

Assignment 5
BS22B009

All the python codes are present here :

<https://colab.research.google.com/drive/1Z1O-hDsg5gOT9KMfOD6BQV2JtBE1v72Z?usp=sharing>

1. Analyze the occurrence of similar proteins in “nr” and SWISS-PROT database for the sequence given below: >1336093|Genbank|Outer membrane integral membrane protein|HrcC

MVEKRELRCRLLGALLMLCATLPAGAQTPADWKEQSYAYSADRTPLSTVLQDFADGHSVD
LHLGNVEDTEVTAKIRAENASAFDLRLALEHHFQWFVYNNTLYVSPQDEQSSERLEISPD
AAPDIKQALSGIGLLDPRFGWGELPDDGVVLVTGPPQYLELVKRFSEQREKKEDRRKVM
FPLRYASVADRTIHYRDQTVVIPGVATMLNELMNGKRAAPASASGIDSTPGGPDTNSMMQ
NTQTLLSRLSSRNKTSNRAGGRDNEIEDVSGRISADVRNNALLIRDDDKRHDEYSQLIAK
IDVPQNLVEIDAVILDIDRTALNRLEANWQATLGGVTGGSSLMSGSGTLFVSDFKRFFAD
IQALEGEGTASIVANPSVLTLENQPAVIDFSQTAYITATGERVADIQPVTAGTSLQVTPR
AVGNEGHSSIQLMIDIEDGHVQTNGDGQATGVKRGTVSTQALISENRALVLGGFHVEESA
DRDRRIPLLGDIPWLGLFSSKRHEISQRQRLFILTPRLIGDQTDPTRYVTADNRQQQLSD
AMGRVERRHSSVNQHVDVVENALRDLAEGQSPAGFQPQTSGTRLSEVCRSTPALLFESTRG
QWYSSSTNGVQLSVGVVRNTSSKPLRFDEANCASKRTLAVAVWPHSALAPGESAEVYLA
M DPSRVLHASRESLLNR

Steps:

Use BLAST (Basic Local Alignment Search Tool):

- Go to NCBI BLASTp (for protein sequences):
<https://blast.ncbi.nlm.nih.gov/Blast.cgi>
- Select "Protein BLAST (BLASTp)".
- Enter the given sequence in the query box.
- Choose the nr (non-redundant) protein database for a broad search.
- Run the BLAST search and analyze the results (sequence similarity, identity, and alignment).

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

blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins

blastn **blastp** blastx tblastn tblastx

BLASTP programs search protein databases using a protein query. more...

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

>1336093|Genbank|Outer membrane integral membrane protein|HrcC
MVEKRELRCRLLGALLMLCATLPAGAQTPADWKEQSYAYSADRTPLSTVLQ
DFADGHSVDLHLGNVEDTEVTAKIRAENASAFDLRLALEHHFQWFVYNNNTLY
VSPQDEQSSERLEISPDAAPIKQALSGIGLLDPRFGWGLPDDGWVLVTGP

Query subrange [?](#)
From
To

Or, upload file No file chosen [?](#)

Job Title
Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Databases ☒ Standard databases (nr etc.): ☐ Experimental databases

Compare ☐ Select to compare standard and experimental database [?](#)

Standard

Database ?

Organism Optional ☒ exclude
Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

Exclude Optional ☐ Models (XM/XP) ☐ Non-redundant RefSeq proteins (WP) ☐ Uncultured/environmental sample sequences

Program Selection

Algorithm ☐ Quick BLASTP (Accelerated protein-protein BLAST)
☒ blastp (protein-protein BLAST)
☐ PSI-BLAST (Position-Specific Iterated BLAST)
☐ PHI-BLAST (Pattern Hit Initiated BLAST)
☐ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)
Choose a BLAST algorithm [?](#)

Search database nr using Blastp (protein-protein BLAST)
☒ Show results in a new window

blast.ncbi.nlm.nih.gov/Blast.cgi

Descriptions Graphic Summary Alignments Taxonomy

Sequences producing significant alignments Download Select columns Show 100

select all 100 sequences selected GenPept Graphics Distance tree of results Multiple alignment MSA Viewer

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
type III secretion system outer membrane ring subunit SctC [Erwinia amylovora]	Erwinia amylovora	1387	1387	100%	0.0	100.00%	676	WP_004155366.1
type III secretion system outer membrane ring subunit SctC [Erwinia amylovora]	Erwinia amylovora	1385	1385	100%	0.0	99.85%	676	WP_168385176.1
type III secretion system outer membrane ring subunit SctC [Erwinia amylovora]	Erwinia amylovora	1385	1385	100%	0.0	99.85%	676	WP_168421624.1
type III secretion system outer membrane ring subunit SctC [Erwinia amylovora]	Erwinia amylovora	1383	1383	100%	0.0	99.85%	676	WP_400167015.1
type III secretion system outer membrane ring subunit SctC [Erwinia amylovora]	Erwinia amylovora	1371	1371	100%	0.0	98.97%	677	WP_004168436.1
Type III secretion system outer membrane pore HrcC [Erwinia amylovora ATCC BAA-2158]	Erwinia amylovo...	1369	1369	100%	0.0	98.82%	677	CBX79367.1
type III secretion system outer membrane ring subunit SctC [Erwinia sp. Eip617]	Erwinia sp. Eip617	1347	1347	100%	0.0	96.89%	676	WP_014543268.1
type III secretion system outer membrane ring subunit SctC [Erwinia pyrifoliae]	Erwinia pyrifoliae	1343	1343	100%	0.0	96.75%	676	WP_259816781.1
type III secretion system outer membrane ring subunit SctC [Erwinia pyrifoliae]	Erwinia pyrifoliae	1339	1339	100%	0.0	96.45%	676	WP_012669302.1
HrcC [Erwinia pyrifoliae]	Erwinia pyrifoliae	1337	1337	100%	0.0	96.30%	676	ABA39798.2
type III secretion system outer membrane ring subunit SctC [Erwinia piriflorinigrans]	Erwinia piriflorini...	1269	1269	100%	0.0	92.46%	676	WP_023653761.1
type III secretion system outer membrane ring subunit SctC [Erwinia tasmaniensis]	Erwinia tasmani...	1242	1242	97%	0.0	93.62%	676	WP_012440288.1
type III secretion system outer membrane ring subunit SctC [Erwinia psidii]	Erwinia psidii	1211	1211	100%	0.0	86.67%	677	WP_124231871.1
type III secretion system outer membrane ring subunit SctC [Erwinia tracheiphila]	Erwinia tracheip...	1184	1184	100%	0.0	83.85%	674	WP_233479868.1
MULTISPECIES: type III secretion system outer membrane ring subunit SctC [Pantoea]	Pantoea	1184	1184	100%	0.0	84.09%	679	WP_275222444.1
type III secretion system outer membrane ring subunit SctC [Erwinia tracheiphila]	Erwinia tracheip...	1183	1183	100%	0.0	83.70%	674	WP_016191026.1
type III secretion system outer membrane ring subunit SctC [Pantoea sp. AS142]	Pantoea sp. AS...	1181	1181	100%	0.0	84.09%	679	WP_337013830.1
type III secretion system outer membrane ring subunit SctC [Pantoea vagans]	Pantoea vagans	1180	1180	100%	0.0	83.95%	679	WP_061060948.1
type III secretion system outer membrane ring subunit SctC [Pantoea vagans]	Pantoea vagans	1177	1177	100%	0.0	83.80%	679	WP_336749516.1
type III secretion system outer membrane ring subunit SctC [Pantoea agglomerans]	Pantoea agglom...	1176	1176	100%	0.0	83.80%	679	WP_323169453.1
type III secretion system outer membrane ring subunit SctC [Pantoea agglomerans]	Pantoea agglom...	1176	1176	100%	0.0	83.65%	679	WP_031590814.1

