

# STATISTICAL DATA ANALYSIS 2

Kamil Łopuszański

## *Phylogenetics of the hemoglobin*

The goal of the project was to conduct phylogenetic analysis on hemoglobin DNA sequences coming from a set of animals.

Libraries used in the project:

- ape
- Biostrings
- phytools
- phangorn
- msa
- httr
- XML
- stringr

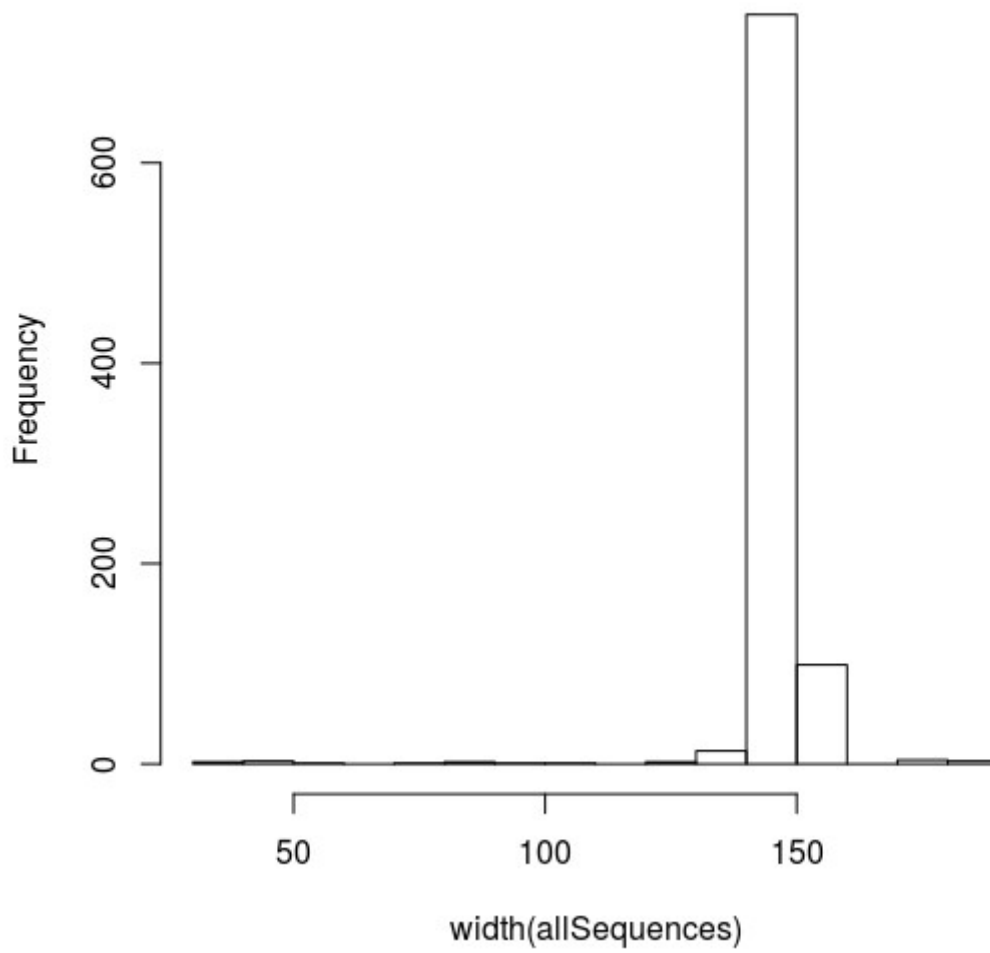
After all sequences were loaded, headers were processed. Uniprot ID of each header was extracted, as substring on positions from 1 to 8. The ID was used to request an XML document using Uniprot API. Status code of the response from the server was always 200 (OK). After the XML document was ready, it was parsed. To get the scientific name, child of the child of the root tag was extracted. The scientific name was split to 2 words and a new string was created (function **toShortName(name)**). Apart from the scientific name, taxons of the organism was extracted. The author checked the class of an organism in the following way. If *Homo* was part of the taxons set, the organism belonged to *Homo sapiens* species. Otherwise, the following conditions (and corresponding classes) were crafted:

- taxons included *Mammalia* (mammals)
- taxons included *Aves* (birds)
- taxons included *Amphibia* (amphibians)
- taxons included one of *Lepidosauria*, *Archelosauria* (reptiles)
- taxons included one of *Actinopterygii*, *Elasmobranchii*, *Coelacanthiformes*, *Hyperoartia*, *Hyperotreti*, *Dipnoi* (fish)

A new file **cleaned-hemoglobin.fasta** was created, which was the file created from the original **hemoglobin.fasta**, on cleaning the data (replacing existing headers with new names).

