

Tracking the Evolution of the Hemoglobin Beta (HBB) Gene Across Species

Project Objective:

The objective of this mini-project is to apply the bioinformatics skills acquired in Module 2 to explore the evolutionary conservation of the Hemoglobin Beta (HBB) gene across six different species. This involves retrieving gene sequences using BLAST, performing pairwise and multiple sequence alignments to assess similarity, generating sequence logos to visualize conserved residues, and constructing a phylogenetic tree to infer evolutionary relationships. Through this analysis, we aim to understand how the HBB gene has been conserved or diversified across species and gain insights into its functional and evolutionary significance.

Project Tasks

1: Sequence Retrieval & BLAST Search

Species name	Accession number	% Identity with human HBB
Pan troglodytes	XM_508242.5	100%
Sus scrofa	X86791.1	81.94%
Bos taurus	AB512624.1	80.76%
Mus musculus	EF605506.1	79.03%

2: Pairwise Sequence Alignment

Homosapien Vs Pan troglodytes	Homosapien Vs Mus musculus
Length: 1619 Identity: 627/1619 (38.7%) Similarity: 627/1619 (38.7%) Gaps: 991/1619 (61.2%) Score: 2622.0	Length: 1703 Identity: 751/1703 (44.1%) Similarity: 751/1703 (44.1%) Gaps: 791/1703 (46.4%) Score: 2279.5

Interpretation:

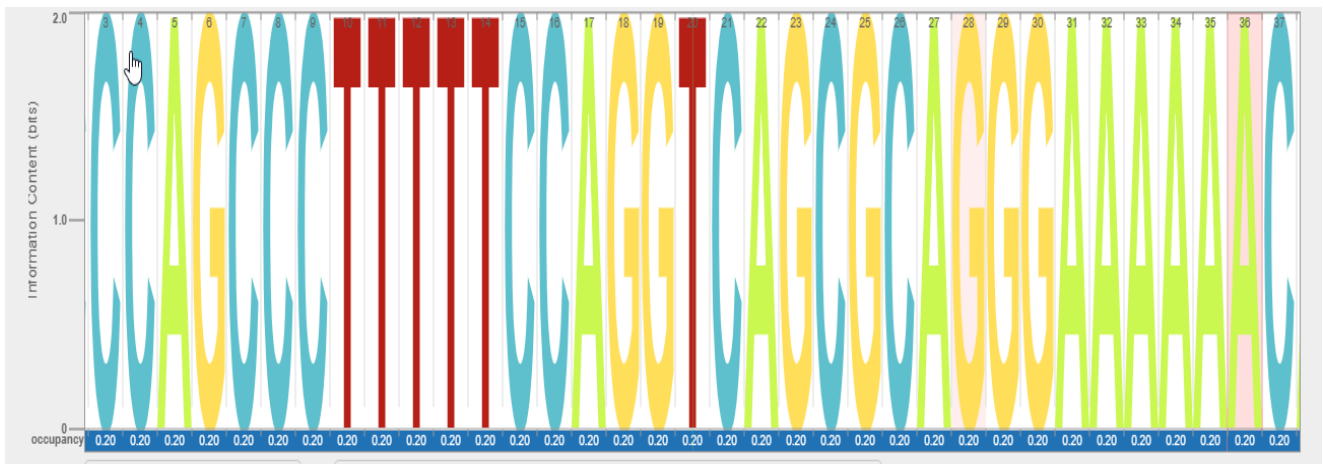
- Higher sequence identity (44.1%) was observed between Human and Mouse HBB genes compared to Human and Chimpanzee (38.7%). However, Chimpanzee alignment has a higher score (2622.0) despite lower identity, likely due to better match regions with fewer substitutions. Gap percentage is significantly higher in Human–Chimpanzee (61.2%) than in Human–Mouse (46.4%), suggesting more alignment interruptions in chimpanzee, which may be due to differences in transcript regions included.

- Human vs. Mouse HBB appears more conserved based on identity percentage and lower gap content. This result is a bit unexpected, since chimpanzees are evolutionarily closer to humans. The discrepancy might be due to differences in transcript lengths or sequence formatting in the alignment files.

