<u>**Task 1 DNA/Protein Sequence Analysis**</u>


<u>**Sequence Similarity Analysis between Streptococcus pyogenes M-Protein and Human Cardiac Myosin**</u>

# Introduction

Protein sequence comparison is one of the most important bioinformatics techniques used to study evolutionary relationships, structural similarity, and potential functional overlap between proteins.

The **M protein** of *Streptococcus pyogenes* is a major virulence factor. It helps the bacteria evade the host immune system by resisting phagocytosis. However, certain regions of this bacterial protein show **molecular mimicry** with human heart proteins, especially **cardiac myosin**. This similarity is believed to play a key role in **autoimmune cross-reactivity**.

To explore this concept at the molecular level, a **protein sequence comparison** was performed between the M protein of *S. pyogenes* and human cardiac myosin using NCBI databases and BLAST tools.
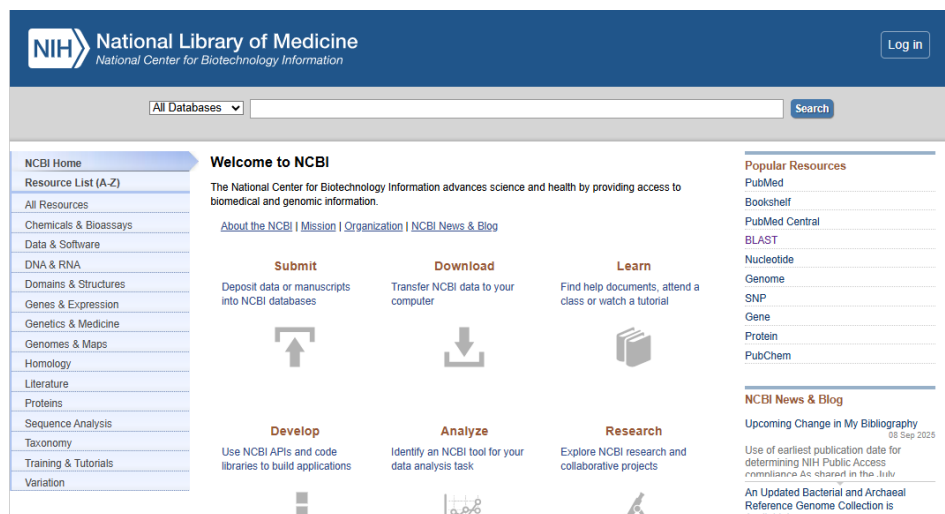
# Background Information

1. **M Protein (Streptococcus pyogenes)**
   o A surface-exposed protein responsible for virulence.
   o Highly variable in sequence, with over 200 known serotypes.
   o Some conserved regions show cross-reactivity with human tissues.
2. **Cardiac Myosin (Human)**
   o A major contractile protein in cardiac muscle.
   o Plays a key role in muscle contraction and heart function.
   o Has both conserved functional domains and variable regions.

# Procedure:
**Steps to Retrieve FASTA sequences from NCBI:**

1. Open the **NCBI** website

2. In the search bar, type the **protein name**



3. From the results, open the **Protein** record of the selected protein.



4. Scroll to the top-right section and click on the **"FASTA"** link.

5. The FASTA format sequence will appear; you may **copy the sequence directly** for use or to download the sequence, click on **"Graphics"** at the top of the record.



6. Select **"Download"** and choose **FASTA format** to save the file.



7. Repeat the same steps for the **second protein** (e.g., Myosin).

**Steps to Perform Multiple Sequence Alignment Using NCBI BLAST**

1. Open the **BLAST** page on NCBI



2. Select **Protein BLAST (blastp)**.
3. Under the query box, click **"Align two or more sequences"**.
4. Upload or paste the **FASTA sequences** of both proteins (Actin and Myosin).
5. Scroll down and click **"BLAST"** to run the alignment.
6. Wait for the results to load; the alignment output will be displayed below.

# Understanding the BLAST Alignment Output

- **Matching letters** → identical amino acids at that position.
- **Dots ( . )** → similar amino acids (conserved substitution).
- **Spaces ( )** → mismatch; no similarity at that position.
- **Query** → first protein sequence you input.
- **Subject** → second protein sequence.
- **Score / Bit score** → strength of alignment (higher = better).
- **E-value** → significance of the match (lower = more significant).