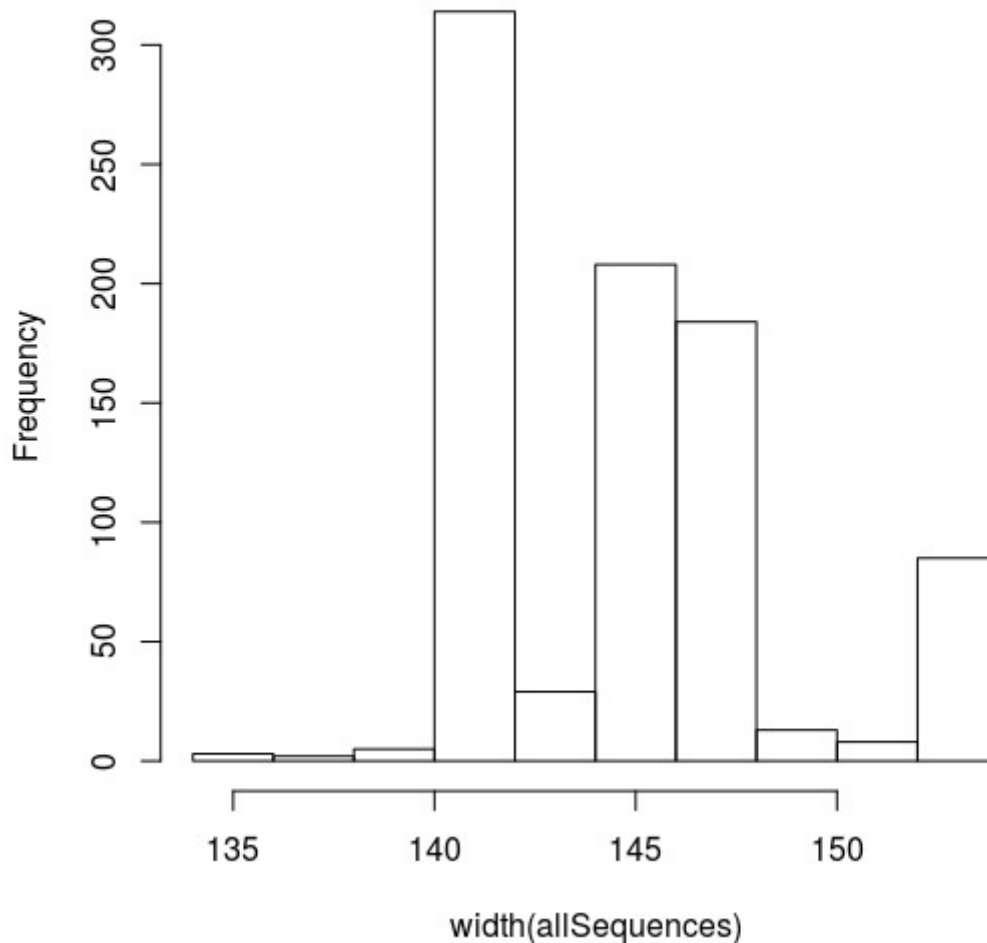
The consequent analysis was done on the cleaned data. Mean and standard deviation of the length of sequences was calculated:  $\mu=144.52, \sigma=10.81$  . Histogram:

### Histogram of width(allSequences)



To remove outliers, all the sequences whose length was either shorter than 134 (about  $\mu - std$ ) or longer than 154 (about  $\mu + std$ ). Histogram after this operation:

## Histogram of width(allSequences)



And the mean was practically unchanged (the new mean equaled 145.22), as well as the number of sequences (originally 880, then 851).

