# Sample Scenario (17, translated)

# Evidence for evolution by comparing proteins in different species

## Summary

## 2 Procedure

2.1 Find the sequences of the CFTR protein in several species in the UniProtKB database
*UniProtKB is a database that contains all protein sequences listed from public data (~230 million proteins).*
*N.B.: to study evolution, protein sequences are often used: they are more relevant from an evolutionary point of view. Here we are working with the CFTR protein, which causes cystic fibrosis when it is faulty, but other proteins can be used. (cf. underlined selected proteins)*

Open [UniProtKB](#) ; in Query, type gene: CFTR.

*You have to search for proteins by their gene name (1)*

UniProt will search for all entries in which the term CFTR *solution* the list is huge.
Select entries "reviewed" *solution (Insulin)* solution CFTR *(as of 23.09.22).*

(1) *N.B.: The name of the protein (cf. list of underlined selected proteins) often differs from the name of the gene. It is best to search for a protein by its gene name, as gene names are 'standardized' by expert committees. UniProtKB is a protein database: the information is focused on protein sequences but the list of gene names is exhaustive. The same gene can sometimes have several names. Example: the CFTR gene has a synonym: abcc7*

### 2.2 Selecting the same protein in several organisms

Sample results obtained by searching for 'INS' as a gene name.

Select (tick) from the list the proteins of the species that interests you.
*A minimum of 5 protein sequences is required to construct a sequence alignment that makes sense. (cf.list of selected proteins).*
*It may be wise to impose a few species on each group to facilitate comparisons (e.g., human, chimpanzee, mouse, rat, bovine, horse, Danio rero (fish), fruit fly) and let them choose a few others at will.*
*[…] The selected sequences are marked with a blue rectangle at the beginning of the line and their number appears in the blue band at the top: e.g. "21 rows selected out of 97".*
*[…]*

## 2.3 Building an alignment

- Click the "Align" button above the strip at the top of the list. A window displays the sequences that will be aligned in FASTA format.
- Click the button at the bottom right Align XX sequences. In principle do not touch anything After a certain time (several seconds, even minutes) a "Tool results" window appears displaying the alignments carried out during the session.

- • Click "Completed" to display the alignment in a new page. <u>Example for CFTR</u>



Alignment of a few CFTR protein sequences from various species - similarity enabled.

Aligned sequences are displayed in table, one sequence (species) per line, 66 per line. (In Overview mode = default)
On a line below, the protein signal is shown in red.
*NB: "-" (dash) means that the alignment program introduced a space (insertion/deletion or "gap") to be able to align sequences of different lengths.*
*To find the name of the species, click on the accession number in blue -> a new window opens with the entry UniProtKB corresponding to this protein and the full name of the species.*

**Nota Bene**

The « Similarity » highlight above the table on the left is enabled by default: similar regions are highlighted (purple) - allowing discussion of those that are more evolutionarily conserved.

Highlight "Physical properties" then "Hydrophobicity": highlights regions likely to be transmembrane

**Additional information and possible questions**

The physico-chemical properties of the different amino acids are described *<u>here</u>*

A correspondence table of 3-letter and 1-letter amino acid and codon codes can be accessed _here_.

**A phylogenetic tree?**

At the top of the page, instead of the "overview" visualization, you can click "Trees" to display trees with different visualizations...

_Note that it is a guided tree, used by the program to construct the alignment and is based solely on observed differences between sequences. It is not a proper phylogenetic tree that is much more complex to establish! See for example scenario. Phylogeny, biodiversity and pizza..._

See also scenario 19 _Find the date of evolutionary divergence of two species_

# 3 _Possible student questions:_

- Find the organisms closest to (human, mouse, cow, ...) for each of the proteins studied by the class.
  Try to find another explanation than the common origin and the independent evolution since the separation.
- Why use "protein" sequences rather than nucleotides in this evolutionary context?
- Are there sequences that are similar in a large number of organisms?
- Observing that the sequences of insulin, EPO, a Histone, CFTR, etc. are very similar in many species: what are the possible explanations for this fundamental similarity?
- Insulin varies in many regions, while Histone (H4) is almost identical between very many organisms. perform its functions in DNA winding and regulation.
- Do similar sequences between a large number of organisms indicate that there are no mutations in those places?
- Ask the students to find the time when the species studied separated and compare with the tree obtained. There is no site offering this, but a search on the internet often makes it possible to find an estimate of the last common ancestor between two species.

_Discussion of these alignments can illustrate the fundamental concepts of common origin and divergence common to various proteins._

Scenario established with scientific advice from Dr. Marie-Claude Blatter of SIB Swiss Institute of Bioinformatics