



NATIONAL TECHNICAL UNIVERSITY OF ATHENS  
SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING  
M.Sc. IN TRANSLATIONAL ENGINEERING IN HEALTH AND MEDICINE

## Translational bioinformatics

*bio1100*

---

## REPORT

**EVANGELOS STAMOS**

**Supervisors:** George Matsopoulos, Ioannis Makris  
Professor NTUA | Teaching and Research Associate

Athens, February 2024

---





NATIONAL TECHNICAL UNIVERSITY OF ATHENS  
SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING  
M.Sc. IN TRANSLATIONAL ENGINEERING IN HEALTH AND MEDICINE

## Translational bioinformatics

*bio1100*

---

REPORT

**EVANGELOS STAMOS**

**Supervisors:** George Matsopoulos, Ioannis Makris  
Professor NTUA | Teaching and Research Associate

Approved by Prof. George Matsopoulos and Dr. Ioannis Makris in 14th February 2024.

(Signature)

.....  
George Matsopoulos, Ioannis Makris  
Professor NTUA | Teaching and Research Associate

Athens, February 2024





NATIONAL TECHNICAL UNIVERSITY OF ATHENS  
SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING  
M.Sc. IN TRANSLATIONAL ENGINEERING IN HEALTH AND MEDICINE

Copyright © - All rights reserved.

Evangelos Stamos, 2024.

The copying, storage and distribution of this report, exall or part of it, is prohibited for commercial purposes. Reprinting, storage and distribution for non - profit, educational or of a research nature is allowed, provided that the source is indicated and that this message is retained.

The content of this assignment does not necessarily reflect the views of the Department, the Supervisor, or the committee that approved it.

**DISCLAIMER ON ACADEMIC ETHICS AND INTELLECTUAL PROPERTY RIGHTS**

Being fully aware of the implications of copyright laws, I expressly state that this report, as well as the electronic files and source codes developed or modified in the course of this assingment, are solely the product of my personal work and do not infringe any rights of intellectual property, personality and personal data of third parties, do not contain work / contributions of third parties for which the permission of the authors / beneficiaries is required and are not a product of partial or complete plagiarism, while the sources used are limited to the bibliographic references only and meet the rules of scientific citing. The points where I have used ideas, text, files and / or sources of other authors are clearly mentioned in the text with the appropriate citation and the relevant complete reference is included in the bibliographic references section. I fully, individually and personally undertake all legal and administrative consequences that may arise in the event that it is proven, in the course of time, that this assignment or part of it does not belong to me because it is a product of plagiarism.

(Signature)

.....

Evangelos Stamos

14th February 2024



# **Abstract**

---

In this assignment, we delve into the realm of translational bioinformatics with a focus on understanding genetic variations and their implications. Through practical exercises, we will employ bioinformatics tools to analyze genomes, identify mutations, and design primers for PCR amplification. Additionally, the assignment involves using BLAST for sequence alignment to elucidate genetic information relevant to health and disease. This project aims to equip us with the skills to bridge the gap between computational biology and clinical applications, highlighting the significance of bioinformatics in advancing personalized medicine.

## **Keywords**

Bioinformatics, Translational, BLAST, Genomes, Mutations, Primers, PCR, Phylogenetic Trees



# Table of Contents

---

<b>Abstract</b>	<b>1</b>
<b>Preface</b>	<b>9</b>
<b>1 Exercise 1</b>	<b>11</b>
<b>2 Exercise 2</b>	<b>21</b>
<b>3 Exercise 3</b>	<b>27</b>
3.1 Sequence . . . . .	30
3.2 Structure . . . . .	30
3.3 Post-Translational Modification (PTM) . . . . .	30
3.4 The 4 amino acids that change and cause beta-thalassemia . . . . .	33
<b>4 Exercise 4</b>	<b>35</b>
4.1 NCBI Blast Neighbor Joining Phylogenetic Tree . . . . .	36
4.2 MEGA Neighbor Joining Phylogenetic Tree . . . . .	38
4.3 Differences between the two constructed phylogenetic trees . . . . .	38
<b>List of Abbreviations</b>	<b>43</b>



# List of Figures

---

1.1	blastn Sequence Result . . . . .	12
1.2	Webcutter 2.0 Settings . . . . .	15
1.3	Webcutter 2.0 Settings for restriction enzymes that cut the sequence in 4 positions . . . . .	16
2.1	NCBI Primer-Blast - Analysis settings . . . . .	23
2.2	NCBI Primer-Blast - PCR template highly similar to Homo sapiens tyrosinase (TYR) mRNA . . . . .	23
2.3	NCBI Primer-Blast - Detailed primer reports I . . . . .	24
2.4	NCBI Primer-Blast - Detailed primer reports II . . . . .	24
2.5	NCBI Primer-Blast - Detailed primer reports II . . . . .	24
3.1	NCBI Gene Information - HEXA . . . . .	28
3.2	NCBI Gene Information - HFE . . . . .	28
3.3	NCBI Gene Information - PKU . . . . .	28
3.4	NCBI Gene Information - HBB . . . . .	29
3.5	NCBI Gene Information - LDLR . . . . .	29
3.6	UniProt - HBB Sequence . . . . .	30
3.7	UniProt - HBB Structure I . . . . .	32
3.8	UniProt - HBB Structure II . . . . .	32
3.9	UniProt - HBB PTM . . . . .	32
3.10	UniProt - HBB The 4 amino acids that change and cause beta-thalassemia . . . . .	33
4.1	HBB Homo sapiens blastp . . . . .	36
4.2	blastp sequences selected I . . . . .	36
4.3	blastp sequences selected II . . . . .	37
4.4	blastp sequences selected III . . . . .	37
4.5	NCBI Neighbor Joining Phylogenetic Tree . . . . .	37
4.6	Obtained sequences files in .fasta format . . . . .	39
4.7	All species sequences in .fasta format . . . . .	39
4.8	All species aligned sequences . . . . .	39
4.9	Neighbor Joining Phylogenetic Tree settings . . . . .	39
4.10	Neighbor Joining Phylogenetic Tree . . . . .	40
4.11	Neighbor Joining Phylogenetic Tree clean . . . . .	40



## List of Tables

---

1	Exercises implementation and short answers . . . . .	9
1.1	Endonucleases that do not cut the sequence . . . . .	16
1.2	Table of restriction enzymes 1 of 3 . . . . .	17
1.3	Table of restriction enzymes 2 of 3 . . . . .	18
1.4	Table of restriction enzymes 3 of 3 . . . . .	19
1.5	Table of restriction enzymes that cut sequence in 4 positions . . . . .	20
2.1	Primer pair and product size . . . . .	25
3.1	Genes mutation and diseases causality . . . . .	30
3.2	Neutron, NMR, EM available HBB structures on UniProt . . . . .	31
3.3	Structural Features of HBB . . . . .	31



**Table 1.** Exercises implementation and short answers

Question	Implementation	Answer(s)
1	LaTeX, Webcutter 2.0, Python	596 bp: BglII   4 cuts: BstX2I, BstYI, MfI and XholI
2	LaTeX, NCBI	Primers in 2.1
3	LaTeX, NCBI, UniProt	HBB: $\beta$ thalassemia, Sequence 3.6, Structure 3.7, 3.8 PTM 3.9 Glutamic acid (E), lysine (K), leucine (L) and proline (P) 3.10
4	LaTeX, NCBI, MEGA11	Neighbor Joining Phylogenetic Tree 4.5, 4.11

## Preface

---

This report was written during the Fall Semester of the academic year 2023 - 2024, in the context of project of the course bio1100 Translational bioinformatics.

In the following table all my personal information is included. For any question I am available on given email.

Category	Personal Data
Name	Evangelos Stamos
Registry Number	<b>03500050</b>
Student Type	Postgraduate student
Program Enrolled	<b>MSc Translational Engineering in Health and Medicine</b>
Email	stamosevangelos AT mail DOT ntua DOT gr

Report was implemented in LaTeX (pdfLaTeX) based on personal template specifically designed for course context. The table 1 gives detailed information about the implementation of each problem.



## Chapter 1

### Exercise 1

---

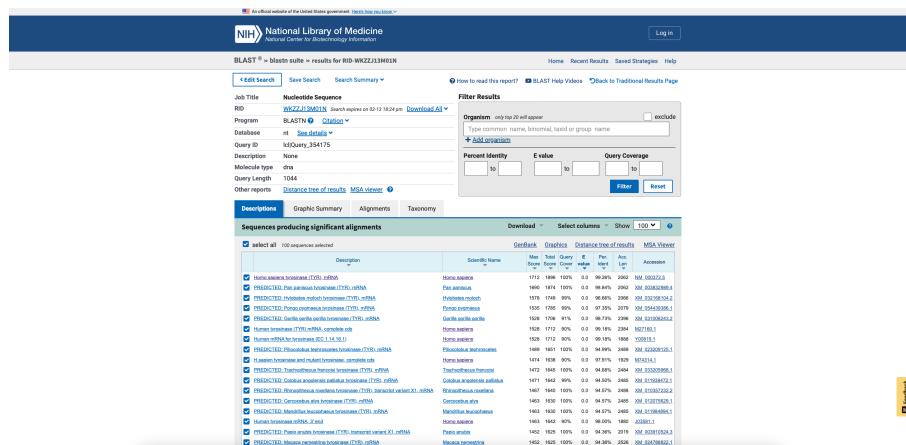
Given the following sequence:

CGAGGGACCTTACGGCGTAATCCTGGAAACCATGACAAATCCAGAACCCAAGGCTCCCTTTCA  
GCTGATGTAGAATTTCCTGAGTTGACCCATGGAAGGATTGCTAGTCCACTACTGGGATAGCGGATG  
CCTCTCAAAGAGCATGCACAATGCCCTGCACATCTATGAATGGAACAATGTCCCAGGCAGGGATCTGCC  
AACGATCCTATCTCCTCTTCACCATGCATTGTTGACAGTATTTGAGCAGTGGCTCCGAAGGCACCG  
TCCTCTCAAGAAGTTATCCAGAAGCCAATGCACCCATTGGACATAACCGGAAATCCTACATGGTCTTA  
TACCACTGTACAGAAATGGTATTCTTATTCATCCAAAGATCTGGCTATGACTATAGCTATCTACAA  
GATTCAAGACCCAGACTCTTCAAGACTACATTAAGTCCTATTGGAACAAGCGAGTCGGATCTGGTCATG  
GCTCCTGGGCGGGATGGTAGGGGCCCTCACTGCCCTGCTGGCGGGCTGTGAGCTGCTGTGCT  
TCACAAGAGAAAGCAGCTCCTGAAGAAAAGCAGCCACTCCTCATGGAGAAAGAGGATTACCACAGCTTGT  
ATCAGAGCCATTATAAAAGGCTTAGGCAATAGAGTAGGGCAAAAAGCCTGACCTCACTCTAACTCAAAG  
TAATGTCCAGGTTCCCAGAGAATATCTGCTGGTATTCTGTAAAGACCATTGCAAAATTGTAACCTAA  
TACAAAGTGTAGCCTCTCCAACTCAGGTAGAACACACCTGTCTTGCTGTCTTCACTCAGCCCT  
TTAACATTCTCCCTAACGCCATATGTCTAACGGAAAGGATGCTATTGGAATGAGGAACGTGTTATTGT  
ATGTGAATTAAAGTGTCTTATTAAAAAATTGAAATAATTGATTTGCCTCTGATTATTAAAGA  
TCTATATATGTTTATTGGCCCCCTCTTATTAAATAAAACAGTGAGAAATCT

With the help of the program Webcutter 2.0 use the appropriate restriction endonuclease to cut the sequence in such a way that you get a product 596 bp. Also, find which restriction enzymes cut the above sequence in 4 positions.