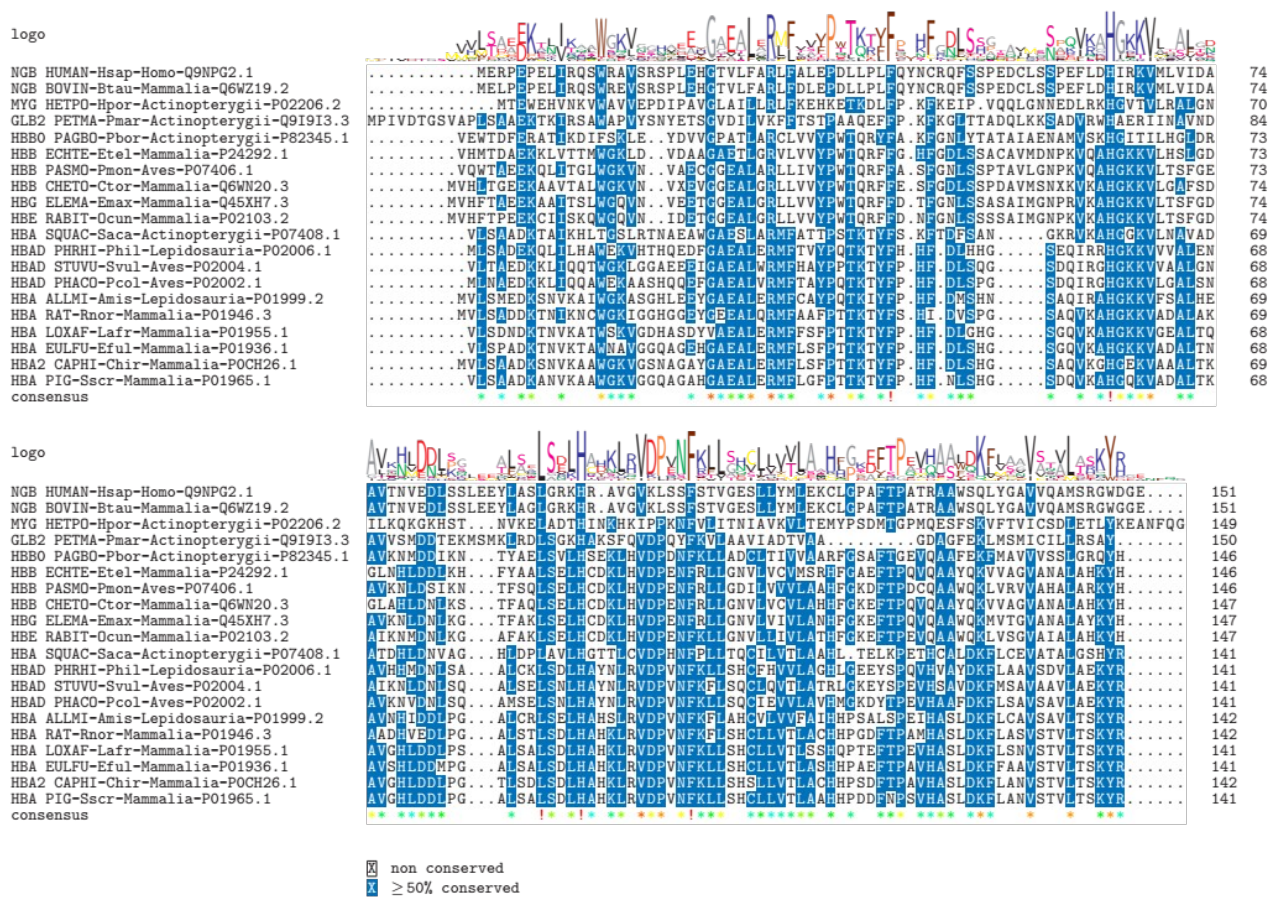
To take a sample consisting of 1 homo sapiens, 4 mammals, 3 birds, 2 reptiles and 3 fish, sets corresponding to each class were created. Then, each sequence was added to the right set. The core (set of essential organism) was created as a union of samples taken on the aforementioned sets. The picked organism were removed from the sets. Remaining organisms (not taken into the core) were gathered, let's call it **R**. In order to create a sample big enough, the sample was taken from **R**. The final sample was created as a union of the core and the sample from **R**. The sequences were written out to files 'subsetA.fasta', 'subsetB.fasta', 'subsetC.fasta'.

