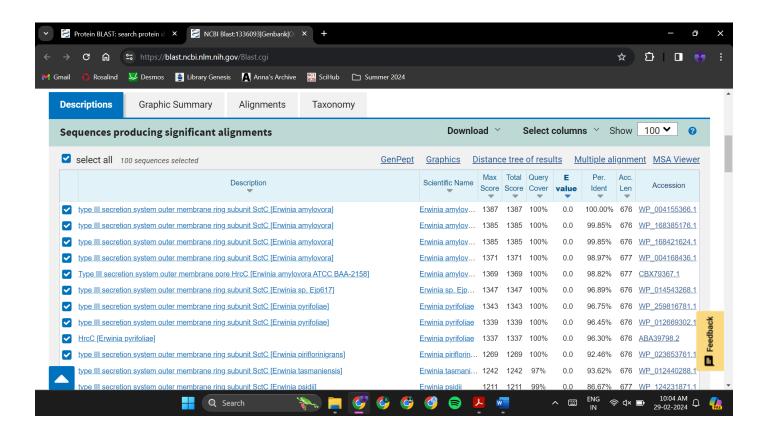Anirudh Rao
BE21B004

# BT3040 – Bioinformatics

## Practical 5

1

After running BLASTP with the database as "nr", the following protein sequences are found to be similar to the given sequence:
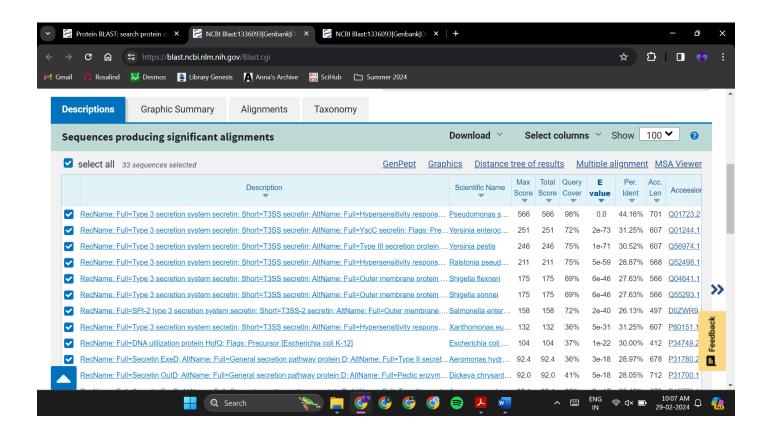


## Analysis

- All 100 similar sequences have an E value of 0
- 29 sequences have 100% query coverage
- The lowest query coverage is 90%
- Only 1 sequence has 100% identity (type III secretion system outer membrane ring subunit SctC [Erwinia amylovora])
- Lowest percentage identity is 66.96% (type III secretion system outer membrane ring subunit SctC [Dickeya lacustris])

Thus, the given sequence is likely to be type III secretion system outer membrane ring subunit SctC [Erwinia amylovora].

After running BLASTP with the database as "swissprot", the following protein sequences are found to be similar to the given sequence:
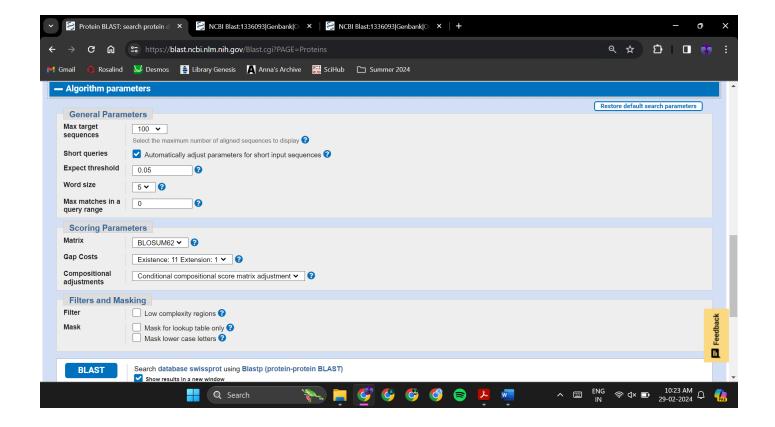