Feedback

Analysis of Nr non reductant protein sequence BLAST Results:

Top Matches with Erwinia amylovora

- The best hits (top 5 sequences) are from Erwinia amylovora, a well-known plant pathogen responsible for fire blight in apples and pears.
- These sequences have 100% Query Coverage and 100% or 99.85% identity, confirming near-exact matches.
- Max Score: 1387-1385
- E-value: 0.0 (Indicating an extremely strong match)

Presence of Other Erwinia Species

- The T3SS outer membrane ring subunit SctC is highly conserved across multiple Erwinia species, including:
 - E. pyrifoliae, E. piriflorinigrans, E. tasmaniensis, E. psidii, E. tracheiphila, E. papayae, and E. mallotivora.
- These hits have 96-86% sequence identity, showing some level of divergence.

Matches with Pantoea Species

- Pantoea species, closely related to Erwinia, also exhibit significant similarity, with scores ranging from 1184 to 1149 and identity between 84–81%.
- Pantoea includes both opportunistic plant pathogens and beneficial endophytes, so this result suggests a common evolutionary origin for the T3SS structure.

All 100 similar sequences have an E value of 0 and 29 sequences have 100% query coverage

The lowest query coverage is 90% while the lowest percentage identity is 66.96% (type III secretion system outer membrane ring subunit SctC [*Dickeya lacustris*])

The non-redundant (NR) protein sequence closely matches the T3SS outer membrane ring subunit SctC from *Erwinia amylovora* with 100% identity.

The sequence is highly conserved across multiple *Erwinia* species and some *Pantoea* species, confirming its essential function.

Given its role in bacterial pathogenesis, this protein could be a potential target for disease control in plant pathology.

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

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per Ident	Acc. Len	Accession
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Hypersensitivity response sec...	<i>Pseudomonas sy...</i>	566	566	99%	0.0	44.16%	701	Q01723.2
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=YscC secretin; Flags: Precurs...	<i>Yersinia enterocol...</i>	251	251	72%	2e-73	31.25%	607	Q01244.1
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Type III secretion protein Ysc...	<i>Yersinia pestis</i>	246	246	76%	1e-71	30.52%	607	Q56974.1
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Hypersensitivity response sec...	<i>Ralstonia pseudo...</i>	211	211	75%	5e-59	28.87%	568	Q52498.1
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Outer membrane protein Mxi...	<i>Shigella flexneri</i>	175	175	70%	6e-46	27.63%	566	Q04641.1
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Outer membrane protein Mxi...	<i>Shigella sonnei</i>	175	175	70%	6e-46	27.63%	566	Q55293.1
RecName: Full=SPI-2 type 3 secretion system secretin; Short=T3SS-2 secretin; AltName: Full=Outer membrane prot...	<i>Salmonella enteri...</i>	158	158	73%	2e-40	26.13%	497	D02W69.1
RecName: Full=Type 3 secretion system secretin; Short=T3SS secretin; AltName: Full=Hypersensitivity response sec...	<i>Xanthomonas eu...</i>	132	132	37%	5e-31	31.25%	607	P80151.1
RecName: Full=DNA utilization protein HofQ; Flags: Precursor [Escherichia coli K-12]	<i>Escherichia coli K...</i>	104	104	37%	1e-22	30.00%	412	P34749.2
RecName: Full=Secretin ExeD; AltName: Full=General secretion pathway protein D; AltName: Full=Type II secretion...	<i>Aeromonas hydro...</i>	92.4	92.4	37%	3e-18	28.97%	678	P31780.2
RecName: Full=Secretin OutD; AltName: Full=General secretion pathway protein D; AltName: Full=Pectic enzymes s...	<i>Dickeya chrysant...</i>	92.0	92.0	42%	5e-18	28.05%	712	P31700.1
RecName: Full=Secretin ExeD; AltName: Full=General secretion pathway protein D; AltName: Full=Type II secretion...	<i>Aeromonas salm...</i>	90.1	90.1	37%	2e-17	28.42%	678	P45778.1
RecName: Full=Secretin OutD; AltName: Full=General secretion pathway protein D; AltName: Full=Pectic enzymes s...	<i>Dickeya dadanti...</i>	89.0	89.0	42%	4e-17	27.66%	710	Q01565.1
RecName: Full=Secretin OutD; AltName: Full=General secretion pathway protein D; AltName: Full=Pectic enzymes s...	<i>Pectobacterium c...</i>	86.7	86.7	84%	2e-16	22.80%	650	P31701.2
RecName: Full=Competence protein E; AltName: Full=DNA transformation protein ComE; Flags: Precursor [Haemop...	<i>Haemophilus influ...</i>	79.7	79.7	37%	2e-14	26.41%	445	P31772.2
RecName: Full=Putative secretin GspD; AltName: Full=Putative general secretion pathway protein D; AltName: Full=...	<i>Escherichia coli K...</i>	79.7	79.7	37%	3e-14	26.28%	650	P45758.2
RecName: Full=Secretin PulD; AltName: Full=General secretion pathway protein D; AltName: Full=Pullulanase secret...	<i>Klebsiella pneum...</i>	79.0	79.0	76%	6e-14	23.45%	660	P15644.1
RecName: Full=Fimbrial assembly protein PilQ; Flags: Precursor [Pseudomonas aeruginosa PAO1]	<i>Pseudomonas ae...</i>	76.3	76.3	41%	4e-13	22.63%	714	P34750.2
RecName: Full=Type IV pilus biogenesis and competence protein PilQ; Flags: Precursor [Neisseria meningitidis Z2491]	<i>Neisseria meningi...</i>	73.2	73.2	42%	4e-12	23.97%	761	Q9JVV4.1
RecName: Full=Type IV pilus biogenesis and competence protein PilQ; Flags: Precursor [Neisseria meningitidis MC58]	<i>Neisseria meningi...</i>	72.8	72.8	42%	5e-12	23.97%	769	Q70M91.2
RecName: Full=Type IV pilus biogenesis and competence protein PilQ; Flags: Precursor [Neisseria meningitidis H44/...	<i>Neisseria meningi...</i>	72.8	72.8	42%	5e-12	23.97%	777	Q9ZHF3.2

