

Article

Fernando Freire

December 2, 2018

Contents

1	Phylogenetics Evaluation	2
1.1	Best-fit substitution model	2
1.1.1	Selected models from jModelTest	2
1.1.2	Selected model for Prottest	3
1.2	Maximum parsimony (MP) tree.	5
1.2.1	16S nucleotides	5
1.2.2	RAG1 nucleotides	5
1.2.3	RAG1 amino acid	6
1.3	Distances (Neighbor Join) trees	7
1.3.1	16S nucleotides	7
1.3.2	RAG1 nucleotides	8
1.3.3	RAG1 amino acid	9
1.4	Maximum likelihood (ML) trees	10
1.4.1	16S nucleotides	10
1.4.2	RAG1 nucleotides	11
1.4.3	RAG1 amino acid	12
1.5	Class AlignSequences	13
1.6	MSA	25
1.7	MSA generic methods	26
1.8	T-COFFEE methods	36
1.9	Main MSA method	39
1.10	T-COFFEE RAG1_AA	43
1.11	CLUSTAL	48

1 Phylogenetics Evaluation

For each of the three matrices (16S, RAG1 nucleotides, RAG1 amino acids) present:

- Best-fit substitution model.
- Maximum parsimony (MP) tree.
- Distances (NJ) tree.
- Maximum likelihood (ML) tree.
- For every tree, show bootstrap values as well.
- Brief discussion of the results: comparison of topologies, branch lengths, support of phylogenetic relationships and resolution.
- Compare the trees of the three matrices among them, and also with the reference phylogeny.

1.1 Best-fit substitution model

1.1.1 Selected models from jModelTest

We use **jModelTest** to obtain the most appropriate model of nucleotide evolution given the sequence data of RAG1 and 16S.

These are the versions involved:

```
Phyml version = 3.0
Phyml binary = PhyML_3.0_macOS_i386
jModelTest version = 2.1
OS:          Mac OS X (10.14.1)
Arch:        x86_64
Java:        1.8.0_144 (Oracle Corporation)
```

16S: GTR+I+G The models selected under AIC and BIC criteria are equivalent:

```
Model    GTR+I+G
partition 012345
-lnL     13398.4723
K        27
freqA    0.3697  R(a)    4.4083
freqC    0.2395  R(b)    6.5698
freqG    0.1687  R(c)    5.3153
freqT    0.2222  R(d)    0.3201
ti/tv    -      R(e)    20.8173
R(f)     1.0000
p-inv    0.1950  gamma  0.9680
```

RAG1: GTR+I+G The models selected under AIC and BIC criteria are equivalent:

```
Model    GTR+I+G
partition 012345
-lnL     10202.7853
K        27
freqA    0.2859  R(a)    2.8475
```

```

freqC  0.2349  R(b)    7.3488
freqG  0.2584  R(c)    3.1158
freqT  0.2208  R(d)    1.1479
ti/tv   -    R(e)   11.5334
R(f)    1.0000
p-inv   0.3840  gamma   1.2830

```

1.1.2 Selected model for Prottest

We use **ProtTest3** to infer the most appropriate model of protein evolution given the sequence data of RAG1 amino acid.

These are the versions employed:

```

Version:  3.4.2 : 8th May 2016
OS:       Mac OS X (10.14.1)
Arch:     x86_64
Java:     1.8.0_144 (Oracle Corporation)

```

RAG1 JTT+I+G Prottest results:

```

*****
Best model according to BIC: JTT+G
Confidence Interval: 100.0
*****
Model          deltaBIC      BIC          BICw          -lnL
-----
JTT+G          0.00          8062.58       0.76          3975.25
JTT+I+G        2.33          8064.91       0.24          3973.30
...
*****
Relative importance of parameters
*****
alpha          (+G):      0.763
p-inv          (+I):      0.000
alpha+p-inv    (+I+G):    0.237
freqs          (+F):      No +F models

*****
Model-averaged estimate of parameters
*****
alpha (+G):      0.379
p-inv (+I):      0.520
alpha (+I+G):    1.190
p-inv (+I+G):    0.380

*****
Best model according to LnL: JTT+I+G

```

Confidence Interval: 100.0

Model	deltaLnL	LnL	LnLw	-lnL
JTT+I+G	0.00	3973.30	0.71	3973.30
JTT+G	1.95	3975.25	0.27	3975.25
LG+I+G	7.29	3980.59	0.02	3980.59

...

Relative importance of parameters

alpha (+G): 0.276
p-inv (+I): 0.000
alpha+p-inv (+I+G): 0.724
freqs (+F): No +F models

Model-averaged estimate of parameters

alpha (+G): 0.379
p-inv (+I): 0.520
alpha (+I+G): 1.183
p-inv (+I+G): 0.379

The models are estimated under Maximum Likelihood LnL information criterion and under Bayesian Information Criterion. Both criterion selected as the best substitution models: JTT+I+G and JTT+G, exchanging the rank between them.

We launch another run, now with the graphical interface, and seek for AIC criterion (the previous run is launched in console with this command line `java -jar prottest-3.4.jar -i HRSV-A_aa.phy -all-matrices -all-distribution`).

The result is:

AKAIKE INFORMATION CRITERION

Best model according to AIC: JTT+I+G

Confidence Interval: 100.0

Model	deltaAIC	AIC	AICw	-lnL
JTT+I+G	0.00	7984.61	0.59	3973.30
JTT+G	1.89	7986.50	0.23	3975.25

...

So, by consensus, we choose the **JTT+I+G** model.

1.2 Maximum parsimony (MP) tree.

I use the MEGA v7.0 program.

1.2.1 16S nucleotides

The bootstrap consensus tree inferred from 1000 replicates. The MP tree was obtained using the Tree-Bisection-Regrafting (TBR) algorithm with search level 1. Analysis Settings:

Analysis	= Phylogeny Reconstruction
Statistical Method	= Maximum Parsimony
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Nucleotide
Gaps/Missing Data Treatment	= Use all sites
No. of Initial Trees (random addition)	= 10
MP Search level	= 1
Max No. of Trees to Retain	= 10

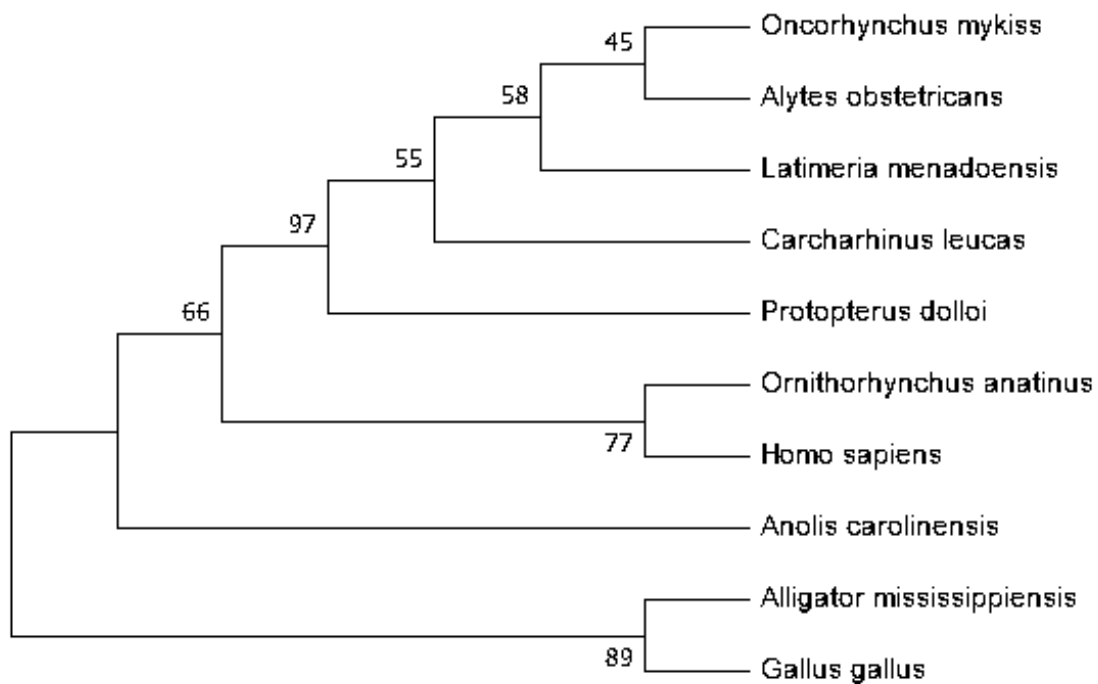


Figure 1: 16S. Maximum parsimony tree

1.2.2 RAG1 nucleotides

The bootstrap consensus tree inferred from 1000 replicates. The MP tree was obtained using the Tree-Bisection-Regrafting (TBR) algorithm with search level 1. Analysis Settings:

Analysis	= Phylogeny Reconstruction
----------	----------------------------

Statistical Method	= Maximum Parsimony
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Nucleotide
Gaps/Missing Data Treatment	= Use all sites
No. of Initial Trees (random addition)	= 10
MP Search level	= 1
Max No. of Trees to Retain	= 100
Codons Included	= 1st+2nd

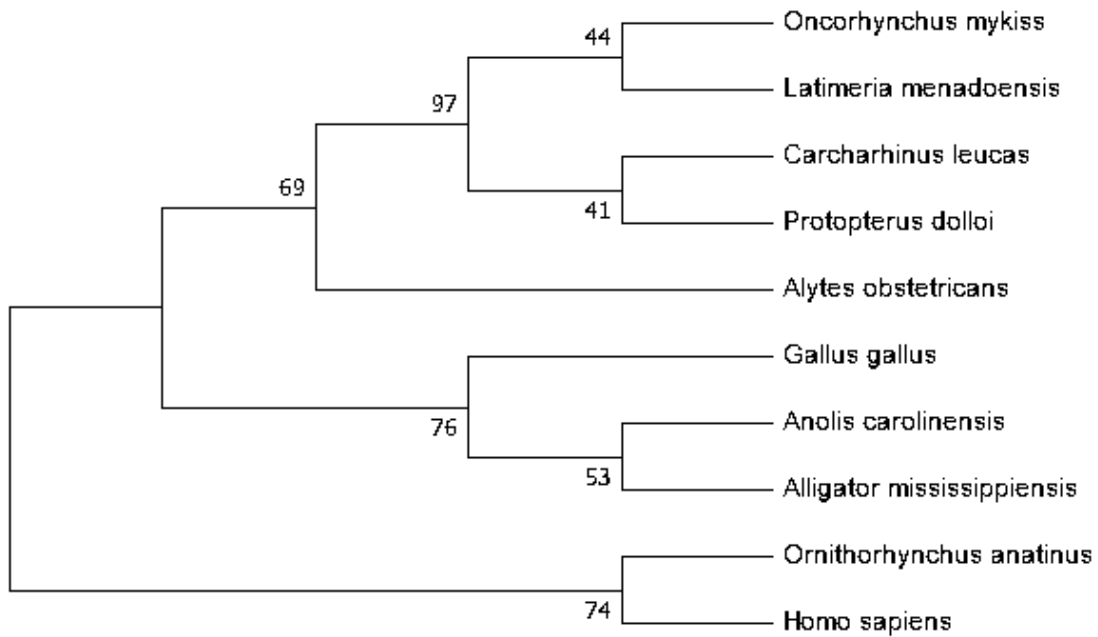


Figure 2: RAG1 nucleotides. Maximum parsimony tree

1.2.3 RAG1 amino acid

The bootstrap consensus tree inferred from 1000 replicates. The MP tree was obtained using the Tree-Bisection-Regrafting (TBR) algorithm with search level 1. Analysis Settings:

Analysis	= Phylogeny Reconstruction
Statistical Method	= Maximum Parsimony
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Amino acid
Gaps/Missing Data Treatment	= Use all sites
No. of Initial Trees (random addition)	= 10
MP Search level	= 1
Max No. of Trees to Retain	= 100

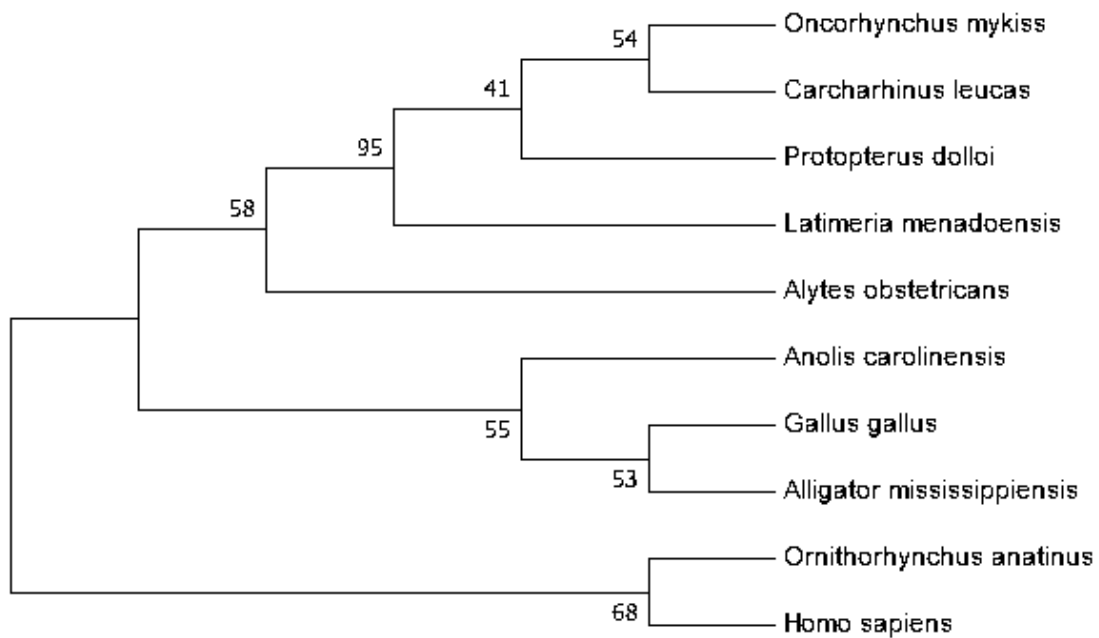


Figure 3: RAG1 AA. Maximum parsimony tree

1.3 Distances (Neighbor Join) trees

I use the MEGA v7.0 program.

