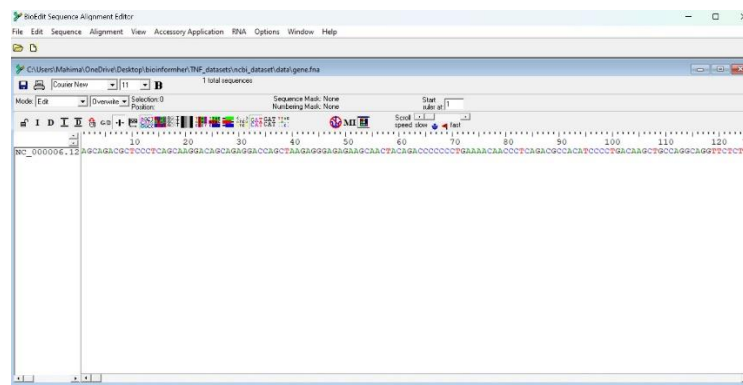# Comprehensive Analysis of the Human Tumor Necrosis Factor (TNF) Gene Sequence

**Objective:** This project aims to perform a comprehensive analysis of the Human Tumor Necrosis Factor (TNF) gene sequence using various bioinformatics tools such as NCBI, BioEdit, PROMO, GENSCAN, and MEME Suite. This mini-project is the first task of the BioinformHER Module 1. BioinformHER is an initiative taken by HackBio aimed at introducing young women in STEM to Bioinformatics using simple and easy-to-follow guidelines.

## Tasks and interpretations

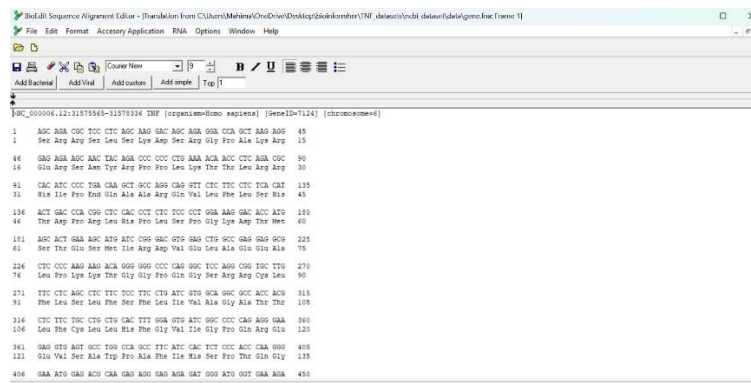### Task 1: Download a Biological Sequence from NCBI and View/Edit It

1. Accessed the NCBI homepage at NCBI
2. Searched for the human TNF gene using the term 'human TNF gene.'
3. Located the correct sequence record (e.g., 'Homo sapiens TNF')
4. There were multiple sequences available for the human TNF gene. Therefore, I selected the most recent sequence [NC_000006.12 (31575565-31578336)]
5. Downloaded the sequence in FASTA format
6. Opened the sequence in BioEdit and viewed it.



Task 1 Output

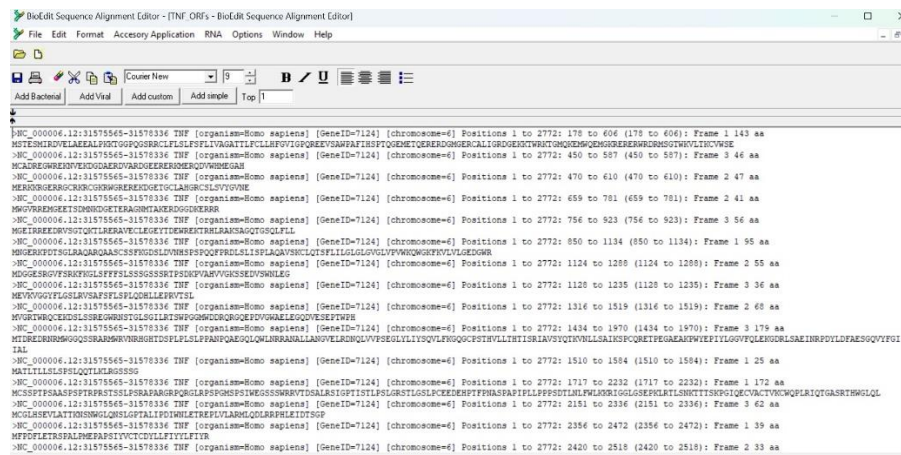### Task 2: Generate a Translation of a DNA or RNA Sequence into Amino Acids