First we need to identify if given sequence is linear or circular. We run the relevant blastn 1.1 and we observe that given sequence has 100% query coverage with **Homo sapiens tyrosinase (TYR), mRNA**. This sequence, identified as Homo sapiens tyrosinase (TYR) mRNA, is **linear**. They are transcribed from DNA in the nucleus and then translated into proteins in the cytoplasm. Tyrosinase is an enzyme critical in the melanin biosynthesis pathway, involved in determining the color of the skin, hair, and eyes. In the context of molecular biology, mRNA molecules are synthesized from DNA templates during transcription. After transcription, mRNA undergoes processing events, including splicing, 5' capping, and 3' polyadenylation, resulting in a mature linear mRNA molecule that exits the nucleus to be translated into protein by ribosomes in the cytoplasm. Therefore, the Homo sapiens tyrosinase (TYR) mRNA, like all other eukaryotic mRNAs, is linear.

To perform analysis, we have to paste the sequence into the provided box. We then select All enzymes in the database on "Please indicate which enzymes to include in the analysis" setting and click the Analyze sequence button as depicted in figure 1.2. We perform a linear sequence analysis. Searching for restriction enzymes that cut the sequence into 4 positions can be performed more effectively if we select in "Please indicate which enzymes to include in the display" setting Enzymes cut exactly 4 times, so we do not have to search along all results depicted in figure 1.3.



**Figure 1.1.** *blastn* Sequence Result

According to performed analysis 1.1, 1.2, 1.3, 1.4 , restriction endonuclease **BgIII** cuts the sequence in such a way resulting in a product of 596 bp. A python script was written based on analysis performed on Heimanlab Webcutter 2.0 tool to identify which restriction endonuclease results in the desired product size base pairs.

---

```

9      TCTTCACCATGCATTGTTGACAGTATTTGAGCAGTGGCTCGAAGGCACCGTCTTCAGAAGTTATCC
10     AGAAGCCAATGCACCCATTGGACATAACCGGAATCCTACATGGTTCTTACCAACTGTACAGAAATGGTATT
11     CTTTATTCATCCAAAGATCTGGCTATGACTATAGCTATACAAGATTCAAGACCCAGACTCTTCAAGACTA
12     CATTAAGTCCTATTTGAAACAAGCGAGTCGGATCTGGTATGGCTCTTGGGCGGCATGGTAGGGGCCGTCT
13     CACTGCCCTGCTGGCGGGCTTGTGAGCTGCTGTGTCACAAGAGAAAAGCAGCTCTGAAGAAAAGCAGCCA
14     CTCCCTCATGGAGAAAGAGGATTACCACAGCTTGATCAGAGCATTATAAAAGGCTTAGGCAATAGAGTAGGGC
15     CAAAAAGCCTGACCTCACTCAAAGTAATGTCAGGTTCCAGAGAATATCTGCTGGTATTTCTGAA
16     AGACCATTTGCAAATTTGTAACCTAACAAAGTGTAGCCTCTTCAACTCAGGTAGAACACACCTGCTTTGT
17     CTTGCTGTTTCACTCAGCCCTTTAACATTTCCCCTAAGCCCATATGTCTAAGGAAAGGATGCTATTGGTAA
18     TGAGGAACGTATTGTATGTGAAATTAAAGTGTCTTATTTAAAAAATTGAAATAATTGATTTGCTTCTC
19     TGATTATTAAAGATCTATATGTTTATTGGCCCTCTTATTTAAATAAAACAGTGAGAAATCT"
20
21
22 # Calculate the length of the DNA sequence
23 sequence_length = len(sequence)
24
25 # Print the length
26 print("Length of given sequence:", sequence_length, "bp")
27
28 # Based on analysis performed on the sequence, the following endonucleases and their cut positions are
29 # provided.
30 # Analysis was performed using the Heimanlab Webcutter 2.0 tool (http://heimanlab.com/cut2.html) with the
31 # following settings:
32 # Type of analysis: Linear sequence analysis
33 # Restriction sites displayed: Map of restriction sites
34 # Table of sites, sorted alphabetically by enzyme name
35 # Displayed enzymes : All enzymes
36 # Enzymes to include in the analysis: All enzymes in the database
37
38 # Given the provided list of endonucleases and their cut positions, identify which one produces a
39 # fragment of 596 bp.
40
41 # Defining the endonucleases with their cut positions
42 endonucleases = {
43     "AccB1I": [274], "AccB7I": [123], "AciI": [135, 507, 543], "AclWI": [205, 217, 485],
44     "AcsI": [77], "AfaI": [360], "AluI": [68, 412, 552, 580, 630], "Alw21I": [936],
45     "AlwI": [205, 217, 485], "ApoI": [77], "AspHI": [936], "AspS9I": [5, 517, 673, 1008],
46     "AsuI": [5, 517, 673, 1008], "AvaII": [5], "BanI": [274], "BbvI": [154],
47     "Bbv12I": [936], "BbvI": [581, 599], "BcnI": [330], "Bfai": [112], "BglI": [537],
48     "BglII": [392, 988], "Bme18I": [5], "BmyI": [936], "BsaJI": [50, 97, 192, 497],
49     "Bsc4I": [101, 122, 197, 502, 538], "BseII": [127], "BseDI": [50, 97, 192, 497],
50     "BseNI": [127], "BseRI": [607], "BshNI": [274], "BsiHKAI": [936], "BsiSI": [329],
51     "BsiYI": [102, 123, 198, 503, 539], "BsI": [102, 123, 198, 503, 539], "BsmFI": [8, 194],
52     "BsoFI": [505, 578, 596], "Bsp1286I": [936], "Bsp1407I": [358], "Bsp143I": [201, 212, 392, 481, 988],
53     "Bsp19I": [97], "BsrGI": [358], "BsrI": [127], "BsrSI": [127], "BssT1I": [50, 97, 497],
54     "Bst2UI": [24, 194, 714], "Bst71I": [581, 599], "BstDEI": [86, 657, 801, 840, 863, 877],
55     "BstDSI": [97], "BstF5I": [138, 389, 890], "BstNI": [24, 194, 714], "BstOI": [24, 194, 714],
56     "BstSFI": [407], "BstX2I": [201, 392, 481, 988], "BstXI": [395], "BstYI": [201, 392, 481, 988],
     "BsRI": [519, 675, 1009], "Cac8I": [152, 538, 554], "Cfr13I": [5, 517, 673, 1008],
     "Csp6I": [359], "CviJI": [55, 68, 266, 306, 400, 412, 494, 519, 544, 552, 580, 598, 630, 642, 656,
       675, 683, 789, 844, 867, 1009],
     "DdeI": [86, 657, 801, 840, 863, 877], "DpnI": [203, 214, 394, 483, 990], "DpnII": [201, 212, 392,
```

```

481, 988],
57 "DraI": [944, 985], "DraII": [5], "DsaI": [97], "Eam1104I": [66, 288], "EarI": [66, 288],
58 "Eco130I": [50, 97, 497], "Eco47I": [5], "Eco57I": [68, 590], "Eco64I": [274], "Eco0109I": [5],
59 "EcoRII": [22, 192, 712], "EcoT14I": [50, 97, 497], "EcoT22I": [239], "ErhI": [50, 97, 497],
60 "Esp1396I": [123], "FauI": [544], "FauNDI": [871],
61 "FokI": [138, 389, 890], "Fsp4HI": [505, 578, 596],
62 "HaeIII": [519, 675, 1009], "HapII": [329], "HgiEI": [5],
63 "HincII": [245], "HindII": [245], "HinfI": [333, 423, 435, 476],
64 "HpaII": [329], "HphI": [234, 373], "Hsp92II": [35, 101, 154, 237, 344, 493, 610],
65 "ItaI": [505, 578, 596], "Ksp632I": [66, 288], "Kzo9I": [201, 212, 392, 481, 988],
66 "MaeI": [112], "MaeIII": [563, 768], "MboI": [201, 212, 392, 481, 988],
67 "MboII": [66, 224, 231, 288, 591, 797], "Mfli": [201, 392, 481, 988],
68 "MnlI": [5, 63, 142, 285, 527, 607, 620, 692, 906], "Mph1103I": [239],
69 "MseI": [454, 850, 927, 943, 984, 1023], "MsI": [173], "MspAII": [68],
70 "MspI": [329], "MspR9I": [24, 194, 330, 714], "MvaI": [24, 194, 714],
71 "MwoI": [272, 537], "NciI": [330], "NcoI": [97], "NdeI": [871],
72 "NdeII": [201, 212, 392, 481, 988], "NlaIII": [35, 101, 154, 237, 344, 493, 610],
73 "NlaIV": [6, 56, 267, 276, 495, 518, 718, 1010], "NsI": [239], "NspBII": [68],
74 "NspI": [154], "PaeI": [154], "PalI": [519, 675, 1009], "PflMI": [123],
75 "PleI": [439, 480], "Ppu10I": [235], "PpuMI": [5], "Psp5II": [5],
76 "PspN4I": [6, 56, 267, 276, 495, 518, 718, 1010], "PvuII": [68], "RsaI": [360],
77 "Sau3AI": [201, 212, 392, 481, 988], "Sau96I": [5, 517, 673, 1008],
78 "ScrFI": [24, 194, 330, 714], "SduI": [936], "SfaNI": [139, 891], "SfcI": [407],
79 "SinI": [5], "SphI": [154], "Sse9I": [77, 764, 924, 948, 957], "SspBI": [358],
80 "StyI": [50, 97, 497], "TfiI": [333, 423], "TruII": [454, 850, 927, 943, 984, 1023],
81 "Tru9I": [454, 850, 927, 943, 984, 1023], "Tsp45I": [563], "Tsp509I": [77, 764, 924, 948, 957],
82 "TspEI": [77, 764, 924, 948, 957], "TspRI": [265, 359, 531, 1036], "XhoII": [201, 392, 481, 988],
83 "Zsp2I": [239]}

84
85

86 # Function to calculate fragment sizes
87 def calculate_fragment_sizes(cut_positions, sequence_length):
88     # Add the start and end of the sequence to the list of cut positions
89     cut_positions = [0] + sorted(cut_positions) + [sequence_length]
90     # Calculate the sizes of the resulting fragments
91     fragment_sizes = [cut_positions[i+1] - cut_positions[i] for i in range(len(cut_positions)-1)]
92     return fragment_sizes
93

94 # Check each endonuclease for a 596 bp fragment
95 enzyme_with_596bp = {}

96

97 for enzyme, positions in endonucleases.items():
98     sizes = calculate_fragment_sizes(positions, sequence_length)
99     if 596 in sizes:
100         enzyme_with_596bp[enzyme] = sizes
101

102 if not enzyme_with_596bp: # Zero endonucleases found that produce a 596 bp fragment
103     print('No endonucleases found that produce a 596 bp fragment.')
104 else: # Print the endonuclease(s) that produce a 596 bp fragment
105     print('Endonuclease(s) that cut the sequence resulting to a product of 596 base pairs:',
106         enzyme_with_596bp)

