**NAME: RACHAEL OLUWAKAMIYE ABOLADE**

## Project Title:

Comprehensive Sequence Analysis of the Human TNF Gene

## Objective:

To apply bioinformatics skills learned in Module 1 to download, analyze, and interpret the sequence of the human TNF gene, which encode a proinflammatory cytokine call TNF.

## Project Overview:

In this mini project, you will perform a series of bioinformatics tasks using the human TNF gene as your sequence of interest. The project will guide you through downloading the sequence, translating it, finding ORFs, analyzing sequence composition, identifying transcription factor binding sites, searching for functional motifs, predicting coding/non-coding regions, and converting sequence file formats.

**Date: AUGUST 20ᵀᴴ 2024**

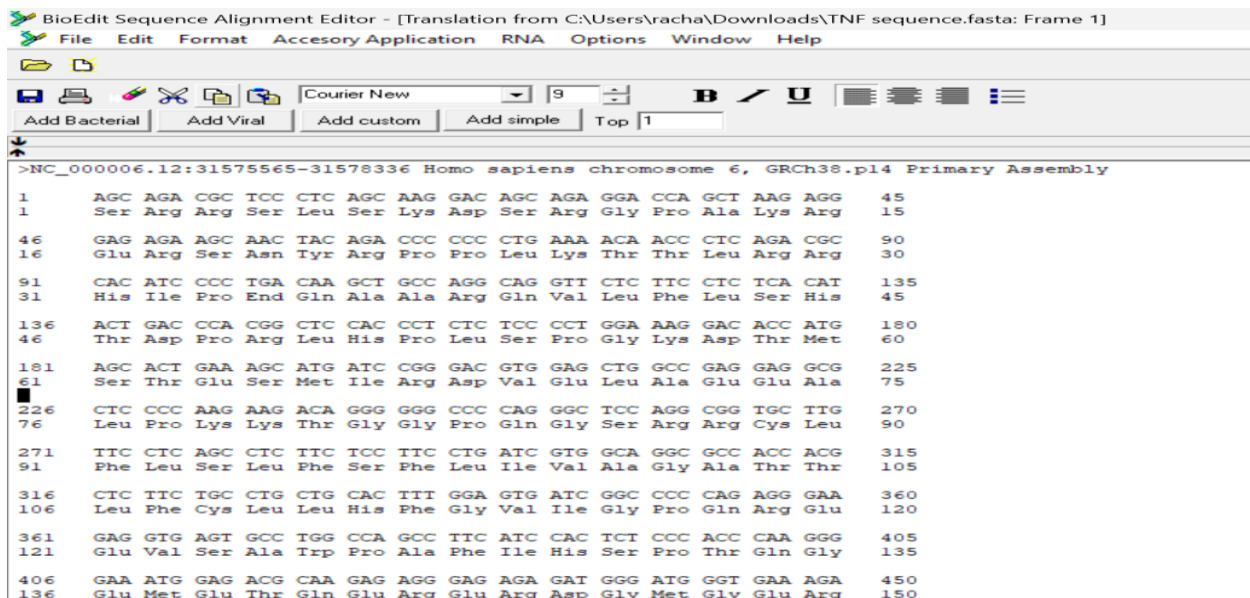## Task 1: Download a Biological Sequence from NCBI and View/Edit It

The human TNF gene sequence was downloaded from NCBI and viewed using BioEdit. Here is a screenshot of the sequence display:



The output shows that the human TNF gene with accession number NC_00006, region: 31575565..31578336 is 2722 bp long.

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

The DNA sequence of the TNF gene was translated into an amino acid sequence. Here is a screenshot of the amino acid sequence generated:

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

Start and stop positions, lengths, and protein translations of the ORFs.