1. Opened the downloaded TNF gene sequence in BioEdit
2. Used the 'Translate' feature in BioEdit to generate the amino acid sequence



Task 2 output

## Task 3: Find ORFs (Open Reading Frames) in a DNA or RNA Sequence

1. Used BioEdit's ORF Finder tool to find ORFs in the TNF gene sequence
2. Recorded the start and stop positions, lengths, and protein translations of the ORFs
3. Selected the longest ORF and analyzed it using SMART BLAST to determine whether it encodes the Human Tumor Necrosis Factor (TNF)
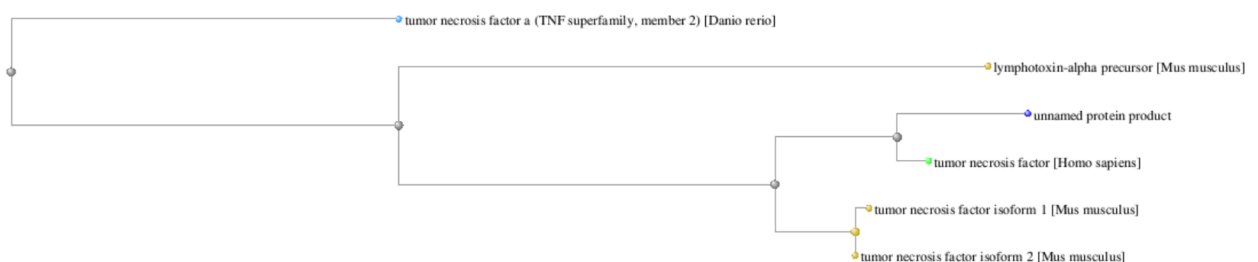


Task 3 output

**Interpretation:** I analyzed the longest ORF using NCBI's SMART BLAST to determine the protein-encoding region of the TNF gene. The ORF with a length of 179 residues showed similarity to TNF (*Homo sapiens*) indicating that it is the protein-encoding region of the TNF gene.

**>NC_000006.12:31575565-31578336 TNF [organism=Homo sapiens] [GeneID=7124] [chromosome=6] Positions 1 to 2772: 1434 to 1970 (1434 to 1970): Frame 3 179 aa**
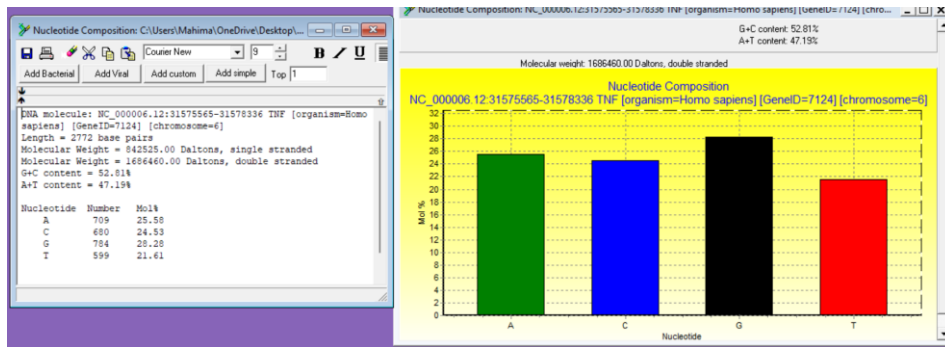
**MTDREDRNRMWGGQSSRARMWRVNRHGHTDSPLPLSLPPANPQAEGQLQWLNRRANALLANGVELR DNQLVVPSEGLYLIYSQVLFKGQGCPSTHVLLTHTISRIAVSYQTKVNLLSAIKSPCQRETPEGAEAKPWYEPIY LGGVFQLEKGDRLSAEINRPDYLDFAESGQVYFGIIAL**

The phylogenetic tree generated for the protein product is as follows:



## Task 4: Analyze Sequence Composition (Nucleotide or Amino Acid Frequencies)

1. Used BioEdit to analyze the sequence composition of the TNF gene
2. Calculated the frequencies of each nucleotide and the overall GC content
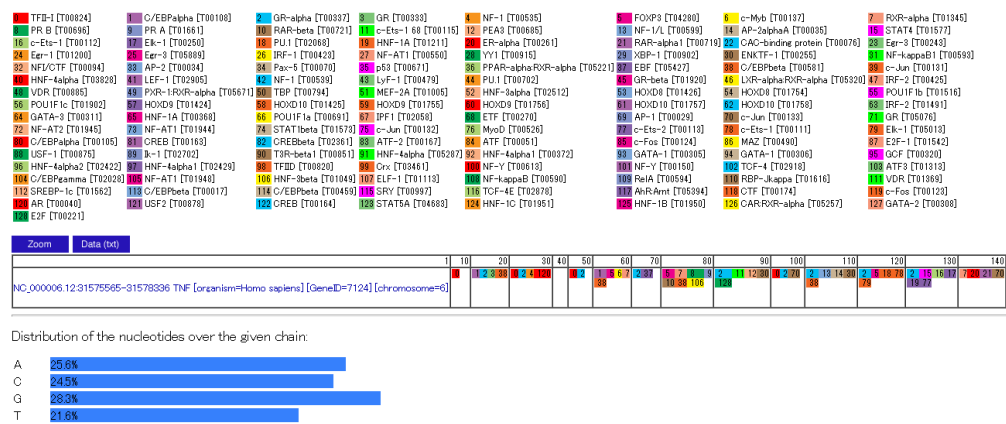3. Interpreted the results and saved the analysis

Task 4 output

**Interpretation:** The sequence has a GC content of 52.81% which is higher than the AT content indicating that it is stable because of stronger H-bonding between G and C.

## Task 5: Identify Transcription Factor Binding Sites Using the PROMO Tool

1. Accessed the PROMO tool at https://alggen.lsi.upc.es/cgi-bin/promo_v3/promo/promoinit.cgi?dirDB=TF_8.3
2. Selected 'Homo sapiens' as the species
3. Input the entire TNF gene sequence
4. Identified potential transcription factor binding sites

PROMO is a virtual laboratory for identifying putative transcription factor binding sites (TFBS) in DNA sequences from a species or groups of species of interest.



Task 5 output

Several transcription factors can bind to the promoter region of the Human TNF gene and initiate its transcription at different locations in the gene.

## Task 6: Search for Functional Motifs in a Genome or Transcriptome Using MEME Suite

1. Accessed the MEME Suite at https://meme-suite.org/meme/
2. Uploaded the TNF gene sequence in FASTA format
3. Used the default settings to search for motifs
4. Interpreted and saved the results of the motif search

Task 6 output

**Interpretation:** MEME (Multiple Em for Motif Elicitation) found the lowest E-value motifs with the highest statistical significance as shown above. Motif 1 has a width of 38 and is found at 5 different sites in the sequence with an E-value of 2.6e+000 which is the lowest, suggesting its significance.

## Task 7: Predict Coding/Non-Coding Regions in a Genome Using GENSCAN

1. Accessed the GENSCAN tool at http://hollywood.mit.edu/GENSCAN.html
2. Input the TNF gene sequence in the appropriate format
3. Ran the analysis to predict coding and non-coding regions
4. Saved and interpreted the results.



Task 7 output

**Interpretation:** GENSCAN predicted 5 exons/coding regions in the positive strand of the gene sequence.

- Exon 1 is an initiator exon that marks the beginning of the coding sequence with the start codon. It starts from residue 218 and ends at residue 403 with a length of 186 residues.
- Exon 2 and Exon 3 are internal exons
- Exon 4 is a terminal exon that includes the stop codon
- Exon 5 is the poly-A signal. It is a sequence in the 3' untranslated region (UTR) of the mRNA that signals for the addition of the poly-A tail to the mRNA

## Task 8: Convert Between Sequence File Formats Using BioEdit (FASTA to PHYLIP)

1. Opened the TNF gene sequence in BioEdit
2. Used the 'Save As...' feature to convert the file to PHYLIP format
3. Verified the conversion by opening the PHYLIP file in a text editor

| gene | ✓ | 17-08-2024 00:14 | PHY File |
|------|---|------------------|----------|

Task 8 output