3: Multiple Sequence Alignment (MSA)

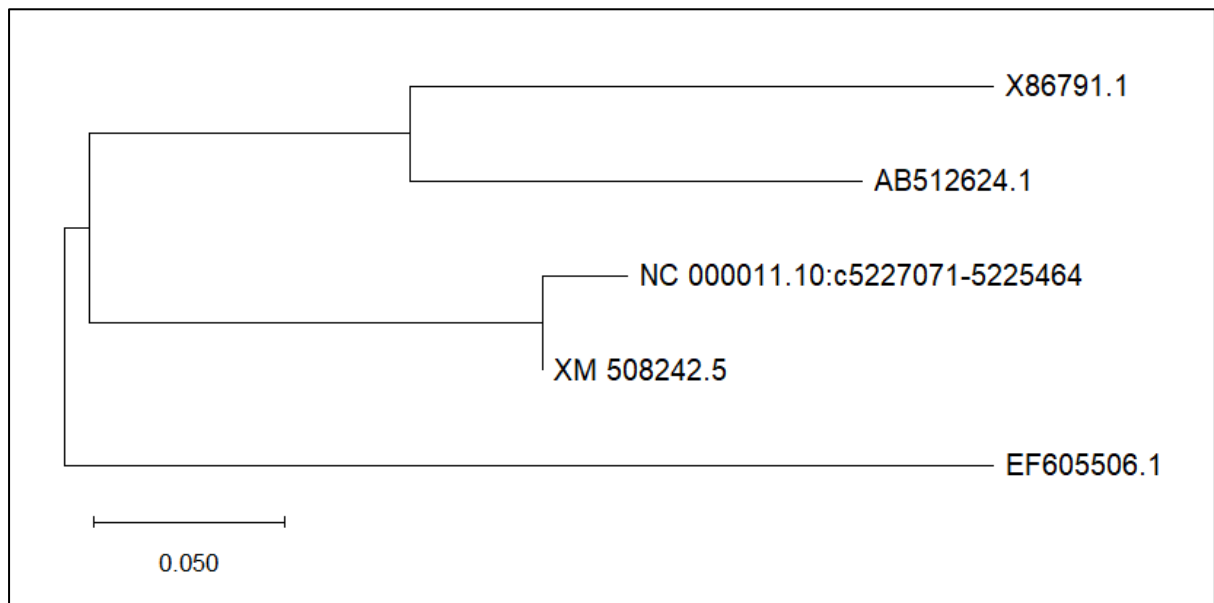
consensus/100%		AGGC GC GG T C ACCC=GGAC=CA
consensus/90%		AGGC GC GG T C ACCC=GGAC=CA
consensus/80%			S CSUUSAGSSUC UUC SSC=SS=CS.S .S.C S S SSC=SSSS.SAGGC GC GG T C ACCC=GGAC=CA
consensus/70%			S CSUUSAGSSUC UUC SSC=SS=CS.S .S.C S S SSC=SSSS.SAGGC GC GG T C ACCC=GGAC=CA
	cov	pid 1121	
1 EF605506.1	100.0%	100.0%	GCGG AC GA AGC GGAGACC A CC C CC C C C A CA GGG AA CCCAAGG GAAGGCCA GGCAAAA
2 X86791.1	97.2%	27.7%	GAGG TC GAGTC G GGGAGCC G CCAA CCGA CCG CA GGGCAA CCCAAGG GAAGGCCA GGCAAG
3 AB512624.1	86.1%	35.0%	GAGG TC GAGTC G GGGAGCC G CCAA CCGA CCG CA GGGCAA CCCAAGG GAAGGCCA GGCAAG
4 NC_000011.10:c5227071-5225464	92.1%	41.6%	GAGG TC GAGTC G GGGAGCC G CCAA CCGA CCG CA GGGCAA CCCAAGG GAAGGCCA GGCAAG
5 XM_508242.5	56.0%	44.0%	GAGG TC GAGTC G GGGAGCC G CCAA CCGA CCG CA GGGCAA CCCAAGG GAAGGCCA GGCAAG
consensus/100%			G=GG S C G=SSS C G=GA=SS U CC=SS S C=SS C C=SU S A G=UUS AA=CC=AAAG GAAGGC=CA=GGCAAU
consensus/90%			G=GG S C G=SSS C G=GA=SS U CC=SS S C=SS C C=SU S A G=UUS AA=CC=AAAG GAAGGC=CA=GGCAAU
consensus/80%			GAGG C GAG CC GGGA=C G CCAC S C=GA C G S A GGCAA=CC=AAAG GAAGGC=CA GGCAAGA
consensus/70%			GAGG C GAG CC GGGA=C G CCAC S C=GA C G S A GGCAA=CC=AAAG GAAGGC=CA GGCAAGA
	cov	pid 1201	
1 EF605506.1	100.0%	100.0%	AGG GA AAC CC AACGAGGCC GAAAAACC GGACAA CAAGGGCACC GCCAGCC CAG GAGC CCAC
2 X86791.1	97.2%	27.7%	AGG GC CCAGTCC C GTGA GGCC GAAAC TC CGACAA CAAGGGCACC GCT AGC GAGC GCAC
3 AB512624.1	86.1%	35.0%	AGG GC GATCC C GTAT GGCA GAAAC TC CGACAA CAAGGGCACC GCTGCGC GAG GAGC GCAC
4 NC_000011.10:c5227071-5225464	92.1%	41.6%	AGG GC CGG CC GTGA GGCC GGCTCACC GGACAA CAAGGGCACC GCCACAC GAG GAGC GCAC
5 XM_508242.5	56.0%	44.0%	AGG GC CGG CC GTGA GGCC GGCTCACC GGACAA CAAGGGCACC GCCACAC GAG GAGC GCAC
consensus/100%			A-G S SSSSS C S USU S GGC S G=SSS S C S G=SU C CAAGGGCACC GC=SU S AG=GAGC S CAC
consensus/90%			A-G S SSSSS C S USU S GGC S G=SSS S C S G=SU C CAAGGGCACC GC=SU S AG=GAGC S CAC
consensus/80%			A-G C SU S CC S AG GA=GGCC G=SSC S C S GACAA CAAGGGCACC GC=SU C GAG GAGC GCAC
consensus/70%			A-G C SU S CC S AG GA=GGCC G=SSC S C S GACAA CAAGGGCACC GC=SU C GAG GAGC GCAC
	cov	pid 1281	
1 EF605506.1	100.0%	100.0%	G GACAAGC GCA G GGA CC GAGAAC CAGGG GAG C GA GGGCACC CC GGG CC CCCC GGC A
2 X86791.1	97.2%	27.7%	G GACCAGC GCA G GGA CC GAGAAC CAGGG GAG C GGGGACCCCTCA-C TTC CCG G-- CTCCTGGG C
3 AB512624.1	86.1%	35.0%	G GATAGC GCA G GGA CC GAGAAC CAGGG GAG T GTGAATCTCAG TTC CT C----- T T
4 NC_000011.10:c5227071-5225464	92.1%	41.6%	G GACAAGC GCA G GGA CC GAGAAC CAGGG GAG C ATGGGACGCT GAT TT CTT CCCC CTT TC
5 XM_508242.5	56.0%	44.0%	G GACAAGC GCA G GGA CC GAGAAC CAG-----
consensus/100%			G GA=AGC GCA= G GGA CC GAGAAC CAU.....
consensus/90%			G GA=AGC GCA= G GGA CC GAGAAC CAU.....
consensus/80%			G GACAAGC GCA= G GGA CC GAGAAC CAGGG GAG S U=SS U=SS SSS S..... SSSSS
consensus/70%			G GACAAGC GCA= G GGA CC GAGAAC CAGGG GAG S U=SS U=SS SSS S..... SSSSS
	cov	pid 1361	
1 EF605506.1	100.0%	100.0%	C GC CAACC CC A CAGAAGGAAGGGGAAGCA C- AGGGAGCAG CSA GA GG G GGA G G-----