For each subset, the alignment was performed. Used algorithms:

- ClustalW
- ClustalOmega
- MUSCLE

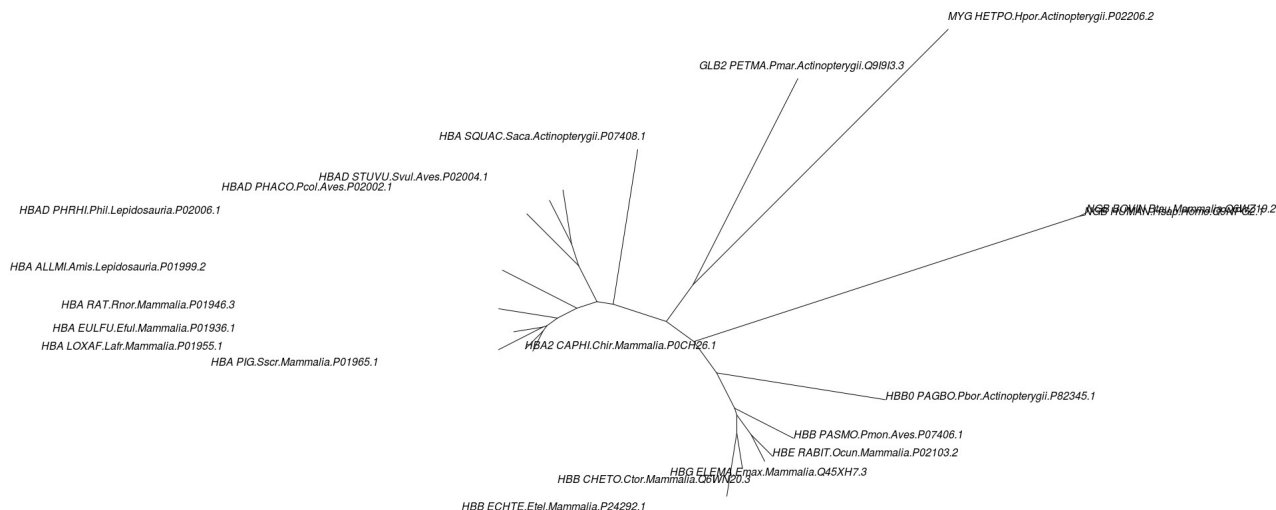
The aligned sequences were written out to the folder `out`, to files 'subsetA clustalW.fasta', 'subsetA clustalOmega.fasta', 'subsetA MUSCLE.fasta', 'subsetB clustalW.fasta', 'subsetB clustalOmega.fasta', 'subsetB MUSCLE.fasta', 'subsetC clustalW.fasta', 'subsetC clustalOmega.fasta', 'subsetC MUSCLE.fasta'.