Analysis of SwissProt BLAST Results:

- All 33 related sequences have extremely low E values, close to zero.
- The maximum query coverage observed is 98%, associated with the Type 3 secretion system secretin from *Pseudomonas syringae* pv. *Syringae*.
- The minimum query coverage recorded is 23%, occurring in three sequences.
- The highest identity (44.16%) is with *Pseudomonas syringae* T3SS secretin (Q01723.2).
- *Pseudomonas syringae* (highest identity) is a well-known plant pathogen.
- *Yersinia*, *Shigella*, and *Salmonella* are human pathogens. *Ralstonia* and *Xanthomonas* are plant pathogens. *Escherichia coli* and *Klebsiella pneumoniae* are opportunistic human pathogens.
- The lowest identity percentage is 21.75%, found in the uncharacterized protein y4xJ from *Sinorhizobium fredii* NGR234.
- The low sequence identity (~30%) in some cases but functional similarity indicates potential structural conservation despite sequence divergence.

- The query protein is highly similar to T3SS secretins, particularly HrpH from *Pseudomonas syringae*.

2. List the algorithm parameters used for the search (Q1).

Default parameters which we used for the first question:

Algorithm parameters

General Parameters

Max target sequences	100 ▼	Select the maximum number of aligned sequences to display ?
Short queries	<input checked="" type="checkbox"/> Automatically adjust parameters for short input sequences ?	
Expect threshold	0.05	?
Word size	5 ▼	?
Max matches in a query range	0	?

Scoring Parameters

Matrix	BLOSUM62 ▼	?
Gap Costs	Existence: 11 Extension: 1 ▼	?
Compositional adjustments	Conditional compositional score matrix adjustment ▼	?

Filters and Masking

Filter	<input type="checkbox"/> Low complexity regions ?
Mask	<input type="checkbox"/> Mask for lookup table only ?
	<input type="checkbox"/> Mask lower case letters ?

3. What is the sequence identity of the query sequence (given in Q1) with AAK81929.1?

Steps:

Use align 2 or more sequences and Run BLASTP using the HrcC protein sequence and AAK81929.1 in subject sequence.

The screenshot displays the NCBI BLASTP web interface. At the top, there are tabs for different BLAST programs: blastn, **blastp**, blastx, tblastn, and tblastx. The title 'Align Sequences Pro' is visible on the right. Below the tabs, a blue banner reads 'BLASTP programs search protein subjects'. The interface is divided into three main sections: 'Enter Query Sequence', 'Enter Subject Sequence', and 'Program Selection'. In the 'Enter Query Sequence' section, the 'Enter accession number(s), gi(s), or FASTA sequence(s)' field contains '1336093'. The 'Query subrange' fields (From and To) are empty. Below this, the 'Or, upload file' section shows a 'Choose File' button and 'No file chosen'. The 'Job Title' field is empty, with a prompt 'Enter a descriptive title for your BLAST search'. A checkbox labeled 'Align two or more sequences' is checked. The 'Enter Subject Sequence' section has the 'Enter accession number(s), gi(s), or FASTA sequence(s)' field containing 'AAK81929.1'. The 'Subject subrange' fields (From and To) are empty. Below this, the 'Or, upload file' section shows a 'Choose File' button and 'No file chosen'. The 'Program Selection' section shows the 'Algorithm' dropdown set to 'blastp (protein-protein BLAST)'. At the bottom, there is a blue 'BLAST' button and a checkbox for 'Show results in a new window'.

blastn **blastp** blastx tblastn tblastx Align Sequences Pro

BLASTP programs search protein subjects

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#) Query subrange [?](#)

1336093 From
To

Or, upload file Choose File No file chosen [?](#)

Job Title
Enter a descriptive title for your BLAST search [?](#)

☒ Align two or more sequences [?](#)

Enter Subject Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#) Subject subrange [?](#)

AAK81929.1 From
To

Or, upload file Choose File No file chosen [?](#)

Program Selection

Algorithm ☒ blastp (protein-protein BLAST)
Choose a BLAST algorithm [?](#)

BLAST Search [protein sequence](#) using **Blastp (protein-protein BLAST)**
☐ Show results in a new window

