Name: Thi Hanh Nguyen Ly (0772489)

**Practical Computing of Bioinformatics – Assignment 1**

**Question 1.**

1. The IDs from **ID_list.txt** is from UniProt database of protein sequences.
2. The IDs are derived from Caspase-3 protein sequence from **CASP3** gene of human (**Homo Sapiens**) lying on chromosome 4.
3. The protein resides in the cytoplasm, nucleus, nucleoplasm, cytosol, plasma membrane of the cell.
   It associates mostly apoptosis processes. More specifically, the protein links to immune system process, autophagy, cell adhesion, programmed cell death, reproductive process, signalling, cell differentiation, protein catabolic process, anatomical structure development, nervous system process, protein maturation, …
4. Information for caspase 3 gene in mouse (Mus musculus) can be found under **EBI Search** in ebi.ac.uk website. Further information about this gene can be reviewed when clicking "View in Ensembl".
   Under "gene-based display" panel, when clicking to the option **Phenotypes,** a list of all phenotypes related to Casp3 gene of mouse can be found. Each phenotype will be associated with its corresponding allele. Thus, we can search for the **cleft palate** phenotype and its corresponding allele. In this case, allele **Casp3_hith2** can be found (Figure 1).
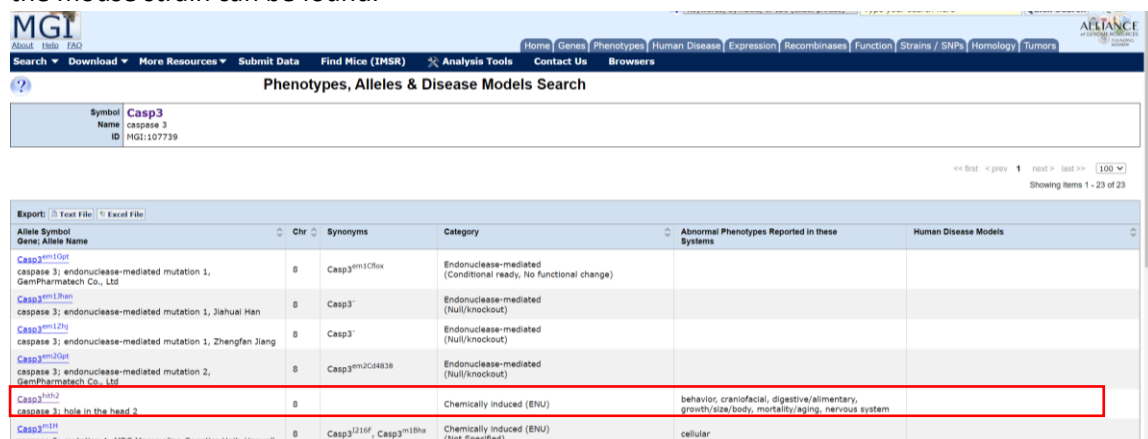


Figure 1. Allele for cleft palate development of Casp3 gene in mouse.

Information in Figure 1 also includes a link to **MGI source**, an online database of mouse genome informatics where all mutations and alleles information can be found. When choosing the **Casp3$^{hith2}$ allele** (Figure 2), an information of how the allele was introduced in the mouse strain can be found.

In figure 3, it is said in the mutation details that **Casp3hith2 allele** is "identified in an ENU screen and mapped to a 4 Mb interval on Chromosome 8 containing Casp3. A mutation two bases upstream of exon 3 results in splicing errors causing exon 3 splicing to be 4 nucleotides into the wild type exon or loss of exon 3 altogether. Both forms result in premature stop codons."



Figure 3. Mutation description of Casp3hith2 allele.

5. Retrieving GO term in BioMart for Casp3 gene of Human:
   Casp3 gene of Human has reference gene id of ENSG00000164305 (Ensembl) (Figure 4)



Figure 4. Retrieving GO-terms for ENSG00000164305

The reference gene id of Casp3 gene is input in the **Filters** part when searching in BioMart. GO-term accession and GO-domain are chosen for the output attributes. In this list, the GO-terms for GO-domain "biological process" are: GO:0006508, GO:0006915, GO:0006508.

**Question 2.**

1.  BLAST algorithm from https://blast.ncbi.nlm.nih.gov/ (BLOSUM-62 matrix) is used to analyse the sequence in the file **unknown_sequence.fsa.** The result is shown to support that the sequence belongs to the **Insulin-like receptor protein** of **Caenorhabditis elegans** (Roundworm species) (Figure 5).



Figure 5. BLAST analysis of unknown_sequence.fsa

2.  When searching for the **Accession** "NP_497650.4" in UniProtKB, the result for this protein sequence is Q968Y9 - INSR_CAEEL. It is shown that this protein sequence is derived from gene **daf-2 (**Figure 6**)**



Figure 6. Genomic information of NP_497650.4

Searching for Homologous genes in NCBI for "Caenorhabditis elegans daf-2" (HomoloGene database) results in **IGF1R (insulin-like growth factor 1 receptor) gene** for Homo Sapiens (Human) and **InR (insulin-like receptor) gene** for D. melanogaster (Fruit fly).