**Output:**

```
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 178 to 606 (178 to 606): Frame 1 143 aa
MSTESMIRDVELAEEALPKKTGGPQGSRRCLFLSLFSFLIVAGATTLFCLLHFGVIGPQREEVSAWPAFIHSPTQGEMETQERERDGMGERCALIGRDGEKKTWRKTGMQKEMWQEMGKRERERWRDRMSGTWKVLTKCVWSE
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 450 to 587 (450 to 587): Frame 3 46 aa
MCADREGWREKNVEKDGDAERDVARDGEERERKMERQDVWHMEGAH
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 470 to 610 (470 to 610): Frame 2 47 aa
MERKKRGERRGCRKRCGKRWGREREKDGETGCLAHGRCSLSVYGVNE
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 659 to 781 (659 to 781): Frame 2 41 aa
MWGVRREMGEETSDMNKDGETERAGNMTAKERDGGDKERRR
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 756 to 923 (756 to 923): Frame 3 56 aa
MGEIRREEDRVSGTQKTLRERAVECLEGEYTDEWREKTRHLRAKSAGQTGSQLFLL
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 850 to 1134 (850 to 1134): Frame 1 95 aa
MNGERKPDTSGLRAQARQAASCSSFKGDSLDVNHSPSPQQFPRDLSLISPLAQAVSKCLQTSFLILGLGLGVGLVPVWKQWGKFKVLVLGEDGWR
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1124 to 1288 (1124 to 1288): Frame 2 55 aa
MDGGESRGVFSRKFKGLSFFFSLSSSGSSSRTPSDKPVAHVVGKSSEDVSWNLEG
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1128 to 1235 (1128 to 1235): Frame 3 36 aa
MEVKVGGYFLGSLRVSAFSFLSPLQDHLLEPRVTSL
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1316 to 1519 (1316 to 1519): Frame 2 68 aa
MVGRTWRQCEKDSLSSREGWRNSTGLSGILRTSWPGGMWDDRQRGQEPDVGWAELEGQDVESEPTWPH
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1434 to 1970 (1434 to 1970): Frame 3 179 aa
MTDREDRNRMWGGQSSRARMWRVNRHGHTDSPLPLSLPPANPQAEGQLQWLNRRANALLANGVELRDNQLVVPSEGLYLIYSQVLFKGQGCPSTHVLLTHTISRIAVSYQTKVNLLSAIKSPCQRETPEGAEAKPWYEPIYLGGVFQLEKGDRLSAEINRPDYLDFAESGQVYFGI
IAL
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1510 to 1584 (1510 to 1584): Frame 1 25 aa
MATLTLLSLSPSLQQTLKLRGSSSG
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1717 to 2232 (1717 to 2232): Frame 1 172 aa
MCSSPTPSAASPSPTRPRSTSSLPSRAPARGRPQRGLRPSPGMSPSIWEGSSSWRRVTDSALRSIGPTISTLPSLGRSTLGSLPCEEDEHPTFPNASPAPIPLLPPPSDTLNLFWLKKRIGGLGSEPKLRTLSNKTTTSKPGIQECVACTVKCWQPLRIQTGASRTHWGLQL
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 2151 to 2336 (2151 to 2336): Frame 3 62 aa
MCGLHSEVLATTKNSNWGLQNSLGPTALIPDIWNLETREPLVLARMLQDLRRPHLEIDTSGP
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 2356 to 2472 (2356 to 2472): Frame 1 39 aa
MFPDFLETRSPALPMEPAPSIYVCTCDYLLFIYYLFIYR
>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 2420 to 2518 (2420 to 2518): Frame 2 33 aa
```

BioEdit searches the TNF sequence for open reading frames (ORFs) in 6 reading frames (3 in either direction). ORFs are sequences that contain start codons and stop codons with a minimum distance between them. BioEdit returns the range of each ORF, along with its protein translation. Each ORFs may potentially be protein coding regions. Protein-coding regions are usually of a significant length and contain multiple amino acids. Short sequences or those that do not align well with known proteins might be non-coding. Potential protein-coding regions from the output may include:

>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1434 to 1970 (1434 to 1970): Frame 3 179 aa

>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 1717 to 2232 (1717 to 2232): Frame 1 172 aa

>NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly Positions 1 to 2722: 178 to 606 (178 to 606): Frame 1 143 aa