43.20% sequence identity

blast.ncbi.nlm.nih.gov/Blast.cgi

NIH National Library of Medicine
National Center for Biotechnology Information

BLAST® » blastp suite-2sequences » results for RID-VY17JB2B114

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 1336093|GENBANK|OUTER MEMBRANE INTEGRAL
 RID [VY17JB2B114](#) Search expires on 02-27 19:11 pm [Download All](#)
 Program Blast 2 sequences [Citation](#)
 Query ID lc|Query_3304695 (amino acid)
 Query Descr 1336093|GENBANK|OUTER MEMBRANE INTEGRAL MEME ...
 Query Length 676
 Subject ID [AAK81929.1](#) (amino acid)
 Subject Descr RscC [Pseudomonas fluorescens]
 Subject Length 713
 Other reports [Multiple alignment](#) [MSA viewer](#)

Filter Results

Percent Identity to E value to Query Coverage to
[Filter](#) [Reset](#)

Descriptions **Graphic Summary** **Alignments** **Dot Plot**

Sequences producing significant alignments Download Select columns Show 100

☒ select all 1 sequences selected [GenPept](#) [Graphics](#) [Multiple alignment](#) [MSA Viewer](#)

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per Ident	Acc. Len	Accession
<input checked="" type="checkbox"/> RscC [Pseudomonas fluorescens]	Pseudomonas fluorescens	530	530	97%	0.0	43.20%	713	AAK81929.1

4. How far are hemoglobin (beta) sequences in humans and chicken similar?

Human hemoglobin beta: P68871.

Chicken hemoglobin beta: P02112.

P68871 · HBB_HUMAN

Proteinⁱ | Hemoglobin subunit beta
 Geneⁱ | HBB
 Statusⁱ | UniProtKB reviewed (Swiss-Prot)
 Organismⁱ | Homo sapiens (Human)

Amino acids | 147 [\(go to sequence\)](#)
 Protein existenceⁱ | Evidence at protein level
 Annotation scoreⁱ | [5/5](#)


[Entry](#) [Variant viewer](#) [794](#) [Feature viewer](#) [Genomic coordinates](#) [Publications](#) [External links](#) [History](#)

[Tools](#) [Download](#) [Add](#) [Community curation](#) (3) [Add a publication](#) [Entry feedback](#)

P02112 · HBB_CHICK

Proteinⁱ | Hemoglobin subunit beta


Geneⁱ | HBB

Statusⁱ |  UniProtKB reviewed (Swiss-Prot)

Organismⁱ | [Gallus gallus \(Chicken\)](#)

Amino acids | 147 ([go to sequence](#))

Protein existenceⁱ | Evidence at protein level

Annotation scoreⁱ | 

Based on the BLASTP alignment results, the similarity between human and chicken hemoglobin beta sequences is as follows:

1. Identity: 69% (102 out of 147 amino acids are exactly the same)
2. Positives: 82% (121 out of 147 amino acids share similar chemical properties, meaning functional similarity)
3. Gaps: 0% (No insertions or deletions, indicating a well-aligned sequence)
4. E-value: $1e-80$ (Highly significant match, indicating evolutionary conservation)

blast.ncbi.nlm.nih.gov/Blast.cgi#alnHdr_P02112

Descriptions Graphic Summary **Alignments** Dot Plot

Alignment view Pairwise [Restore defaults](#)

1 sequences selected

[Download](#) [GenPept](#) [Graphics](#)

hemoglobin subunit beta [Gallus gallus]

Sequence ID: [NP_990820.1](#) Length: 147 Number of Matches: 1

[See 5 more title\(s\)](#) [See all Identical Proteins\(IPG\)](#)

Range 1: 1 to 147 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

Score	Expect	Method	Identities	Positives	Gaps
221 bits(564)	1e-80	Compositional matrix adjust.	102/147(69%)	121/147(82%)	0/147(0%)
Query 1	MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVYPWTQRFESFGDLSTPDAMGNPK	60			
Sbjct 1	MVH T EEK +T LWGKVVN E G EAL RLL+VYPWTQRF SFG+LS+P A++GNP	60			
Query 61	VKAHGKKVLGAFSDGLAHLNDNLKGTFTLSELHCDKLHVDPENFRLLGNVLCVLAHFG	120			
Sbjct 61	V+AHGKKVL +F D + +LDN+K TF+ LSELHCDKLHVDPENFRLLG++L+ VLA HF	120			
Query 121	KEFTPPVQAAYQKVVAGVANALAHKYH	147			
Sbjct 121	K+FTP QAA+QK+V VA+ALA KYH	147			

5. Write a program to list all the matching pentapeptides (which occur in both the sequences) and their frequency of occurrence in given sequences.

Human hemoglobin beta: P68871.

Chicken hemoglobin beta: P02112.

Got these protein sequence from uniprot:

```
>sp|P02112|HBB_CHICK Hemoglobin subunit beta OS=Gallus gallus OX=9031
GN=HBB PE=1 SV=2
MVHWTAEKQLITGLWGKVNVAECGAEALARLLIVYPWTQRFASFGNLSPTAILGNPM
VRAHGKKVLTSGDAVKNLNDIKNTFSQSELHCDKLHVDPENFRLLGDILIVLAHFS
KDFTPECQAAWQKLVRVVAHALARKYH
```

```
>sp|P68871|HBB_HUMAN Hemoglobin subunit beta OS=Homo sapiens OX=9606
GN=HBB PE=1 SV=2
```

MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK
K
VKAHGKKVLGAFSDGLAHLNLIKGTATLSELHCDKLHVDPENFRLLGNVLCVLAHHFG
KEFTPPVQAAYQKVVAGVANALAHKYH

This code identifies common pentapeptides (five-amino-acid sequences) between human and chicken hemoglobin sequences. It does this by generating all possible pentapeptides from both sequences and counting their occurrences using Python's `Counter` from the `collections` module.