1.3.1 16S nucleotides

The evolutionary distances were computed using the Kimura 2-parameter method. Analysis Settings:

Analysis	= Phylogeny Reconstruction
Scope	= All Selected Taxa
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Nucleotide
Substitutions to Include	= d: Transitions + Transversions
Rates among Sites	= Gamma Distributed (G)
Gamma Parameter	= 0.817
Pattern among Lineages	= Same (Homogeneous)
Gaps/Missing Data Treatment	= Complete deletion

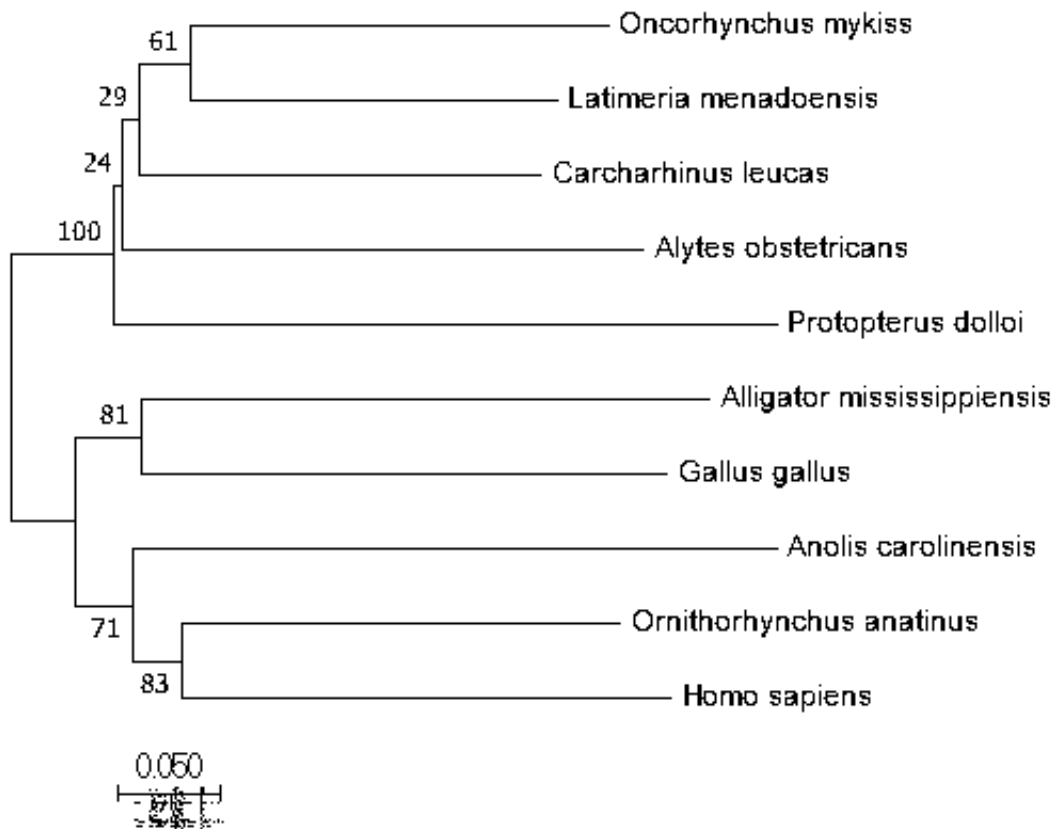


Figure 4: 16S. Neighbor join optimal tree. Sum of branch length = 2.66249603

1.3.2 RAG1 nucleotides

The evolutionary distances were computed using the Kimura 2-parameter method. Analysis Settings:

Analysis	= Phylogeny Reconstruction
Scope	= All Selected Taxa
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Nucleotide
Substitutions to Include	= d: Transitions + Transversions
Rates among Sites	= Gamma Distributed (G)
Gamma Parameter	= 0.817
Pattern among Lineages	= Same (Homogeneous)
Gaps/Missing Data Treatment	= Complete deletion

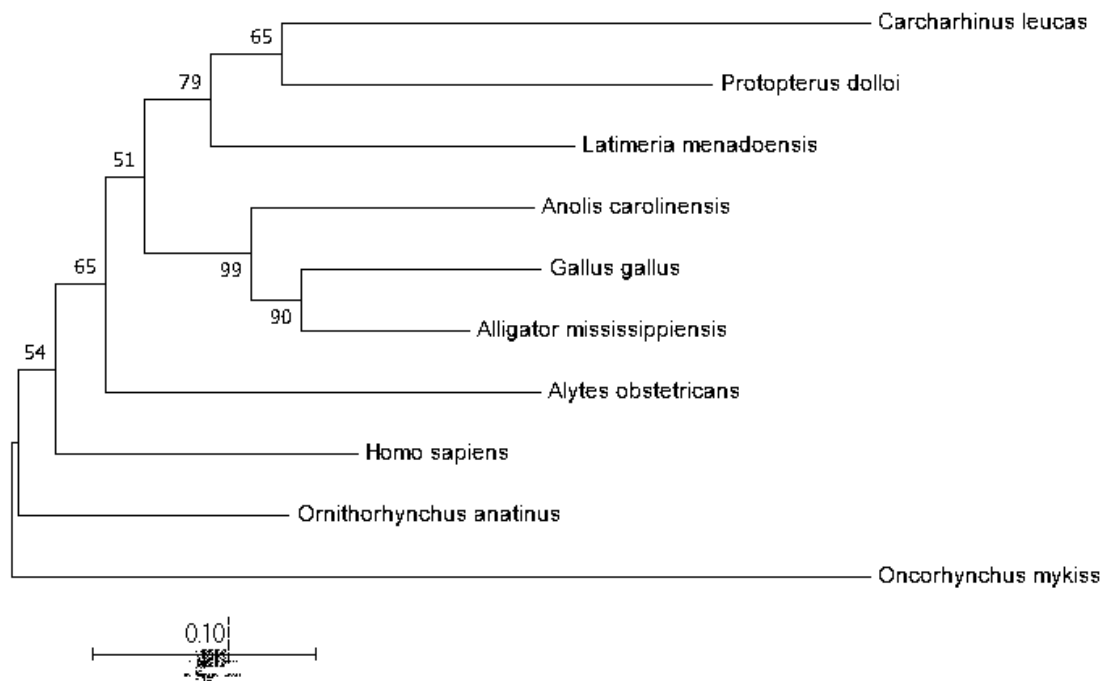


Figure 5: RAG1 nucleotides. Neighbor join optimal tree. Sum of branch length = 1.95639526

1.3.3 RAG1 amino acid

Analysis Settings:

Analysis	= Phylogeny Reconstruction
Scope	= All Selected Taxa
Test of Phylogeny	= Bootstrap method
No. of Bootstrap Replications	= 1000
Substitutions Type	= Amino acid
Model/Method	= Poisson model
Rates among Sites	= Gamma Distributed (G)
Gamma Parameter	= 0.817
Pattern among Lineages	= Same (Homogeneous)
Gaps/Missing Data Treatment	= Complete deletion

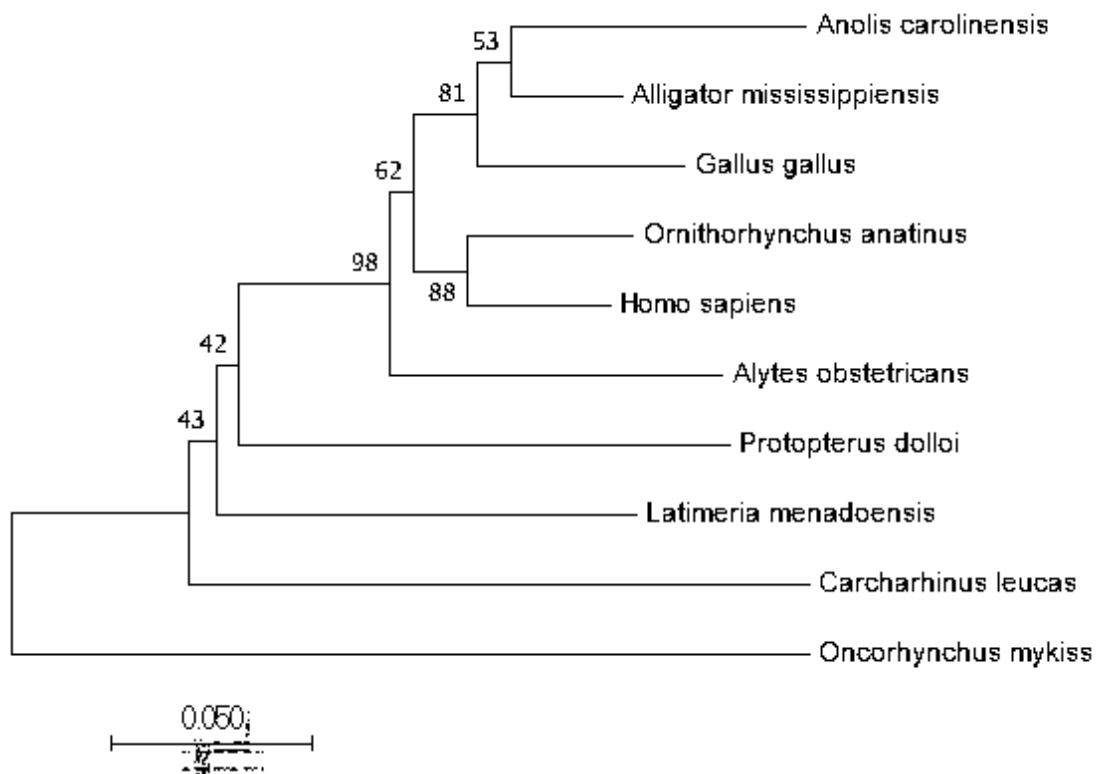


Figure 6: RAG1 AA. Neighbor join optimal tree. Sum of branch length = 1.0316

1.4 Maximum likelihood (ML) trees

I use the MEGA v7.0 program.

1.4.1 16S nucleotides

The bootstrap consensus tree inferred from 1000 replicates. The tree is computed by Maximum Likelihood method based on the Tamura-Nei. Initial tree(s) for the heuristic search were obtained automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using the Maximum Composite Likelihood (MCL) approach, and then selecting the topology with superior log likelihood value.

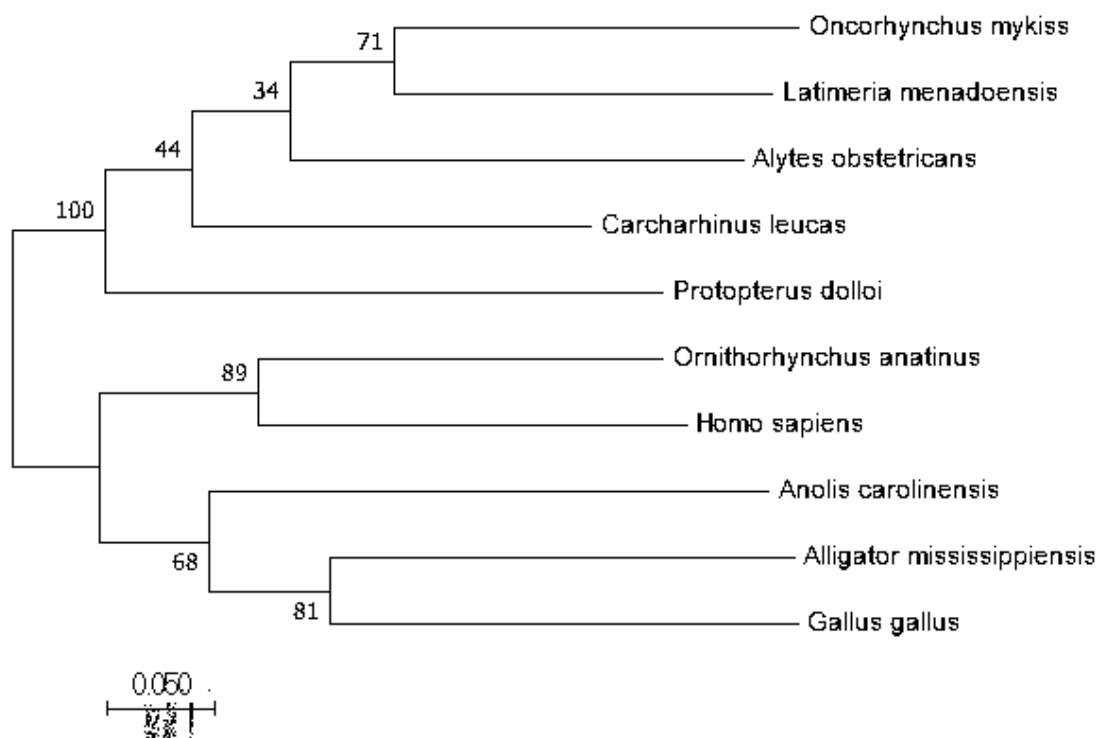


Figure 7: 16S. Maximum likelihood tree

1.4.2 RAG1 nucleotides

The bootstrap consensus tree inferred from 1000 replicates. The tree is computed by Maximum Likelihood method based on the Tamura-Nei. Initial tree(s) for the heuristic search were obtained automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using the Maximum Composite Likelihood (MCL) approach, and then selecting the topology with superior log likelihood value.

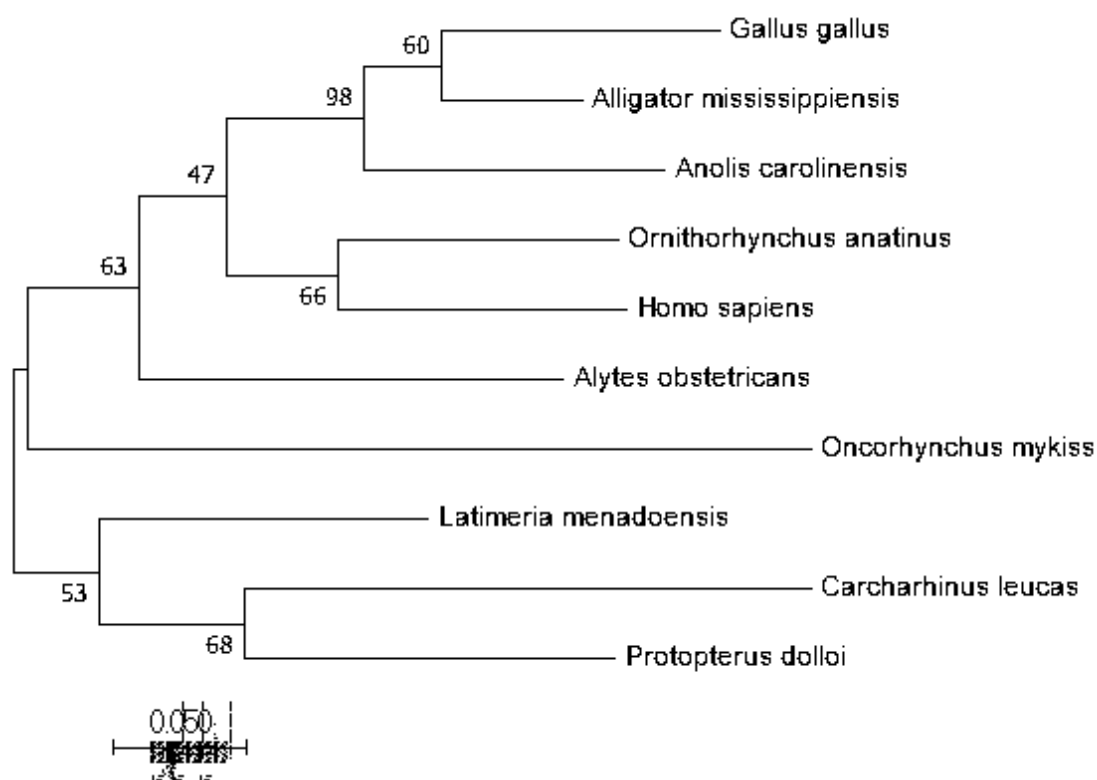


Figure 8: RAG1 nucleotides. Maximum likelihood tree

1.4.3 RAG1 amino acid

The bootstrap consensus tree inferred from 1000 replicates. The tree is computed by Maximum Likelihood method based on Equal Input model. Initial tree(s) for the heuristic search were obtained automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using the Maximum Composite Likelihood (MCL) approach, and then selecting the topology with superior log likelihood value.

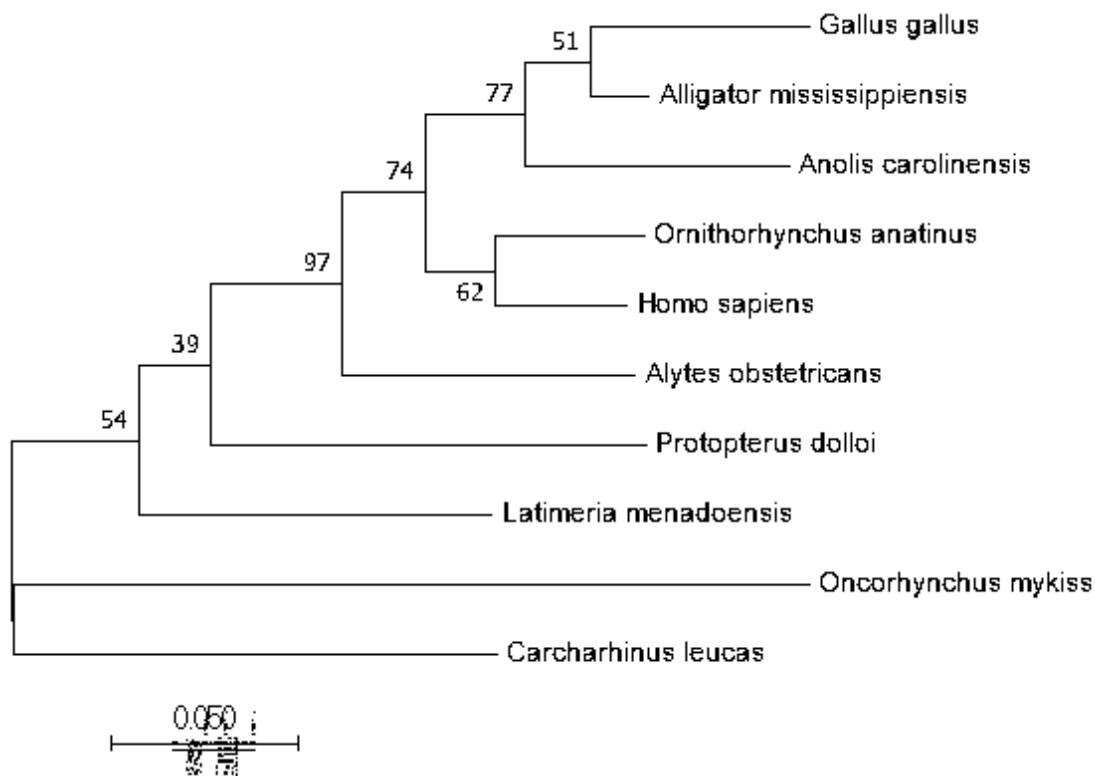


Figure 9: RAG1 AA. Maximum likelihood tree

1.5 Class AlignSequences

This class implements recursively three alignment algorithms: 1. Global alignment (Needleman-Wunsch inspired).

