

# Bioinformatics Project 3 - Clustering and Phylogeny

09.12.2024

The goal of the project is to conduct a comprehensive phylogenetic analysis and submit a report with a detailed description of the methods applied, the datasets used, and the results obtained.

## 1. Clustering:

- (a) Choose an algorithm to work with.
- (b) Create your own dataset:
  - i. choose 8 different human proteins,
  - ii. run BLAST (online version) on these proteins, and for each protein, download the sequences for seven different organisms. Be careful not to select identical sequences—refer to the tutorial below for guidance,
  - iii. remember to include the identifiers of the sequences used in your report, for example, in the form of a table,
- (c) use BLAST (local version) or another comparison method (e.g., MSA) to obtain input data (similarity between sequences). Select the appropriate input format for your chosen clustering method,
- (d) cluster the provided sequences.
- (e) Check the results. Do the clusters correspond to similar proteins identified by BLAST?

## 2. Phylogenetics:

- (a) use [this](#) module to build trees in three different ways:
  - i. separate tree for each "group" of proteins
  - ii. separate tree for each cluster
  - iii. one common tree for all downloaded sequences
- (b) create consensus trees from your two approaches (clusters vs. the groups of proteins downloaded together).
- (c) create different visualizations (from the two consensus trees and the single tree generated from all sequences) to compare the results and draw conclusions.
  - i. Color tree branches based on the organism
  - ii. Color tree branches based on the protein "group"
- (d) What is your observation? Which approach seems to work best?
- (e) Are the trees similar to each other? Does the evolution look exactly the same in different trees?

### 3. Report

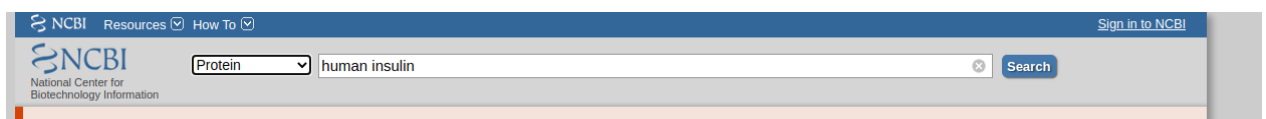
- (a) The grade will be mainly based on the description of the research tasks and results.
- (b) Prepare a detailed report of your work.
- (c) Include the following points:
  - i. Description of your dataset (organisms and proteins that you chose).
  - ii. Short description of the methods and algorithms that you used.
  - iii. Description of the steps of your analysis.
  - iv. Description of your results and conclusions from the project.
  - v. Figures showing consensus trees and examples of trees from your groups and clusters (see 2e).
  - vi. Tree generated from all the sequences (see 2e).
  - vii. Answers to questions 2f and 2g.



### 4. Grading

- (a) max 10 points
- (b) Dataset preparation 2p
- (c) Clustering 3p
- (d) Constructing trees inside groups and clusters 3p
- (e) Consensus trees 2p
- (f) All of the above will be graded based on both your description in the report and included code.
- (g) five points can be deducted for a careless and inaccurate report

### 5. Deadline: 20/12/2024 23:59

#### 1. Search for chosen protein on NCBI:



GENE Was this helpful?  

**INS – insulin**

[Homo sapiens \(human\)](#)

Also known as: IDDM, IDDM1, IDDM2, ILPR, IRDN, MODY10, PNDM4

Gene ID: 3630

[RefSeq transcripts \(4\)](#) [RefSeq proteins \(4\)](#) [RefSeqGene \(2\)](#) [PubMed \(952\)](#)

[Orthologs](#) [Genome Data Viewer](#) [BLAST](#) [Download](#)

**RefSeq Sequences** +

**Items: 1 to 20 of 12411**

<< First < Prev Page 1 of 621 Next > Last >>

- ☐ [insulin, isoform 2 precursor \[Homo sapiens\]](#)
- 1. **200 aa protein**  
Accession: NP\_001035835.1 GI: 109148522  
[BioProject](#) [Nucleotide](#) [PubMed](#) [Taxonomy](#)  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [RecName: Full=Insulin; Contains: RecName: Full=Insulin B chain; Contains: RecName: Full=Insulin A chain; Flags: Precursor](#)
- 2. **110 aa protein**  
Accession: P01308.1 GI: 124617  
[PubMed](#) [Taxonomy](#)  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [insulin \[Homo sapiens\]](#)
- 3. **110 aa protein** ← →  
Accession: AAA59172.1 GI: 386828  
[Nucleotide](#) [PubMed](#) [Taxonomy](#)  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)

**Find related data**

Database: Select

[Find items](#)

**Search details**





human insulin[Protein Name] OR (human[All Fields] AND insulin[All Fields])

[Search](#) [See more...](#)

**Recent activity**


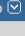
[Turn Off](#) [Clear](#)

human insulin (12411) Protein


-  Promoter Capture Hi-C: High-resolution, Genome-wide Profiling of Promoter Intera...
-  Automatic analysis and 3D-modelling of Hi-C data using TADbit reveals structural...
-  Comparison of computational methods for Hi-C data analysis
-  Modeling chromosomes: beyond pretty pictures

[See more...](#)

## 2. Run BLAST on chosen protein.

NCBI Resources  How To  Sign in to NCBI

Protein Protein Advanced [Search](#) [Help](#)

**COVID-19 Information** 


[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#) | [Español](#)

GenPept Send to: Change region shown Customize view

**insulin [Homo sapiens]**

GenBank: AAA59172.1

[Identical Proteins](#) [FASTA](#) [Graphics](#)

[Go to:](#) 

LOCUS AAA59172 110 aa linear PRI 10-JUN-2016

DEFINITION insulin [Homo sapiens].