```

```

107 # Produced result
108 # 'Length of given sequence: 1044 bp'
109 #'Endonuclease(s) that cut the sequence resulting to a product of 596 base pairs: {'BglII': [392, 596,
      56]}

```

Listing 1.1: Exercise 1 - Python Code

Restriction enzymes **AspS9I**, **AsuI**, **BsaJI**, **BseDI**, **BstX2I**, **BstYI**, **Cfr13I**, **HinfI**, **MfI**, **MspR9I**, **Sau96I**, **ScrFI**, **TspRI** and **XhoII** cut the given sequence in 4 positions, with further information available in table 1.5.

The screenshot shows the Webcutter 2.0 user interface. At the top right is the logo "web cutter". Below it is a text area for entering a sequence, with instructions: "Please enter a title for this sequence:" and "Exercise 1 Sequence". A large text box contains the DNA sequence: CGAAGGACCTTACGGCCTAAATCCCTGAAAACCGATGAAACATTCAAGACCGAACGTCCTTGACGTTGAAATTTGGCTTA... (the sequence continues for several lines). Below the sequence entry is a section titled "Please select the type of analysis you would like": 

- Linear sequence analysis
- Circular sequence analysis
- Find sites which may be introduced by silent mutagenesis

Next is a section titled "Please indicate how you would like the restriction sites displayed":

- Map of restriction sites
- Table of sites, sorted alphabetically by enzyme name
- Table of sites, sorted sequentially by base pair number

Then is a section titled "Please indicate which enzymes to include in the display":

- All enzymes
- Enzymes not cutting
- Enzymes cutting once
- Enzymes cutting exactly  times
- Enzymes cutting at least  times, and at most  times
- Rainbow  highlights for enzymes from the Standard  polylinker

Finally, a section titled "Please indicate which enzymes to include in the analysis":

- All enzymes in the database
- Only enzymes with recognition sites equal to or greater than  bases long
- Only the following enzymes:  AspI  AccI  AccI13I  AccI16I  AccI95I

Figure 1.2. Webcutter 2.0 Settings

**Figure 1.3.** Webcutter 2.0 Settings for restriction enzymes that cut the sequence in 4 positions

**Table 1.1.** Endonucleases that do not cut the sequence

<b>Endonucleases</b>	<b>Endonucleases (continue)</b>
AatI, AatII, Acc113I, Acc16I, Acc65I, AccBSI, AccI, AccII, AccIII, AclNI, AcyI, AfeI, AflII, AflIII, AgeI, AhdI, Alw26I, Alw44I, AlwNI, Ama87I, AocI, Aor51HI, ApaI, ApaLI, Ascl, AseI, AsnI, Asp700I, Asp718I, AspEI, AspI, AspLEI, AtsI, AvaI, AvII, AvrII, BbaI, BamHI, BanII, BanIII, BbeI, BbiII, BbrPI, BbsI, Bbv16II, BcgI, BclI, BcoI, BfrI, BlnI, BlpI, BpiI, BpmI, Bpu1102I, Bpu14I, BpuAI, Bsa29I, BsaAI, BsaBI, BsaHI, BsaI, BsaMI, BsaOI, BsaWI, BscI, Bse118I, Bse21I, Bse8I, BseAI, BseCI, BsePI, BsgI, Bsh1236I, Bsh1285I, Bsh1365I, BsiEI, BsiI, BsiMI, BsiWI, BsmAI, BsmBI, BsmI, BsoBI, Bsp106I, Bsp119I, Bsp120I, Bsp13I, Bsp143II, Bsp1720I, Bsp68I, BspCI, BspDI, BspEI, BspHI, BspLU11I, BspMI, BspTI, BspXI, BsrBI, BsrBRI, BsrDI, BsrFI, BssAI, BssHII, BssSI, Bst1107I, Bst98I, BstBI, BstD102I, BstEII, BstH2I, BstI, BstMCI, BstPI, BstSNI, BstUI, BstZI, Bsu15I, Bsu36I, CciNI, CellII, CfoI, Cfr10I, Cfr42I, Cfr9I, CfrI, ClaI, CpoI, Csp45I, CspI, CvnI, DralII, DrdI	EaeI, EagI, Eam1105I, Ecl136II, EclHKI, EclXI, Eco105I, Eco147I, Eco24I, Eco255I, Eco31I, Eco32I, Eco47III, Eco52I, Eco72I, Eco81I, Eco88I, Eco91I, EcoICRI, EcoNI, EcoO65I, EcoRI, EcoRV, EheI, Esp3I, FbaI, FriOI, FseI, FspI, GsuI, HaeII, HgaI, Hhal, Hin1I, Hin6I, HinP1I, HindIII, HpaI, Hsp92I, HspAI, KasI, Kpn2I, KpnI, Ksp22I, KspI, LspI, MaeII, Mami, MfeI, MluI, MluNI, MroI, MroNI, MscI, Msp17I, MspCI, MunI, Mva1269I, MvnI, NaeI, NarI, NgoAIV, NgoMI, NheI, NotI, NruI, NspV, PacI, PaeR7I, Pfl23II, PinAI, Ple19I, PmaCI, Pme55I, PmeI, PmlI, PshAI, PshBI, Psp124BI, Psp1406I, PspAI, PspALI, PspEI, PspLI, PspOMI, PstI, PstNHI, PvUI, RcaI, RsrII, SacI, SacII, Sall, Sapi, SbfI, ScaI, SexAI, SfiI, Sfr274I, Sfr303I, SfuI, Sgfl, SgrAI, SmaI, Smil, SnabI, SpeI, SpII, SrfI, Sse8387I, SseBI, SspI, SstI, SstII, StuI, SunI, Swal, TaqI, ThaI, Tth111I, TthHB8I, Vha464I, VneI, VspI, XbaI, XcmI, XhoI, XmaI, XmaIII, Xmni

**Table 1.2.** Table of restriction enzymes 1 of 3

Enzyme name	No. of cuts	Positions of sites	Recognition sequence
AccB1I	1	274	g/gyrcc
AccB7I	1	123	ccannnn/ntgg
AciI	3	135 507 543	ccgc
AclWI	3	205 217 485	ggatc
AcsI	1	77	r/aatty
AfaI	1	360	gt/ac
AluI	5	68 412 552 580 630	ag/ct
Alw21I	1	936	gwgcw/c
AlwI	3	205 217 485	ggate
ApoI	1	77	r/aatty
AspHI	1	936	gwgcw/c
AspS9I	4	5 517 673 1008	g/gncc
AsuI	4	5 517 673 1008	g/gncc
AvaII	1	5	g/gwcc
BanI	1	274	g/gyrcc
BbuI	1	154	gcatg/c
Bbv12I	1	936	gwgcw/c
BbvI	2	581 599	gcagc
BcnI	1	330	cc/nggg
BfaI	1	112	c/tag
BglII	1	537	gccnnnn/nggg
<b>BglIII</b>	<b>2</b>	<b>392 988</b>	<b>a/gatct</b>
Bme18I	1	5	g/gwcc
BmyI	1	936	gdgch/c
BsaJI	4	50 97 192 497	c/cnngg
Bsc4I	5	101 122 197 502 538	ccnnnn/nnnngg
BseII	1	127	actgg
BseDI	4	50 97 192 497	c/cnngg
BseNI	1	127	actgg
BseRI	1	607	gaggag
BshNI	1	274	g/gyrcc
BsiHKAI	1	936	gwgcw/c
BsiSI	1	329	c/cgg
BsiYI	5	102 123 198 503 539	ccnnnnnn/nnggg
BsII	5	102 123 198 503 539	ccnnnnnn/nnggg
BsmFI	2	8 194	gggac
BsoFI	3	505 578 596	gc/ngc
Bsp1286I	1	936	gdgch/c
Bsp1407I	1	358	t/gtaca
Bsp143I	5	201 212 392 481 988	/gatc
Bsp19I	1	97	c/catgg
BsrGI	1	358	t/gtaca
BsrI	1	127	actgg
BsrSI	1	127	actgg
BssT1I	3	50 97 497	c/cwwgg
Bst2UI	3	24 194 714	cc/wgg
Bst71I	2	581 599	gcagc
BstDEI	6	86 657 801 840 863 877	c/tnag
BstDSI	1	97	c/crygg
BstF5I	3	138 389 890	ggatg
BstNI	3	24 194 714	cc/wgg
BstOI	3	24 194 714	cc/wgg
BstSFI	1	407	c/tryag
BstX2I	4	201 392 481 988	r/gatcy
BstXI	1	395	ccannnnnn/ntgg
BstYI	4	201 392 481 988	r/gatcy