2. Local alignment (Smith-Waterman inspired).

3. Longest common substring. (the search for the longest common sequence can also be considered a type of alignment).

Script 1.5.1 (python)

```
1  """This module shows alternative recursive implementations of global sequence alignments:
2      Global alignment (Needleman-Wunsch based)
3      Local (Smith-Waterman based)
4      Finding of the longest common substring.
5  Todo:
6      * Return all the solutions of the alignments. Now it only returns one solution
7      * Control of errors
8      * Implement multi-alignments
9      * Implement heuristic algorithms
10 """
11 import re
12 from Bio import pairwise2
```

```

13 from Bio.pairwise2 import format_alignment
14 from Bio.SubsMat import MatrixInfo
15 import time
16 import sys
17
18 MIN = -sys.maxsize - 1
19 COMPAC = 100000
20 """int: Constant to compact max score."""
21 SCORE_MATCH = 2
22 """int: Default match score."""
23 SCORE_NO_MATCH = -3
24 """int: Default no match score."""
25 SCORE_GAP_INI = -10
26 """int: Default gap ini in affine gap penalty."""
27 SCORE_GAP_CONT = -2
28 """int: Default gap continuation in affine gap penalty"""
29 DEFAULT_SUBST_MATRIX = {('A', 'A'): 0, ('A', 'C'): 1, ('A', 'G'): 1, ('A', 'T'): 1, ('C',
→ 'A'): 1, ('C', 'C'): 0, ('C', 'G'): 1, ('C', 'T'): 1, ('G', 'A'): 1, ('G', 'C'): 1,
→ ('G', 'G'): 0, ('G', 'T'): 1, ('T', 'A'): 1, ('T', 'C'): 1, ('T', 'G'): 1, ('T', 'T'): 0}
30 """dict: Default substitution matrix (for "ACGT" common nucleotide alphabet)"""
31
32 sys.setrecursionlimit(5000)
33
34 class AlignSequences:
35     """Recursive implementation of global, local and long substring alignments methods.
36
37     Attributes:
38         sequences (list of str): Contains the two sequences to align. The first
39         one (index 0) is the query sequence (BLAST concept) or bottom sequence on
→ alignment prints
40         or vertical sequence in the common graphical representation of score matrix.
41         len_seq0 (int): Sequence 0 length.
42         len_seq1 (int): Sequence 1 length.
43         mode (str): Computation mode:
44             'GLOBAL'           Global Alignment
45             'LOCAL'           Local Alignment
46             'LONG_SUBSTRING'   Obtain long common substring
47         score_match (int): Score of match characters.
48         score_no_match (int): Score of no match characters.
49         score_gap_ini (int): Score of gap init.
50         score_gap_cont (int): Score of gap continuation.
51         score (int): Score of last computed alignment.
52         gaps (int): Number of gaps of the last computed alignment.
53         matches (int): Number of matches of the last computed alignment.
54         unmatches (int): Number of unmatches of the last computed alignment.
55         align_seq0 (str): Sequence 0 with the gaps necessary for the alignment.
56         align_seq1 (str): Sequence 1 with the gaps necessary for the alignment.
57         matching (str): Printable line with the align relations ('|', '.', ' ') between
58         both align_seq, necessary for printing the alignment.
59         ini_time (int): Initial time of computation, for profiling purposes
60         finish_time (int): Final time of computation, for profiling purposes

```

```

61         score_store (dict of tuple int): Store of scores, for each calculated cell with
↳ tuple (i,j,g)
62         where i is the coordinate of the bottom sequence, j the coordinate of the top
↳ sequence
63         and g has the value 1 if the cell is a gap init cell and 0 if it's a gap
↳ continuation.
64         For a explanation of calculared cell see align method.
65         matches_store (dict of tuple int): Store of the number of matches in the calculated
↳ cell
66         gaps_store (dict of tuple int): Store of the number of gaps in the calculated cell
67         max_score_index (tuple of int): Cell coordinate tuple of the cell with the maximum
↳ score
68         max_score (int): maximum computed score
69         forward_arrow (dict of str): Store of the optimal displacements accomplished at a
↳ cell
70         to guarantee an optimal score: 'v' vertical (down), 'h' horizontal (righth),
71         'd' diagonal.
72         stacks (list of list of str): Stacks of sequences related to principal sequences in
↳ a msa
73         stacks_indexes (list of str): Indexes of the sequences of stack relatives to
↳ original sequences
74         stacks_refs (list of dict): References of the char in sequence os stack relatives to
75         char positions on original sequences
76         matrix_mode (str): If "SUBST" it's a substitution matrix, if not it's a weight matrix
77         with the keys
78             i = position of first sequence in stack
79             j = position of second sequence on stack
80             pos_i = coordinate of char on first sequence
81             pos_j = coordinate of char on second sequence
82             and the value is the weight to score this position
83             if not match, the score is 0.
84     """
85
86     def __init__(self, sequences, mode="ALIGN", score_match=SCORE_MATCH,
↳ score_no_match=SCORE_NO_MATCH,\
87                 score_gap_ini=SCORE_GAP_INI, score_gap_cont=SCORE_GAP_CONT, subst_matrix={})
↳ ):
88         """Init parameters of alignment"""
89         self.set_sequences(sequences)
90         self.set_stacks()
91         self.len_seq0 = len(self.sequences[0])
92         self.len_seq1 = len(self.sequences[1])
93         self.init_stores()
94         self.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
95         self.set_mode(mode)
96         self.score = 0
97         self.matches = 0
98         self.unmatches = 0
99         self.gaps = 0
100         self.align_seq0 = ""
101         self.align_seq1 = ""
102         self.matching = ""

```

```

103     self.ini_time = 0
104     self.finish_time = 0
105     self.set_subst_matrix(subst_matrix)
106     self.set_matrix_mode()
107
108     def init_stores(self):
109         """Init dictionary that store temp data of the alignment"""
110         self.score_store = {}
111         self.matches_store = {}
112         self.gaps_store = {}
113         self.max_score_index = (0, 0, 0)
114         self.max_score = 0
115         self.forward_arrow = {}
116
117     def set_sequences(self, sequences):
118         """Update the target sequences of the alignment"""
119         self.sequences = sequences
120
121     def set_stacks(self, stack_0=[], stack_1=[],\
122                   stack_0_indexes=[], stack_1_indexes=[], stack_0_refs=[], stack_1_refs=[]):
123         """Update the stacks for msa"""
124         self.stacks = [stack_0, stack_1]
125         self.stacks_indexes = [stack_0_indexes, stack_1_indexes]
126         self.stacks_refs = [stack_0_refs, stack_1_refs]
127
128     def set_matrix_mode(self, mode="SUBST"):
129         """Update matrix mode"""
130         self.matrix_mode = mode
131
132     def set_subst_matrix(self, subst_matrix={}):
133         """Update the score matrix"""
134         self.subst_matrix = subst_matrix
135
136     def set_scores(self, score_match=SCORE_MATCH, score_no_match=SCORE_NO_MATCH,\
137                   score_gap_ini=SCORE_GAP_INI, score_gap_cont=SCORE_GAP_CONT):
138         """Update the weight scores of the alignment"""
139         self.score_match = score_match
140         self.score_no_match = score_no_match
141         self.score_gap_ini = score_gap_ini
142         self.score_gap_cont = score_gap_cont
143
144     def set_mode(self, mode="ALIGN"):
145         """Set computation mode"""
146         self.mode = mode
147
148     def forward_track(self, index):
149         """Calc alignments in forward direction.
150
151         The alignment strings are calculated from init cell (0,0) in global
152         alignments or maximum score cell in local alignments.
153

```



```

154         In local mode it's necessary to extend the alignments (local) to the total
155     → length of
156         the sequences to show the location of the alignment, and in order to compare with
157         BioPython outputs.
158
159     Args:
160         index (tuple of int): Cell coordinates of the starting cell
161
162     Returns:
163         string: align sequence 0 (bottom) for printing purposes
164         string: align sequence 1 (top) for printing purposes
165         tuple of int: Coordinates of the last cell
166
167     """
168     ret_align_seq0, ret_align_seq1 = "", ""
169     (i, j, gap_ini) = index
170     ret_final_pos = (self.len_seq0, self.len_seq1)
171     while i < self.len_seq0 or j < self.len_seq1:
172         if self.mode == "LOCAL" and self.score_store[(i, j, gap_ini)] == 0:
173             ret_final_pos = (i, j)
174             break
175         arrow = self.forward_arrow[(i, j, gap_ini)]
176         if arrow == "d":
177             ret_align_seq0 += self.sequences[0][i]
178             ret_align_seq1 += self.sequences[1][j]
179             i, j, gap_ini = i + 1, j + 1, 1
180         elif arrow == "h":
181             ret_align_seq0 += "-"
182             ret_align_seq1 += self.sequences[1][j]
183             i, j, gap_ini = i, j + 1, 0
184         elif arrow == "v":
185             ret_align_seq0 += self.sequences[0][i]
186             ret_align_seq1 += "-"
187             i, j, gap_ini = i + 1, j, 0
188     #compute the complete align in local mode
189     if self.mode == "LOCAL":
190         ret_align_seq0 = self.sequences[0][0:index[0]] + \
191             ret_align_seq0 + self.sequences[0][ret_final_pos[0]:]
192         ret_align_seq1 = self.sequences[1][0:index[1]] + \
193             ret_align_seq1 + self.sequences[1][ret_final_pos[1]:]
194         diff_pos_ini = index[1] - index[0]
195         if diff_pos_ini > 0:
196             ret_align_seq0 = '-' * diff_pos_ini + ret_align_seq0
197         else:
198             ret_align_seq1 = '-' * -diff_pos_ini + ret_align_seq1
199         diff_len = len(ret_align_seq1) - len(ret_align_seq0)
200         if diff_len > 0:
201             ret_align_seq0 += '-' * diff_len
202         else:
203             ret_align_seq1 += '-' * -diff_len
204     return ret_align_seq0, ret_align_seq1, ret_final_pos
205
206 def calc_matching(self, align_seq0, align_seq1, ini_pos=(), final_pos=()):

```

```

205 """Calc matching string
206
207 The matching string is the string line to print between the top and
208 bottom alignment strings. It contains the match (!), no match (.) and
209 gap ( ) indicators.
210
211 Args:
212 align_seq0 (string): Bottom sequence
213 align_seq1 (string): Top sequence
214 ini_pos (tuple of int): Initial cell coordinates
215 final_pos (tuple of int): Final cell coordinates
216
217 Returns:
218 string: Matching string
219
220 """
221 count = 0
222 ret_matching = ""
223 diff_pos_ini = ini_pos[1] - ini_pos[0]
224 if diff_pos_ini > 0:
225     delta_pos = diff_pos_ini
226 else:
227     delta_pos = 0
228 for n, (i, j) in enumerate(zip(align_seq0, align_seq1)):
229     if self.mode == "LOCAL" and not (n >= ini_pos[0] + delta_pos and n <
230     ↪ final_pos[0] + delta_pos):
231         ret_matching += ' '
232     else:
233         if i == j: ret_matching += '|'
234         elif i != j and i != '-' and j != '-': ret_matching += '.'
235         else: ret_matching += ' '
236     count += 1
237 return ret_matching
238
239 def store(self, key, score, matches, gaps):
240     """Store info related to a computed cell
241     The maximum score is computed having into account the number of matches, if there are
242     most than one solution. If the score are equal, the path with more matches is
243     ↪ selected.
244     Args:
245     key (tuple of int): Cell coordinates
246     score (int): Cell score
247     matches (int): Cell matches
248     gaps (int): Cell gaps
249     """
250     self.score_store[key] = score
251     super_score = score * COMPAC + 10 * matches
252     if super_score > self.max_score:
253         self.max_score_index = key
254         self.max_score = super_score
255     self.matches_store[key] = matches
256     self.gaps_store[key] = gaps

```

```

255 def calc_score_binary(self, seq_0, seq_1, i, j, seq_0_index=0, seq_1_index=1, pos_0=0,
256 ↪ pos_1=0):
257     """Compute alignment scores for two sequences
258     If there are a substitution matrix (actually dictionary) defined,
259     the scores are computed from the dictionary.
260     Args:
261         seq_0 (int): Sequence 0
262         seq_1 (int): Sequence 1
263         i (int): Sequence 0 char index
264         j (int): Sequence 1 char index
265         seq_0_index (int): Sequence 0 index on original sequences (MSA)
266         seq_1_index (int): Sequence 1 index on original sequences (MSA)
267         pos_0 (int): Sequence 0 index on stack 0
268         pos_1 (int): Sequence 1 index on stack 0
269     """
270     if self.subst_matrix:
271         if self.matrix_mode == "SUBST":
272             #print("PAIR", i, j, seq_0[i], seq_1[j])
273             subst_matrix_index = (seq_0[i], seq_1[j])
274             subst_matrix_index_swap = (seq_1[j], seq_0[i])
275             if subst_matrix_index in self.subst_matrix:
276                 matrix_score = self.subst_matrix[subst_matrix_index]
277             elif subst_matrix_index_swap in self.subst_matrix:
278                 matrix_score = self.subst_matrix[subst_matrix_index_swap]
279             else: #weight matrix
280                 if pos_0 in self.stacks_refs and i in self.stacks_refs[pos_0]:
281                     i_orig = self.stacks_refs[pos_0][i]
282                 else:
283                     i_orig = i
284                 if pos_1 in self.stacks_refs and j in self.stacks_refs[pos_1]:
285                     j_orig = self.stacks_refs[pos_1][j]
286                 else:
287                     j_orig = j
288                 if i_orig in self.subst_matrix[seq_0_index][seq_1_index] and \
289                     j_orig in self.subst_matrix[seq_0_index][seq_1_index][i_orig]:
290                     matrix_score = self.subst_matrix[seq_0_index][seq_1_index][i_orig][j_ori
291
292             else:
293                 matrix_score = 0
294
295             # Gaps in almost one of the sequences. This case only arises in MSA
296             # There is no matrix related entry. If matrix is a weight matrix we compute
297             # as zero (as defined in T-Coffee)
298             if seq_0[i] == "-" or seq_1[j] == '-':
299                 inc_matches = 0
300                 if self.subst_matrix:
301                     if self.matrix_mode == "SUBST":
302                         inc_score = self.score_gap_cont
303                     else:
304                         inc_score = 0
305                 else:
306                     inc_score = self.score_gap_cont

```

g]

```

305     else:
306         if seq_0[i] == seq_1[j]:
307             if self.subst_matrix:
308                 inc_score = matrix_score
309             else:
310                 inc_score = self.score_match
311             inc_matches = 1
312         else:
313             if self.subst_matrix:
314                 inc_score = matrix_score
315             else:
316                 inc_score = self.score_no_match
317             inc_matches = 0
318
319     return inc_score, inc_matches
320
321 def calc_score(self, i, j):
322     """Compute alignment scores.
323     If there are stacks associated with the sequence, we compute the score weighing the
324     scores of the stacks (SOP: Score of Pairs). Stacks contains also the guiding
325     ↪ sequences.
326     Args:
327         i (int): Sequence 0 index
328         j (int): Sequence 1 index
329     """
330     if self.stacks == [[], []]:
331         return self.calc_score_binary(self.sequences[0], self.sequences[1], i, j, 0, 1,
332         ↪ 0, 0)
333     else:
334         computed_score = 0
335         computed_matches = 0
336         nvalues = 0
337         for pos_0, (seq_0, index_0) in enumerate(zip(self.stacks[0],
338         ↪ self.stacks_indexes[0])):
339             for pos_1, (seq_1, index_1) in enumerate(zip(self.stacks[1],
340             ↪ self.stacks_indexes[1])):
341                 score, matches = self.calc_score_binary(seq_0, seq_1, i, j, index_0,
342                 ↪ index_1, pos_0, pos_1)
343                 computed_score += score
344                 computed_matches += matches
345                 nvalues += 1
346         ret_score = computed_score / nvalues
347         ret_matches = computed_matches / nvalues
348         return ret_score, ret_matches
349
350 def align(self, i=0, j=0, ini_gap=1):
351     """Recursive align of sequences
352     For each cell, which coordinates are (i, j, ini_gap), calc the maximum score path
353     ↪ from
354     three alternative displacements:
355     """

```

1) To $(i + 1, j + 1, 1)$, that is, matching or no matching the $seq0(i)$ and $seq1(i)$ characters.
 This is a diagonal displacement.
 2) To $(i, j + 1, 0)$, that is, setting a gap in $seq0$ and advance $seq1$. Horizontal displacement.
 3) To $(i + 1, j, 0)$, that is, setting a gap in $seq1$ and advance $seq0$. Vertical displacement.