The program then finds overlapping pentapeptides and displays their frequencies in each sequence. This analysis helps compare structural similarities between the two proteins, revealing conserved regions that might be functionally important.

```
p5_BS22B009.ipynb ☆ ☁
File Edit View Insert Runtime Tools Help
+ Code + Text

#Question5
from collections import Counter

def get_pentapeptide_counts(sequence):
    pentapeptides = [sequence[i:i+5] for i in range(len(sequence) - 4)]
    return Counter(pentapeptides)

human_seq = "MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK"
human_seq += "VKAHGKKVLGAFSDGLAHLNLIKGTATLSELHCDKLHVDPENFRLLGNVLCVLAHHFG"
human_seq += "KEFTPPVQAAYQKVVAGVANALAHKYH"

chicken_seq = "MVHTAEKQLITGLWGKVVNVAECGAEALARLLIVYPWTQRFFASFGNLSPTAILGNPM"
chicken_seq += "VRAHGKKVLTSFGDAVKNLNLIKNTFSQLSELHCDKLHVDPENFRLLGDILIIIVLAHFS"
chicken_seq += "KDFTPECQAAWQKLVRAHALARKYH"

human_counts = get_pentapeptide_counts(human_seq)
chicken_counts = get_pentapeptide_counts(chicken_seq)

common_pentapeptides = set(human_counts.keys()) & set(chicken_counts.keys())

print("Matching Pentapeptides and Their Frequency:")
for peptide in sorted(common_pentapeptides):
    human_count = human_counts[peptide]
    chicken_count = chicken_counts[peptide]
    total_count = human_count + chicken_count
    print(f"{peptide} occurs {total_count} time(s) in both peptides "
          f"({human_count} time(s) in Seq. 1 and {chicken_count} time(s) in Seq. 2)")
```

Result for the code:

Matching Pentapeptides and Their Frequency:

```
AHGKK occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
CDKLH occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
DKLHV occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
DPENF occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
ELHCD occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
ENFRL occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
FRLLG occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
GKKVL occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
GKVVN occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
HCDKL occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
HGKKV occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
HVDPE occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
KLHVD occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
LHCDK occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
LHVDP occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
LSELH occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
LWGKV occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
NFRLL occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
PENFR occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
PWTQR occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
SELHC occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
TQRFF occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
VDPEN occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
VYPWT occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
WGKVN occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
WTQRF occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
YPWTQ occurs 2 time(s) in both peptides (1 time(s) in Seq. 1 and 1 time(s) in Seq. 2)
```

6. Write a program to compute sequence identity, similarity, query coverage and gap percentage from the alignment of human and chicken hemoglobin sequences (refer Q4).

Here's my Python program to compute sequence identity, similarity, query coverage, and gap percentage between human and chicken hemoglobin sequences. I used a character-by-character comparison, considering exact matches for identity and grouping similar amino acids for similarity. Query coverage is calculated based on the aligned portion, and gaps are accounted for separately.

It calculates four key metrics:

1. Sequence Identity – The percentage of exact matches between the two sequences.
2. Sequence Similarity – The percentage of residues that either match exactly or belong to the same biochemical similarity group.

3. Query Coverage – The proportion of the query sequence (human hemoglobin) that aligns with the target sequence (chicken hemoglobin).
4. Gap Percentage – The fraction of positions where either sequence has a gap ('-').

To determine similarity, I grouped amino acids based on their properties, such as hydrophobicity and charge. Then, I iterated through the sequences to count matches, similar residues, and gaps. Finally, I calculated and printed the results.

```
p5_BS22B009.ipynb ☆
File Edit View Insert Runtime Tools Help
Code + Text
+ human_seq = ("MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMVGNPK"
               "VKAHGKKVLGAFTSDGLAHLNLTGKTFATLSSEILHCDKLVDPENFRLLGNVLVCVLAHHFG"
               "KEFTTPVQAAAYQKVAVGAVANALAHKYH")
)
+ chicken_seq = ("MVHLTAEKQLITGLWGKVNVAECGAEALRLLIVYPWTQRFFASFGNLSPTAILGNPM"
                 "VRAHGKKVLTSFGDAVKNLNLTGKTFATLSSEILHCDKLVDPENFRLLGDILIVLAHFS"
                 "KDFTEPCQAAYQKLVAVAHALARKYH")
)
+ def calculate_metrics(seq1, seq2):
    matches = 0
    similar = 0
    gaps = 0
    total_length = min(len(seq1), len(seq2))
    similarity_groups = [
        set("NDEQ"), set("MILV"), set("FYW"), set("KRH"), set("ST"), set("AGP")
    ]
    for i in range(total_length):
        if seq1[i] == seq2[i]:
            matches += 1
            similar += 1
        elif seq1[i] == '-' or seq2[i] == '-':
            gaps += 1
        else:
            for group in similarity_groups:
                if seq1[i] in group and seq2[i] in group:
                    similar += 1
                    break
    identity = (matches / total_length) * 100
    similarity = (similar / total_length) * 100
    query_coverage = (total_length / len(seq1)) * 100
    gap_percentage = (gaps / total_length) * 100
```

```

    break
    identity = (matches / total_length) * 100
    similarity = (similar / total_length) * 100
    query_coverage = (total_length / len(seq1)) * 100
    gap_percentage = (gaps / total_length) * 100
    return identity, similarity, query_coverage, gap_percentage

identity, similarity, query_coverage, gap_percentage = calculate_metrics(human_seq, chicken_seq)

print(f"Sequence Identity: {identity:.2f}%")
print(f"Sequence Similarity: {similarity:.2f}%")
print(f"Query Coverage: {query_coverage:.2f}%")
print(f"Gap Percentage: {gap_percentage:.2f}%")
```

Result:

```
Sequence Identity: 69.39%  
Sequence Similarity: 83.67%  
Query Coverage: 100.00%  
Gap Percentage: 0.00%
```

7. Obtain the multiple sequence alignment for TIM barrel proteins from different organisms (select 20 proteins, for example). Compare the results obtained with Clustal Omega, MAFFT, and MUSCLE. List 5 residue positions which are aligned differently in these three methods.