Analysis

- All 33 similar sequences have E values close to 0
- The highest query coverage is 98% (Type 3 secretion system secretin [Pseudomonas syringae pv. syringae])
- The lowest query coverage is 23%, which is the coverage for 3 sequences
- The highest percentage identity is 44.16% (Type 3 secretion system secretin [Pseudomonas syringae pv. syringae])
- The lowest percentage identity is 21.75% (Uncharacterized protein y4xJ [Sinorhizobium fredii NGR234])
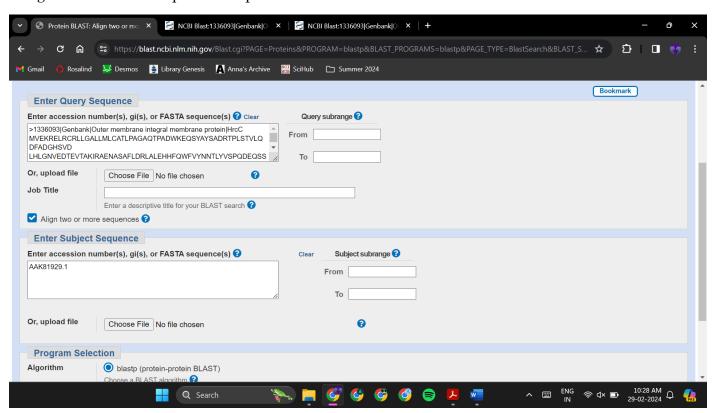
2

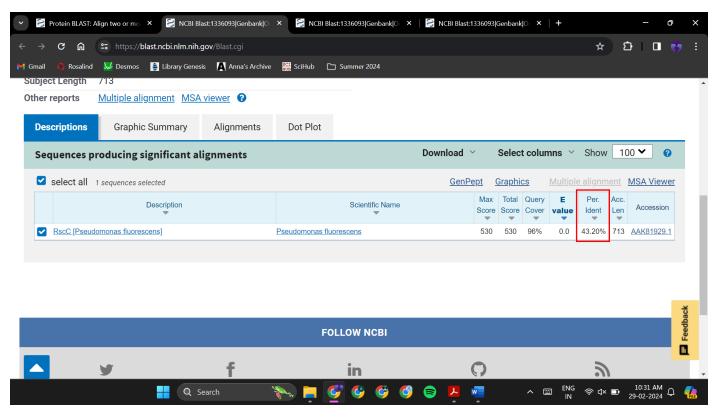The algorithm parameters used in the previous question are shown below:

The sequency identity of the query sequence with AAK81929.1 was found by selecting the "Align two or more sequences" option in BLASTP:
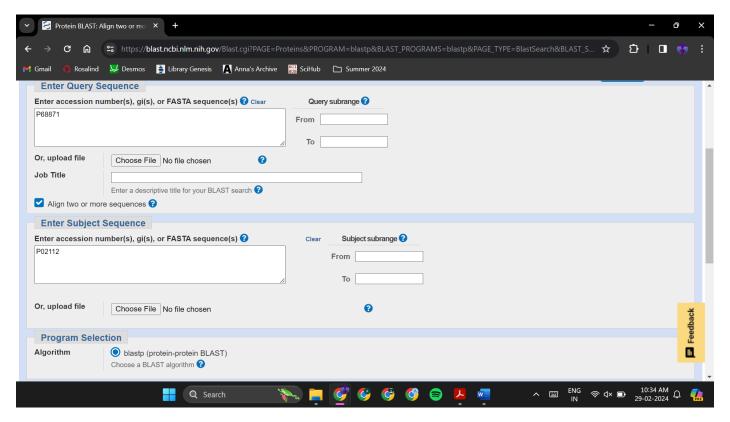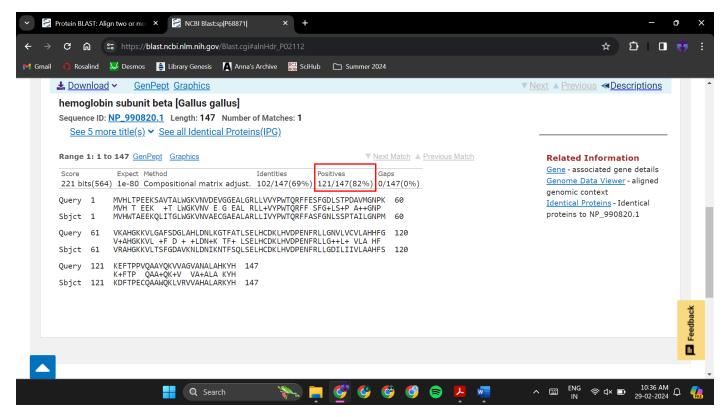
The percentage identity was found to be 43.20%.