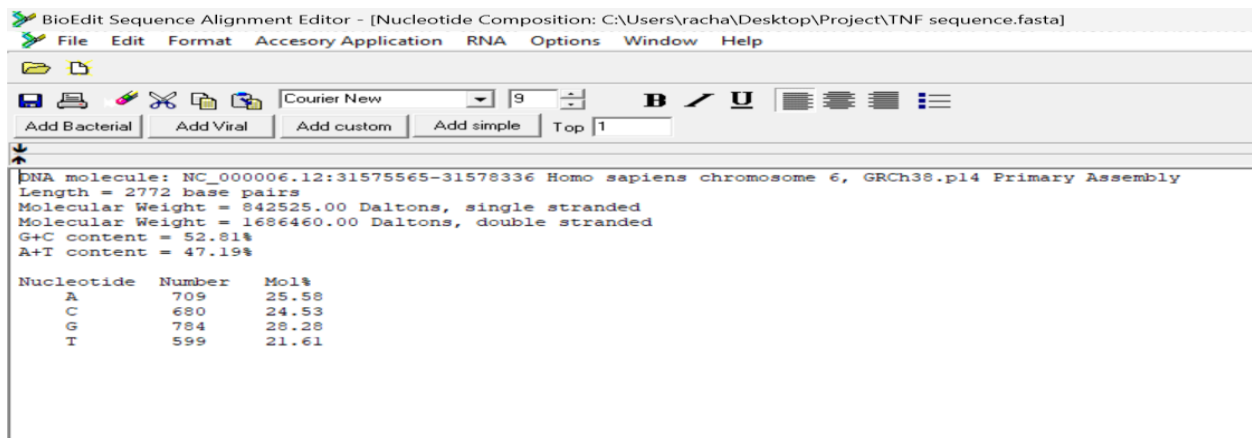**Task 4: Analyze Sequence Composition (Nucleotide or Amino Acid Frequencies)**

The nucleotide composition of the TNF gene sequence was analyzed:
The output shows that the Human TNF gene sequence (2722bp) contains 25.58%, 24.53%, 28.28%, 21.61% of Nucleotide A, C, G, T respectively.
The sequence has a high overall G+C content of 52.81%. A higher GC content often correlates with higher DNA stability due to the three hydrogen bonds between G and C compared to the two hydrogen bonds between adenine (A) and thymine (T).
This information is useful for various applications, such as understanding the genetic content of the region, designing primers for PCR, or studying the DNA's structural properties.
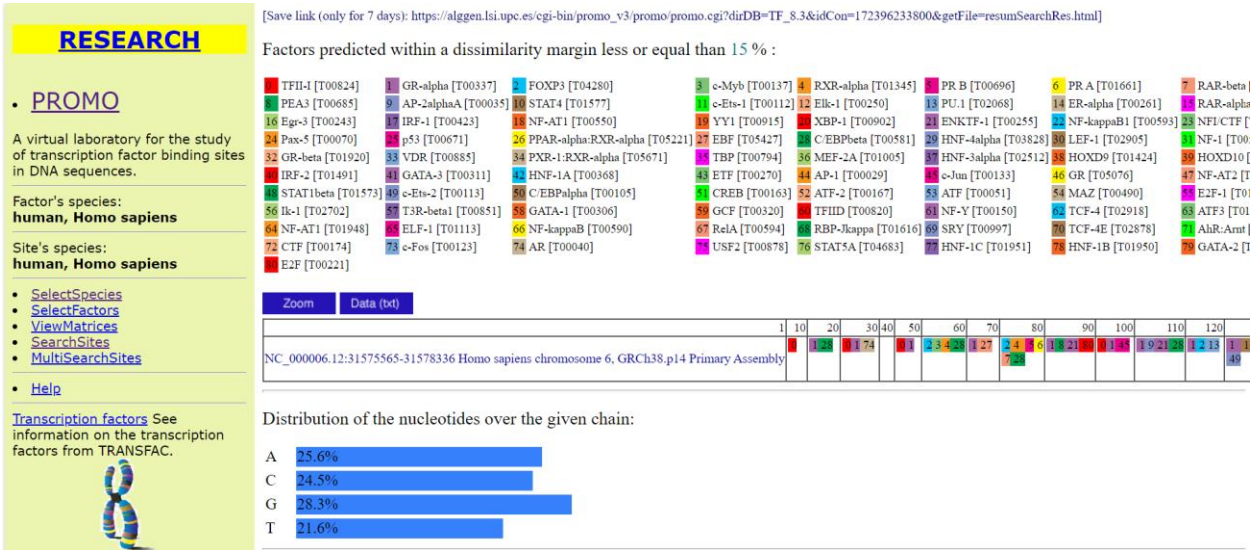
**Output:**