Steps:

1. Collect 20 TIM barrel protein sequences from UniProt.
2. Use three alignment tools:
 - Clustal Omega (<https://www.ebi.ac.uk/Tools/msa/clustalo/>)
 - MAFFT (<https://mafft.cbrc.jp/alignment/server/>)
 - MUSCLE (<https://www.ebi.ac.uk/Tools/msa/muscle/>)
3. Compare alignments and list 5 residue positions with differences.

20 sequences download in FASTA format from Uniprot :

← → ↺ 🌐 mafft.cbrc.jp/alignment/server/spool/_ho.25030421147746201TEgar7ekGWAMHvCXihlsfnormal.html

[Clustal format](#) | [Fasta format](#) | [MAFFT result](#) | [View](#) | [Tree](#) | [Refine dataset](#) | [Return to home](#)

[View](#)

[Reformat](#) to GCG, PHYLIP, MSF, NEXUS, uppercase/lowercase, etc. with Readseq

[GUIDANCE2](#) computes the residue-wise confidence scores and extracts well-aligned residues.

[Refine dataset](#)

[Phylogenetic tree](#)

MAFFT-**L-INS-i** Result

```
CLUSTAL format alignment by MAFFT (v7.511)

sp|O14744|MA-----AMAVGG-----AGGSRVSSGRDLNCVPEI
sp|P46580|MSNRTYADNLFPOQVAEQHEEQMSGSSPKSNSPSRSISVVEANSRIHIQW-----M
sp|P38274|-----MHSNVFVGVRPGFNHKKHKKSRFLEN
sp|P0A6C1|-----
sp|P22936|-----
sp|P24215|-----
sp|P30147|-----
sp|Q81RQ4|-----
sp|P45541|-----
sp|P35914|-----
sp|Q8XKU1|-----
sp|Q92520|-----
sp|Q9RUP5|-----
sp|P21826|-----
sp|P30952|-----
sp|P37330|MSQITQSRRLRIDANFRFVDEEVLPGTGLDAAAFWRNFDEIVHDLAPENROLLAERDR
sp|P9WK16|MTDRVSVGNLRIRARVLYDFVNNELPGTDIDPDSFWAGVDRKVADLTQONQALLNARDEL
sp|P32170|-----
sp|P54802|-----
sp|P4062|-----

sp|O14744|ADTLGAVAKQGFDFLCMPVFHPRFKREFIQEP-----AKNRPGPQTRSD
sp|P46580|ATTLDAENLDRHVATFCTRLGEFKYFVYVPIGGVVRAF--WTPNGSAENHPVIDLPD
sp|P38274|VSSHSPELPSNYDYVLLPITTPRYKEIVGQVFRDQFQSQIQNWKPLQIPEPOLQDICI
sp|P22936|-----
sp|P24215|-----
sp|P30147|-----
sp|Q81RQ4|-----
sp|P45541|-----
sp|P35914|-----
sp|Q8XKU1|-----
sp|Q92520|-----
sp|Q9RUP5|-----
sp|P21826|-----
```

Multiple sequence alignment (MSA) with MUSCLE:

← → ↺ 🌐 ebi.ac.uk/jdispatcher/msa/muscle/summary?jobid=muscle-120250304-121957-0456-57805862-p1m8js=pass

Results for Job ID muscle-120250304-121957-0456-57805862-p1m [Copy](#) [Resubmission](#)

Tool Output **Alignments** Phylogenetic Tree Results Viewers Result Files Submission Details

Nightingale

COLOR SCHEME LEGEND

clustal2 ARND CQEGHILKMF PSTWYVBXZ

20 sequences

100 200 300 400 500 600 700 800 900 1,000

1

SPIQ81RQ4|ASBF_BACAN
SPIP04062|GBA1_HUMAN
SPIP30147|HYL_ECOLI
SPIP45541|FRLC_ECOLI
SPIP32170|RHAA_ECOLI
SPIP54802|IANAG_HUMAN
SPIP0A6C1|END4_ECOLI
SPIP22936|APN1_YEAST
SPIP30952|MLS2_YEAST
SPIP24215|UXUA_ECOLI
SPIP38274|HSL7_YEAST
SPIO14744|ANM5_HUMAN
SPIP46580|ANM5_CAEL
SPIP35914|HMGCL_HUMAN
SPIP37330|IMASZ_ECOLI
SPIP9WK16|IMASZ_MYCTO
SPIQ9RUP5|ITPIS_DEIRA
SPIQ8XKU1|ITPIS_CLOPE
SPIQ8XKU1|ITPIS_CLOPE

-----WKYSLCTISFRHQLISFTDI-----VQFAYEN-----
-----VEFSSPSRECPKPLSRVSIWAGSLTGL-----
-----VLRFSANLSWLFGEYDFLA-----
-----WKTGVTCTGHORLP1EH-----
-----MTTOLEQANIELAKQRFAAVGIIDVEEALROL-----
MEAVAVAAAVGVLLLAGAGGAAGDEAREAAVRLVARLLGPGPAADFSVVERALAAKPGLDTYSLGGGGAARVVRGSGTGAAAAAGHRYLRDFCG
-----WKYIGAHVSAAGGLANAAIRAA-----
-----WPSPTSPFVRSAYSKYKFG-----
MVKISLDNTALYADITTPQFEPKSTTVADI1TKDALEFI-----VLLHRTFINSTRKOLLANRSLQSKLDSGEYRF-----
MVKVSLDNVLKLVVDVDEKPEFFKPSSTTVADI1TKDALEFI-----VLLHRTFINSTRKOLLANRQIVQKLLDSGSYHL-----
-----VEQTWRVYGPNDPVSADVRQAGAT-----
-----MHSNVFVGVRPGFNHKKHKKSRFLEN-----ENVSSHSPELPSNYDYVLLPITTPRYKEIVGQVIF-----
-----MAAVAVGGAGGSRVSSGRDL-----NCVPEIADTLGAVAKQGFDFLCMPVFH-----
-----VSNRTYADNLFPOQVAEQHEEQMSGSSPKSNSPSRSISVVEANSRIHIQWATTLDVAENLDRHVATFCTRLGEFKYFVYVPIG-----
-----MAVRKALPRLVGLASLRAVSTSSN-----
-----VTSQITQSRRLRIDANFRFVDEEVLPGTGLDAAAFWRNFDEIVHDLAPENROLLAERDRIQAAALDENHRSNPGPVKD-----
-----MTDRVSVGNLRIRARVLYDFVNNELPGTDIDPDSFWAGVDRKVADLTQONQALLNARDELQAIQKIHRRRVIEPID-----
-----VOTLLALNKKWKTPTTEAR-----DKVYADLTQONQALLNARDELQAIQKIHRRRVIEPID-----
-----WRTPIIAGNNKNNHYTIDEAV-----

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

1. Residue Position 1:

- Clustal Omega aligns M (Methionine) in P46580, while the other sequences have gaps.
- MAFFT aligns M in Q9Z520, with gaps in other sequences.
- MUSCLE aligns M in P54802, while others have gaps.