Using UniProtKB, the accession number of human hemoglobin beta was found to be P68871, while that of chicken hemoglobin beta was found to be P02112.

These accession numbers were used along with the "Align two or more sequences" option in BLASTP:

The percentage similarity was found to be 82%.

The Python code to list all the matching pentapeptides (which occur in both the human and chicken hemoglobin beta sequences) and their frequency of occurrence in given sequences is given below:

```python
def find_pentapeptides(seq1, seq2):
    seq1_pentapeptides = {}

    for i in range(len(seq1) - 5 + 1):
        pentapeptide = seq1[i:i + 5]
        if pentapeptide not in seq1_pentapeptides.keys():
            seq1_pentapeptides[pentapeptide] = 1
        else:
            seq1_pentapeptides[pentapeptide] += 1

    matched_pentapeptides = {}
    for i in range(len(seq2) - 5 + 1):
        pentapeptide = seq2[i:i + 5]
        if pentapeptide in seq1_pentapeptides.keys():
            if pentapeptide not in matched_pentapeptides.keys():
                matched_pentapeptides[pentapeptide] = seq1_pentapeptides[pentapeptide]
+ 1
            else:
                matched_pentapeptides[pentapeptide] += 1

    for pentapeptide, frequency in matched_pentapeptides.items():
        print(
            f"{pentapeptide} occurs {frequency} time(s) in both peptides
({seq1_pentapeptides[pentapeptide]} time(s) in Seq. 1 and {frequency -
seq1_pentapeptides[pentapeptide]} time(s) in Seq. 2)")
```
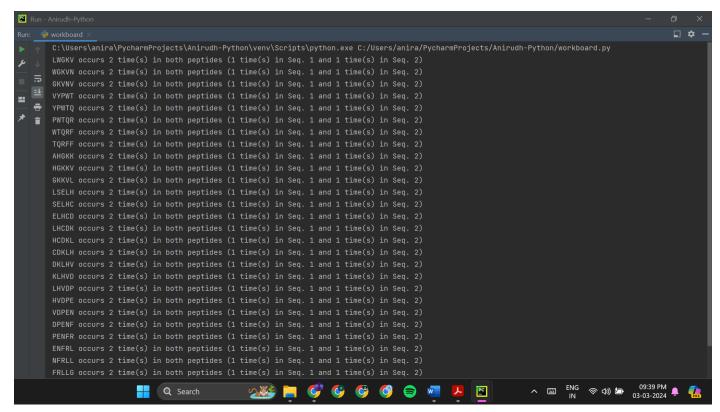
```
human =
"MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTF
ATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH"
chicken =
"MVHWTAEEKQLITGLWGKVNVAECGAEALARLLIVYPWTQRFFASFGNLSSPTAILGNPMVRAHGKKVLTSFGDAVKNLDNIKNTF
SQLSELHCDKLHVDPENFRLLGDILIIVLAAHFSKDFTPECQAAWQKLVRVVAHALARKYH"
find_pentapeptides(human, chicken)
```

The output looks like:



**6**

The Python code to compute sequence identity, similarity, query coverage and gap percentage from the alignment of human and chicken hemoglobin sequences is given below:

```python
def alignment_statistics(seq1, seq2, align):
    identity = 0
    similarity = 0
    gaps = 0

    for char in align:

        if char.isalpha():
            identity += 1
            similarity += 1
        if char == "+":
            similarity += 1
        if char == "-":
            gaps += 1

    alignment_identity = 100 * identity / len(align)
```

```python
    alignment_similarity = 100 * similarity / len(align)
    query_coverage = 100 * len(align) / len(seq2)
    gap_percentage = 100 * gaps / len(align)

    print(f"The sequence identity is {alignment_identity:.2f}%.")
    print(f"The sequence similarity is {alignment_similarity:.2f}%.")
    print(f"The query coverage is {query_coverage:.2f}%.")
    print(f"The gap percentage is {gap_percentage:.2f}%.")


human =
"MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTF
ATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH"
chicken =
"MVHWTAEEKQLITGLWGKVNVAECGAEALARLLIVYPWTQRFFASFGNLSSPTAILGNPMVRAHGKKVLTSFGDAVKNLDNIKNTF
SQLSELHCDKLHVDPENFRLLGDILIIVLAAHFSKDFTPECQAAWQKLVRVVAHALARKYH"
alignment = "MVH T EEK  +T LWGKVNV E G EAL RLL+VYPWTQRFF SFG+LS+P A++GNP V+AHGKKVL +F D
+ +LDN+K TF+ LSELHCDKLHVDPENFRLLG++L+ VLA HF K+FTP  QAA+QK+V  VA+ALA KYH"

alignment_statistics(human, chicken, alignment)
```
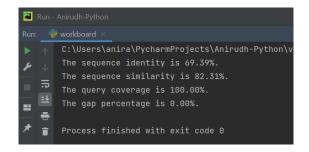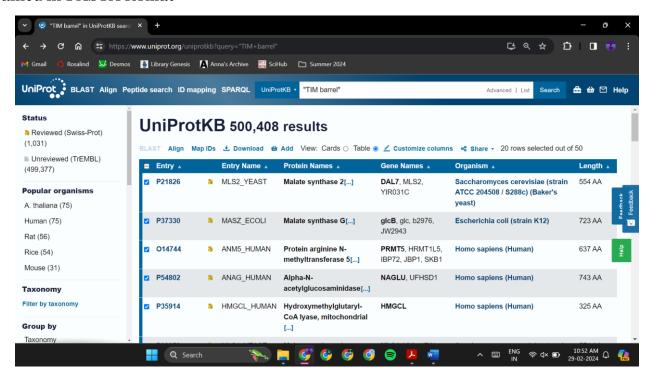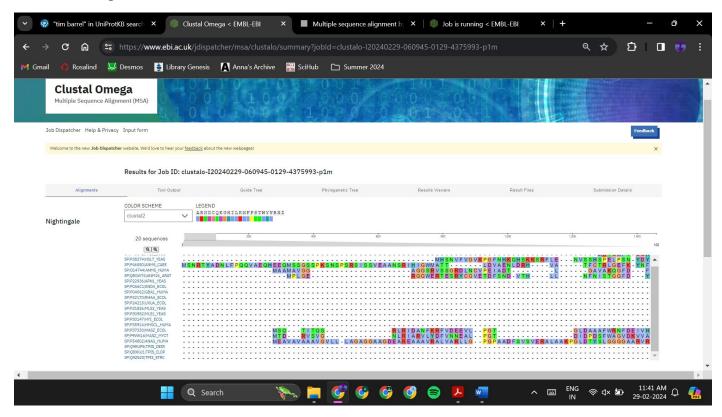
The output looks like:

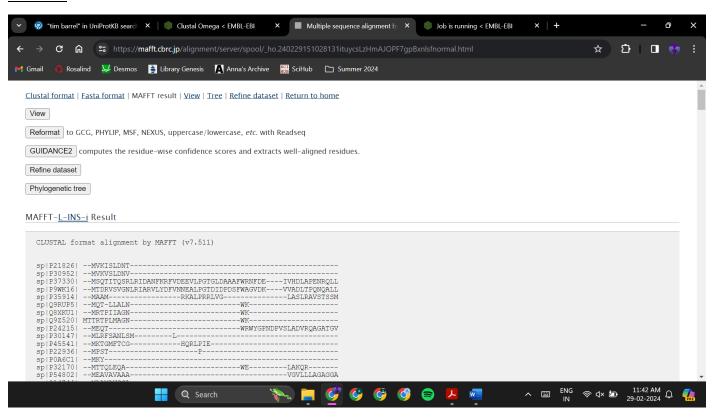20 sequences with a TIM barrel domain were selected using UniProt and their sequences were obtained in FASTA format:

