr 'cluster' of identical sequences

# UniProtKB entries at NCBI…

# A UniProtKB/Swiss-Prot entry with the NCBI look

RecName: Full=Carbonic anhydrase 2; AltName: Full=Carbonic anhydrase II; Short=CA-II; AltName: Full=Carbonate dehydratase II; AltName: Full=Carbonic anhydrase C; Short=CAC

Swiss-Prot: P00918.2

FASTA    Graphics

Comment    Features    Sequence

```
LOCUS       CAH2_HUMAN                  260 aa            linear   PRI 10-AUG-2010
DEFINITION  RecName: Full=Carbonic anhydrase 2; AltName: Full=Carbonic
            anhydrase II; Short=CA-II; AltName: Full=Carbonate dehydratase II;
            AltName: Full=Carbonic anhydrase C; Short=CAC.
ACCESSION   P00918
VERSION     P00918.2  GI:115456
DBSOURCE    UniProtKB: locus CAH2_HUMAN, accession P00918;
            class: standard.
            extra accessions:B2R7G8,Q6FI12,Q96ET9
            created: Jul 21, 1986.
            sequence updated: Jan 23, 2007.
            annotation updated: Aug 10, 2010.
            xrefs: M77181.1, AAA51909.1, M77176.1, M77177.1, M77178.1,
            M77179.1, M77180.1, Y00339.1, CAA68426.1, X03251.1, CAA27012.1,
            J03037.1, AAA51908.1, CR536526.1, CAG38763.1, CR541875.1,
            CAG46673.1, AK312978.1, BAG35815.1, CH471068.1, EAW87136.1,
            BC011949.1, AAH11949.1, M36532.1, AAA51911.1, CRHU2, NP_000058.1,
            12CA_A, 1A42_A, 1AM6_A, 1AVN_A, 1BCD_A, 1BIC_A, 1BN1_A, 1BN3_A,
            1BN4_A, 1BNM_A, 1BNN_A, 1BNQ_A, 1BNT_A, 1BNU_A, 1BNV_A, 1BNW_A,
            1BV3_A, 1CA2_A, 1CA3_A, 1CAH_A, 1CAI_A, 1CAJ_A, 1CAK_A, 1CAL_A,
            1CAM_A, 1CAN_A, 1CAO_A, 1CAY_A, 1CAZ_A, 1CCS_A, 1CCT_A, 1CCU_A,
            1CIL_A, 1CIM_A, 1CIN_A, 1CNB_A, 1CNC_A, 1CNG_A, 1CNH_A, 1CNI_A,
            1CNJ_A, 1CNK_A, 1CNW_A, 1CNX_A, 1CNY_A, 1CRA_A, 1CVA_A, 1CVB_A,
            1CVC_A, 1CVD_A, 1CVE_A, 1CVF_A, 1CVH_A, 1DCA_A, 1DCB_A, 1EOU_A,
            1F2W_A, 1FQL_A, 1FQM_A, 1FQN_A, 1FQR_A, 1FR4_A, 1FR7_A, 1FR7_B,
            1FSN_A, 1FSN_B, 1FSQ_A, 1FSQ_B, 1FSR_A, 1FSR_B, 1G0E_A, 1G0F_A,
            1G1D_A, 1G3Z_A, 1G45_A, 1G46_A, 1G48_A, 1G4J_A, 1G4O_A, 1G52_A,
            1G53_A, 1G54_A, 1H4N_A, 1H9N_A, 1H9Q_A, 1HCA_A, 1HEA_A, 1HEB_A,
            1HEC_A, 1HED_A, 1HVA_A, 1I8Z_A, 1I90_A, 1I91_A, 1I9L_A, 1I9M_A,
            1I9N_A, 1I9O_A, 1I9P_A, 1I9Q_A, 1IF4_A, 1IF5_A, 1IF6_A, 1IF7_A,
            1IF8_A, 1IF9_A, 1KWQ_A, 1KWR_A, 1LG5_A, 1LG6_A, 1LGD_A, 1LUG_A,
            1L7V_A, 1MOO_A, 1MUA_A, 1OKL_A, 1OKM_A, 1OKN_A, 1OO5_A, 1RAY_A,
```

UniProt

SIB
Swiss Institute of
Bioinformatics

**Important remarks concerning the datasets**

**Different servers…**

**UniProtKB/TrEMBL entries are not available at NCBI**
The same protein sequence might be present,
but not with the UniProtKB/TrEMBL AC (**with some exceptions**)
(not the case for UniProtKB/Swiss-Prot entries)

# UniProtKB/TrEMBL entries are not available at NCBInr with some exceptions…

# ID/AC mapping

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

# http://www.ebi.ac.uk/Tools/picr/

**Understanding protein function** is critical to research in many areas of science such as biology, medicine and biotechnology.

**Keeping up with all of this information** is a daunting task for most researchers.

UniProt helps with this in the following ways:

- it provides an **up-to-date**, comprehensive body of protein information at a single site;
- it aids scientific discovery by collecting, **interpreting and organising** this information so that it is easy to access and use;
- it saves researchers countless hours of work in monitoring and **collecting** this information themselves;
- it provides tools to help with **protein sequence analysis**;
- it **provides links** to related information in more than 150 other biological databases to help you access additional information in more specialised collections.

- https://www.ebi.ac.uk/training/online/course/uniprot-exploring-protein-sequence-and-functional/why-do-we-need-uniprot

# Thank you !

Thanks to Emmanuel Boutet and Marc Feuermann
for some of the slides !

Thanks to Diana Marek & Geoff Fucile
for the organisation of this course

**Nucleic acid sequence databases**
              **INSDC, RefSeq, Ensembl**

**Protein sequence databases**

**UniProtKB**

        **UniProtKB/Swiss-Prot**
              **Protein sequences**
              **Biological knowledge**

        **UniProtKB/TrEMBL**
              **Protein sequences**
              **Biological knowledge**

        **GO annotation**

        **UniProt Proteomes**

**NCBI protein (RefSeq)**

**Practicals**