3.  The three protein sequences extracted from UniProt are aligned with **t-coffee algorithm** in ebi.ac.uk website. The result is then viewed in MView (Figure 7)

t-coffee result: https://www.ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=tcoffee-I20231112-172932-0756-99822783-p1m&analysis=alignments

MView result: https://www.ebi.ac.uk/Tools/services/rest/mview/result/mview-I20231112-173008-0791-33417150-p1m/aln-html

```
Reference sequence (1): sp|P08069|IGF1R_HUMAN
Identities normalised by aligned length.
Colored by: identity

                                cov     pid   1 [        .         .         .         .         :         .         .         . 80
1 sp|P08069|IGF1R_HUMAN 100.0% 100.0%     MKSGS-----------------------------------------------------------------------
2 sp|P09208|INSR_DROME   97.0%  24.2%     MFNMPRGVTKSKSKRGKIKMENDMAAAATTTACTLGHICVLCRQEMLLDTCCCRQAVEAVDSPASSEEAYSSSNSSSCQA
3 sp|Q968Y9|INSR_CAEEL   97.1%  23.0%     MTRMNIV-----------------------------------------RCRRRHKILENLEEENLGPSCSSTTSTTA
  consensus/100%                          Mhphs.......................................................................
  consensus/90%                           Mhphs.......................................................................
  consensus/80%                           Mhphs.......................................................................
  consensus/70%                           Mhphs.......................................................................

                                cov     pid  81          .         1         .         .         .         .         :         . 160
1 sp|P08069|IGF1R_HUMAN 100.0% 100.0%     ---------------------------------------------------------------------------
2 sp|P09208|INSR_DROME   97.0%  24.2%     SSEISAEEVWFLSHDDIVLCRRPKFDEVETTGKKRDVKCSGHQCSNECDDGSTKNNRQQRENFNIFSNCHNILRTLQSLL
3 sp|Q968Y9|INSR_CAEEL   97.1%  23.0%     -----AT---------------------------------------EALGTTTEDMRLKQQRSSSRATEHDIVD------
  consensus/100%                          ...........................................................................
  consensus/90%                           ...........................................................................
  consensus/80%                           ...........................................................................
  consensus/70%                           ...........................................................................

                                cov     pid 161          .         .         .         2         .         .         .         . 240
1 sp|P08069|IGF1R_HUMAN 100.0% 100.0%     ---------------------------------------------------------------------------
2 sp|P09208|INSR_DROME   97.0%  24.2%     LLMFNCGIFNKRRRRQHQQQHHHHYQHHHQQHHQQHHQRQQANVSYTKFLLLLQTLAAATTRLSLSPKNYKQQQQLQHNQ
3 sp|Q968Y9|INSR_CAEEL   97.1%  23.0%     -----------------------GNHHDDEHITMRR---------------------------------LRLVKNSR
  consensus/100%                          ...........................................................................