Here is also a graphical illustration of the TNF Nucleotide composition

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

The potential transcription factor binding sites in the TNF gene promoter region were Identified.
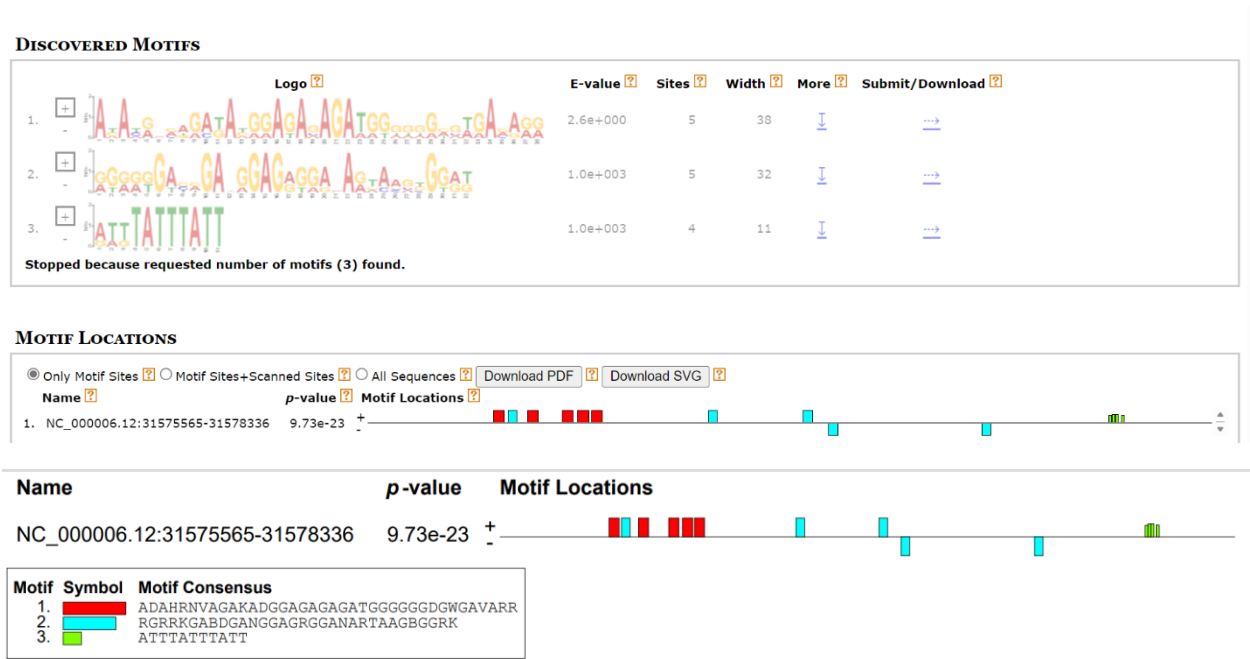
**Output:**

[Save link (only for 7 days): https://alggen.lsi.upc.es/cgi-bin/promo_v3/promo/promo.cgi?dirDB=TF_8.3&idCon=172396233800&getFile=resumSearchRes.html]

Factors predicted within a dissimilarity margin less or equal than 15 % :

| | | |
|---|---|---|
| TFII-I [T00824] | 1 GR-alpha [T00337] | 2 FOXP3 [T04280] |
| 8 PEA3 [T00685] | 9 AP-2alphaA [T00035] | 10 STAT4 [T01577] |
| 16 Egr-3 [T00243] | 17 IRF-1 [T00423] | 18 NF-AT1 [T00550] |
| 24 Pax-5 [T00070] | 25 p53 [T00671] | 26 PPAR-alpha:RXR-alpha [T05221] |
| 32 GR-beta [T01920] | 33 VDR [T00885] | 34 PXR-1:RXR-alpha [T05671] |
| IRF-2 [T01491] | 41 GATA-3 [T00311] | 42 HNF-1A [T00368] |
| STAT1beta [T01573] | 49 c-Ets-2 [T00113] | 50 C/EBPalpha [T00105] |
| 56 Ik-1 [T02702] | 57 T3R-beta1 [T00851] | 58 GATA-1 [T00306] |
| 64 NF-AT1 [T01948] | 65 ELF-1 [T01113] | 66 NF-kappaB [T00590] |
| 72 CTF [T00174] | 73 c-Fos [T00123] | 74 AR [T00040] |
| E2F [T00221] | | |