**Table 1.3.** Table of restriction enzymes 2 of 3

<b>Enzyme name</b>	<b>No. of cuts</b>	<b>Positions of sites</b>	<b>Recognition sequence</b>
BsuRI	3	519 675 1009	gg/cc
Cac8I	3	152 538 554	gcn/ngc
Cfr13I	4	5 517 673 1008	g/gncc
Csp6I	1	359	g/tac
CviJI	21	55 68 266 306 400 412 494 519 544 552 580 598 630 642 656 675 683 789 844 867 1009	rg/cy
DdeI	6	86 657 801 840 863 877	c/tnag
DpnI	5	203 214 394 483 990	ga/tc
DpnII	5	201 212 392 481 988	/gatc
DraI	2	944 985	ttt/aaa
DraII	1	5	rg/gnccy
DsaI	1	97	c/crygg
Eam1104I	2	66 288	ctcttc
EarI	2	66 288	ctcttc
Eco130I	3	50 97 497	c/cwwgg
Eco47I	1	5	g/gwcc
Eco57I	2	68 590	ctgaag
Eco64I	1	274	g/gyrcc
EcoO109I	1	5	rg/gnccy
EcoRII	3	22 192 712	/ccwgg
EcoT14I	3	50 97 497	c/cwwgg
EcoT22I	1	239	atgea/t
ErhI	3	50 97 497	c/cwwgg
Esp1396I	1	123	ccannnn/ntgg
FauI	1	544	cccg
FauNDI	1	871	ca/tatg
FokI	3	138 389 890	ggatg
Fsp4HI	3	505 578 596	gc/ngc
HaeIII	3	519 675 1009	gg/cc
HapII	1	329	c/cgg
HgiEI	1	5	g/gwcc
HincII	1	245	gtt/rac
HindII	1	245	gtt/rac
HinfI	4	333 423 435 476	g/antc
HpaII	1	329	c/cgg
HphI	2	234 373	ggta
Hsp92II	7	35 101 154 237 344 493 610	catg/
ItaI	3	505 578 596	gc/ngc
Ksp632I	2	66 288	ctcttc
Kzo9I	5	201 212 392 481 988	/gatc
MaeI	1	112	c/tag
MaeIII	2	563 768	/gtac
MboI	5	201 212 392 481 988	/gatc
MboII	6	66 224 231 288 591 797	gaaga
MfII	4	201 392 481 988	r/gatcy
MnlI	9	5 63 142 285 527 607 620 692 906	cctc
Mph1103I	1	239	atgea/t
MseI	6	454 850 927 943 984 1023	t/taa
MslI	1	173	caynn/nmrng
MspA1I	1	68	cmg/ckg
MspI	1	329	c/cgg
MspR9I	4	24 194 330 714	cc/ngg
MvaI	3	24 194 714	cc/wgg
MwoI	2	272 537	gcnnnnn/nngc

**Table 1.4.** Table of restriction enzymes 3 of 3

Enzyme name	No. of cuts	Positions of sites	Recognition sequence
NciI	1	330	cc/sgg
NcoI	1	97	c/catgg
NdeI	1	871	ca/tatg
NdeII	5	201 212 392 481 988	/gatc
NlaIII	7	35 101 154 237 344 493 610	catg/
NlaIV	8	6 56 267 276 495 518 718 1010	ggn/ncc
NsiI	1	239	atgca/t
NspBII	1	68	cmg/ckg
NspI	1	154	rcatg/y
PaeI	1	154	gcatg/c
PaiI	3	519 675 1009	gg/cc
PflMI	1	123	ccannnn/ntgg
PleI	2	439 480	gagtc
Ppu10I	1	235	a/tgcat
PpuMI	1	5	rg/gwccy
Psp5II	1	5	rg/gwccy
PspN4I	8	6 56 267 276 495 518 718 1010	ggn/ncc
PvuII	1	68	cag/ctg
RsaI	1	360	gt/ac
Sau3AI	5	201 212 392 481 988	/gatc
Sau96I	4	5 517 673 1008	g/gncc
ScrFI	4	24 194 330 714	cc/ngg
SduI	1	936	gdgch/c
SfaNI	2	139 891	gcatc
SfcI	1	407	c/tryag
SinI	1	5	g/gwcc
SphI	1	154	gcatg/c
Sse9I	5	77 764 924 948 957	/aatt
SspBI	1	358	t/gtaca
StyI	3	50 97 497	c/cwwgg
TfiI	2	333 423	g/awtc
Tru1I	6	454 850 927 943 984 1023	t/taa
Tru9I	6	454 850 927 943 984 1023	t/taa
Tsp45I	1	563	/gtsac
Tsp509I	5	77 764 924 948 957	/aatt
TspEI	5	77 764 924 948 957	/aatt
TspRI	4	265 359 531 1036	cagtg
Van91I	1	123	ccannnn/ntgg
XbaII	4	201 392 481 988	r/gatcy
Zsp2I	1	239	atgca/t

**Table 1.5.** Table of restriction enzymes that cut sequence in 4 positions

Enzyme name	No. of cuts	Position site	Recognition sequence
AspS9I	4	5 517 673 1008	g/gncc
AsuI	4	5 517 673 1008	g/gncc
BsaJI	4	50 97 192 497	c/cnngg
BseDI	4	50 97 192 497	c/cnngg
BstX2I	4	201 392 481 988	r/gatcy
BstYI	4	201 392 481 988	r/gatcy
Cfr13I	4	5 517 673 1008	g/gncc
HinfI	4	333 423 435 476	g/antc
MfII	4	201 392 481 988	r/gatcy
MspR9I	4	24 194 330 714	cc/ngg
Sau96I	4	5 517 673 1008	g/gncc
ScrFI	4	24 194 330 714	cc/ngg
TspRI	4	265 359 531 1036	cagtg
XhoII	4	201 392 481 988	r/gatcy

## **Chapter 2**

### **Exercise 2**

---

For the 596 bp product you got in Exercise 1, select primers using any program you want to carry out the polymerase chain reaction (PCR). Select the pair of primers that gives you the largest product.

In Exercise 1, we observed that using restrictive endonuclease **BglIII** results in a product of 596 base pairs. Given that BglII cuts the linear sequence at positions 392 and 988, and considering the length of our sequence is 1044 base pairs, the digestion with BglII produces the following fragments:

- Fragment 1: This fragment is from the start of the sequence to the first cut site at position 392, 392 base pairs long.
- Fragment 2: This fragment is between the two cut sites, from position 392 to 988, 596 base pairs long.
- Fragment 3: This fragment is from the second cut site at position 988 to the end of the sequence at position 1.044, 56 base pairs long.

We are interested in **Fragment 2**, which starts from position 392 of the given sequence to position 988.

An alternative way of finding the desired cut sequence is implemented in a Python script.

```
1 # Author: Stamos Evangelos
2 # Date: 2024-02-02
3 # Title: bio1100 | Translational bioinformatics | Assignment | Part 2
4 # Description: This script calculates the fragment of a sequence given cut positions.
5
6 sequence = "CGAGGGACCTTACGGCGTAATCCTGAAACCATGACAATCCAGAACCCCAGGCTCCCTTCAGCTGATGT
7 AGAATTTGCCTGAGTTGACCAGTGGAAAGGATTGCTAGTCCACTTACTGGGATCGGGATGCCCTCAAAGAG
8 CATGCACAATGCCCTGCAATCTATGAATGAAACAATGTCCCAGGCGAGGATCTGCCAACGATCCTATCTTCC
9 TCTTCACCATGCATTTGTTGACAGTATTTGAGCAGTGGCTCGAAGGCACCGTCTTCAAGAAGTTATCC
10 AGAAGCCAATGCACCCATTGGACATAACCGGAATCCTACATGGTTTACACTGTACAGAAATGGTGATT
11 CTTTATTCATCCAAGATCTGGCTATGACTATAGCTATCTACAAGATTGACCCAGACTCTTCAAGACTA
12 CATTAAAGTCCTATTGGAACAAGCGAGTCGGATCTGGTCATGGCTCCTGGGGCGCGATGGTAGGGGCCGTCC
13 CACTGCCCTGCTGGCGGCTTGTGAGCTGCTGTGTCGTACAAGAGAAAGCAGCTCCTGAAGAAAAGCAGCCA
14 CTCCCTATGGAGAAAGAGGATTACCACAGCTGTATCAGAGCCATTATAAAGGCTTAGGCAATAGAGTAGGGC
15 CAAAAAGCCTGACCTCACTCTAAAGTAATGTCCAGGTTCCAGAGAATATCTGCTGGTATTTCTGTAA
16 AGACCAATTGCAAAATTGTAACCTAATACAAAGTGTAGCCTTCCAECTCAGGTAGAACACACCTGTCTTGT
17 CTTGCTGTTTCACTCAGCCCTTTAACATTTCCCTAAGCCCATATGTCTAAGGAAAGGATGCTATTGGTAA
18 TGAGGAACTGTTATTGTATGTAATTAAAGTGTCTTATTTAAAAAATTGAAATAATTGATTTGATTTGCCCTC
19 TGATTATTTAAAGATCTATATGTTTATTGGCCCTTCTTATTTAAATAAAACAGTGAGAAATCT"
20
21 # Note 1st position is 0
22 fragment_596bp = sequence[391:987]
23
24 print("596 bp fragment:", fragment_596bp)
25
26 # Output
27 # 596 bp fragment:AGATCTGGCTATGACTATAGCTATCTACAAGATTGACCCAGACTCTTCAAGACTACATTAAGTCC
28 TATTGGAACAAGCGAGTCGGATCTGGTCATGGCTCCTGGGGCGCGATGGTAGGGCCGTCTACTGCCCTGCTGGCGGGCTTG
29 TGAGCTGCTGTGTCACAAGAGAAAGCAGCTCCTGAAGAAAAGCAGCCACTCCTCATGGAGAAAGAGGATTACACAGCTTGT
30 ATCAGAGCCATTATAAAGGCTAGGCAATAGAGTAGGGCAAAAAGCCTGACCTCACTCTAAAGTAATGTCCAGGTTCCC
31 AGAGAATATCTGCTGGTATTTCTGAAAGACCATTGCAAAATTGTAACCTAACAAAGTGTAGCCTTCCAECTCAGGTAG
32 AACACACCTGTCTTGCTGCTTCACTCAGCCCTTTAACATTTCCCTAAGCCCATATGTCTAAGGAAAGGATGCTATT
```

```

33 GGTAAATGAGGAACCTGTTATTTGATGTGAATTAAAGTGCTTATTTAAAAAATTGAAATAATTTGATTTGCCTCTGATTAT
34 TTAA

```

### Listing 2.1: Exercise 2 - Python Code

Since we have obtained our sequence, we may proceed to find primers for PCR and select the pair of primers that results in the biggest longer product, using NCBI Primer-BLAST. We set Forward prime range From 392 and Reverse primer range to 988 to input the initial sequence 2.1. Alternatively, we can input the cut sequence produced by our Python script and not select any range for primers. As we can see in detailed primer reports figures 2.3, 2.4, 2.5 and at table 2.1 **Primer pair 7 with Forward primer TAGCTATCTACAAGATTTCAGACCCA and Reverse primer AGGGCTGAGTGAAAACAGCA results in the longer product of 439 base pairs length.**