The scores of these displacements are calculated adding the score of the target cells (that are computed recursively) and the matrix, default of gap scores in each case.

The score, matches, gaps and forward_arrow are stored at related dictionary entry based on coordinates (i, j, ini_gap) , all of them associated to the maximum score of the three possible paths starting from the cell, avoiding recomputation of the cell if it's called from another recursive path.

Each cell has a third score coordinate, because a cell could be called from a cell with yet has a gap (only from horizontal or vertical prior displacement) or from a cell with has a match/no match.

Then we need to store two scores, matches, gaps and forward_arrows related to the two possible cell incarnations at coordinates $(i, j, 0)$ and $(i, j, 1)$.

We store matches and gaps in order to have one additional criterion to tiebreaker if some of the scores are equal. We are using this approach in local alignment computation. If two scores are equal we choose the solution with the greatest number of matches.

We store the displacement directions in forward_arrow dict to compute the alignment. It's possible to avoid this, using only the score information, but we have let this approach as proof of concept and for clarity in the algorithm.

In this scenario we observe that the differences between the global, local and long substring algorithms are minimal.

Local algorithm:

Starting from the global algorithm, which would be the most general, the local algorithm only changes two aspects:

1. Rejection of the roads with negative values of the score, equaling these values to 0, that is, not letting previous alignments of poor quality affect the final result.
2. Use as cell of beginning of the alignment the one with the highest scores. In our implementation we also take into account the number of matches, as we have already mentioned.

390 Finally, but outside the algorithm of alignment itself (at forward_track and
 391 → matching methods)
 392 it only remains to extend the alignment obtained to show its location within the
 393 → chains to be aligned.

394 Search algorithm of the long common substring:

395 Modify the global algorithm in the following aspects:
 396 1. Only computes matches between characters or gaps in one or another
 397 → initial sequence.

398 Args:

399 i (int): Sequence 0 index
 400 j (int): Sequence 1 index
 401 ini_gap (int): 1 if gap initiation, 0 if gap continuation

```
402 """
403 score_diag, score_hor, score_ver = MIN, MIN, MIN
404 matches_diag, matches_hor, matches_ver = MIN, MIN, MIN
405 gaps_diag, gaps_hor, gaps_ver = MIN, MIN, MIN
406 #align and advance seq0 and seq1
407 #in long_substring mode only matches are processed
408 if i < self.len_seq0 and j < self.len_seq1 and\
409 (self.mode != "LONG_SUBSTRING" or self.sequences[0][i] == self.sequences[1][j]):
410     inc_score, inc_matches = self.calc_score(i, j)
411     key = (i + 1, j + 1, 1)
412     if key in self.score_store:
413         score_diag, matches_diag, gaps_diag = \
414             self.score_store[key] + inc_score, self.matches_store[key] + inc_matches,
415             → self.gaps_store[key]
416     else:
417         score, matches, gaps = self.align(i + 1, j + 1, 1)
418         self.store(key, score, matches, gaps)
419         score_diag, matches_diag, gaps_diag = score + inc_score, matches +
420         → inc_matches, gaps
421 #don't align and gap in seq0 (advance seq1)
422 if j < self.len_seq1:
423     gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
424     key = (i, j + 1, 0)
425     if key in self.score_store:
426         score_hor, matches_hor, gaps_hor = self.score_store[key] + gap_score,\
427         self.matches_store[key], self.gaps_store[key] + 1
428     else:
429         score, matches, gaps = self.align(i, j + 1, 0)
430         self.store(key, score, matches, gaps)
431         score_hor, matches_hor, gaps_hor = score + gap_score, matches, gaps + 1
432 #don't align and gap in seq1 (advance seq0)
433 if i < self.len_seq0:
434     gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
435     key = (i + 1, j, 0)
436     if key in self.score_store:
437         score_ver, matches_ver, gaps_ver = \
```

```

436         self.score_store[key] + gap_score, self.matches_store[key],
437         ↪ self.gaps_store[key] + 1
438     else:
439         score, matches, gaps = self.align(i + 1, j, 0)
440         self.store(key, score, matches, gaps)
441         score_ver, matches_ver, gaps_ver = score + gap_score, matches, gaps + 1
442     #choose the high score path
443     matcher_diag, matcher_hor, matcher_ver = score_diag, score_hor, score_ver
444     if i < self.len_seq0 or j < self.len_seq1:
445         if self.mode == "LOCAL" and matcher_diag < 0 and matcher_hor < 0 and matcher_ver
446         ↪ < 0:
447             score_diag, score_hor, score_ver = 0, 0, 0
448             #matcher_diag, matcher_hor, matcher_ver = 0, 0, 0
449         if matcher_diag > matcher_hor and matcher_diag > matcher_ver:
450             ret_score, ret_matches, ret_gaps, ret_arrow = \
451             score_diag, matches_diag, gaps_diag, "d"
452         elif matcher_hor > matcher_ver:
453             ret_score, ret_matches, ret_gaps, ret_arrow = \
454             score_hor, matches_hor, gaps_hor, "h"
455         else:
456             ret_score, ret_matches, ret_gaps, ret_arrow = \
457             score_ver, matches_ver, gaps_ver, "v"
458     else:
459         ret_score, ret_matches, ret_gaps, ret_arrow = \
460         0, 0, 0, ""
461     self.forward_arrow[(i, j, ini_gap)] = ret_arrow
462     if i == 0 and j == 0:
463         self.store((0, 0, 1), ret_score, ret_matches, ret_gaps)
464         if self.mode in ["GLOBAL", "LONG_SUBSTRING"]: self.max_score_index = (0, 0, 1)
465         else: ret_score = self.max_score // COMPAC
466         ret_matches = self.matches_store[self.max_score_index]
467         ret_gaps = self.gaps_store[self.max_score_index]
468
469     return ret_score, ret_matches, ret_gaps
470
471 def compute(self, mode="LOCAL", silent=False):
472     """Calc alignment
473     Args:
474     mode (str): Type of algorithm (local, global or long substring)
475     silent (bool): If true don't show alignment output
476     """
477     self.ini_time = time.time()
478     self.init_stores()
479     self.set_mode(mode)
480     self.score, self.matches, self.gaps = self.align()
481     self.align_seq0, self.align_seq1, final_pos = self.forward_track(self.max_score_index)
482
483     self.matching = self.calc_matching(self.align_seq0, self.align_seq1,
484     ↪ self.max_score_index, final_pos)
485     self.unmatches = self.matching.count('.')
486     self.gaps = self.matching.count(' ')
487     self.finish_time = time.time()

```

x)

```

484         if not silent:
485             self.view()
486
487     def get_len_long_common_substring(self):
488         """Getter for the len of the common substring
489         That is equal to the number of matches of the alignment
490         """
491         return self.matches
492
493     def get_long_common_substring(self):
494         """Returns the longest common substring
495         whitout alignment (positional) information
496         """
497         long_common_substring = ""
498         for (char, match_char) in zip(self.align_seq1, self.matching):
499             if match_char == '|':
500                 long_common_substring += char
501         return long_common_substring
502
503     def view(self):
504         """Prints the alignment data"""
505         #unmatches = self.matching.count('.')
506         #gaps = self.matching.count(' ')
507         if self.matching:
508             gap_groups = self.matching.count('| ') + self.matching.count('. ') +
509                 ↪ self.matching[0].count(' ')
510         else:
511             gap_groups = 0
512         print(" ")
513         if self.mode == "LOCAL":
514             print("### AlignSequences. Local alignment (Smith-Waterman)")
515         elif self.mode == "LONG_SUBSTRING":
516             print("### AlignSequences. Long substring finder")
517         else:
518             print("### AlignSequences. Global alignment (Needleman-Wunsch)")
519         if self.subst_matrix:
520             print("\tUsing score matrix with matrix mode", self.matrix_mode)
521         print(self.align_seq1)
522         print(self.matching)
523         print(self.align_seq0)
524         print("\tScore:", self.score)
525         print("\tSimilarity (wo gaps):", self.matches / (self.matches + self.unmatches))
526         print("\tDistance (wo gaps):", self.unmatches / (self.matches + self.unmatches))
527         print("\tDistance:", self.unmatches / (self.matches + self.unmatches + self.gaps))
528         print("\tInit index:", self.max_score_index)
529         print("\tMatches:", self.matches, " Unmatches:", self.unmatches, " Gaps:",
530             ↪ self.gaps, " Gap groups:", gap_groups)
531         #simple scoring verification todo: apply to matrix
532         if not self.subst_matrix:
533             print("\tScore verified:", self.matches * self.score_match + self.unmatches *
534                 ↪ self.score_no_match \
535                 + self.gaps * self.score_gap_cont + gap_groups * self.score_gap_ini)

```



```

533     print("\tFinish. Execution milliseconds:", round((self.finish_time - self.ini_time)
534           ↪ * 1000))
535     print("\tScore Dictionary Size", len(list(self.score_store.keys())))
536
537 def edit_distance(self, score_match=0, score_no_match=-1, score_gap_ini=0,
538     ↪ score_gap_cont=-1):
539     """Calculates an edit distance as requested in questions 1 and 3
540       It's the same computation as a global alignment with -1 penalties applied to
541       score_gap_cont and score_nomatch and 0 in score_match and score_gap_ini
542
543       Args:
544         score_match (int): Score of match characters.
545         score_no_match (int): Score of no match characters.
546         score_gap_ini (int): Score of gap init.
547         score_gap_cont (int): Score of gap continuation.
548     """
549     self.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
550     self.compute("GLOBAL", True)
551     return abs(self.score)

```

1.6 MSA

In these functions, what is necessary to perform a multiple alignment of sequences is developed.

The method of progressive alignment based on a guide tree is used.

The guide tree can be obtained in two alternative ways: **Unweighted Pair Group Method with Arithmetic Mean (UPGMA)** and **Neighbor Join (NJ)**, the same options present in *CLUSTAL* software.

The alignment has three known phases. These are the particularities of my implementation on each phase:

1. Perform pairwise alignments between all the sequences involved and assign them a score.

In the case of **UPGMA** we use the proportion (in percentages) between matches and matches plus no matches (without taking gaps into account). That is, we use a measure of the identity between the two sequences involved. You can also use the distance, which would be the complement to 100 of identity, but we wanted to do so to be able to compare with the information that *CLUSTAL* throws at the beginning of his output. It does not affect the result, we simply have to look for maximum identities to build the guide tree, instead of minimum distances. In the case of **NJ**, we have chosen to use distances, computed also in percentages. Also not taking into account the number of gaps in the denominators.

2. Build the guide tree. As I said, we can do it using UPGMA or NJ. The NJ method generates an unrooted tree. As we need a root, for purposes of the subsequent alignment, we have chosen to root it by clustering the two nodes that have no relation. We do not know if it is the method used by *CLUSTAL* or similar programs, but it seems to work.
3. Multiple alignment. It is, as we know, a progressive alignment following the order indicated by the guide tree. In each step we need to align two groups of sequences, of length $n \geq 1$ and $m \geq 1$. Each group is aligned as a whole, in the sense that the gaps entered in one of the sequences of the group must be introduced in the same positions in the rest of the sequences of their group.

To assign a score to a position, the combined score of all the residuals of that position is used. To do this we produce the Cartesian product $n \times m$ of all the characters of that position and calculate the average of scores:

$$\frac{\sum_{\substack{0 \leq i < n \\ 0 \leq j < m}} \text{matrix}(i, j)}{nm}$$

If in any of the positions we have a gap, we have chosen to penalize it as the sum of penalties assigned to the start of the gap plus gap continuation penalty. It is a criterion, *CLUSTAL* we know that it uses another one.

If we already have a pairwise development, as it was my case, it would be easy to extend it to address MSA?. The answer is affirmative. With slight modifications in the class **AlignSequences**, we have managed to address an MSA, in the following way:

1. Generalize the one-position scoring algorithm to take into account all the sequences of both groups, averaging the scores as indicated above.
2. Take a sequence from each group (the first) to perform a simple pairwise alignment (but with the scores calculated as indicated in 1).
3. Compare the sequences resulting from the pairwise alignment with their originals from each group, compute where the gap is introduced and introduce the gap at the same positions in the rest of the sequences of each group.

1.7 MSA generic methods

Script 1.7.1 (python)

```

1  """This methods shows alternative implementations of multiple sequence alignments, CLUSTAL,
   →  T-COFFEE.
2
3  TODO:
4      * Many more tests. Create a test battery.
5
6      * Achieve that the results obtained are more similar to those of CLUSTAL (if they have
   →  to be).
7      Given the lack of detailed information it will be necessary to resort
8      to the sources (in C++) of CLUSTAL.
9
10     * Allow to configure the initial alignment and the final multialignments with different
   →  parameters.
11
12     * Draw the alignments in a more standard way.
13
14     * Draw the phylogenetic trees.
15
16     * Include all new methods in AlignSequences or in another class.
17  """
18  from ete3 import Tree, TreeStyle
19  MIN_SCORE = 0
20
21  def draw_guide_tree(tree):
22      """
23      Draw guide tree with ETE library
24      """
25      t = Tree(tree + ";")
26      ts = TreeStyle()

```

```

27     ts.show_leaf_name = True
28     ts.show_branch_length = False
29     ts.show_branch_support = False
30     ts.scale = 160
31     ts.branch_vertical_margin = 40
32     print(t)
33     return t, ts
34
35 def readFasta(file):
36     """
37     Reads all sequences of a FASTA file
38     Args:
39         file (str): name of the input FASTA file
40     Returns:
41         dict of str, str: sequences readed
42     """
43     ret_seqs = {}
44     seq = ""
45     key_found = False
46     with open(file, 'r') as f:
47         key = ""
48         for line in f:
49             line = line.replace('\n', '')
50             if len(line) > 0:
51                 if line[0] == ">":
52                     if key_found:
53                         ret_seqs[key] = seq
54                     key_found = True
55                     key = line[1:].split(" ")[0]
56                     seq = ""
57                 elif key_found:
58                     seq += line
59         if key_found:
60             ret_seqs[key] = seq
61     return ret_seqs
62
63 def pairwise_align(s1, s2, matrix, matrix_mode, mode, score_gap_ini=0, score_gap_cont=-8,\
64                   score_match=3, score_no_match=-2):
65     """
66     Performs initial pairwise alignments against the class AlignSequences
67     returning the %identity.
68
69     Args:
70         s1 (str): First sequence to compare.
71         s2 (str): Second sequence to compare.
72         matrix (dict of tuples of int): Substitution matrix, Biopython format
73         matrix_mode (str): Type of matrix
74             'SUBST'          Substitution matrix
75             'WEIGHT'         Weight matrix
76         mode (str): Computation mode:
77             'GLOBAL'         Global Alignment
78             'LOCAL'          Local Alignment