4: Sequence Logo Generation



The sequence logo generated from the multiple sequence alignment (MSA) of the HBB gene across selected species using Skylign visually represents the conservation of amino acid residues at each position in the alignment. In the logo, taller stacks of letters indicate positions with higher conservation, while shorter or more mixed stacks reflect greater variability. Several positions exhibit tall, single-letter stacks, signifying highly conserved residues. These conserved amino acids are likely critical for the structural integrity and biological function of the hemoglobin beta protein. Such regions may correspond to essential domains involved in oxygen binding, heme interaction, or subunit folding. Their evolutionary conservation suggests that mutations in these sites could disrupt protein function and are therefore likely to be deleterious, which explains their preservation across species. This highlights the biological importance of these conserved residues in maintaining the functional role of hemoglobin in oxygen transport.

5: Phylogenetic Tree Construction



Interpretation:

The phylogenetic tree constructed from the HBB gene sequences shows the evolutionary relationships among five species based on sequence similarity. In the tree, NC000011.10:c5227071-5225464 and XM 508242.5 are grouped closely together, indicating that these two species (likely human and chimpanzee or another primate) share a recent common ancestor and have highly similar HBB sequences. This suggests strong evolutionary conservation of the gene in these species. In contrast, EF605506.1 forms the most distant branch, suggesting it is the least closely related among the group—possibly representing a non-primate mammal like a rodent. The tree topology aligns with expected evolutionary relationships, where primates are more closely related to each other, and other mammals like rodents or distant species show greater divergence in their HBB gene sequences.