**Figure 2.1.** NCBI Primer-Blast - Analysis settings

**Figure 2.2.** NCBI Primer-Blast - PCR template highly similar to *Homo sapiens tyrosinase (TYR)* mRNA

**Figure 2.3.** NCBI Primer-Blast - Detailed primer reports I

**Figure 2.4.** NCBI Primer-Blast - Detailed primer reports II

**Figure 2.5.** NCBI Primer-Blast - Detailed primer reports II

**Table 2.1.** Primer pair and product size

Forward Primer	Reverse Primer	Product size
AGCTATCTACAAGATTCAAGACCCA	GACACAGCAAGCTCACAAGC	153
GCTATCTACAAGATTCAAGACCCAGA	CTTGTGACGACACAGCAAGC	160
AGCTATCTACAAGATTCAAGACCCAG	TCCGACTCGCTTGTCCAAA	73
GCTATCTACAAGATTCAAGACCCAG	CCAGATCCGACTCGCTTGT	77
GCTATCTACAAGATTCAAGACCCA	GATCCGACTCGCTTGTCCA	74
AGCTATCTACAAGATTCAAGACCC	AGATCCGACTCGCTTGTCC	76
<b>TAGCTATCTACAAGATTCAAGACCCA</b>	<b>AGGGCTGAGTGAAACAGCA</b>	<b>439</b>
TAGCTATCTACAAGATTCAAGACCC	CGACACAGCAAGCTCACAAG	155
CTATCTACAAGATTCAAGACCCAGA	GTGAGGTCAAGCTTTGGC	282
ATAGCTATCTACAAGATTCAAGACCC	GCCATGACCAGATCCGACTC	87



## Chapter **3**

### **Exercise 3**

---

The mutations of the genes HEXA, HFE, PKU, PKU, HBB, LDLR cause respectively 5 well-known genetic diseases in humans. Which gene mutation causes  $\beta$  thalassemia?

Find for this protein:

- Its sequence
- Its structure
- The post-translational modification
- The 4 amino acids that change and cause beta-thalassemia

Based on NCBI Gene search information obtained the following results are displayed in 3.1, 3.2, 3.3, 3.4, 3.5.

**HEXA hexosaminidase subunit alpha [ Homo sapiens (human) ]**

Gene ID: 3073, updated on 25-Jan-2024

**Summary**

**Official Symbol:** HEXA [provided by HGNC]  
**Official Full Name:** hexosaminidase subunit alpha [provided by HGNC]  
**Primary source:** HGNC/HGNC-4878  
**Secondary source:** EntrezGene/ENSG00000213614 MM-808889; AllianceGenome/HGNC-4878  
**Gene type:** protein coding  
**RefSeq status:** REVIEWED  
**Organism:** Homo sapiens  
**Lineage:** Eukaryote; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominoidea; Homo  
**Also known as:** Summary  
TSD  

This gene encodes a member of the glycosidase hydrolase 25 family of proteins. The encoded preprotein is proteolytically processed to generate the alpha subunit of the lysosomal enzyme beta-hexosaminidase. This enzyme, together with the cofactor GM2 activator protein, catalyzes the degradation of the ganglioside GM2, and other molecules containing terminal N-acetyl hexosamines. Mutations in this gene lead to an accumulation of GM2 ganglioside in neurons, the underlying cause of neurodegenerative disorders termed the GM2 gangliosidoses, including Tay-Sachs disease (GM2-gangliosidosis type I). Alternative splicing results in multiple transcript variants, at least one of which encodes a preprotein that is proteolytically processed. [provided by RefSeq, Jan 2016]

**Expression:** Tissue expression in placenta (RPKM 36.4), thyroid (RPKM 34.8) and 29 other tissues [See more]  
**Ortholog:** mouse, rat  
[Try the new Gene table](#)  
[Try the new Transcript table](#)

**Figure 3.1. NCBI Gene Information - HEXA**

**HFE homeostatic iron regulator [ Homo sapiens (human) ]**

Gene ID: 3077, updated on 31-Jan-2024

**Summary**

**Official Symbol:** HFE [provided by HGNC]  
**Official Full Name:** homeostatic iron-regulator [provided by HGNC]  
**Primary source:** HGNC/HGNC-4898  
**Secondary source:** EntrezGene/ENSG0000019704 MM-815009; AllianceGenome/HGNC-4898  
**Gene type:** protein coding  
**RefSeq status:** REVIEWED  
**Organism:** Homo sapiens  
**Lineage:** Eukaryote; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominoidea; Homo  
**Also known as:** Summary  
HFE; HFE1; HLA-H; MCD7; TFR2  

This gene encodes a membrane protein that is similar to MHC class I-type protein and associates with beta2-microglobulin (b2M). It is thought that this protein functions to regulate iron absorption by regulating the interaction of the transferrin receptor with transferrin. The iron storage disorder, hemochromatosis, is a recessive genetic disorder that results from defects in this gene. [provided by RefSeq, May 2022]

**Expression:** Tissue expression in liver (RPKM 5.0), gall bladder (RPKM 4.6) and 24 other tissues [See more]  
**Ortholog:** mouse, rat  
[Try the new Gene table](#)  
[Try the new Transcript table](#)

**Figure 3.2. NCBI Gene Information - HFE**

**PAH phenylalanine hydroxylase [ Homo sapiens (human) ]**

Gene ID: 5053, updated on 23-Nov-2023

**Summary**

**Official Symbol:** PAH [provided by HGNC]  
**Official Full Name:** phenylalanine hydroxylase [provided by HGNC]  
**Primary source:** HGNC/HGNC-8582  
**Secondary source:** EntrezGene/ENSG00000171798 MM-81249; AllianceGenome/HGNC-8582  
**Gene type:** protein coding  
**RefSeq status:** REVIEWED  
**Organism:** Homo sapiens  
**Lineage:** Eukaryote; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominoidea; Homo  
**Also known as:** Summary  
PH; PROK; PKU1  

This gene encodes a member of the tyrosine-dependent aromatic amino acid hydroxylase protein family. The encoded phenylalanine hydroxylase enzyme hydroxylates phenylalanine to tyrosine and is the rate-limiting step in phenylalanine catabolism. Deficiency of this enzyme activity results in the autosomal recessive disorder phenylketonuria. [provided by RefSeq, Aug 2017]

**Expression:** Tissue expression in liver (RPKM 237.1), kidney (RPKM 196.4) and 1 other tissue [See more]  
**Ortholog:** mouse, rat  
[Try the new Gene table](#)  
[Try the new Transcript table](#)

**Figure 3.3. NCBI Gene Information - PKU**

The mutations in the genes mentioned above are associated with the following diseases:

**HEXA:** Mutations in this gene cause Tay-Sachs disease, a rare inherited disorder that results in progressive destruction of nerve cells in the brain and spinal cord.

**HFE:** Mutations in the HFE gene cause Hemochromatosis, a condition of excess iron in the body that can lead to serious conditions such as diabetes, heart problems, and liver disease.

**PKU:** This is likely a reference to the PAH gene, which when mutated, causes Phenylketonuria (PKU). PKU is a genetic disorder that causes increased levels of phenylalanine (an

**Figure 3.4.** NCBI Gene Information - HBB

**Figure 3.5.** NCBI Gene Information - LDLR

amino acid) in the body.

**HBB: Mutations in the HBB gene cause  $\beta$  thalassemia.** The HBB gene provides instructions for making a protein called beta-globin, which is a component of hemoglobin. Hemoglobin is the protein in red blood cells that carries oxygen throughout the body. Beta thalassemia is a blood disorder that reduces the production of hemoglobin. Mutations in the HBB gene can lead to reduced ( $\beta^+$ ) or absent ( $\beta^0$ ) production of beta-globin. This reduction or absence disrupts the normal balance of alpha-globin to beta-globin, leading to the formation of abnormal hemoglobin molecules that can destroy red blood cells, causing anemia and other related health issues. The severity of beta thalassemia can vary depending on the specific mutations in the HBB gene and how much they affect beta-globin production. The condition is inherited in an autosomal recessive pattern, meaning that two copies of the mutated gene (one from each parent) are necessary for a child to be affected by the disorder.

**LDLR: Mutations in the LDLR gene cause Familial Hypercholesterolemia, a form of high cholesterol that can lead to serious health conditions like coronary artery disease.**

Having identified that mutation on gene HBB causes  $\beta$ -thalassemia we proceed further on analysis using UniProt.

**Table 3.1.** Genes mutation and diseases causality

Gene Mutation	Caused Disease
HEXA	Tay-Sachs Disease
HFE	Hemochromatosis
PKU	Phenylketonuria
<b>HBB</b>	<b><math>\beta</math> thalassemia</b>
LDLR	Hypercholesterolemia

## 3.1 Sequence

The sequence of HBB has length 147 and Mass (Da) 15,998, as depicted in 3.6 is the following:

The screenshot shows the UniProt entry page for HBB. The top navigation bar includes links for BLAST, Align, Peptide search, ID mapping, SPARQL, and UniProt. The main content area is titled 'Entry' and shows the UniProt ID P04813. It displays the protein name HBB, a complete sequence status, and a length of 147 amino acids with a mass of 15,998 Da. The sequence itself is shown in a grey box with several highlighted regions. Below the sequence, it says 'Computationally mapped potential isoform sequences'. A table lists three isoforms: HBB1 (111 aa), FBWSP1 (90 aa), and AGADIPYWKK (55 aa). The bottom section, titled 'Features', shows a sequence logo and a table of features with columns for ID, POSITION, and DESCRIPTION. One feature is highlighted in blue.

**Figure 3.6.** UniProt - HBB Sequence

MVHLTPEEKSAVTALWGKVNVDVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPK  
VKAHGKKVLGAFSDGLAHDNLKGTFATLSELHCDKLHVDPENFRLGNVLVCVLAHHFG  
KEFTPPVQAAYQKVVAGVANALAHKYH

## 3.2 Structure

The structure of HBB is displayed in figure 3.7 and 3.8. It should be noted that displayed HBB structure is sourced by Protein Data Bank (PDB), with identifier 1A00, X-ray method, resolution 2.00 Å, chain B/D and positions 1 - 147. Many more structures are available on Uni-Prot, most of them of X-ray method while there are also 2 Neutron, 2 Nuclear Magnetic Resonance (NMR), 8 Electron Microscopy (EM) method structures from PDB and 1 predicted structure from AlphaFold 3.2. Structural Features of HBB are included in table 3.3.