| | | | |
|---|---|---|---|
| 3 c-Myb [T00137] | 4 RXR-alpha [T01345] | PR B [T00696] | 6 PR A [T01661] | 7 RAR-beta |
| 11 c-Ets-1 [T00112] | 12 Elk-1 [T00250] | 13 PU.1 [T02068] | 14 ER-alpha [T00261] | 15 RAR-alpha |
| 19 YY1 [T00915] | XBP-1 [T00902] | 21 ENKTF-1 [T00255] | 22 NF-kappaB1 [T00593] | 23 NFI/CTF |
| 27 EBF [T05427] | 28 C/EBPbeta [T00581] | 29 HNF-4alpha [T03828] | 30 LEF-1 [T02905] | 31 NF-1 |
| TBP [T00794] | 36 MEF-2A [T01005] | 37 HNF-3alpha [T02512] | 38 HOXD9 [T01424] | 39 HOXD10 |
| 43 ETF [T00270] | 44 AP-1 [T00029] | c-Jun [T00133] | 46 GR [T05076] | 47 NF-AT2 |
| 51 CREB [T00163] | 52 ATF-2 [T00167] | 53 ATF [T00051] | 54 MAZ [T00490] | 55 E2F-1 |
| 59 GCF [T00320] | TFIID [T00820] | 61 NF-Y [T00150] | 62 TCF-4 [T02918] | 63 ATF3 |
| 67 RelA [T00594] | 68 RBP-Jkappa [T01616] | 69 SRY [T00997] | 70 TCF-4E [T02878] | 71 AhR:Arnt |
| 75 USF2 [T00878] | 76 STAT5A [T04683] | 77 HNF-1C [T01951] | 78 HNF-1B [T01950] | 79 GATA-2 |

Zoom   Data (txt)

NC_000006.12:31575565-31578336 Homo sapiens chromosome 6, GRCh38.p14 Primary Assembly

Distribution of the nucleotides over the given chain:

A 25.6%
C 24.5%
G 28.3%
T 21.6%

The output shows eighty (80) potential Transcription factors within a dissimilarity margin less or equal than 15%. These Transcription factors play a critical role in regulating the TNF gene expression by binding to specific sequence in the Nucleotide and either activating or repressing transcription.

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

Functional motifs in the TNF gene sequence were identified using MEME Suite.

**Output:**



The output shows three consensus motifs spanning the human TNF gene sequence.

Sequence motifs are short recurring patterns in DNA that are presumed to have a biological function. The above screenshot indicates which nucleotides or amino acids are most frequently observed at each position in the alignment. This is shown in a logo format, where each position shows the most common base or amino acid.

The overall p-value of 9.73e-23 of the motifs indicates that the motifs are highly statistically significant, suggesting a very low probability that the observed pattern is due to chance. In practical terms, this p-value provides strong evidence that the motif is likely to be biologically relevant.

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

The coding and non-coding regions within the TNF gene sequence were predicted using GENSCAN.

**Output:**

Here is a screenshot of the output:



Coding regions are parts of the DNA sequence that contain the instructions for building proteins, while non-coding regions do not contain these instructions.
GENSCAN output helps in identifying and mapping coding regions (exons) and non-coding regions (introns and intergenic regions) in a DNA sequence. It provides important information for gene prediction and annotation, which is crucial for understanding gene structure and function.

GENSCAN has predicted five exons on the positive strand on the gene. One initial exon, two internal exons, a terminal exon and a polyadenylation signal within the TNF DNA sequence. The confidence levels for these predictions vary, with some predictions showing high probability and others more moderate indicating strong confidence in the gene predictions. The output also includes a predicted peptide sequence 233 aa in length derived from the gene model which includes various motifs and domains that may be functionally significant, which is essential for functional annotation and further analysis.

```
Predicted peptide sequence(s):



>/tmp/08_21_24-13:29:49.fasta|GENSCAN_predicted_peptide_1|233_aa

MSTESMIRDVELAEEALPKKTGGPQGSRRCLFLSLFSFLIVAGATTLFCLLHFGVIGPQR

EEFPRDLSLISPLAQAVRSSSRTPSDKPVAHVVANPQAEGQLQWLNRRANALLANGVELR

DNQLVVPSEGLYLIYSQVLFKGQGCPSTHVLLTHTISRIAVSYQTKVNLLSAIKSPCQRE

TPEGAEAKPWYEPIYLGGVFQLEKGDRLSAEINRPDYLDFAESGQVYFGIIAL
```

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

Convert the TNF gene sequence from FASTA format to PHYLIP format.

**Output**

The Fasta and Phylip format of the sequence is attached to this report.