ACCESSION AAA59172

VERSION AAA59172.1

DBSOURCE accession [AH002844.2](#)

KEYWORDS .

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 110)

AUTHORS Bell, G.I., Pictet, R. and Rutter, W.J.

TITLE Analysis of the regions flanking the human insulin gene and sequence of an Alu family member

JOURNAL Nucleic Acids Res. 8 (18), 4091-4109 (1980)

PUBMED [6253909](#)

COMMENT On Aug 28, 1993 this sequence version replaced gi:186432.

METHOD Method: conceptual translation.

**Analyze this sequence**


**Run BLAST** ← →

[Identify Conserved Domains](#)

[Highlight Sequence Features](#)

[Find in this Sequence](#)

**Protein 3D Structure**

 Monoclinic human insulin in complex with p-coumaric acid

PDB: 6TC2

Source: Homo sapiens

Method: X-ray Diffraction

Resolution: 1.36 Å

[See all 258 structures...](#)

## 3. Filter the results of Blast. We don't want to have identical results and our goal is to search for sequences from different organisms than the original.



Descriptions Graphic Summary Alignments Taxonomy

**Sequences producing significant alignments**

☐ select all 6 sequences selected

GenPept

Download New Select columns Show 100

FASTA (complete sequence)  
FASTA (aligned sequences)  
GenBank (complete sequence)  
Hit Table (text)  
Hit Table (CSV)  
Text  
Descriptions Table (CSV)  
XML  
ASN.1

Multiple alignment New MSA Viewer

Query cover	E value	Per. Ident	Acc. Len	Accession
100%	9e-71	96.36%	110	<a href="#">XP_011896559.1</a>
100%	1e-70	96.36%	110	<a href="#">XP_023039285.1</a>
100%	4e-70	96.36%	110	<a href="#">XP_025211013.1</a>
100%	5e-70	96.36%	147	<a href="#">XP_023039287.1</a>
100%	2e-69	96.36%	147	<a href="#">XP_025211011.1</a>
100%	4e-69	92.98%	114	<a href="#">AA072172.1</a>
100%	2e-61	89.09%	98	<a href="#">XP_018891846.1</a>
100%	1e-60	88.18%	98	<a href="#">XP_024110668.1</a>

insulin (Tupaia chinensis)

insulin isoform X1 [Ptilocobus tephrosceles]  
insulin isoform X4 [Theropithecus gelada]  
insulin isoform X2 [Ptilocobus tephrosceles]  
☒ insulin isoform X2 [Theropithecus gelada]  
☐ synthetic preproinsulin [synthetic construct]  
☒ insulin isoform X2 [Gorilla gorilla gorilla]  
☒ insulin isoform X3 [Pongo abelii]

```

>XP_025211011.1 insulin isoform X2 [Theropithecus gelada]
MAGLLKRLGVSPGAPGQGTWPSAGLRPACLPGHCP SAMALWMRLPLLALLALWGPDSPAFVNQHL CGSHLVEALYLVC
GERGFFYTPKTRREAEDPQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSI CSLYQLENYCN
>XP_018891846.1 insulin isoform X2 [Gorilla gorilla gorilla]
MALWMRLLPLLALLALWGPDPAAAFVNQHL CGSHLVEALYLVCGERGFFYTPKTRREAEDLQGS LQPLALEGSLQKRGIV
EQCCTSI CSLYQLENYCN
>XP_024110668.1 insulin isoform X3 [Pongo abelii]
MALWMRLLPLLALLALWGPDPAAAFVNQHL CGSHLVEALYLVCGERGFFYTPKTRREAEDLQGS LQPLALEGSLQKRGIV
EQCCTSI CSLYQLENYCN
>XP_032748073.1 insulin-1 [Rattus rattus]
MALWMRFLPLLALLVWEPKPAQAFVKQHL CGPHLVEALYLVCGERGFFYTPKSRREVEDPQVPQLELGGGPEAGDLQTL
ALEVARQKRGIVDQCCTSI CSLYQLENYCN
>ABB89749.1 preproinsulin 2 [Mus caroli]
MALWMRFLPLVALLFLWESHPTQAFVKQHL CGSHLVEALYLVCGERGFFYTPMSRREVEDPQVAQLELGGGPGAGDLQTL
ALEVAQQKRGIVDQCCTSI CSLYQLENYCN
>XP_032747262.1 insulin-2 [Rattus rattus]
MALWIRFLPLLALLVLWEPRPAQAFVKQHL CGSHLVEALYLVCGERGFFYTPVSRREVEDPQVAQLELGGGPGAGDLQTL
ALEVARQKRGIVDQCCTSI CSLYQLENYCN

```

- Remember to also download the reference sequence (Homo sapiens insulin in this case)
- Repeat the process with different proteins. This time search only for sequences in previously chosen organisms. This way you will get 8 different proteins (or similar proteins) from the same 7 organisms.