## 3.3 Post-Translational Modification (PTM)

Post-translational modification (PTM) is the following:

<b>Source</b>	<b>Identifier</b>	<b>Method</b>	<b>Chain</b>	<b>Positions</b>
PDB	2DXM	Neutron	B/D	2-147
PDB	3KMF	Neutron	C/G	2-147
PDB	2H35	NMR	B/D	2-147
PDB	2M6Z	NMR	B/D	2-147
PDB	5NI1	EM	B/D	2-147
PDB	6NBC	EM	B/D	2-144
PDB	6NBD	EM	B/D	2-144
PDB	7PCF	EM	B	2-147
PDB	7PCH	EM	B/D	2-147
PDB	7PCQ	EM	B/D	2-147
PDB	7VDE	EM	B/D	1-147
PDB	7XGY	EM	B/D	1-147
AlphaFold	AF-P68871-F1	Predicted		1-147

**Table 3.2.** Neutron, NMR, EM available HBB structures on UniProt

<b>Type</b>	<b>Positions</b>	<b>Sequence</b>
Helix	6-17	PEEKSAVTALWG
Helix	21-35	VDEVGGEALGRLVV
Helix	37-42	PWTQRF
Helix	44-46	ESF
Helix	52-57	PDAVMG
Helix	59-75	PKVKAHGKKVLGAFSDG
Turn	78-80	HLD
Helix	82-95	LKGTFATLSELHCD
Helix	102-119	ENFRLLGNVLVCVLAHHF
Helix	120-122	GKE
Helix	125-143	PPVQAAYQKVVAGVANALA
Helix	144-146	HKY

**Table 3.3.** Structural Features of HBB

Glucose reacts non-enzymatically with the N-terminus of the beta chain to form a stable ketoamine linkage. This takes place slowly and continuously throughout the 120-day life span of the red blood cell. The rate of glycation is increased in patients with diabetes mellitus.

S-nitrosylated; a nitric oxide group is first bound to  $Fe^{2+}$  and then transferred to Cys-94 to allow capture of  $O_2$ .

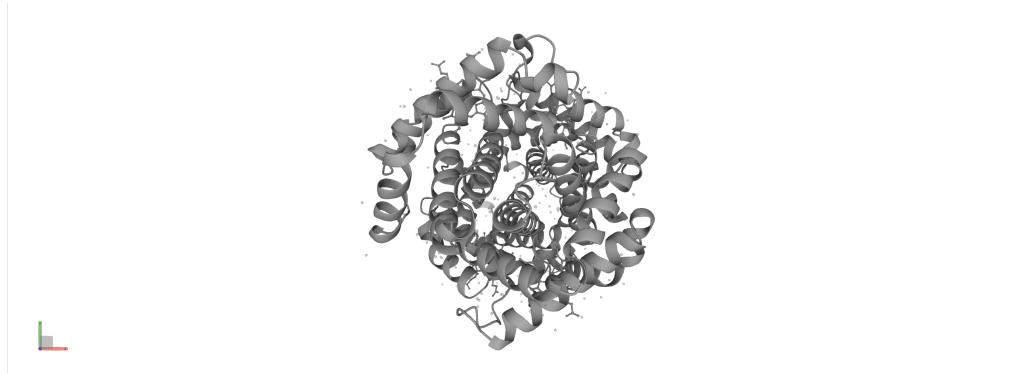
Acetylated on Lys-60, Lys-83 and Lys-145 upon aspirin exposure.

This screenshot shows the UniProt entry for HBB. The top navigation bar includes links for BLAST, Align, Peptide search, ID mapping, SPARQL, UniProtKB, and HBB. The main content area is titled 'Structure' and displays a 3D ribbon model of the HBB protein. Below the structure is a table listing PDB entries:

SOURCE	IDENTIFIER	METHOD	RESOLUTION	CHAIN	POSITIONS	LINKS
PDB	1A00	X-ray	2.00 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1A01	X-ray	1.80 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1A0U	X-ray	2.14 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1A0Z	X-ray	2.00 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1A3N	X-ray	1.80 Å	B/D	2-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1A3O	X-ray	1.80 Å	B/D	2-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1ARW	X-ray	2.00 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1ABY	X-ray	2.60 Å	B/D	3-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1AJP	X-ray	2.20 Å	B	2-147	PDBe RCSB-POB-POBj-PDBsum
PDB	1B86	X-ray	2.50 Å	B/D	2-147	PDBe RCSB-POB-POBj-PDBsum

Below the table, a section titled 'Features' shows 'Showing features for helix<sup>1</sup>, turn<sup>1</sup>'.

**Figure 3.7. UniProt - HBB Structure I**



**Figure 3.8. UniProt - HBB Structure II**

This screenshot shows the UniProt entry for HBB, specifically the PTM/Processing section. The top navigation bar includes links for BLAST, Align, Peptide search, ID mapping, SPARQL, UniProtKB, and HBB. The main content area is titled 'PTM/Processing' and displays a sequence alignment and a table of post-translational modifications:

TYPE	ID	POSITION(S)	SOURCE	DESCRIPTION
► Initiator methionine		1	Uniprot	Removed [By Similarity] [1 Publication]
► Modified residue		2	Uniprot	N-acetylvaline [By Similarity]
► Modified residue		2	Uniprot	N-pyruvate 2-imino-valine; in Hb A1b
► Glycosylation		2	Uniprot	N-linked (Glc) (glycation) valine; in Hb A1c [1 Publication]
► Chain	PRO_0000052976	2-147	Uniprot	Hemoglobin subunit beta
► Modified residue (large scale data)		5	PRIDE	Phosphotyrosine [Combined Sources]
► Glycosylation		9	Uniprot	N-linked (Glc) (glycation) lysine [1 Publication]
► Modified residue		10	Uniprot	Phosphoserine [Combined Sources]
► Modified residue (large scale data)		10	PRIDE	Phosphoserine [Combined Sources]

Below the table, a section titled 'Post-translational modification' provides details about the modifications, mentioning glucose reacts non-enzymatically with the N-terminus of the beta chain to form a stable ketoamine linkage. It also notes that glycation is increased in patients with diabetes mellitus. Another section discusses S-nitroylation, mentioning it is a nitrile oxide group first bound to Fe<sup>2+</sup> and then transferred to Cys-94 to allow capture of O<sub>2</sub>.

**Figure 3.9. UniProt - HBB PTM**

### 3.4 The 4 amino acids that change and cause beta-thalassemia

The four (4) amino acids that change and cause beta-thalassemia are **glutamic acid (E)**, **lysine (K)**, **leucine (L)** and **proline (P)**, as depicted in figure 3.10.

1. A mutation that causes a change from **glutamic acid (E)** to **lysine (K)** at position 27.
2. A mutation that leads to a change from **leucine (L)** to **proline (P)** at position 115.
3. A mutation that results in a change from **alanine (A)** to **aspartic acid (D)** at position 116.
4. A mutation that changes **valine (V)** to **glycine (G)** at position 127.

MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRF FESFGDLSTPDAMGNPK  
VKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLG NVLVCV**L**AHHFG  
KEFTP **P**VQAAYQKV VAGVANALAHKYH

The screenshot shows the UniProt entry for the HBB gene. The top navigation bar includes links for BLAST, Align, Peptide search, ID mapping, SPARQL, UniProtKB, and HBB. The left sidebar lists various protein properties: Function, Names & Taxonomy, Subcellular Location, Disease & Variants, PTM/Processing, Expression, Interaction, Structure, Family & Domains, Sequence, and Similar Proteins. The main content area is titled 'Beta-thalassemia (B-THAL)' and includes a 'Variants' section. The 'Note' section states: 'The disease is caused by variants affecting the gene represented in this entry. A form of thalassemia. Thalassemias are common monogenic diseases occurring mainly in Mediterranean and Southeast Asian populations. The hallmark of beta-thalassemia is a failure in globin-chain production in the adult Hb molecule. Absence or partial chain length polymorphisms, while reduced amounts of detectable hemoglobin raises beta-H thalassemia. In the severe form of beta-thalassemia, the excess alpha-globin chains accumulate in the developing erythroid precursors in the marrow. Their degradation leads to a vast increase in erythroid apoptosis that in turn causes ineffective erythropoiesis and severe microcytic hypochromic anemia. Clinically, beta-thalassemia is divided into thalassemia major which is transfusion dependent, thalassemia intermedia (of intermediate severity), and thalassemia minor that is asymptomatic.' Below this is a 'Description' section. A 'See also' link points to MIM:613985. The 'Natural variants in B-THAL' table lists four variants:

Natural variants in B-THAL	VARIANT ID	POSITION(S)	CHANGE	DESCRIPTION
	VAR_002907	27	E>K	in B-THAL; Hb E; confers resistance to severe malaria; dbSNP rs33950507
	VAR_010145	115	L>P	in B-THAL; Dohaen-NC; dbSNP rs3601591
	VAR_003037	116	A>D	in B-THAL; Haderer-Kralow; unstable; dbSNP rs3348509
	VAR_003056	127	V>G	in B-THAL; Dhonburi/Hepatic; unstable; dbSNP rs33925391

**Figure 3.10.** UniProt - HBB The 4 amino acids that change and cause beta-thalassemia



## **Chapter 4**

### **Exercise 4**

---

Using the beta thalassemia protein sequence, perform the relevant BLAST. Select the sequences *Homo sapiens*, *Gorilla gorilla gorilla*, *Pongo abelii*, *Nomascus leukogenys*, *Trachypithecus francois*, *Pongo pygmaeus* and using whichever program you wish, construct a Neighbor Joining phylogenetic tree. What do you observe?

## 4.1 NCBI Blast Neighbor Joining Phylogenetic Tree

Having obtained beta thalassaemia protein sequence (HBB Homo sapiens), we perform the relevant Protein BLAST blastp 4.1. We select the species for which we want to create a neighbor joining phylogenetic tree as shown in figures 4.2, 4.3, 4.4 and we save their sequences to perform a separate analysis with MEGA11, and we select 'Distance tree of results' to display it, choosing 'Neighbor Joining' in Tree Method option. We select 'Taxonomic Name (Sequence ID)' in Sequence Label option for a more clear visualized result. 4.5

The screenshot shows the NCBI BLAST suite interface. In the 'Enter Query Sequence' section, a protein sequence (HBB Homo sapiens) is entered. In the 'Choose Search Set' section, 'Standard' databases are selected. Under 'Program Selection', 'blastp (protein-protein BLAST)' is chosen. At the bottom, the 'BLAST' button is highlighted.