  consensus/90%                           ...........................................................................
  consensus/80%                           ...........................................................................
  consensus/70%                           ...........................................................................

                                cov     pid 241          :         .         .         .         .         3         .         . 320
1 sp|P08069|IGF1R_HUMAN 100.0% 100.0%     ---------------------------GGGSPTSLYGLLFLSAALSLWPTS------------------------
2 sp|P09208|INSR_DROME   97.0%  24.2%     QLPRATPQQKQQEK-----DRHKCFHYKHNYSYSPGISLLLFILLANTLAIQAVVLPAHQQHLLHNDIADGLDKTALSVS
3 sp|Q968Y9|INSR_CAEEL   97.1%  23.0%     TRRRTTPDSSMDCYEENPPSQKTSINYSWISKKSSMTSLMLLLFAFVQPCASIV-------------------------
  consensus/100%                          .............................thus.hSLhhhlhhu.s.sh.shs.......................
  consensus/90%                           .............................thus.hSLhhhlhhu.s.sh.shs.......................
  consensus/80%                           .............................thus.hSLhhhlhhu.s.sh.shs.......................
  consensus/70%                           .............................thus.hSLhhhlhhu.s.sh.shs.......................

                                cov     pid 321          .         .         .         :         .         .         .         . 4 400
1 sp|P08069|IGF1R_HUMAN 100.0% 100.0%     ---------------------GEICGPGIDIRNDYQQL----------------KRLENCIVIEGYLHILLISKAE----
2 sp|P09208|INSR_DROME   97.0%  24.2%     GTQSRWTRSESNPTMRLSQNVKPCKS-MDIRNMVSHF---------------NQLENCIVIEGFLLIDLINDAS----
3 sp|Q968Y9|INSR_CAEEL   97.1%  23.0%     --------------------EKRCGP-IDIRNRPWDIKPQWSKLGDPNEKDLAGQRMVCIVVEGSLTISFVLKHKTKAQ
  consensus/100%                          ...............c.Cts.hDIRN...ph.............pph.NCIVIEG.LhI.hl.ctp....
  consensus/90%                           ...............c.Cts.hDIRN...ph.............pph.NCIVIEG.LhI.hl.ctp....
  consensus/80%                           ...............c.Cts.hDIRN...ph.............pph.NCIVIEG.LhI.hl.ctp....
  consensus/70%                           ...............c.Cts.hDIRN...ph.............pph.NCIVIEG.LhI.hl.ctp....
```
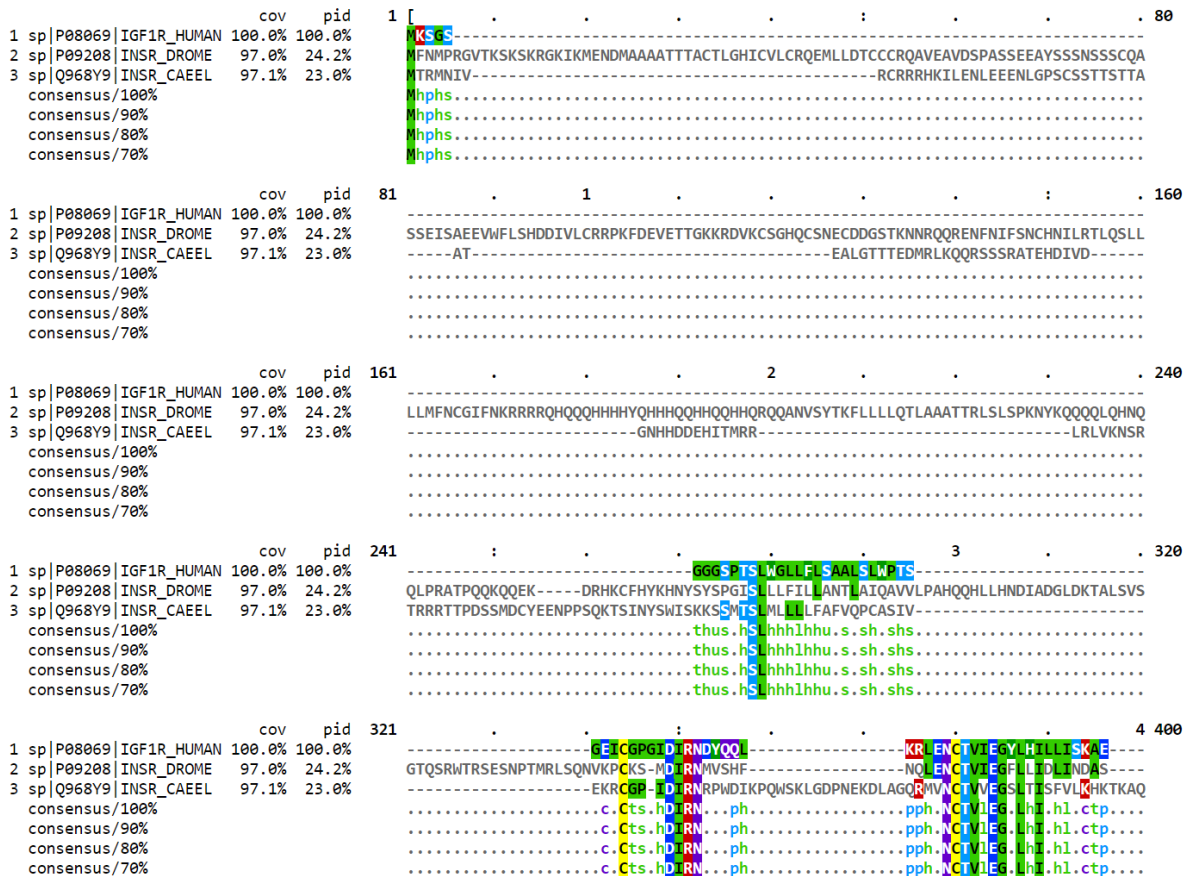
Figure 7. Partial MView result from multiple alignment.

According to Figure 7, the conservation percentage reflects the degree of homology in the three protein sequences. Thus, if we take the sequence of human as a reference, then fruit fly and roundworm have around 97% similarity in their protein sequence.

## Question 3 - Linux

The result from Linux assignment is store in Linux folder. The folder contains analyze_GFF_features.sh, analyzeGFF_job.slurm, the two error, output log files of the slurm job and the gff3 file of chromosome 1. There is also another folder call "text files result" containing the text file output from running the code of the bash file.