```

```

79         'LONG_SUBSTRING'    Obtain long common substring
80     score_gap_ini (int): Score of gap init.
81     score_gap_cont (int): Score of gap continuation.
82     score_match (int): Score of match characters (used if no matrix informed)
83     score_no_match (int): Score of no match characters (used if no matrix informed)
84
85     Returns:
86         (int): % identity between sequences
87     """
88     align = AlignSequences([s1, s2])
89     align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
90     align.set_subst_matrix(matrix)
91     align.set_matrix_mode(matrix_mode)
92     align.compute(mode.upper(), silent = True)
93     return round((align.matches + 1) * 100 / (align.matches + align.unmatches + 2))
94
95 def pairwise_align_distance(s1, s2, matrix, matrix_mode, mode, score_gap_ini=0,
96 → score_gap_cont=-8):
97     """
98     Performs initial pairwise alignments against the class AlignSequences
99     returning the distance between 0 and 100.
100
101     Args:
102         s1 (str): First sequence to compare.
103         s2 (str): Second sequence to compare.
104         matrix (dict of tuples of int): Substitution matrix, Biopython format.
105         matrix_mode (str): Type of matrix
106             'SUBST'          Substitution matrix
107             'WEIGHT'         Weight matrix
108         mode (str): Computation mode:
109             'GLOBAL'        Global Alignment
110             'LOCAL'         Local Alignment
111             'LONG_SUBSTRING' Obtain long common substring
112         score_gap_ini (int): Score of gap init.
113         score_gap_cont (int): Score of gap continuation.
114         score_match (int): Score of match characters (used if no matrix informed)
115         score_no_match (int): Score of no match characters (used if no matrix informed)
116
117     Returns:
118         (int): distance between sequences
119     """
120     align = AlignSequences([s1, s2])
121     align.set_scores(0, 0, score_gap_ini, score_gap_cont)
122     align.set_subst_matrix(matrix)
123     align.set_matrix_mode(matrix_mode)
124     align.compute(mode.upper(), silent = True)
125     identity = round((align.matches + 1) * 100 / (align.matches + align.unmatches + 2))
126     return 100 - identity
127
128 def guide_tree_UPGMA(sequences, matrix, matrix_mode, mode,\
129 → score_gap_ini, score_gap_cont,\
130 → score_match, score_no_match):

```

```

130 """
131 Performs initial pairwise alignments against the class AlignSequences
132 returning the guide_tree derived from UPGMA method.
133
134 Args:
135     sequences (lit of str): Sequences to align
136     matrix (dict of tuples of int): Substitution matrix, Biopython format.
137     matrix_mode (str): Type of matrix
138         'SUBST'          Substitution matrix
139         'WEIGHT'         Weight matrix
140     mode (str): Computation mode:
141         'GLOBAL'         Global Alignment
142         'LOCAL'          Local Alignment
143         'LONG_SUBSTRING' Obtain long common substring
144     score_gap_ini (int): Score of gap init.
145     score_gap_cont (int): Score of gap continuation.
146     score_match (int): Score of match characters (used if no matrix informed)
147     score_no_match (int): Score of no match characters (used if no matrix informed)
148
149 Returns:
150     (list of 3-tuples of int): guide three, the third position of the tuple contains the
151     → root
152     of the other two nodes.
153     (dict of int, boolean = True): contains all nodes
154
155 """
156 tree = {} #initial tree
157 guide_tree = [] #guided tree, pairs to align in sequence
158 max_score = MIN_SCORE
159 max_score_position = ()
160 for i in range(0, len(sequences)):
161     for j in range(0, i):
162         if (i,j) not in tree:
163             score = pairwise_align(sequences[i], sequences[j], matrix, matrix_mode, mode,\
164                                     score_gap_ini, score_gap_cont, score_match, score_no_match)
165             tree[(i,j)] = score
166             if score >= max_score:
167                 max_score = score
168                 max_score_position = (i,j)
169
170 print(tree)
171 len_tree = len(sequences)
172 guide_tree_nodes = {}
173 # Generate guide tree. At every step we compute another row averaging the
174 # most closer rows and removing all their row coordinates from the tree
175 while len(tree) > 0:
176     (imax, jmax) = max_score_position
177     guide_tree.append((imax, jmax, len_tree))
178     guide_tree_nodes[imax] = True
179     guide_tree_nodes[jmax] = True
180     guide_tree_nodes[len_tree] = False

```

```

181     # Average scores from i,j rows into new row in new_row_pos
182     for j in range(0, len_tree):
183         if j in [imax, jmax]:
184             continue
185         nscores = 0.0;
186         for coordinate in [(imax, j), (j, imax), (jmax, j), (j, jmax)]:
187             if coordinate in tree:
188                 score = tree[coordinate]
189                 nscores += 1
190                 if (len_tree, j) not in tree:
191                     tree[(len_tree, j)] = score
192             else:
193                 tree[(len_tree, j)] += score
194         if nscores > 0:
195             tree[(len_tree, j)] = tree[(len_tree, j)] / nscores
196
197     # Tree cleaning and calc max score
198     max_score = MIN_SCORE
199     max_score_position = ()
200     for i in range(0, len_tree + 1):
201         for j in range(0, len_tree + 1):
202             if (i,j) in tree:
203                 if i == imax or i == jmax or j == imax or j == jmax:
204                     del(tree[(i,j)])
205             else:
206                 if tree[(i,j)] >= max_score:
207                     max_score = tree[(i,j)]
208                     max_score_position = (i,j)
209
210     len_tree += 1
211
212     return guide_tree, guide_tree_nodes
213
214 def q(i, j, nseq, n, dmatrix):
215     """
216     NJ method: calculate element of intermediate Q matrix.
217     """
218     d = (nseq - 2) * dmatrix[(i,j)]
219     for k in range(0, n):
220         if (i,k) in dmatrix:
221             d -= dmatrix[(i,k)]
222         if (j,k) in dmatrix:
223             d -= dmatrix[(j,k)]
224     return d
225
226 def calc_qmatrix(nseq, n, dmatrix):
227     """
228     NJ method: calculate intermediate Q matrix.
229     """
230     qmatrix = {}
231     for (i,j) in dmatrix:
232         qmatrix[(i,j)] = q(i, j, nseq, n, dmatrix)

```

```

233     return qmatrix
234
235 def smallest_q(qmatrix):
236     """
237     NJ method: returns the coordinates of the minimum score in intermediate Q matrix.
238     """
239     sq = ()
240     min_sq = - MIN
241     for key in qmatrix.keys():
242         if qmatrix[key] < min_sq:
243             min_sq = qmatrix[key]
244             sq = key
245     return sq
246
247 def djoin(joined_pair, nseq, n, dmatrix):
248     """
249     NJ method: returns distances of joined nodes to the rooted node, so, it returns the
    → branch lengths
250     """
251     (i, j) = joined_pair
252     d_i_1 = dmatrix[(i,j)] / 2.0
253     d_i_2 = 0
254     for k in range(0, n):
255         if (i,k) in dmatrix:
256             d_i_2 += dmatrix[(i,k)]
257         if (j,k) in dmatrix:
258             d_i_2 -= dmatrix[(j,k)]
259     d_i = d_i_1 - d_i_2 / (2*(nseq - 2))
260     d_j = dmatrix[(i,j)] - d_i
261     return d_i, d_j
262
263 def dnjoin(k, joined_pair, dmatrix):
264     """
265     NJ method: returns distance of sequence k to the new node that routes the joined_pair.
266     The distance is the mean of the distances from k to each of nodes joined.
267     """
268     (i, j) = joined_pair
269     d_k = 0
270     if (i,k) in dmatrix:
271         d_k += dmatrix[(i,k)]
272     if (j,k) in dmatrix:
273         d_k += dmatrix[(j,k)]
274     d_k = (d_k - dmatrix[(i,j)]) / 2.0
275     return d_k
276
277 def recalc_dmatrix(joined_pair, n, dmatrix):
278     """
279     NJ method: recalc distance matrix taking into account the joined pair
280     """
281     (i, j) = joined_pair
282     # Recalculate distances
283     for k in range(0, n):

```

```

284         if (i,k) in dmatrix and (j,k) in dmatrix:
285             dmatrix[(n + 1, k)] = dnjoin(k, joined_pair, dmatrix)
286             dmatrix[(k, n + 1)] = dmatrix[(n + 1, k)]
287     # Remove joined rows from dmatrix
288     for k in range(0, n + 1):
289         for l in range(0, n + 1):
290             if k == i or k == j or l == i or l == j:
291                 if (k, l) in dmatrix:
292                     del(dmatrix[(k, l)])
293     return
294
295 def guide_tree_NJ(sequences, matrix, matrix_mode, mode,\
296                  score_gap_ini, score_gap_cont,\
297                  score_match, score_no_match):
298     """
299     Performs initial pairwise alignments against the class AlignSequences
300     returning the guide_tree derived from NJ method.
301
302     Args:
303         sequences (list of str): Sequences to align
304         matrix (dict of tuples of int): Substitution matrix, Biopython format.
305         matrix_mode (str): Type of matrix
306             'SUBST'          Substitution matrix
307             'WEIGHT'         Weight matrix
308         mode (str): Computation mode:
309             'GLOBAL'         Global Alignment
310             'LOCAL'          Local Alignment
311             'LONG_SUBSTRING' Obtain long common substring
312         score_gap_ini (int): Score of gap init.
313         score_gap_cont (int): Score of gap continuation.
314         score_match (int): Score of match characters (used if no matrix informed)
315         score_no_match (int): Score of no match characters (used if no matrix informed)
316
317     Returns:
318         (list of 3-tuples of int): guide tree, the third position of the tuple contains the
319         ↪ root
320         ↪ of the other two nodes.
321         (dict of int, boolean = True): contains all nodes
322
323     """
324     dmatrix = {} #initial distance matrix
325     n = len(sequences)
326     for i in range(0, n):
327         for j in range(0, i):
328             if (i,j) not in dmatrix:
329                 distance = pairwise_align_distance(sequences[i], sequences[j], matrix,
330                 ↪ matrix_mode,\
331                 ↪ mode, score_gap_ini, score_gap_cont)
332                 dmatrix[(i,j)] = distance
333                 dmatrix[(j,i)] = distance
334
335     nseq = n
336     new_nodes = n - 2

```



```

334 guide_tree = [] #guided tree, pairs to align in sequence
335 guide_tree_nodes = {} #guided tree rooted nodes to complete
336 for i in range(0, n):
337     guide_tree_nodes[i] = False
338 while new_nodes > 0:
339     qmatrix = calc_qmatrix(nseq, n, dmatrix)
340     (joined_i, joined_j) = smallest_q(qmatrix)
341     #print("JOIN:", (joined_i, joined_j))
342     guide_tree.append((joined_i, joined_j, n + 1))
343     guide_tree_nodes[joined_i] = True
344     guide_tree_nodes[joined_j] = True
345     guide_tree_nodes[n + 1] = False
346     recalc_dmatrix((joined_i, joined_j), n, dmatrix)
347     n += 1
348     nseq -= 1
349     new_nodes -= 1
350 # Root the tree
351 #print("DMATRIX:", dmatrix)
352 rooting_tuple = []
353 for node in guide_tree_nodes:
354     if not guide_tree_nodes[node]:
355         rooting_tuple.append(node)
356 rooting_tuple.append(n + 1)
357 guide_tree_nodes[n + 1] = True
358 #print("Rooting tuple:", rooting_tuple)
359 if len(rooting_tuple) == 3:
360     guide_tree.append(tuple(rooting_tuple))
361 assert len(rooting_tuple) == 3
362 return guide_tree, guide_tree_nodes
363
364 def gapeator(a, a_gapped, b_stack, b_stack_refs):
365     """
366     Introduces gaps in all the sequences of b_stack taking into account the positions
367     and the gaps introduced in sequence a to obtain sequence a_gapped
368     Args:
369         a (str): template sequence not gapped
370         a_gapped (str): template sequence gapped
371         b_stack (list of str): stack of b sequences ungapped
372         b_stack_refs (list of dict): stack of references to original positions
373     Returns:
374         list of str: stack b gapped as a does
375         list of dict: stack b coordinates referred to original sequence
376     """
377     ini_a_gapped = a_gapped
378     b_gapped_stack = []
379     b_references_stack = []
380     len_a_gapped = len(a_gapped)
381     for b, b_refs in zip(b_stack, b_stack_refs):
382         b_gapped = ""
383         b_gapped_references = {}
384         a_gapped = ini_a_gapped
385         base_ref = 0

```

```

386     for k, (i, j) in enumerate(zip(a, b)):
387         index = a_gapped.index(i)
388         a_gapped = a_gapped[index + 1:]
389         #print("a_gapped", a_gapped )
390         #print(i, j, index)
391         b_gapped += "-" * index + j
392         if k in b_refs:
393             b_gapped_references[base_ref + k + index] = b_refs[k]
394             base_ref += index
395             #print("b_gapped", b_gapped )
396         b_gapped += b[k+1:]
397         remaining_gaps = "-" * (len_a_gapped - len(b_gapped))
398         b_gapped += remaining_gaps
399         b_gapped_stack.append(b_gapped)
400         b_references_stack.append(b_gapped_references)
401     return b_gapped_stack, b_references_stack
402
403 def pairwise_align_msa_step(stack_0, stack_1, sequences, matrix, matrix_mode,\
404                             mode, stack_0_indexes, stack_1_indexes, stack_0_refs,
405                             ↪ stack_1_refs,\
406                             score_match, score_no_match, score_gap_ini, score_gap_cont):
407     """
408     Performs msa alignment of sequence stack 0 and 1.
409
410     Args:
411         stack_0 (list of str): First stack of sequences to align.
412         stack_1 (list of str): Second stack of sequences to align.
413         sequences (list of str): Sequences to align.
414         matrix (dict of tuples of int): Substitution matrix, Biopython format.
415         matrix_mode (str): Type of matrix
416             'SUBST'          Substitution matrix
417             'WEIGHT'         Weight matrix
418         mode (str): Computation mode:
419             'GLOBAL'        Global Alignment
420             'LOCAL'         Local Alignment
421             'LONG_SUBSTRING' Obtain long common substring
422         stack_0_indexes(list of int) : Indexes of initial sequences related to stack
423 ↪ sequences 0
424         stack_1_indexes(list of int) : Indexes of initial sequences related to stack
425 ↪ sequences 1
426         stack_0_refs(list of dict) : stack_0 references to original sequences
427         stack_1_refs(list of dict) : stack_1 references to original sequences
428         score_match (int): Score of match characters (used if no matrix informed)
429         score_no_match (int): Score of no match characters (used if no matrix informed)
430         score_gap_ini (int): Score of gap init.
431         score_gap_cont (int): Score of gap continuation.
432
433     Returns:
434         list of str: stack_0 gapped (with the gaps necessary for the alignment)
435         list of str: stack_1 gapped (with the gaps necessary for the alignment)
436         list of dict: stack_0 references to original sequences
437         list of dict: stack_1 references to original sequences

```

```

435 """
436 align = AlignSequences([stack_0[0], stack_1[0]])
437 align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
438 align.set_subst_matrix(matrix)
439 align.set_matrix_mode(matrix_mode)
440 align.set_stacks(stack_0, stack_1, stack_0_indexes, stack_1_indexes, stack_0_refs,
441 ↪ stack_1_refs)
441 align.compute(mode.upper(), silent = True)
442 # align_seq0 align_seq1 are the seq0 and seq1 alignments
443 # we need to deduce the rest of alignments.
444 # what we do is perform the same gap insertions, if any, as the first sequence of the
445 ↪ stacks
446 # the gap insertions where performed taken into account the initial sequence
447 # to compute the references to initial sequence in order to employ a weight matrix if
448 ↪ informed
447 stack_0_gapped, stack_0_references = gapeator(stack_0[0], align.align_seq0, stack_0,
448 ↪ stack_0_refs)
448 stack_1_gapped, stack_1_references = gapeator(stack_1[0], align.align_seq1, stack_1,
449 ↪ stack_1_refs)
449 return stack_0_gapped, stack_1_gapped, stack_0_references, stack_1_references
450
451 def get_name(index, sequence_names):
452 """
453 Obtain sequence name from index
454 """
455 name = ""
456 if index < len(sequence_names):
457     name = sequence_names[index]
458 else:
459     name = str(index)
460 return name
461
462 def to_newick(tree, sequence_names):
463 """
464 Obtain guide tree in newick format
465 """
466 # Change format to intermediate roots
467 roots = {}
468 newick_tree = ""
469 for branch in tree:
470     (i, j, k) = branch
471     name_i = get_name(i, sequence_names)
472     name_j = get_name(j, sequence_names)
473     name_k = get_name(k, sequence_names)
474     if name_i in roots:
475         new_root_i = roots[name_i]
476     else:
477         new_root_i = name_i
478     if name_j in roots:
479         new_root_j = roots[name_j]
480     else:
481         new_root_j = name_j