**Figure 4.1.** HBB Homo sapiens blastp

The screenshot shows the NCBI BLAST suite interface displaying search results for the query HBB Homo sapiens. The results table includes columns for Description, Scientific Name, Max Score, Total Score, Query Cover, E value, Per cent Ident, Acc Len, and Accession. Several entries are listed, including synthetic construct, Homo sapiens, and Gorilla gorilla gorilla.

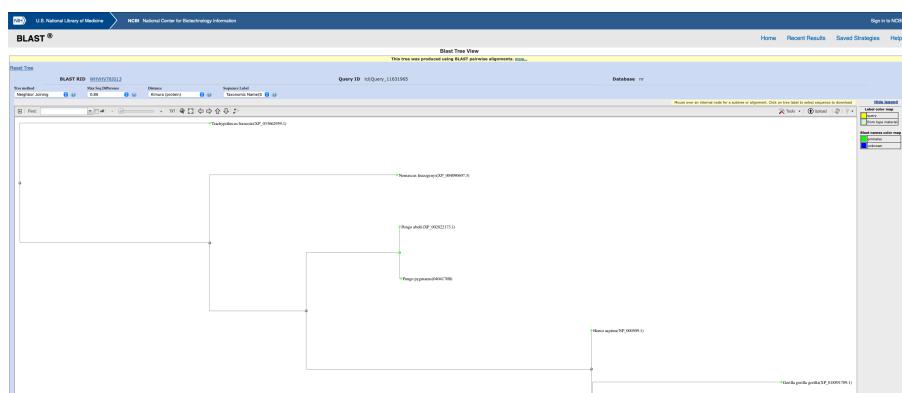
**Figure 4.2.** blastp sequences selected I

Based on the constructed phylogenetic tree the following we can make the following observations:

<input checked="" type="checkbox"/> Hemoglobin subunit beta [Gorilla gorilla gorilla]	Gorilla gorilla gorilla	301	301	100%	Se-103	100.00%	147	NP_000508.1
<input type="checkbox"/> beta subunit chain variant [Homo sapiens]	Homo sapiens	299	299	100%	Se-102	99.32%	147	XP_1989120.1
<input type="checkbox"/> beta subunit [Homo sapiens]	Homo sapiens	299	299	100%	Se-102	99.32%	147	AAB4548.1
<input type="checkbox"/> hemoglobin beta chain [Homo sapiens]	Homo sapiens	299	299	100%	Se-102	99.32%	147	ACU6684.1
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	299	299	100%	Se-102	99.32%	147	AAQ18696.1
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	298	298	99%	Se-102	100.00%	146	TANJ_B
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	298	298	100%	Se-102	99.32%	147	QB168124.1
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	298	298	100%	Se-102	99.32%	147	AAAB8254.1
<input type="checkbox"/> HBB [Intrinsic construct]								
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	298	298	100%	Se-101	99.32%	148	7KAB_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	298	298	99%	Se-101	99.32%	146	ZY5E_B
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	297	297	99%	Se-101	99.32%	146	TADQ_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	297	297	99%	Se-101	99.32%	146	AMGD_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (DEOXY) BETA-V87T [Homo sapiens]	Homo sapiens	297	297	99%	Se-101	99.32%	146	IHDQ_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (DEOXY) (BETA-CHAIN) [Homo sapiens]	Homo sapiens	297	297	98%	Se-101	100.00%	146	IDWZ_B
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	297	297	98%	Se-101	100.00%	145	SEED_B
<input type="checkbox"/> mutant beta-chain [Homo sapiens]	Homo sapiens	297	297	100%	Se-101	99.32%	147	AAJ6878.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	297	297	99%	Se-101	99.32%	146	INQD_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	297	297	99%	Se-101	99.32%	146	IXYK_B
<input type="checkbox"/> hemoglobin beta variant Hb S-Wake [Homo sapiens]	Homo sapiens	296	296	100%	Se-101	98.84%	147	AAN11320.1
<input checked="" type="checkbox"/> hemoglobin subunit beta [Pongo abelii]	Pongo abelii	296	296	100%	Se-101	98.84%	147	XP_022622123.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	296	296	98%	Se-101	100.00%	145	YB99_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	296	296	99%	Se-101	98.83%	146	1Y2Z_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	296	296	99%	Se-101	98.83%	146	1O1D_B
<input type="checkbox"/> sickle beta-hemoglobin [Homo sapiens]	Homo sapiens	296	296	100%	Se-101	98.84%	147	AAJ2996.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	296	296	99%	Se-101	98.83%	146	Y3YD_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	296	296	99%	Se-101	98.83%	146	1AOI_B
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	296	296	99%	Se-101	98.83%	146	Y3YF_B

**Figure 4.3.** blastp sequences selected II

<input checked="" type="checkbox"/> Hemoglobin subunit beta [Nomascus leucogenys]	Nomascus leucogenys	295	295	100%	Se-100	97.96%	147	XP_05099687.3
<input type="checkbox"/> Chain B. HEMOGLOBIN [Homo sapiens]	Homo sapiens	294	294	99%	Se-100	99.32%	146	1LGR_B
<input checked="" type="checkbox"/> Hemoglobin subunit beta [Macaca fasciata fasciata]	Macaca fasciata fasciata	294	294	100%	Se-100	99.32%	147	XP_03362895.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	294	294	99%	Se-100	98.83%	146	1Y2A_B
<input type="checkbox"/> Redshank. Full-hemoglobin subunit beta [African: Full-Beta-globin: African: Full-Hemoglobin beta chain]	Holothuria leucospilus	294	294	99%	Se-100	98.83%	146	P20282.1
<input type="checkbox"/> Chain B. HEMOGLOBIN (BETA-CHAIN) [Homo sapiens]	Homo sapiens	294	294	99%	Se-100	98.83%	146	1LZL_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	294	294	99%	Se-100	97.96%	146	1UYV_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	294	294	99%	Se-100	98.83%	146	1R5A_B
<input checked="" type="checkbox"/> Hemoglobin beta chain [Pongo pygmaeus]	Pongo pygmaeus	293	293	99%	Se-100	98.83%	146	XP_011270708.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	293	293	99%	Se-100	98.83%	146	1YQD_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	293	293	99%	Se-100	98.83%	146	1YQF_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	293	293	99%	Se-100	98.83%	146	1YQH_B
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	293	293	100%	Se-100	97.96%	147	AAJ2998.1
<input type="checkbox"/> Chain B. Hemoglobin beta chain [Homo sapiens]	Homo sapiens	293	293	99%	Se-100	97.96%	147	1O1E_B
<input type="checkbox"/> PREDICTED: hemoglobin subunit beta [Rhinoceros hecatenatus]	Rhinoceros hecatenatus	293	293	100%	Se-99	96.60%	147	XP_017717110.1
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Pongo pygmaeus]	Pongo pygmaeus	291	291	97%	Se-99	100.00%	143	ENCL_B
<input type="checkbox"/> Chain B. HEMOGLOBIN [Homo sapiens]	Homo sapiens	291	291	99%	Se-99	98.83%	146	1OBL_B
<input type="checkbox"/> Redshank. Full-Hemoglobin subunit beta [African: Full-Beta-globin: African: Full-Hemoglobin beta chain]	Semnopithecus obreyi	291	291	99%	Se-99	97.96%	146	P02032.1
<input type="checkbox"/> Hemoglobin subunit beta [Chlorocebus sabaeus]	Chlorocebus sabaeus	291	291	100%	Se-99	95.92%	147	XP_00118847.1
<input type="checkbox"/> PREDICTED: hemoglobin subunit beta [Coturnix japonica pallida]	Coturnix japonica pallida	291	291	100%	Se-99	95.92%	147	XP_01189011.1
<input type="checkbox"/> hemoglobin subunit beta [Rhinoceros hecatenatus]	Rhinoceros hecatenatus	290	290	100%	Se-99	95.92%	147	XP_01951461.1
<input type="checkbox"/> hemoglobin subunit beta [Equus caballus]	Equus caballus	290	290	100%	Se-99	95.92%	147	XP_003833.1
<input type="checkbox"/> Redshank. Full-hemoglobin subunit beta [African: Full-Beta-globin: African: Full-Hemoglobin beta chain]	Alotes australis	290	290	100%	Se-99	95.92%	147	QW9N23.3
<input type="checkbox"/> hemoglobin subunit beta [Macacus fasciatus]	Macacus fasciatus	290	290	100%	Se-99	95.92%	147	AVO8937.1
<input type="checkbox"/> Redshank. Full-hemoglobin subunit beta [African: Full-Beta-globin: African: Full-Hemoglobin beta chain]	Macacus sinicus	289	289	100%	Se-99	95.92%	147	XP_003232.2
<input type="checkbox"/> hemoglobin subunit beta [Cercopithecus aethiops]	Cercopithecus aethiops	289	289	100%	Se-99	95.92%	147	NP_001292889.1
<input type="checkbox"/> hemoglobin subunit beta [Cercopithecus aethiops]	Cercopithecus aethiops	289	289	100%	Se-99	95.92%	147	QW9N23.3
<input type="checkbox"/> Redshank. Full-Hemoglobin subunit beta [African: Full-Beta-globin: African: Full-Hemoglobin beta chain]	Layrinya leptocheila	289	289	100%	Se-98	95.92%	147	QW9N23.3
<input type="checkbox"/> hemoglobin subunit beta [Calithrix jacchus]	Calithrix jacchus	289	289	99%	Se-98	95.80%	146	PO028.1
<input type="checkbox"/> hemoglobin subunit beta [Loxodonta africana]	Loxodonta africana	288	288	100%	Se-97	95.24%	147	XP_00254937.1
<input type="checkbox"/> Chain B. Hemoglobin subunit beta [Homo sapiens]	Homo sapiens	288	288	99%	Se-97	100.00%	141	SHL8_B

**Figure 4.4.** blastp sequences selected III**Figure 4.5.** NCBI Neighbor Joining Phylogenetic Tree

- Homo sapiens and Gorilla gorilla gorilla sequences are closely related, indicating a close evolutionary relationship between these two species
- Both Pongo species (Pongo abelii and Pongo pygmaeus) branch off together in the phylogenetic tree, again indicating their close relationship with each other and their separation from the African apes (humans and gorillas).

- The *Nomascus leucogenys* (gibbon) and *Trachypithecus francoisi* (François' leaf monkey) are also placed on separate branches from the great apes, which is consistent with the evolutionary distance between these species.
- Both species of orangutan, *Pongo abelii* and *Pongo pygmaeus*, are grouped together, but they form a separate branch from the rest of the apes, indicating a divergence from the common ancestor they share with the African apes (gorillas, humans) and the Asian apes (*Nomascus leucogenys* and *Trachypithecus francoisi*).