# The output of alignment for subset C, using ClustalOmega:

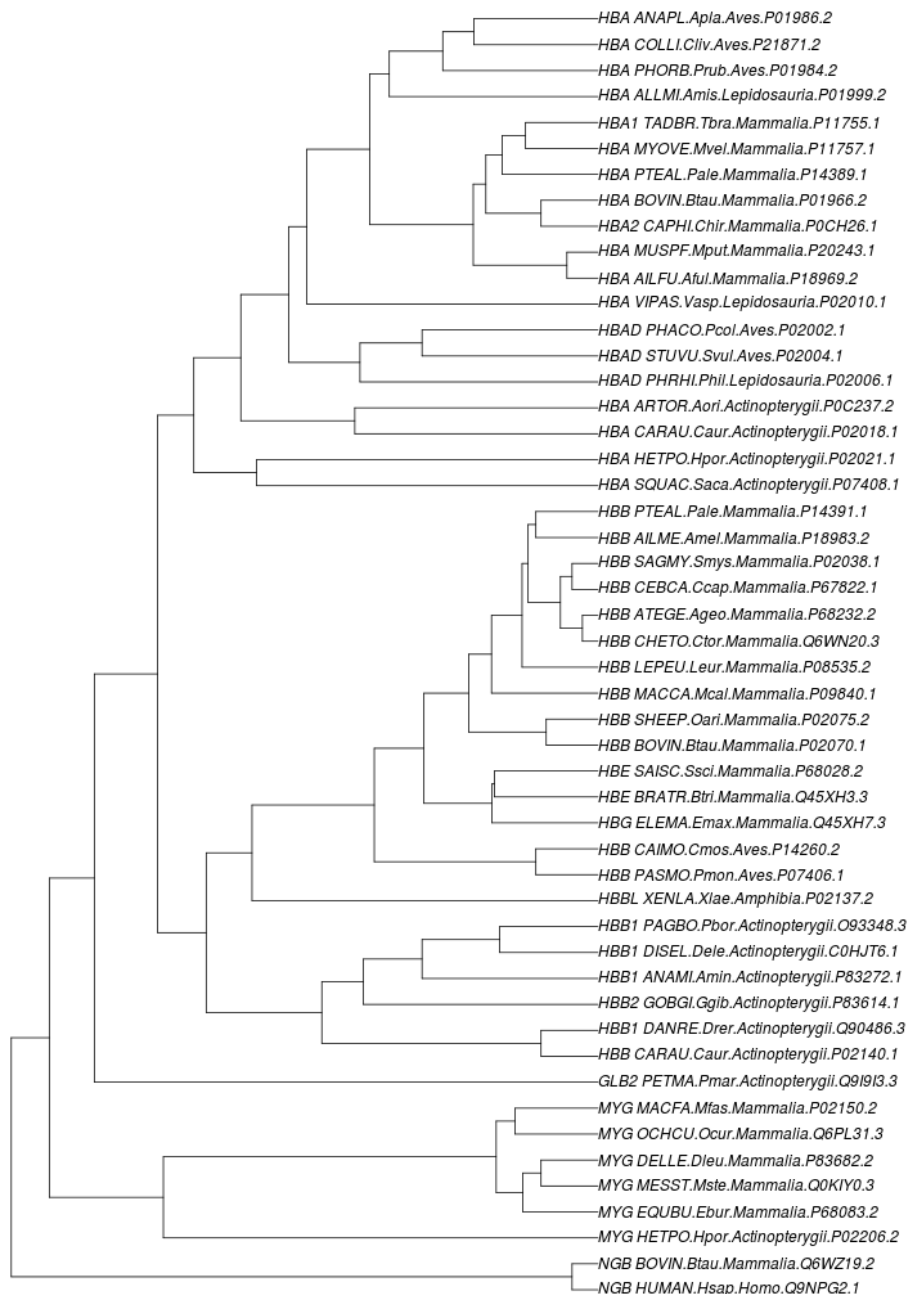


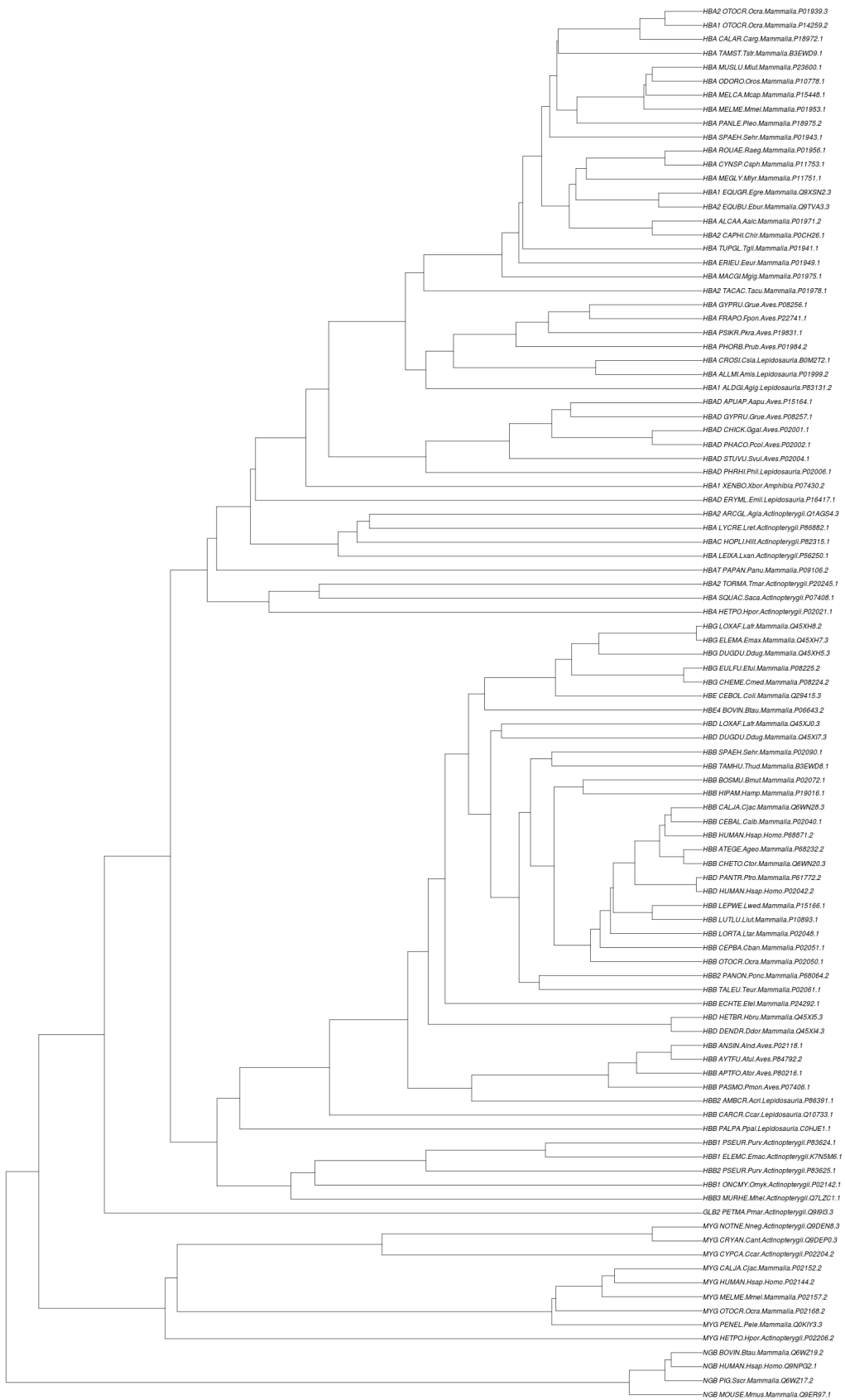
The aligned sequences were loaded and transformed to objects using functions `read.aa` and `phyDat`. Afterwards, distance matrix were calculated, using Hamming and ML metrics. The result was written out to the folder `out/distance`, to files 'distance A clustalW hamming.txt', 'distance A clustalW ml.txt', 'distance A clustalOmega hamming.txt', 'distance A clustalOmega ml.txt', 'distance A MUSCLE hamming.txt', 'distance A MUSCLE ml.txt' etc.

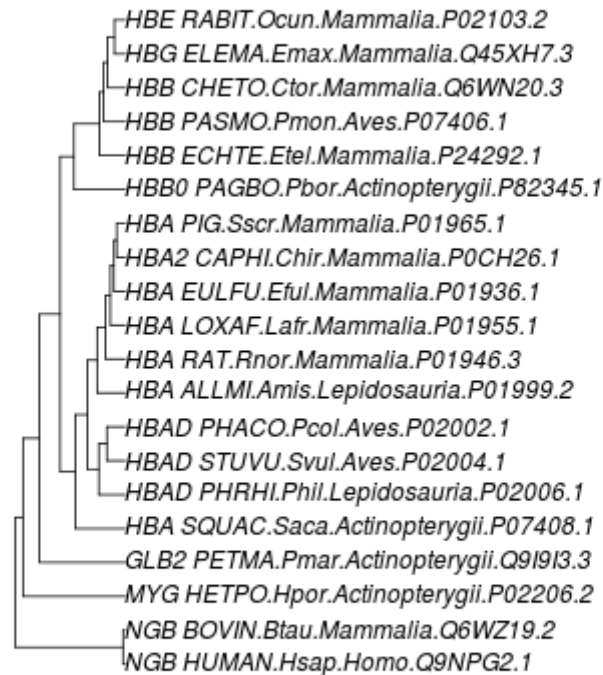
Trees were constructed, based on Neighbor joining algorithm, using all the before calculated distance matrices. Unrooted plot of tree, distance ML, alignment MUSCLE:



The parsimony scores of all trees were computed and compared. ClustalW and ClustalOmega performed better (depending on the set). UPGMA algorithm performed even better (to reduce complexity, only best matrices were considered for UPGMA).  
 The figures show trees produced using UPGMA:







UPGMA performed better than NJ (bigger parsimony score). Best trees were saved to the folders [out/nexus](#) and [out/newick](#), in nexus and newick format respectively.