```

```

482     roots[name_k] = [new_root_i, new_root_j]
483
484     for root in roots.values():
485         s_root = str(root)
486         if len(s_root) > len(newick_tree):
487             newick_tree = s_root.replace("[", "(").replace("]", ")").replace("'", "")
488
489     return newick_tree

```

1.8 T-COFFEE methods

Script 1.8.1 (python)

```

1  # T-COFFEE specific methods
2  def pairwise_align_coffee(s1, s2, matrix, mode, score_gap_ini=0, score_gap_cont=-8, \
3      score_match=3, score_no_match=-2):
4      """
5      Performs initial pairwise alignments against the class AlignSequences
6      returning the %identity and the alignments to construct the primary library
7
8      Args:
9          s1 (str): First sequence to compare.
10         s2 (str): Second sequence to compare.
11         matrix (dict of tuples of int): Substitution matrix, Biopython format.
12         mode (str): Computation mode:
13             'GLOBAL'          Global Alignment
14             'LOCAL'           Local Alignment
15             'LONG_SUBSTRING'  Obtain long common substring
16         score_gap_ini (int): Score of gap init.
17         score_gap_cont (int): Score of gap continuation.
18         score_match (int): Score of match characters (used if no matrix informed)
19         score_no_match (int): Score of no match characters (used if no matrix informed)
20
21     Returns:
22         int: % identity between sequences
23         str: sequence 1 aligned
24         str: sequence 2 aligned
25     """
26     align = AlignSequences([s1, s2])
27     align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
28     align.set_subst_matrix(matrix)
29     align.compute(mode.upper(), silent = True)
30     return round((align.matches + 1) * 100 / (align.matches + align.unmatches + 2)), \
31         align.align_seq0, align.align_seq1
32
33 def get_pos(k, seq_i, align_i):
34     """
35     Obtain position of a character in the original sequence given the
36     position in the alignment(k), the original sequence (seq_i)
37     and the align_i (gapped) sequence

```

```

38     """
39     char = align_i[k]
40     count_char = align_i[0:k+1].count(char)
41     index = -1;
42     for _ in range(0, count_char):
43         index = seq_i.find(char, index + 1)
44     return index
45
46 def update_weight_at_pos(weight_library, i, j, pos_i, pos_j, identity):
47     """
48     Update weight at pos i , j, pos_i, pos_j
49     """
50     if i not in weight_library:
51         weight_library[i] = {}
52     w_i = weight_library[i]
53     if j not in w_i:
54         w_i[j] = {}
55     w_i_j = w_i[j]
56     if pos_i not in w_i_j:
57         w_i_j[pos_i] = {}
58     w_i_j_pi = w_i_j[pos_i]
59     if pos_j not in w_i_j_pi:
60         w_i_j_pi[pos_j] = identity
61     else:
62         w_i_j_pi[pos_j] += identity
63
64 def update_weight_library(weight_library, i, j, identity,\
65                           seq_i, seq_j, align_i, align_j):
66     """
67     Update weights library from alignments and %identity
68     """
69     for k, (c_i, c_j) in enumerate(zip(align_i, align_j)):
70         if c_i != "-" and c_j != "-":
71             pos_i = get_pos(k, seq_i, align_i)
72             pos_j = get_pos(k, seq_j, align_j)
73             #print("Position:", pos_i, pos_j)
74             update_weight_at_pos(weight_library, i, j, pos_i, pos_j, identity)
75             update_weight_at_pos(weight_library, j, i, pos_j, pos_i, identity)
76
77 def compute_library(sequences, matrix={}, weight_library={}, mode="GLOBAL",\
78                   score_gap_ini=0, score_gap_cont=-8,\
79                   score_match=3, score_no_match=-2):
80     """
81     Compute initial library of identities based on scores of PA
82     Args:
83         sequences (list of str): Sequences to compare.
84         matrix (dict of tuples of int): Substitution matrix, Biopython format.
85         mode (str): Computation mode:
86             'GLOBAL'          Global Alignment
87             'LOCAL'           Local Alignment
88             'LONG_SUBSTRING'   Long substring alignment
89         score_gap_ini (int): Score of gap init.

```

```

90     score_gap_cont (int): Score of gap continuation.
91     score_match (int): Score of match characters (used if no matrix informed)
92     score_no_match (int): Score of no match characters (used if no matrix informed)
93
94     Returns:
95         list of str: primary library of alignments
96     """
97     primary_library = {}
98     for i in range(0, len(sequences)):
99         for j in range(0, i):
100             if (i,j) not in primary_library:
101                 identity, align_i, align_j = pairwise_align_coffee(sequences[i],
102                     → sequences[j],\
103                         matrix, mode, score_gap_ini, score_gap_cont,\
104                         score_match, score_no_match)
105                 update_weight_library(weight_library, i, j, identity,\
106                     sequences[i], sequences[j], align_i, align_j)
107                 primary_library[(i,j)] = (align_i, align_j, identity)
108     #print(weight_library)
109     return primary_library
110
111 def extend_weights(weight_library, i, k, j):
112     """
113     Extend weights for pair of sequences (i,j) at pos (pos_i_posj)
114     taken into account the routes using k as
115     intermediate, by means of the alignments (i, k) and (k, j).
116     """
117     for pos_i, pos_i_j in weight_library[i][j].items():
118         for pos_j in pos_i_j.keys():
119             if pos_j in weight_library[j][k].keys():
120                 for pos_k in weight_library[j][k][pos_j].keys():
121                     if pos_k in weight_library[k][i].keys():
122                         for pos_i_new in weight_library[k][i][pos_k].keys():
123                             if pos_i_new == pos_i:
124                                 #print("Extension", pos_i, pos_j, pos_k, weight_library[i][k]
125                                     → ][pos_i][pos_k],weight_library[j][k][pos_j][pos_k])
126                                 m = min(\
127                                     weight_library[i][k][pos_i][pos_k],\
128                                     weight_library[j][k][pos_j][pos_k])
129                                 #print("++", m, weight_library[i][j][pos_i][pos_j])
130                                 weight_library[i][j][pos_i][pos_j] += m
131
132 def extend_library_weights(sequences, weight_library):
133     """
134     Extend library for all triplets of sequences
135     Taken into account the simetry i -> k -> j
136     """
137     len_sequences = len(sequences)
138     for i in range(0, len_sequences):
139         for k in range(0, len_sequences):
140             if k != i:
141                 for j in range(0, len_sequences):

```

```

140         if j != k and j != i:
141             #print("Triplet:", i, k, j)
142             extend_weights(weight_library, i, k, j)
143
144
145 def compute_libraries(sequences, matrix,\
146                     score_gap_ini, score_gap_cont, score_match, score_no_match):
147     """
148     Compute initial library of identities based on scores of pairwise alignments
149     Args:
150         sequences (list of str): Sequences to compare.
151         matrix (dict of tuples of int): Substitution matrix, Biopython format.
152         score_gap_ini (int): Score of gap init.
153         score_gap_cont (int): Score of gap continuation.
154         score_match (int): Score of match characters (used if no matrix informed)
155         score_no_match (int): Score of no match characters (used if no matrix informed)
156
157     Returns:
158         dict : primary library of alignments
159         dict : weight matrix
160     """
161     weight_library = {}
162     primary_library = compute_library(sequences, matrix,\
163                                     weight_library, "GLOBAL", score_gap_ini, score_gap_cont,\
164                                     score_match, score_no_match)
165
166     _ = compute_library(sequences, matrix,\
167                       weight_library, "LOCAL", score_gap_ini, score_gap_cont,\
168                       score_match, score_no_match)
169
170     # _ = compute_library(sequences, matrix,\
171     #                   weight_library, "LONG_SUBSTRING", score_gap_ini, score_gap_cont,\
172     #                   score_match, score_no_match)
173
174     extend_library_weights(sequences, weight_library)
175
176     return primary_library, weight_library

```

1.9 Main MSA method

Script 1.9.1 (python)

```

1 # Generic MSA method
2 def do_msa_from_fasta(file, main_alg="CLUSTAL", method="NJ", matrix={}, matrix_mode="SUBST",\
3                       mode="GLOBAL", score_gap_ini=-10, score_gap_cont=-5, score_match=3,\
4                       score_no_match=-2, verbose=False):
5     """
6     Performs MSA alignments from fasta file
7     Args:
8         file (str): Name of the FASTA file.

```

```

9      main_alg (str): Main algorithm:
10          "CLUSTAL"      Clustal like
11          "T-COFFEE"     T-COFFEE like
12      method (str): NJ neighbor join / UPGMA
13      matrix (dict of tuples of int): Substitution matrix, Biopython format.
14      matrix_mode (str): Type of matrix
15          'SUBST'        Substitution matrix
16          'WEIGHT'        Weight matrix
17      mode (str): Computation mode:
18          'GLOBAL'        Global Alignment
19          'LOCAL'         Local Alignment
20          'LONG_SUBSTRING' Obtain long common substring
21      score_gap_ini (int): Score of gap init.
22      score_gap_cont (int): Score of gap continuation.
23      score_match (int): Score of match characters (used if no matrix informed)
24      score_no_match (int): Score of no match characters (used if no matrix informed)
25      verbose (bool): If True prints verbose info
26
27      Returns:
28          list of 3-tuples of int: guide three, the third position of the tuple contains the
→      root
29          of the other two nodes.
30      list of str: alignments
31      list of str: sequence_names
32      list of int: sequence indexes relating strings in alignment to original sequences
33      """
34      seq_fasta = readFasta(file)
35      sequences = list(seq_fasta.values())
36      sequence_names = list(seq_fasta.keys())
37      return do_msa(sequences, sequence_names,\
38                  main_alg, method, matrix, matrix_mode,\
39                  mode, score_gap_ini, score_gap_cont, score_match,\
40                  score_no_match, verbose)
41
42  def do_msa(sequences, sequence_names, main_alg="CLUSTAL", method="NJ", matrix={},
→      matrix_mode="SUBST",\
43          mode="GLOBAL", score_gap_ini=-10, score_gap_cont=-5, score_match=3,\
44          score_no_match=-2, verbose=False):
45      """
46      Performs MSA alignments from sequences
47      Args:
48          sequences (list of str): Sequences
49          sequence_names (list of str): Names of sequences
50          main_alg (str): Main algorithm:
51              "CLUSTAL"      Clustal like
52              "T-COFFEE"     T-COFFEE like
53          method (str): NJ neighbor join / UPGMA
54          matrix (dict of tuples of int): Substitution matrix, Biopython format.
55          matrix_mode (str): Type of matrix
56              'SUBST'        Substitution matrix
57              'WEIGHT'        Weight matrix
58          mode (str): Computation mode:

```



```

59         'GLOBAL'           Global Alignment
60         'LOCAL'           Local Alignment
61         'LONG_SUBSTRING'   Obtain long common substring
62     score_gap_ini (int): Score of gap init.
63     score_gap_cont (int): Score of gap continuation.
64     score_match (int): Score of match characters (used if no matrix informed)
65     score_no_match (int): Score of no match characters (used if no matrix informed)
66     verbose (bool): If True prints verbose info
67
68     Returns:
69         list of 3-tuples of int: guide three, the third position of the tuple contains the
→     root
70         of the other two nodes.
71     list of str: alignments
72     list of str: sequence_names
73     list of int: sequence indexes relating strings in alignment to original sequences
74     """
75     if main_alg == "T-COFFEE":
76         primary_library, weight_library = compute_libraries(sequences, matrix,\
77             score_gap_ini, score_gap_cont, score_match, score_no_match)
78         #print("Primary library", primary_library)
79         #print("Weight library", weight_library)
80         # Matrix mode and other MSA parameters
81         matrix = weight_library
82         matrix_mode = "WEIGHT"
83     #     score_gap_ini = 0
84     #     score_gap_cont = 0
85     # From here only we need is to compute a MSA with weight matrix as reference.
86     else:
87         matrix_mode = "SUBST"
88     if method == "NJ":
89         guide_tree, guide_tree_nodes = \
90             guide_tree_NJ(sequences, matrix, matrix_mode,\
91                 mode, score_gap_ini, score_gap_cont,\
92                 score_match, score_no_match)
93     else: #UPGMA
94         guide_tree, guide_tree_nodes = \
95             guide_tree_UPGMA(sequences, matrix, matrix_mode,\
96                 mode, score_gap_ini, score_gap_cont,\
97                 score_match, score_no_match)
98     sequences_store = {}
99     sequences_store_indexes = {}
100     sequences_store_refs = {}
101     print("Guide Tree",guide_tree, sequence_names)
102     #return guide_tree, "", sequence_names
103     # Create MSA
104     for i in guide_tree_nodes.keys():
105         if i < len(sequences):
106             sequences_store[i] = [sequences[i]]
107             sequences_store_indexes[i] = [i]
108             sequences_store_refs[i] = []
109             autorefs = {}

```

```

110         for k in range(0, len(sequences[i])):
111             autorefs[k] = k
112             sequences_store_refs[i].append(autorefs)
113
114     if verbose: print(sequences_store)
115     for (i ,j ,k) in guide_tree:
116         stack_i = sequences_store[i]
117         stack_j = sequences_store[j]
118         stack_i_indexes = sequences_store_indexes[i]
119         stack_j_indexes = sequences_store_indexes[j]
120         stack_i_references = sequences_store_refs[i]
121         stack_j_references = sequences_store_refs[j]
122         if verbose: print("Stack i", i, stack_i)
123         if verbose: print("Stack j", j, stack_j)
124         if verbose: print("Stack i_indexes", i, stack_i_indexes)
125         if verbose: print("Stack j_indexes", j, stack_j_indexes)
126         stack_0, stack_1, stack_0_references, stack_1_references = \
127             pairwise_align_msa_step(stack_i, stack_j, sequences, matrix, matrix_mode, \
128                                     mode, stack_i_indexes, stack_j_indexes, \
129                                     stack_i_references, stack_j_references, \
130                                     score_match, score_no_match, score_gap_ini,
131                                     ↪ score_gap_cont)
132
133         sequences_store[k] = []
134         sequences_store_indexes[k] = []
135         sequences_store_refs[k] = []
136         if verbose: print("=====")
137         for s in stack_0:
138             if verbose: print(s)
139             sequences_store[k].append(s)
140         for s in stack_1:
141             if verbose: print(s)
142             sequences_store[k].append(s)
143         for seq_index in stack_i_indexes:
144             sequences_store_indexes[k].append(seq_index)
145         for seq_index in stack_j_indexes:
146             sequences_store_indexes[k].append(seq_index)
147         for seq_refs in stack_0_references:
148             sequences_store_refs[k].append(seq_refs)
149         for seq_refs in stack_1_references:
150             sequences_store_refs[k].append(seq_refs)
151         if verbose: print("=====")
152         if verbose: print("Sequences store indexes",sequences_store_indexes[k])
153         if verbose: print("Sequences store references",sequences_store_refs[k])
154         if verbose: print("New stack:", k, sequences_store[k])
155     alignment = sequences_store[k]
156     newick_tree = to_newick(guide_tree, sequence_names)
157     return newick_tree, alignment, sequence_names, sequences_store_indexes[k]
158
159 def score(alignment, matrix, score_gap_ini=0, score_gap_cont=0):
160     """
161     Score based on sum of pair scores (SOP) taking into account substitution matrix
162     Derived from the objective score of MUSCLE refinement stage