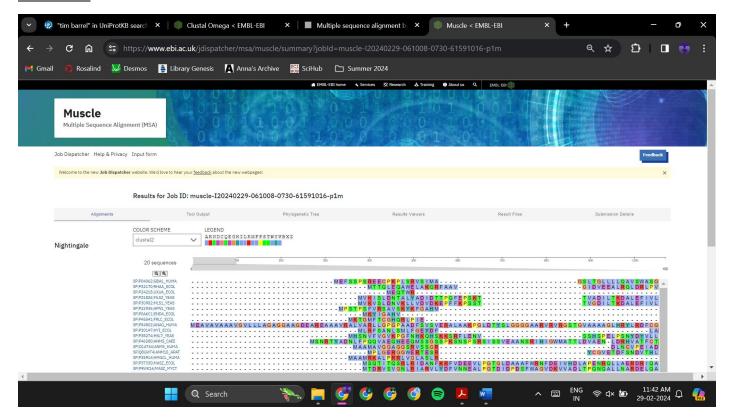These were then aligned using Clustal Omega, MAFFT, and MUSCLE:
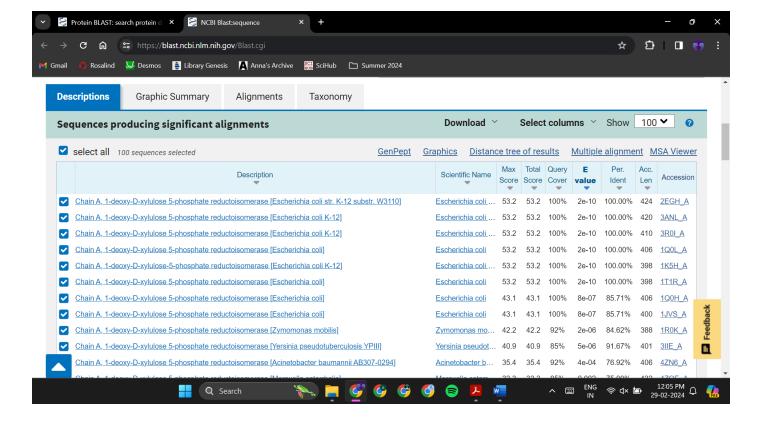
## Clustal Omega



## MAFFT

## MUSCLE



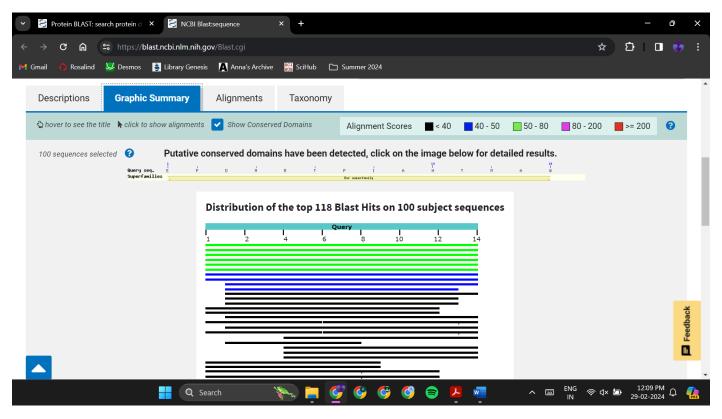5 residue positions that are aligned differently in these three methods are 1, 2, 3, 4, 5.

|   | Clustal Omega | MAFFT | MUSCLE |
|---|---|---|---|
| 1 | M in P46580, gaps in others | M in Q9Z520, gaps in others | M in P54802, gaps in others |
| 2 | S in P46580, gaps in others | T in Q9Z520, gaps in others | E in P54802, gaps in others |
| 3 | N in P46580, gaps in others | T in Q9Z520, M in others | A in P54802, gaps in others |
| 4 | R in P46580, gaps in others | Different amino acids in different sequences | V in P54802, gaps in others |
| 5 | T in P46580, gaps in others | Different amino acids in different sequences | A in P54802, gaps in others |

BLASTP was run using the given short sequence with "pdb" as the database:

- Many of the results have 100% query coverage, high percentage identity, and E values very close to 0.
- The query appears to closely match subsequences in Chain A of 1-deoxy-D-xylulose 5-phosphate reductoisomerase in *Escherichia coli* str. K-12.
- From the graphic summary, we can see that the given query sequence is very closely related to the Dxr superfamily:



Just by using a short sequence, we were able to successfully identify a family of related proteins that share the same motif in their sequences.