2. Residue Position 2:

- Clustal Omega places S (Serine) in P46580, with gaps in the other sequences.
- MAFFT aligns T (Threonine) in Q9Z520, with gaps in the rest.
- MUSCLE aligns E (Glutamic Acid) in P54802, while others have gaps.

3. Residue Position 3:

- Clustal Omega aligns N (Asparagine) in P46580, with gaps elsewhere.
- MAFFT assigns T (Threonine) in Q9Z520, while M (Methionine) is present in other sequences.
- MUSCLE places A (Alanine) in P54802, with gaps in the remaining sequences.

4. Residue Position 4:

- Clustal Omega aligns R (Arginine) in P46580, with gaps in the other sequences.
- MAFFT shows different amino acids across different sequences, meaning no consistent alignment.
- MUSCLE aligns V (Valine) in P54802, while others have gaps.

5. Residue Position 5:

- Clustal Omega aligns T (Threonine) in P46580, with gaps in the other sequences.
- MAFFT presents different amino acids across different sequences, showing variation.
- MUSCLE places A (Alanine) in P54802, with gaps in the rest.

Conclusion:

- Clustal Omega, MAFFT, and MUSCLE show significant differences in alignment at these five positions.

- Clustal Omega and MUSCLE tend to align specific residues while keeping gaps in other sequences, whereas MAFFT exhibits more variation in residue choices.
 - MAFFT shows the highest variability, while Clustal Omega and MUSCLE align specific amino acids with more gaps in others.
 - These differences highlight the varying sensitivity and algorithmic approaches of multiple sequence alignment (MSA) tools.
-

8. Blast the below sequence 'EPDMRTPIAHTMAW' against the PDB database. Analyze the results and discuss the significance of the results.

Steps:

Go to NCBI BLAST → Select BLASTP.

2. Enter sequence 'EPDMRTPIAHTMAW'.

3. Select PDB (Protein Data Bank) database.

4. Run the search and analyze:

- Top hits
- E-value, identity, and coverage

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

EPDMRTPIAHTMAW

Query subrange [?](#)

From

To

Or, upload file

Choose File

No file chosen [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Databases

☒ Standard databases (nr etc.): ☐ Experimental databases

Compare

☐ Select to compare standard and experimental database [?](#)

Standard

Database

Protein Data Bank proteins(pdb) [?](#)

Organism

Optional

☐ exclude [Add organism](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

Exclude

Optional

☐ Models (XM/XP) ☐ Non-redundant RefSeq proteins (WP) ☐ Uncultured/environmental sample sequences

Program Selection

Algorithm

☒ blastp (protein-protein BLAST)

☐ PSI-BLAST (Position-Specific Iterated BLAST)

☐ PHI-BLAST (Pattern Hit Initiated BLAST)

☐ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm [?](#)

BLAST

Search database **pdb** using **Blastp (protein-protein BLAST)**

☒ Show results in a new window

Blast results:

blast.ncbi.nlm.nih.gov/Blast.cgi

Other reports [Distance tree of results](#) [Multiple alignment](#) [MSA viewer](#) [Filter](#) [Reset](#)

Descriptions [Graphic Summary](#) [Alignments](#) [Taxonomy](#)

Sequences producing significant alignments Download Select columns Show 100

☒ select all 100 sequences selected [GenPept](#) [Graphics](#) [Distance tree of results](#) [Multiple alignment](#) [MSA Viewer](#)

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli str. K-12 substr. W3110]	Escherichia coli s...	53.2	53.2	100%	2e-10	100.00%	424	2EGH_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli K-12]	Escherichia coli ...	53.2	53.2	100%	2e-10	100.00%	420	3ANL_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli K-12]	Escherichia coli ...	53.2	53.2	100%	2e-10	100.00%	410	3R0L_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli]	Escherichia coli	53.2	53.2	100%	2e-10	100.00%	406	1Q0L_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli K-12]	Escherichia coli ...	53.2	53.2	100%	2e-10	100.00%	398	1K9H_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli]	Escherichia coli	53.2	53.2	100%	2e-10	100.00%	398	1T1R_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli]	Escherichia coli	43.1	43.1	100%	9e-07	85.71%	406	1Q0H_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Escherichia coli]	Escherichia coli	43.1	43.1	100%	9e-07	85.71%	400	1JVS_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Zymomonas mobilis]	Zymomonas mob...	42.2	42.2	93%	2e-06	84.62%	388	1R0K_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Yersinia pseudotuberculosis YPIII]	Yersinia pseudot...	40.9	40.9	86%	5e-06	91.67%	401	3IIE_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Acinetobacter baumannii AB307-0294]	Acinetobacter ba...	35.4	35.4	93%	4e-04	76.92%	406	4ZN6_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Moraxella catarrhalis]	Moraxella catarrh...	33.3	33.3	86%	0.003	75.00%	432	4ZQE_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Moraxella catarrhalis]	Moraxella catarrh...	33.3	33.3	86%	0.003	75.00%	415	4ZQG_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Vibrio vulnificus CMCP6]	Vibrio vulnificus ...	28.2	28.2	86%	0.17	75.00%	427	5KQO_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Vibrio vulnificus CMCP6]	Vibrio vulnificus ...	28.2	28.2	86%	0.17	75.00%	405	5KRR_A
<input checked="" type="checkbox"/>	Chain A_Glutamate dehydrogenase 2 [Arabidopsis thaliana]	Arabidopsis thali...	26.1	39.0	93%	0.95	69.23%	414	8QWM_A
<input checked="" type="checkbox"/>	Chain A_1-deoxy-D-xylulose 5-phosphate reductoisomerase [Toxoplasma gondii]	Toxoplasma gondii	24.8	24.8	93%	2.7	53.85%	472	8S65_A
<input checked="" type="checkbox"/>	Chain A_Acyl-CoA dehydrogenase [Streptomyces ficellus]	Streptomyces fic...	24.4	24.4	43%	3.8	100.00%	384	8HK0_A
<input checked="" type="checkbox"/>	Chain A_Glutamate dehydrogenase 1 [Arabidopsis thaliana]	Arabidopsis thali...	24.4	37.3	93%	3.8	64.29%	414	6YEH_A
<input checked="" type="checkbox"/>	Chain Bh_30S ribosomal protein S9, mitochondrial [Arabidopsis thaliana]	Arabidopsis thali...	24.4	24.4	43%	3.8	100.00%	430	6XYW_Bh