```

```

161     """
162     msa_score = 0
163     for k in range(0, len(alignment[0])): #columns of msa
164         score_column_k = 0
165         nvalues = 0
166         for i in range(0, len(alignment)):
167             for j in range(i + 1, len(alignment)):
168                 if alignment[i][k] == "-" and alignment[j][k] != "-":
169                     score_column_k += score_gap_cont
170                     if k == 0 or alignment[i][k-1] != "-":
171                         score_column_k += score_gap_ini
172                 if alignment[j][k] == "-" and alignment[i][k] != "-":
173                     score_column_k += score_gap_cont
174                     if k == 0 or alignment[j][k-1] != "-":
175                         score_column_k += score_gap_ini
176                 elif (alignment[i][k], alignment[j][k]) in matrix:
177                     score_column_k += matrix[(alignment[i][k], alignment[j][k])]
178                     nvalues += 1
179                 elif (alignment[j][k], alignment[i][k]) in matrix:
180                     score_column_k += matrix[(alignment[j][k], alignment[i][k])]
181                     nvalues += 1
182             if nvalues > 0:
183                 #score += score_column_k / nvalues
184                 msa_score += score_column_k
185     return msa_score
186
187 def score_from_fasta(file, matrix, score_gap_ini=0, score_gap_cont=0):
188     seq_fasta = readFasta(file)
189     sequences = list(seq_fasta.values())
190     return score(sequences, matrix, score_gap_ini, score_gap_cont)

```

1.10 T-COFFEE RAG1_AA

Script 1.10.1 (python)

```

1 matrix = MatrixInfo.blosum80
2 file = "files/RAG1_AA.fas"
3 quality_score_gap_ini = -10
4 quality_score_gap_cont = -5
5
6 guide_tree_upgma, align, sequence_names, indexes = do_msa_from_fasta(file,\
7     main_alg = "T-COFFEE", method = "UPGMA", \
8     matrix = matrix, matrix_mode = "SUBST",\
9     mode = "GLOBAL", score_gap_ini = -10,\
10    score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
11
12 print("# Guide Tree:", guide_tree_upgma)
13 t_upgma, ts_upgma = draw_guide_tree(guide_tree_upgma)
14 print("# Alignment:")
15 for i, s in enumerate(align):

```

```

16     print(">" + sequence_names[indexes[i]])
17     print(s)
18 print()
19 print("Score", score(align, MatrixInfo.blosum62, quality_score_gap_ini,
    → quality_score_gap_cont))
20 print()
21
22 guide_tree_nj, align, sequence_names, indexes = do_msa_from_fasta(file,\
23     main_alg = "T-COFFEE", method = "NJ", \
24     matrix = matrix, matrix_mode = "SUBST",\
25     mode = "GLOBAL", score_gap_ini = -10,\
26     score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
27 print("# Guide Tree:", guide_tree_nj)
28 t_nj, ts_nj = draw_guide_tree(guide_tree_nj)
29 print("# Alignment:")
30 for i, s in enumerate(align):
31     print(">" + sequence_names[indexes[i]])
32     print(s)
33 print()
34 print("Score", score(align, matrix, quality_score_gap_ini, quality_score_gap_cont))

```

Output

```

{(1, 0): 72, (2, 0): 75, (2, 1): 79, (3, 0): 74, (3, 1): 77, (3, 2): 81, (4, 0): 73, (4, 1):
→ 79, (4, 2): 82, (4, 3): 80, (5, 0): 73, (5, 1): 79, (5, 2): 81, (5, 3): 79, (5, 4): 85,
→ (6, 0): 73, (6, 1): 78, (6, 2): 82, (6, 3): 81, (6, 4): 88, (6, 5): 87, (7, 0): 73, (7,
→ 1): 79, (7, 2): 83, (7, 3): 82, (7, 4): 87, (7, 5): 91, (7, 6): 93, (8, 0): 75, (8, 1):
→ 79, (8, 2): 83, (8, 3): 82, (8, 4): 87, (8, 5): 87, (8, 6): 89, (8, 7): 91, (9, 0): 74,
→ (9, 1): 80, (9, 2): 83, (9, 3): 84, (9, 4): 88, (9, 5): 87, (9, 6): 89, (9, 7): 91, (9,
→ 8): 93}
Guide Tree [(9, 8, 10), (7, 6, 11), (11, 10, 12), (12, 5, 13), (13, 4, 14), (14, 2, 15), (15,
→ 3, 16), (16, 1, 17), (17, 0, 18)] ['Oncorhynchus_mykiss', 'Carcharhinus_leucas',
→ 'Latimeria_menadoensis', 'Protopterus_dolloi', 'Alytes_obstetricans',
→ 'Anolis_carolinensis', 'Gallus_gallus', 'Alligator_mississippiensis',
→ 'Ornithorhynchus_anatinus', 'Homo_sapiens']
# Guide Tree: (((((((Alligator_mississippiensis, Gallus_gallus), (Homo_sapiens,
→ Ornithorhynchus_anatinus)), Anolis_carolinensis), Alytes_obstetricans),
→ Latimeria_menadoensis), Protopterus_dolloi), Carcharhinus_leucas), Oncorhynchus_mykiss)

          /-Alligator_mississippiensis
        /-|
        |  \-Gallus_gallus
      /-|
      |  |  /-Homo_sapiens
    /-|  \-|
    |  |  \-Ornithorhynchus_anatinus
  /-|  |
  |  |  \-Anolis_carolinensis
/-|  |
|  |  \-Alytes_obstetricans
/-|  |

```

```

    | | \-Latimeria_menadoensis
  /-| |
| | \-Protopterus_dolloi
--| |
| | \-Carcharhinus_leucas
|
| \-Oncorhynchus_mykiss
# Alignment:
>Alligator_mississippiensis
RTVKAVTGRQIFQPLHALRTAEKALLPGYHSFEWKPLKNVSANTEVGIIIDGLSGLPHTVDDYPIDTIAKFRYDAALVSALMDMEEDILEGM
↪ KAHDLDYYLNG-PFTVVVKESCDGMGDVSEKHGCGPAVPEKAVRFSFTVMTIAI--THGNTNVRIFEEVKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERETMKNSVLLLEMGGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVLHSITRSHTEN
↪ LERYEVWRSNPYHESVDELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERKRWQSALDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSFNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Gallus_gallus
RTVKAVTGRQIFQPLHALRTAEKALLPGYHPFEWKPLKNVSTNTEVGIIIDGLSGLPLSIDDYPIDTIAKFRYDTALVSALKDMEEEEILEGM
↪ KAKNLDDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNIAI--DHENERIRIFEEVKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKNSSELLLEIGGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVHFSITRSHAEN
↪ LERYEIWRSNPYHESVDELDRVRKGVSAKPF IETVPSIDALHCDIGNATEFYRIFQMEIGEYVKNPDATKEERKRWQLTLDKHLRKKMK
↪ LKPMMRMSGNFARKLMSKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSYNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Homo_sapiens
RTVKAITGRQIFQPLHALRNAEKVLLPGYHHFEWQPPLKNVSSSTDVGIIDGLSGLSSSVDDYPVDTIAKFRYDSALVSALMDMEEDILEGM
↪ RSQDLDDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPVPEKAVRFSFTIMKITI--AHSSQNVKVFEAKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKSELMLLEGGILRTFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVHFSITRSHAEN
↪ LERYEVWRSNPYHESVEELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQLEIGEYVKNPNASKEERKRWQATLDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVDAVCELIPSEERHEALRELMMDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Ornithorhynchus_anatinus
RTVKAITGRQIFQPLHSLRTAEKALLPGYHPFEWDPPLKNVSANTEVGIMDGLSGLPVSVDDYPVDTIAKFRYDAALVSALMDMEEDILEGM
↪ KSQDLDDYYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNITL--AYEQENVKIFEEAKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKDSELKLEMGGILRSFRFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVLHSITRSHAEN
↪ LERYEVWRSNPFHESVEELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQLEIGEAYKNPNASKEERKRWQATLDKHLRKKMK
↪ LKPIMRMNGNFARKLMTKETVEAVCELVHCEERHEALRELMMDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Anolis_carolinensis
RTVKAVTGRQIFQPLHALRTAEKALLPGYHQFEWKPLKNVSSNTEVGIIIDGLSGIQLVDDYPVDTIAKFRYDAALASALMDMEEDILEGL
↪ KRQDLDDYFYG-PFTVVIKESCDGMGDVSEKHGCGPAVPEKAVRFSFTLMTISV--THGNASIRVFEECKPNSELCKPLCLMLADESD
↪ HETLTAILSPLVAERAMKDSVLILDMAGIPRTFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLILHSITRSHAEN
↪ LERYELWRNPNYHETVDELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERRRQSTLDKHLRKKMN
↪ LKPMTRMNGNFARKLMTKETVEAVCELIKSEERHEALRELMMDLYLKMKPVWRSSCPTKECPPELVQYSFNSQRFAELLATKFRYRYAGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Alytes_obstetricans
RTVKATTGRQIFQPLHALRNAEKALLPGYHPFEWKPLKNVSTCTDTGIIDGLSGLNRSIDEYPVEAISKFRYDTALVSALKDMEEDILEGL
↪ RSHDMDDYYLNG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMYISV--PNNNECVRIFDETKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKTSELMLEMGGILRNFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDSTRLEASQNLVHFSITRCHTEN
↪ LQRYETWRANPHESVDELDRVRKGVSAKPF IETLPSIDALHCDIGNAAEFYRLFQLEIGEYVKNPNATKEERKRWQSTLDKHLRKKMN
↪ LKPIMRMNGNFARKLMSKETVEAVCELVHSEERQEILRELMMDLYLKMKPVWRSSCPAKECEPELLYQYSFHSQRFAELLSTKFKYRYAGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Latimeria_menadoensis

```

```

RTVKATTGKQIFQPLHSLRNAEKALLPGFHPFEWQPPLKNVSSTTEVGIIDGMSGMTQFVDEYPLDTISKFRYDAALVSALKDLEEELKGL
↳ IEEDLEDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRYSFTIMTISV-ANSHNENVTFEEGKPNSELCKPLCLMLADESD
↳ HETLTAILGPVIAEREAAMKNSSELFLEMGILRSFKFIFRGTYDEKLIRDVEGLEASGSSYICTLCDSTRSEASQNFILHSITRSHKEN
↳ LERYEIWRSNPYQEPVEELRDRVKGVSAPFIETLPSIDALHCDIGNATEFYKIFQDEIGEIKYKNPNPSREEKKRWSVLDKHLRKNMN
↳ LKPVMRMNGNYARKLMTKETVNAVCELIPSEERQEALKELVDLYLKMKPVWRSTCPAKECEPELLCQYSFHSQRFAELLSTMYRYRYEGK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKYIELEDVL
>Protopterus_dolloi
RTVKAATGRQIFQPLHALRSAEKALLPGYHPFEWQPPLVGSSSTDVGIINGLSGLTSSVDEYPVEALAKFRYDAALVSALKDIEENILEGM
↳ KQNGLEDYLSG-PFTVVIKESCDGMGDVSEKHGSGPPVPEKAVRFSFTIMTISV-AMSDSENVQIFEEFKPNSELCKPLCLMLADESD
↳ HETLTAILGPVIAEREAAMKTSLELMLELGGILRTFKFFFRGTGYDEKLVRVEGLEASGSHYICTLCDATRQEASRNVLHSITRSHAEN
↳ LERYEVWRSNPYNESVDLDRVKGVSAPFIETRPCIDALQCDIGNATEFYKIFQDEVGEVYKRNPNPSKEDRKRWHMTLDKHLRKKLS
↳ LKPVMRMNGNFARKLITKEAVDAVCELIPSEERRAAIRDLVHLYMLMKPVWRSTYPAKECEPELLCQYSFNSQRFAELLSTKFQRYRYEAK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYELEDVL
>Carcharhinus_leucas
KTVKAITGKQIFQPLHALRNAEKTLLPGYYSFEWQPPLANISTNTRVGIIDGLSGWVQCVDYPMETISRRLSYDVALASAVKEMEDDILEGL
↳ RSQNVDEFVSG-PFTVVIKESCDGMGDVSEKHGCGPTVPEKAVRYSFTIMTISV-MNENNEKVKVFEEMKPNSELCCRPLCLMLADESD
↳ RETLTAILGPVIAERQSMKTSDLIVEIGDLYRSFQIFRGTYDEKLVRVEGLEASGSIYICTLCDSTRSEASKNMVLHSITRNHAEN
↳ LERYEIWRSNPYHETADELRDRVKGVSAPFIETQPSIDALHCDIGNATEFYRIFQDEIGEYKNSNSKEERKRWQSMLDKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVEAVCELIPSEERREILRELMHLYLLMKPVWRSTFPTTECPDLLCQYSFNSQRFAELLHTEFSHRYEGK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKSYELEDIL
>Oncorhynchus_mykiss
RTVKATSGRQIFQPLHLTLRTAEKELLPGYHPFEWQPALKSVSTSCHVGIIDGLSGWIASVDDSPADTVTRFRYDVALVSALKDLEEDIMEGL
↳ RERGLEDSACTSGFSVMIKESCDGMGDVSEKHGGPPVPEKPVRFSTIMSVSIQAEGEDEAITIFREPKNSEMSCKPLSLMFVDES
↳ HETLTGVLGPVVAERNAMKHSRLILSVGLSRSFRFHRGTGYDEKVMREMEGLEASGSTYICTLCDSTRAEASQNMVLSVTRSHDEN
↳ LERYELWRTNPHSESAEELRDRVKGVSAPFMETQPTLDALHCDIGNATEFYKIFQDEIGEYVHKANPSREQRSSWRAALDKQLRKKMK
↳ LKPVMRMNGNYARKLMTREAVEAVCELCSEERQEALRELMGLYIQMKPVWRSTCPAKECPDELCRYFSNSQRFAELLSTVFKYRYDGK
↳ ITNYLHKTLAHVPEIVERDGSIGAWASEGNESGNKLFRRFRKMNRQSKTFELEDVL

```

Score 99825

```

Guide Tree [(7, 6, 11), (11, 5, 12), (9, 8, 13), (12, 4, 14), (13, 3, 15), (14, 2, 16), (16,
↳ 1, 17), (15, 0, 18), (17, 18, 19)] ['Oncorhynchus_mykiss', 'Carcharhinus_leucas',
↳ 'Latimeria_menadoensis', 'Protopterus_dolloi', 'Alytes_obstetricans',
↳ 'Anolis_carolinensis', 'Gallus_gallus', 'Alligator_mississippiensis',
↳ 'Ornithorhynchus_anatinus', 'Homo_sapiens']
# Guide Tree: ((((((Alligator_mississippiensis, Gallus_gallus), Anolis_carolinensis),
↳ Alytes_obstetricans), Latimeria_menadoensis), Carcharhinus_leucas), (((Homo_sapiens,
↳ Ornithorhynchus_anatinus), Protopterus_dolloi), Oncorhynchus_mykiss))

```

```

                /-Alligator_mississippiensis
              /-|
            /-|  \-Gallus_gallus
          |  |
        /-|  \-Anolis_carolinensis
      |  |
    /-|  \-Alytes_obstetricans
  |  |
/-|  \-Latimeria_menadoensis
|  |
|  \-Carcharhinus_leucas
--|
    /-Homo_sapiens