---

**NIH National Library of Medicine** — National Center for Biotechnology Information  

Log in

BLAST® » blastp suite-2sequences » results for RID-JPDWFM8T114   Home   Recent Results   Saved Strategies   Help

‹ Edit Search   Save Search   Search Summary ▾

How to read this report?   BLAST Help Videos   Back to Traditional Results Page

| | |
|---|---|
| Job Title | gb|RXH49094.1|:1-579 M protein [Streptococcus |
| RID | JPDWFM8T114  Search expires on 12-01 01:12 am  Download All ▾ |
| Program | Blast 2 sequences  Citation ▾ |
| Query ID | lcl|Query_5120997 (amino acid) |
| Query Descr | gb|RXH49094.1|:1-579 M protein [Streptococcus pyogenes] |
| Query Length | 579 |
| Subject ID | lcl|Query_5120999 (amino acid) |
| Subject Descr | emb|CAA86293.1|:1-1937 Myosin [Homo sapiens] |
| Subject Length | 1937 |
| Other reports | Multiple alignment  MSA viewer |

**Filter Results**

Percent Identity [  ] to [  ]   E value [  ] to [  ]   Query Coverage [  ] to [  ]   Filter   Reset

Clusters | Graphic Summary | Alignments | Dot Plot

**Clusters producing significant alignments**   Download ▾   Select columns ▾   Show 100 ▾

☑ select all  1 clusters selected   Graphics  Multiple alignment  MSA Viewer

| ☑ | Cluster Representative Sequence | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|---|---|---|---|---|---|---|---|---|
| ☑ | emb|CAA86293.1|:1-1937 Myosin [Homo sapiens] | 30.0 | 56.6 | 58% | 0.001 | 24.63% | 1937 | Query_5120999 |

Feedback

---

⬇ Download ▾   Graphics   Sort by: E value ▾   ▾ Next ▲ Previous ◀ Descriptions

**emb|CAA86293.1|:1-1937 Myosin [Homo sapiens]**
Sequence ID: **Query_5120999**  Length: **1937**  Number of Matches: **2**

**Range 1: 1427 to 1745 Graphics**   ▾ Next Match ▲ Previous Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 30.0 bits(66) | 0.001 | Compositional matrix adjust. | 84/341(25%) | 151/341(44%) | 29/341(8%) |

```
Query  149   EYQDLNDDFDLAKQGYASSDKRHQQELEEKEKKVTEATAKVDQISKELETAKQKNESTKQ   208
             E +DL  D + +    A+ DK+ +    +K ++E    K ++    ELE +++++ S
Sbjct  1427  EVEDLMLDVERSNAACAALDKKQRN----FDKVLSEWKQKYEETQAELEASQKESRSLST   1482

Query  209   DLTEKQNRVSELEQELATTK-ENAKKDFELAALGHQLADKEYNAKIAELESKLADAKKDF   267
             +L + +N   E   +L T + EN    E++ L  Q+A  E   +I ELE      KK
Sbjct  1483  ELFKVKNVYEESLDQLETLRRENKNLQQEISDLTEQIA--EGGKQIHELEK----IKKQV   1536

Query  268   ELAALGHQHAHNEYQAKLAEKDGQI--KQLEEQKQILDASRKGTARDLEAVRQAKKATEA   325
             E     Q A E +A L  ++G+I    QLE  +    RK   +D
Sbjct  1537  EQEKCEIQAALEEAEASLEHEEGKILRIQLELNQVKSEVDRKIAEKD-----------E   1584

Query  326   ELNNLKAELAKVTEQKQ-ILDASRKGTARDLEAVRKAKAQVEAALKQLEEQNRISEASRK   384
             E++ LK   +V E Q  LDA  +      L   +K + +    QL   NR++   S +
Sbjct  1585  EIDQLKRNHTRVVETMQSTLDAEIRSRNDALRVKKKMEGDLNEMEIQLNHANRLAAESLR   1644

Query  385   GLRRDLDASREAKKQVEKDL---ANLTAELDKVKEEKQISDASRQGLRRDLDASREAKKQ   441
             R      +E +  ++ L     +L +L V+     + A + L  L+ +  ++K
Sbjct  1645  NYRNTQGILKETQLHLDDALRGQEDLKEQLAIVERRANLLQAEIEELWATLEQTERSRKI   1704

Query  442   VEKALEEANSKLAALEKLNKELEESKKLTEKEKAELQAKLE   482
             E+ L +A+ ++   L   N  L +KK  E + ++LQ+++E
Sbjct  1705  AEQELLDASERVQLLHTQNTSLINTKKKLENDVSQLQSEVE   1745
```

**Range 2: 1807 to 1909 Graphics**   ▾ Next Match ▲ Previous Match ▲ First Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 26.6 bits(57) | 0.008 | Compositional matrix adjust. | 38/119(32%) | 57/119(47%) | 25/119(21%) |

```
Query  357   AVRKAKAQV---EAALKQLE-----EQNRISEASRKGLRRDLDASREAKKQVEKDLANLT   408
             A++  K Q+    EA +++LE    EQ R +EA  KGLR+       +E   Q E+D  N+
Sbjct  1807  ALKGGKKQIQKLEARVRELEGEVENEQKRNAEAV-KGLRKHERRVKELTYQTEEDRKNVL   1865

Query  409   AELDKVKEEKQISDASRQGLRRDLDA-SREAKKQVEKALEEANSKLAALEKLNKELEES   466
                    Q   L   LA  +  K+Q E+A  E++N+ L+   KL  ELEE+
Sbjct  1866  --------------RLQDLVDKLQAKVKSYKRQAEEAEEQSNANLSKFRKLQHELEEA   1909
```