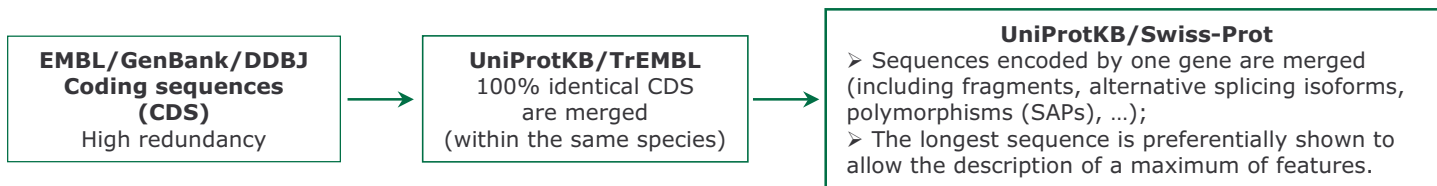


One of the main goals of UniProtKB/Swiss-Prot is to provide the scientific community with the most accurate protein sequence and the description of all its variations.



Manual sequence annotation

1. In the merging process, **UniProtKB/TrEMBL accession** numbers are kept, as well as **references** (publications or submissions) and **cross-references** to EMBL/GenBank/DBJ; **Protein sequencing data** is integrated, if available;
2. Discrepancies (SAPs, alternative promoter usage, alternative splicing, RNA editing, alternative translation initiation or unknown reasons) are analyzed and documented at **single amino acid level**;
3. Sequence comparison between **related organisms** (e.g. human-mouse-rat) allows confirmation of sequence length, sometimes alternative splicing isoforms or SAPs;
4. Once the sequence corrected, high-performance bioinformatics tools are run and predicted domains, post-translational modifications (PTMs), etc. are critically reviewed before integration.

```
ISOFORM_1      VLAAGLVLSVFVAIGFEFYIKSQKNNDIEQAFCHF-----YGLQCKQHTPTNSTSGTTLST
ISOFORM_2      VLAAGLVLSVFVAIGFEFYIKSQKNNDIEQLSFNAIMEELGISLKNQKIKKKSR7TKGKS
                *****..*          *:.*: :.:.*.:.:
ISOFORM_1      DLECGKLIREERGIRKQSSVHTV
ISOFORM_2      SFTSILTCHQRRTRQRTKETVA---
                :.:. :. : * **:.:
                :.:. :. : * **:.:
```

Literature, cDNA vs genome alignment, intron/exon junction check (BLAT)

- **ALTERNATIVE PRODUCTS:** 2 named isoforms [FASTA] produced by alternative splicing. Additional isoforms seem to exist.

VAR_SEQ 870 918 AFCFFYGLQCKQTHPTNSTSGTTLSTDLECGKLIREERGI RKQSSVHTV ->
CLSFNAIMEELGISLKNQKKIKKKSRTKGKSSFTSILTCH QRRTRKETVA (in
isoform 2).

RNA editing annotated in:

References: *RP* line, e.g. *RNA EDITING OF POSITIONS 636*.

Comments: topic *RNA EDITING*.

Keywords: *RNA editing.*

Sequence features: e.g. FT VARIANT 636 636 Q -> R (in RNA edited version). /FTId=VAR_000304

Sequence: the edited sequence is shown in the entry, when fully edited.

Alternative sequences annotated in:

References: *RP* line, e.g. NUCLEOTIDE SEQUENCE [MRNA] (ISOFORM 2).

Comments: topic *ALTERNATIVE PRODUCTS*; contains information about alternative promoter usage, alternative splicing and alternative translation initiation.

Keywords: *Alternative splicing, Alternative promote usage, Alternative initiation*

Sequence features: topic *FT VAR_SEQ* characterized by a unique identifier (/FTId=VSP_000128).

GENOMIC_SEQ	PYEYNNPHPCNPDSVDVNNFTLLNSFWFGVGALMQQSGELMPK
EDITED_MRNA	PYEYNNPHPCNPDSVDVNNFTLLNSFWFGVGALMRQSGELMPK

Literature and cDNA vs genome alignment

- **RNA EDITING:** Modified positions=636; Note=Partially edited.

VARIANT	636	636	1	Q -> R (in RNA edited version).
---------	-----	-----	---	---------------------------------

```

GENOMIC_SEQ      VLAAGL VLSVFVAIGEFYIKSRKNNDIEQAFCCFFYGLQCKQTHPTN
SEQ_1            VLAAGL VLSVFVAIGEFYIKSQKNNDIEQVFCFFYGLQCKQTHPTN
SEQ_2            VLAAGL VLSVFVAIGEFYIKSRKNNDIEQAFCCFFYGLQCKQTHPTN
                  *****

```

Literature and validated SNP from other databases including dbSNP

VARIANT	862	862	1	R -> Q.
VARIANT	870	870	1	A -> V (in dbSNP:363503) [NCBI/Ensembl].

Polymorphisms annotated in:

References: *RP* line, e.g. *VARIANT VAL-98*.

Comments: topic *POLYMORPHISM* and/or *DISEASE*.

Keywords: Polymorphism and/or Disease mutation.

Sequence features: topic *FT VARIANT* characterized by a unique identifier (FTId=VAR_123456). Cross-references to dbSNP.

CAA62631 MAAFSVGTAMNASSYSAEMTEPKSVCVSVDEVVSSNMEATEETDLLNGH
AAK28839 -----MT EPKSVCVSVDEVVTSNMEATEETDLLNGH
AAK28841 MAAFSVGTAMNASSYSAEMTEPKSVCVSVDEVVTSNMEATEETDLLNGH

CONFLICT 38 38 S -> T (in Ref. 5; [AAK28839/AAK28841](#)).

Conflicts

Discrepancies of unknown origin, ranging from sequencing errors to yet uncharacterized polymorphisms.

Cross-references to nucleic acid sequence databases

Various tags in DR EMBL/GenBank/DBJ lines allow the description of discrepancies between the submitted CDS and the UniProtKB/Swiss-Prot sequence:

ALT_INIT, ALT_TERM: erroneous start or stop;

ALT_FRAME: presence of frameshift(s) in the submitted nucleotide sequence. Linked to comment CAUTION: e.g. *Ref.X (PID) sequence differs from that shown due to a frameshift in position Y;*

ALT_SEQ: erroneous gene/CDS prediction or combination of several problems. Linked to comment CAUTION: e.g. *Ref.x (PID) sequence differs from that shown due to erroneous gene model prediction.*