Then, maximum parsimony algorithm was performed.

```
tree.max.pars.A <- optim.parsimony(tree.upgma.A, aligned.A.seqs.clustalOmega)
tree.max.pars.B <- optim.parsimony(tree.upgma.B, aligned.B.seqs.clustalW)
tree.max.pars.C <- optim.parsimony(tree.upgma.C, aligned.C.seqs.clustalOmega)
```

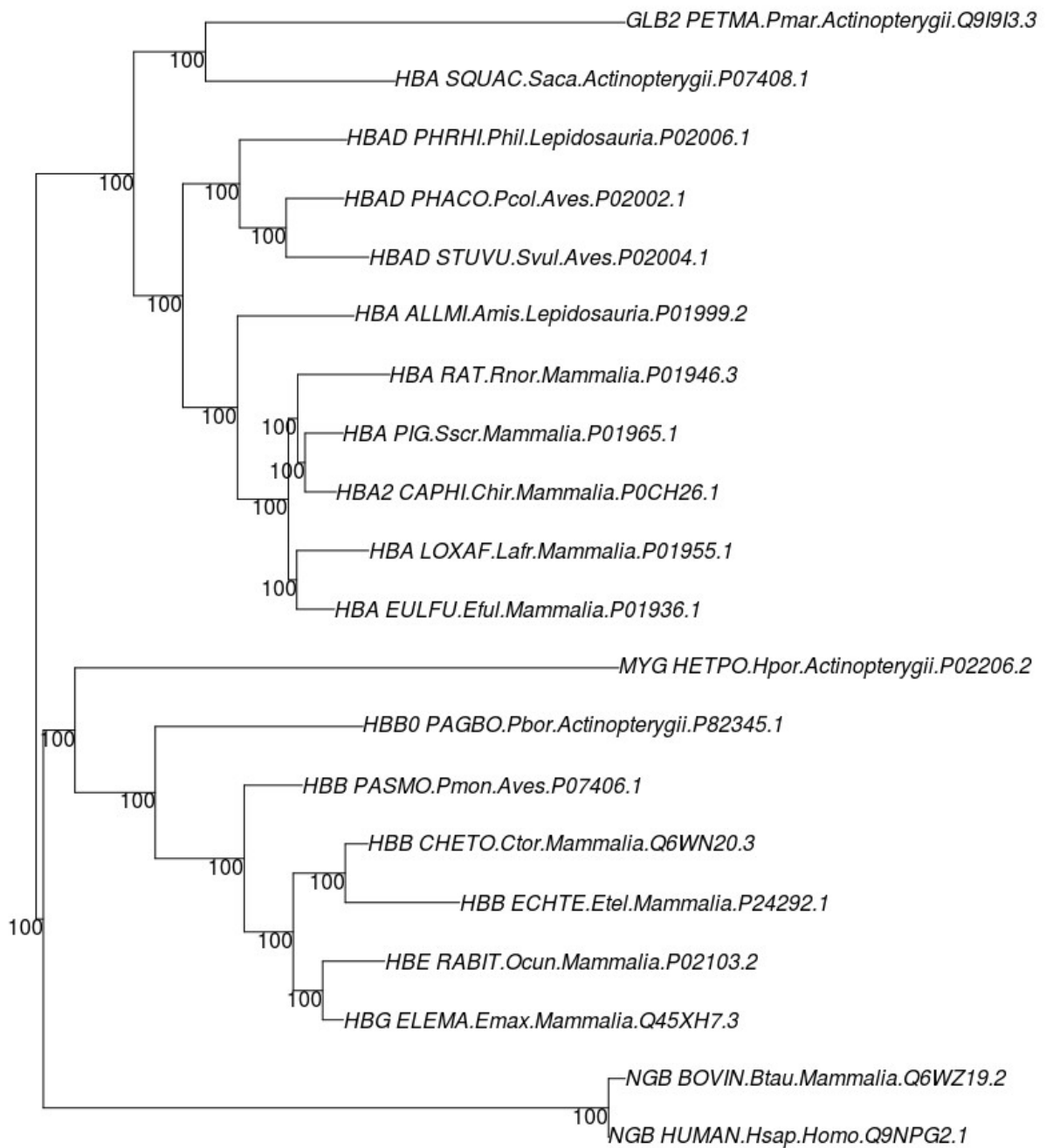
Using maximum likelihood algorithm, one can decrease the loglikelihood:

```
fit.A <- pml(tree.nj.hamming.A.clustalOmega, data = aligned.A.seqs.clustalOmega)
fit.B <- pml(tree.nj.hamming.B.clustalW, data = aligned.B.seqs.clustalW)
fit.C <- pml(tree.nj.ml.C.clustalOmega, data = aligned.C.seqs.clustalOmega)
```

```
fit.A.JC <- optim.pml(fit.A, TRUE)
fit.B.JC <- optim.pml(fit.B, TRUE)
fit.C.JC <- optim.pml(fit.C, TRUE)
```