## 4.2 MEGA Neighbor Joining Phylogenetic Tree

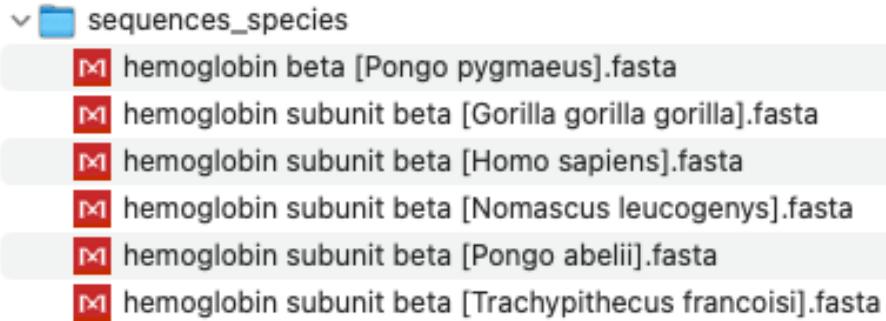
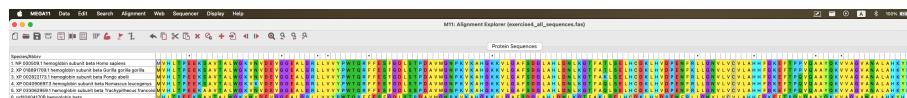
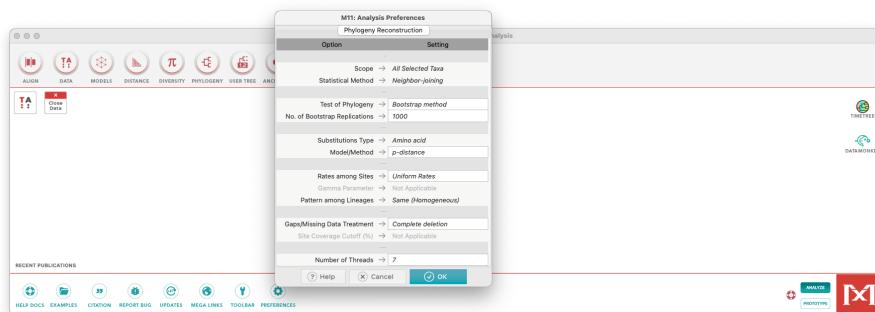
Having obtained all HBB protein sequences for all species, we import them to MEGA 4.7 and we perform multiple alignment 4.8. After alignment is successfully performed, we export aligned sequences in .meg file we open .meg file and we select the 'Construct/Test Neighbor-Joining tree with the appropriate settings 4.9. After clicking Ok we get the results 4.10, 4.11.

Based on the constructed phylogenetic tree the following we can make the following observations:

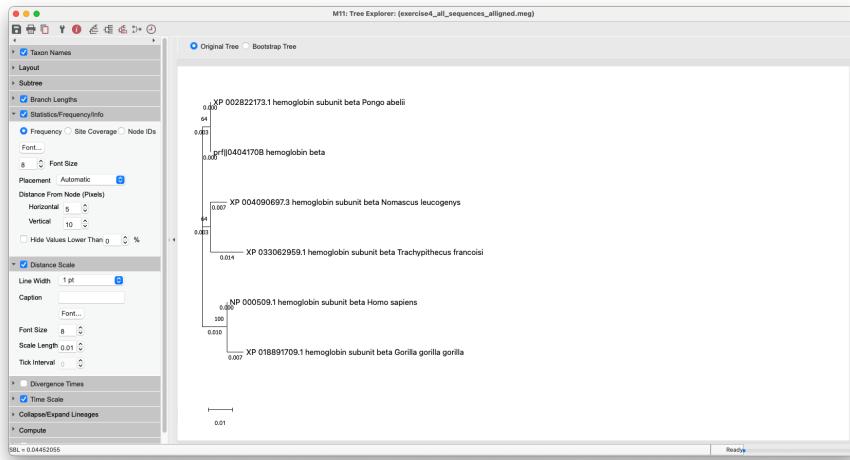
- *Homo sapiens* is shown to be most closely related to *Gorilla gorilla gorilla*, which is consistent with current understanding of hominid evolutionary relationships.
- The tree indicates that all the listed species have a common ancestor, with subsequent branching representing evolutionary divergence.
- The branch lengths represent the amount of genetic change, with the numbers presumably representing substitutions per site. Humans (*Homo sapiens*) and gorillas (*Gorilla gorilla gorilla*) have the shortest branches, suggesting fewer genetic changes since their divergence from a common ancestor compared to the other species.
- Both species of orangutan, *Pongo abelii* and *Pongo pygmaeus*, are grouped together, but they form a separate branch from the rest of the apes, indicating a divergence from the common ancestor they share with the African apes (gorillas, humans) and the Asian apes (*Nomascus leucogenys* and *Trachypithecus francoisi*).
- *Nomascus leucogenys* (gibbon) and *Trachypithecus francoisi* (François' leaf monkey) are grouped on the same branch, which is consistent with them both being Asian primates, but they diverged significantly from both the African apes and orangutans.

## 4.3 Differences between the two constructed phylogenetic trees

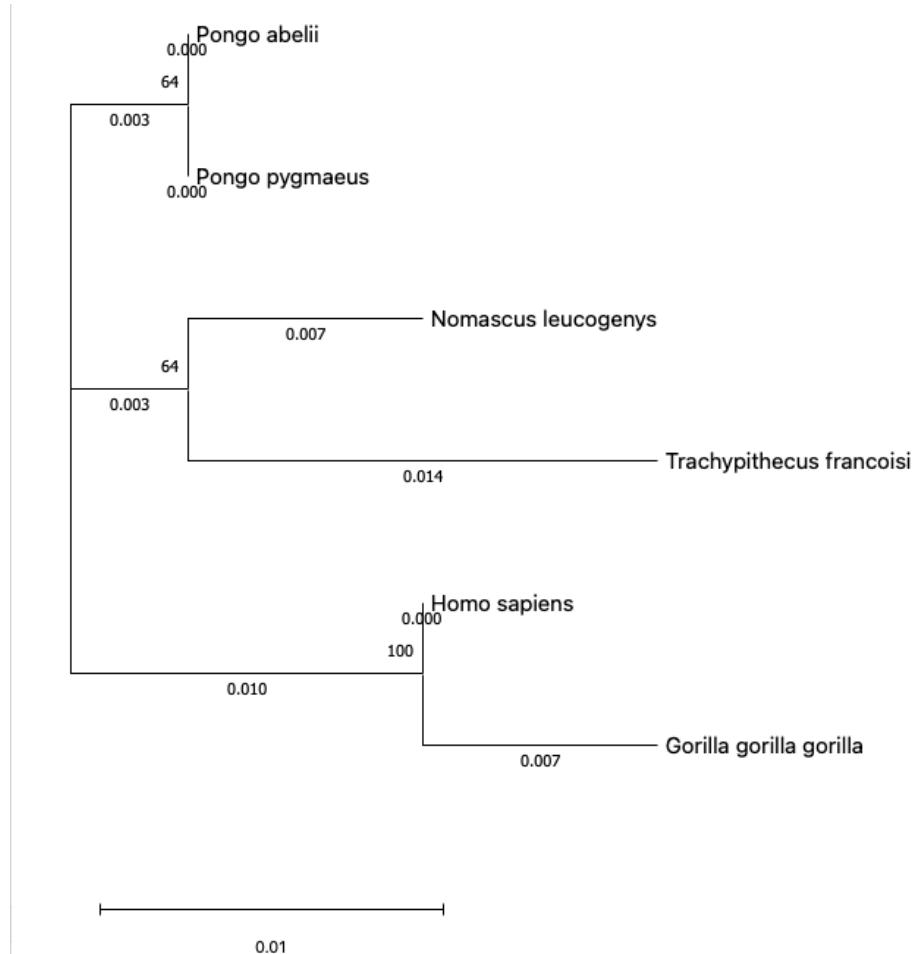
There are several differences between the two phylogenetic trees, both in terms of structure and the information they present. For brevity we will refer to the phylogenetic tree produced by MEGA11 as 'MEGA tree' and for phylogenetic tree produced by NCBI as 'NCBI tree'. Some of the differences are presented below:

**Figure 4.6.** Obtained sequences files in .fasta format**Figure 4.7.** All species sequences in .fasta format**Figure 4.8.** All species aligned sequences**Figure 4.9.** Neighbor Joining Phylogenetic Tree settings

- The MEGA tree has branch lengths with numerical values that represent genetic distances or the number of substitutions per site. This detail gives a sense of how much genetic change has occurred since the divergence from a common ancestor, while the tree, taken from the BLAST interface, does not show these numerical branch length values, making it harder to assess the exact amount of genetic change or distance.
- In MEGA tree, Homo sapiens is placed closer to Gorilla gorilla gorilla, while the two Pongo species (orangutans) are on a separate branch, and the gibbon (Nomascus leucogenys) and leaf monkey (Trachypithecus francoisi) are on another. In NCBI produced tree, the relationships seem to be consistent with the first in that humans are closer to gorillas, and the two orangutan species are grouped together. However, without the branch length data or bootstrap values, the strength of these relationships is not quantified.



**Figure 4.10.** Neighbor Joining Phylogenetic Tree



**Figure 4.11.** Neighbor Joining Phylogenetic Tree clean

- The MEGA tree includes a bootstrap value (64) at the node separating the orangutan species from the other primates, providing some measure of confidence in that grouping, while NCBI blastp produced tree does not have visible bootstrap values or other measures of statistical support.

- While both trees were constructed using the Neighbor-Joining method, the parameters such as the model of evolution (e.g., Kimura protein model in NCBI) and the exact settings differ, which affect the outcome of the tree.
- MEGA tree has a scale bar indicating the genetic distance, which is absent from the NCBI tree. This scale bar is crucial for interpreting the lengths of the branches in evolutionary terms.
- The MEGA tree is a more traditional representation, with clear delineation of branches and distances. The second tree is more schematic.



## **List of Abbreviations**

---

BLAST	Basic Local Alignment Search Tool
blastp	Protein BLAST
bp	base pairs
EM	Electron Microscopy
NCBI	National Center for Biotechnology Information
NMR	Nuclear Magnetic Resonance
PCR	Polymerase Chain Reaction
PDB	Protein Data Bank