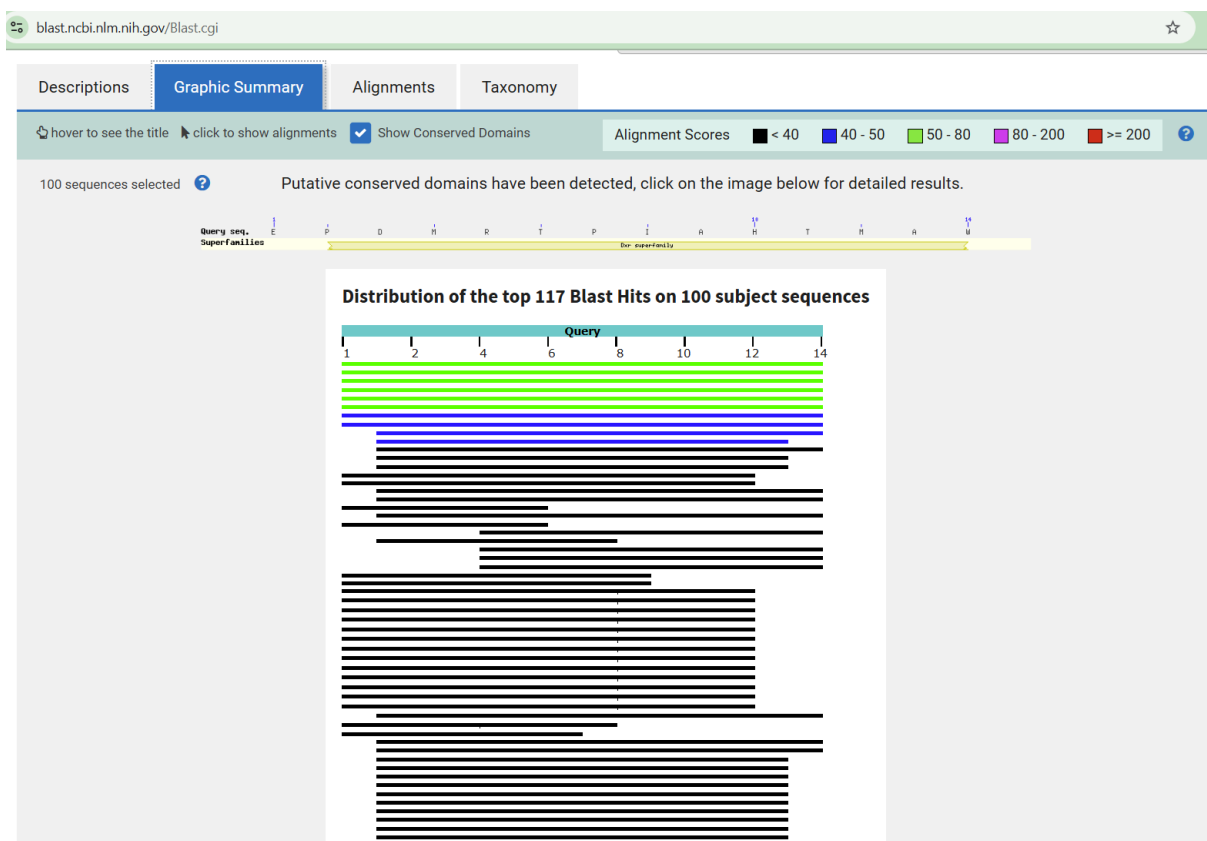
The BLAST results indicate that the query sequence "EPDMRTPIAHTMAW" has 100% query coverage and high sequence identity with 1-deoxy-D-xylulose 5-phosphate reductoisomerase (DXR) in *Escherichia coli* and other bacterial species.

Sequences with lower scores, such as *Zymomonas mobilis*, *Yersinia pseudotuberculosis*, and *Moraxella catarrhalis*, still show high sequence identity (ranging from 85.71% to 100%), indicating a strong conservation of this sequence across different bacterial species.

The Query Cover is 100% for the top sequences, which means that the entire query sequence (all 14 amino acids) aligns completely with the matched sequences in the database.

Some of the lower-ranked hits (e.g., *Arabidopsis thaliana*, *Toxoplasma gondii*) show slightly lower query coverage (~93% or lower), meaning only part of the sequence aligns with their respective proteins.

The E-value (Expect value) for the top matches is 2e-10, which is extremely low. This indicates a highly significant match, meaning the alignment is not due to random chance.



The graphic summary shows that the sequence aligns well with members of the DXR superfamily, highlighting its evolutionary conservation. Given its high sequence identity and conservation, this sequence could be essential for enzyme function or stability.

This suggests the query sequence is part of a highly conserved functional motif in DXR, a key enzyme in the non-mevalonate pathway of isoprenoid biosynthesis. Identifying such a conserved motif with just a short sequence demonstrates the power of sequence alignment in detecting functionally related proteins, which could be useful for drug targeting in bacterial infections.