Bootstrapping gives 100% significance on all branches, unfortunately:

```
bs.C = bootstrap.pml(fit.C.JC, bs=1000, multicore = TRUE, mc.cores = 6)
plotBS(midpoint(fit.A.JC$tree), bs.A, p = 90, type = 'p')
```



At the end, `modelTest()` analysis was performed, using JTT model. Result for subset A:

Model	df	logLik	AIC	AICw	AICc	AICcw	BIC
JTT	197	-12417.97	25229.93	1.181493e-174	NaN	NaN	25845.35
JTT+I	198	-12379.87	25155.74	1.523547e-158	NaN	NaN	25774.28
JTT+G	198	-12019.27	24434.53	6.176710e-02	NaN	NaN	25053.08
JTT+G+I	199	-12015.55	24429.09	9.382329e-01	NaN	NaN	25050.76



Result for subset B:

Model	df	logLik	AIC	AICw	AICc	AICcw	BIC
JTT	97	-7140.691	14475.38	7.539854e-80	14739.44	1.240583e-77	14779.55
JTT+I	98	-7123.369	14442.74	9.240354e-73	14716.03	1.497833e-72	14750.05
JTT+G	98	-6957.993	14111.99	6.136733e-01	14385.28	9.947458e-01	14419.29
JTT+G+I	99	-6957.455	14112.91	3.863267e-01	14395.77	5.254178e-03	14423.35

Result for subset C:

Model	df	logLik	AIC	AICw	AICc	AICcw	BIC
JTT	37	-4084.124	8242.248	3.607868e-33	8264.047	8.687541e-33	8357.614
JTT+I	38	-4069.933	8215.867	1.931479e-27	8239.023	2.358814e-27	8334.351
JTT+G	38	-4008.865	8093.730	6.420081e-01	8116.886	7.840507e-01	8212.214
JTT+G+I	39	-4008.449	8094.898	3.579919e-01	8119.465	2.159493e-01	8216.500

## Conclustions

In the opinion of the author, the best algorithm for creating phylogenetic tree is UPGMA, because it had the biggest parsimony score. Bootstrapping wasn't successful, because all branches were rendered in 100% of cases.

The alignments made using different algorithms differed. The difference wasn't big, though. In distance matrices, the typical difference was about 0.01 or 2%. From the diagram made by **msa** package, we can see that the sequence coding hemoglobin is conserved well, although it depends on the species. For homo sapiens, only a dozen (about 15%) amino acids were conserved. The synteny depends how far apart the species in the phylogenetic tree are.

From the phylogenetic trees, we can see that the results agree with taxonomy of the organism quite well. Mammals are clustered together, homo sapiens is on the neighbor branch with another mammal, organisms of the same class have often the LCA quite close. It seems not ideal, though: for example, Lepidosauria specimen should be closer on the phylogenetic tree, providing that the same class implies the phylogenetic affinity.