```

```

|      /-|
|      /-|      \-Ornithorhynchus_anatinus
|      |      |
|      \-|      \-Protopterus_dolloi
|
|      \-Oncorhynchus_mykiss
# Alignment:
>Alligator_mississippiensis
RTVKAVTGRQIFQPLHALRTAEKALLPGYHSFEWKPLKNVSANTEVGIIDGLSGLPHTVDDYPIDTIAKFRYDAALVSALMDMEEDILEGM
↪ KAHDLDYDLNG-PFTVVVKESCDGMGDVSEKHGCGPAVPEKAVRFSFTVMTIAI--THGNTNVRIFEVVKPNSSELCKPLCLMLADESD
↪ HETLTAILSPLIAERETMKNSVLLLEMGGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVLHSITRSHTEN
↪ LERYEVWRSNPYHESVDELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERKRWQSALDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECPELLCQYSFNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Gallus_gallus
RTVKAVTGRQIFQPLHALRTAEKALLPGYHPFEWKPLKNVSTNTEVGIIDGLSGLPLSIDDYPIDTIAKFRYDTALVSALKDMEEEILEGM
↪ KAKNLDDYDLNG-PFTVVVKECCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNIAI--DHENERIRIFEVVKPNSSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKNSSELLLEIGGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVHFSITRSHAEN
↪ LERYEIWRSNPYHESVDELDRVRKGVSAKPF IETVPSIDALHCDIGNATEFYRIFQMEIGEYVKNPDATKEERKRWQLTLDKHLRKKMK
↪ LKPMRMMSGNFARKLMSKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECPELLCQYSYNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Anolis_carolinensis
RTVKAVTGRQIFQPLHALRTAEKALLPGYHQFEWKPLKNVSSNTEVGIIDGLSGIQLHVDYDPIAKFRYDAALASALMDMEEDILEGL
↪ KRQDLDDYDFKG-PFTVVIKESCDGMGDVSEKHGCGPAVPEKAVRFSFTLMTISV--THGNASIRVFEECKPNSSELCKPLCLMLADESD
↪ HETLTAILSPLVAERAMKDSVLILDMAGIPRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLILHSITRSHAEN
↪ LERYELWRTNPYHETVDELDRVRKGVSAKPF IETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERRRQSTLDKHLRKKMN
↪ LKPMTRMNGNFARKLMTKETVEAVCELIKSEERHEALRELMMDLYLKMKPVWRSSCPTKECPPELVQYSFNSQRFAELLATKFRYRYAGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Alytes_obstetricans
RTVKATTGRQIFQPLHALRNAEKALLPGYHPFEWKPLKNVSTCTDTGIIDGLSGLNRSIDYEPVEAISKFRYDTALVSALKDMEEDILEGL
↪ RSHDMDYDLNG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMTISV--PNNNECVRIFDETKPNSSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKTSSELMEGILRNFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDSTRLEASQNLVHFSITRCHTEN
↪ LQRYETWRANPHHESVDELDRVRKGVSAKPF IETLPSIDALHCDIGNAAEFYRLFQLEIGEYVKNPNATKEERKRWQSTLDKHLRKKMN
↪ LKPIMRMNGNFARKLMSKETVEAVCELVHSEERQEILRELMMDLYLKMKPVWRSSCPAKECPELLYQYSFHSQRFAELLSTKFKYRYAGK
↪ ITNYFHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Latimeria_menadoensis
RTVKATTGKQIFQPLHSLRNAEKALLPGFHPFEWQPPLKNVSSSTTEVGIIDGMSGMTQFVDEYPLDTISKFRYDAALVSALKDLEEELLKGL
↪ IEEDLEDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRYSFTIMTISV-ANSHNENVTIFEEGKPNSELCKPLCLMLADESD
↪ HETLTAILGPVIAERAMKNSSELFLEMGGILRSFKFIFRGTYDEKLIRDVEGLEASGSSYICTLCDSTRSEASQNFILHSITRSHKEN
↪ LERYEIWRSNPYQEPVEELDRVRKGVSAKPF IETLPSIDALHCDIGNATEFYKIFQDEIGEYVKNPNPSREEKKRWHSVLDKHLRKNMN
↪ LKPVMRMNGNYARKLMTKETVNAVCELPSEERQEALKELVDLYLKMKPVWRSTCPAKECPELLCQYSFHSQRFAELLSTMYRYRYEGK
↪ ITNYLHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKYYELEDVL
>Carcharhinus_leucas
KTVKAITGKQIFQPLHALRNAEKTLLPGYYSFEWQPPLANISTNTRVGIIDGLSGWVCVDDYPMETISRRLSYDVALASAVKEMEDDILEGL
↪ RSQNVDEFVSG-PFTVVIKESCDGMGDVSEKHGCGPTVPEKAVRYSFTIMTISV-MNENNEKVKVFEEMKPNSELCCRPLCLMLADESD
↪ RETLTAILGPVIAERQSMKTSDLIVEIGDLRSFQFIFRGTYDEKLVREVEGLEASGSYIYICTLCDSTRSEASKNMVLHSITRNHAEN
↪ LERYEIWRSNPYHETADELDRVRKGVSAKPF IETQPSIDALHCDIGNATEFYRIFQDEIGEYVKNPNSSKEERKRWQSMKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELPSEERREILRELMHLYLLMKPVWRSTFPTTECPDLLCQYSFNSQRFAELLHTEFSHRYEGK
↪ ITNYLHKTLAHVPEI IERDGSIGAWASEGNESGNKLFRRFRKMNRQSKSYELEDIL
>Homo_sapiens

```



```

RTVKAITGRQIFQPLHALRNAEKVLLPGYHHFEWQPPLKNVSSSTDVGIIDGLSGLSSSVDDYPVDITIAKFRYDSALVSALMDMEEDILEGM
↳ RSQDLDDYLSG-PFTVVVKESCDGMGDVSEKHGSGPVPEKAVRFSFTIMKITI--AHSSQNVKVFEEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKSSLMLELGGILRTFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVHFSITRSHAEN
↳ LERYEVWRSNPYHESVEELRDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEYKPNPNASKEERKRWQATLDKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVDVAVCELIPSEERHEALREMLDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Ornithorhynchus_anatinus
RTVKAITGRQIFQPLHSLRTAEKVLLPGYHPFEWDPPLKNVSANTEVGIMDGLSGLPVSVDDYPVDITIAKFRYDAALVSALMDMEEDILEGM
↳ KSQDLDDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNITL--AYEQENVKIFEEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKDSELKLEMGILRSFRFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVLHSITRSHAEN
↳ LERYEVWRSNPFHESVEELRDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEAYKPNPNASKEERKRWQATLDKHLRKKMK
↳ LKPIMRMNGNFARKLMTKETVEAVCELVHCEERHEALREMLDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Protopterus_dolloi
RTVKAATGRQIFQPLHALRSAEKALLPGYHPFEWQPPLVGVSSSTDVGIINGLSGLTSSVDEYPVEALAKFRYDAALVSALKDIEENILEGM
↳ KQNGLDEYLSG-PFTVVIKESCDGMGDVSEKHGSGPPVPEKAVRFSFTIMSISV-AMSDSENVQIFEEFKPNSELCKPLCLMLADESD
↳ HETLTIVILGPVIAERAMKTSELMLELGGILRTFKFFFRGTGYDEKLVREVEGLEASGSHYICTLCDATRQEASQNLVLHSITRSHAEN
↳ LERYEVWRSNPNYESVDELDRVKGVSAPFIIETRPCIDALQCDIGNATEFYKIFQDEVGEVYKPNPNASKEDRKRWHMTLDKHLRKKLS
↳ LKPVMRMNGNFARKLITKEAVDAVCELIPSEERRAAIRDLVHLYMLMKPVWRSTYPAKECPPELLCQYSFNSQRFAELLSTKFQYRYEAK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYELEDVL
>Oncorhynchus_mykiss
RTVKATSGRQIFQPLHLTLRTAEKELLPGYHPFEWQPALKSVSTSCHVGIIDGLSGWIASVDDSPADTVTRRFYDVALVSALKDLEEDIMEGL
↳ RERGLEDSACTSGFSVMIKESCDGMGDVSEKHGGPPVPEKPVRFSTIMSVSIIAQEGEDEAITIFREPKNSEMCKPLSLMFVDES
↳ HETLTGVLGPVVAERNAMKHSRLILSVGLSRSFRFHFRGTGYDEKVMREMEGLEASGSTYICTLCDSTRAEASQNMTLHSVTRSHDEN
↳ LERYELWRTNPHSESAEELRDRVKGVSAPFMETQPTLDALHCDIGNATEFYKIFQDEIGEYVHKANPSREQRWSRAALDKQLRKKMK
↳ LKPVMRMNGNYARKLMTREAVEAVCELCSEERQEALREMLGLYIQMKPVWRSTCPAKECPDELCRYFSFNSQRFAELLSTVFKYRYDYGK
↳ ITNYLHKTLAHVPEIVERDGSIGAWASEGNESGNKLFRRFRKMNRQSKTFELEDVL

```

Score 107036

1.11 CLUSTAL

Script 1.11.1 (python)

```

1 matrix = {}
2 file = "files/16S.fas"
3 quality_score_gap_ini = -10
4 quality_score_gap_cont = -5
5 guide_tree_upgma, align, sequence_names, indexes = do_msa_from_fasta(file,\
6     main_alg = "CLUSTAL", method = "UPGMA", \
7     matrix = matrix, matrix_mode = "SUBST",\
8     mode = "GLOBAL", score_gap_ini = -10,\
9     score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
10 print("# Guide Tree:", guide_tree_upgma)
11 t_upgma, ts_upgma = draw_guide_tree(guide_tree_upgma)
12 print("# Alignment:")
13 for i, s in enumerate(align):
14     print(">" + sequence_names[indexes[i]])
15     print(s)
16 print()

```



```

17 #print("Score", score(aligned, MatrixInfo.blosum62, quality_score_gap_ini,
    ↪ quality_score_gap_cont))
18
19 # print()
20 # guide_tree_nj, aligned, sequence_names, indexes = do_msa_from_fasta(file, \
21 #     main_alg = "CLUSTAL", method = "NJ", \
22 #     matrix = matrix, matrix_mode = "SUBST", \
23 #     mode = "GLOBAL", score_gap_ini = -10, \
24 #     score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
25 # print("# Guide Tree:", guide_tree_nj)
26 # t_nj, ts_nj = draw_guide_tree(guide_tree_nj)
27 # print("# Alignment:")
28 # for i, s in enumerate(aligned):
29 #     print(">" + sequence_names[indexes[i]])
30 #     print(s)
31 # print()
32 # print("Score", score(aligned, matrix, quality_score_gap_ini, quality_score_gap_cont))

```

Output

```

{(1, 0): 70, (2, 0): 71, (2, 1): 72, (3, 0): 69, (3, 1): 68, (3, 2): 69, (4, 0): 71, (4, 1):
↪ 71, (4, 2): 70, (4, 3): 69, (5, 0): 68, (5, 1): 64, (5, 2): 67, (5, 3): 64, (5, 4): 65,
↪ (6, 0): 66, (6, 1): 65, (6, 2): 67, (6, 3): 65, (6, 4): 67, (6, 5): 67, (7, 0): 67, (7,
↪ 1): 66, (7, 2): 68, (7, 3): 65, (7, 4): 66, (7, 5): 67, (7, 6): 69, (8, 0): 69, (8, 1):
↪ 65, (8, 2): 67, (8, 3): 67, (8, 4): 67, (8, 5): 70, (8, 6): 67, (8, 7): 67, (9, 0): 66,
↪ (9, 1): 65, (9, 2): 68, (9, 3): 64, (9, 4): 67, (9, 5): 67, (9, 6): 67, (9, 7): 66, (9,
↪ 8): 71}
Guide Tree [(2, 1, 10), (9, 8, 11), (4, 0, 12), (12, 10, 13), (7, 6, 14), (13, 3, 15), (11,
↪ 5, 16), (16, 14, 17), (17, 15, 18)] ['Carcharhinus_leucas', 'Oncorhynchus_mykiss',
↪ 'Latimeria_menadoensis', 'Protopterus_dolloi', 'Alytes_obstetricans',
↪ 'Anolis_carolinensis', 'Alligator_mississippiensis', 'Gallus_gallus',
↪ 'Ornithorhynchus_anatinus', 'Homo_sapiens']
# Guide Tree: (((Homo_sapiens, Ornithorhynchus_anatinus), Anolis_carolinensis),
↪ (Gallus_gallus, Alligator_mississippiensis)), (((Alytes_obstetricans,
↪ Carcharhinus_leucas), (Latimeria_menadoensis, Oncorhynchus_mykiss)), Protopterus_dolloi))

        /-Homo_sapiens
        /-|
        /-| \-Ornithorhynchus_anatinus
        | |
        /-| \-Anolis_carolinensis
        | |
        | | /-Gallus_gallus
        | \-|
        | \-Alligator_mississippiensis
--|
        | /-Alytes_obstetricans
        | /-|
        | | \-Carcharhinus_leucas
        | /-|
        | | | /-Latimeria_menadoensis

```

```

\ - |    \ - |
|        \-Oncorhynchus_mykiss
|
\ -Protopterus_dolloi

# Alignment:
>Homo_sapiens
GCT-----AAACCTAGCCCCAAACCCACTCCACCTTA-----CTACCAGACAACCTTA-----GCCAAACCA
↳ TTTAC----CCAAATAAAGTATAGGCGATAGAAATTG-AAACCTGGCGCAATAGA----TATAGTACCGCAAGGGAAAGA-TGAAAAAT
↳ TA-----TAACCAAGCATAATATAGCAAGGACTAACCCCTATACCTTCTGCATAATGAATTAAGTAGAAATAAAGTTGCAAGG
↳ AG-AGCCAAAGCTAAGACCCCCGAAA-CCAGACGAGCTACCTAAGAACAGCTAA--AAGAGCACACCCGTCTATGTAGCAAAATAGTG
↳ GGAAGATTTATAGGTAGAGGCGACAAACCTACCGAGCCTGGTGATAGCTGGTTGTCCAAGATA-GAATCTTAGTTCAACTTTAAATTTG
↳ CCC---ACAGAACCCTCTA-AATCCCCT-----TGTAATTTAACTGTTAGTCCAAAGAGGAACAGCTCTTTGGACACTAGG
↳ AAAAAACCTTGTAGAG-AGAGTAAAAAATTTA-----ACACCCATAGTAGGCCTAAAAAGCAGCCCAATTAAGAAA-G
↳ CGTTCAAGCTCAAC--ACCCACTACCTAAAAATCCCAACATA-TAACTGAACTCCTCAC-ACCCAATTGGACCAATCTATCACCCCTA
↳ TAGAAGAACTAATGTTAGTATAAGT-AACATGAAAAC-----ATTCTCCTCCGCATAAGCCTGCGTCAGATTAAAAACTGAACT
↳ GACAATTAACAGCCCAATATCTAC-----AATCAACCAACAAGTCATTATTACC-CTCACTGTC-AACCCAACA
↳ CAGGCATGCTCATAAGGAAAGGTTAAAAAAGTAAAGGAACTCGGCAAA--TCTTACCCCGCCTGTTTACCAAAAAACATCACCTCTAG
↳ CA----TCACCAGTATTAGAGGCACCGCTGCCAGTGACAC---ATGTTTAAACGGCCGCGGTAC-CCTAACCGTGCAAAGGTAGCATA
↳ ATCACTTGTTCTCTAAATAGGGACCTGTATGAATGGCTCCACGAGGGTTCAGCTGTCTCTTACTTTTAAACAGTGAAATTGACCTGCC
↳ GTGAAGAGGCGGGCATAACACAGCAAGACGAGAAGACCCTATGGAGCTTTAATTTATTAATGCAACAG-----TACCT
↳ ACAAACCCACAGTCTAACTACCAAACCTGCATTAAAAATTCGGTTGGGGCGACCTCGGAGCAGAACCCAACTCCGA--GCAGT
↳ ACA-----TG-CTAAGACTTCACCAGTCAAAGCGAACTACTATACT---CAATTGATCCAA-----TAACCTGACC
↳ AACGGAACAAGTTACCTAGGGATAACAGCGCAATCCTATTCTAGAGTCCATATCAACAATAGGGTTTACGACCTCGATGTTGGATCAG
↳ GACATCCCGATGTTGTCAGCCGCTATTAAGGTTTCGTTTGTTCACGATTAA-AGTCCTACGTGATCTGAGTTTCAGACCGGAGTAATCCA
↳ GGTGCGTTTCTATCTACNTTCAAA-TTCCTCCCTGTAC-GAAAGGACAAGAGAAATAAGGCCTACTTCACAAAGCGCCTTCC-----C
↳ CCGTAAATGATATCATCTCAACTTAGTATTATACCCACACCCACC-----CAAGAACAGGGTTT---

>Ornithorhynchus_anatinus
GCC-----AACTCTAGCCCAAATTAATCTTAATACACATATAAA-----CTCTAATAA-----ACTAAAACA
↳ TTTTAAAGCTAACCCCTAGTATAGGAGATAGAAAAGG--AAATTGGAGCCTTAGA----CATAGTACTGTGAAGGAAAAA-TGAAAGAT
↳ TAG-----CTTAAAGCACAAAAAGCCAGGATTTAACCTGTACCTTTTGCATAATGGTTTAGCTAGAAA-ACTTATATACAA
↳ AG-AATTTTAATATAAAATCCCGAAA-CTAGATGAGCTACTATCGAGCAATTTAT--TAGAATCAACCCGTCTATGTTGCAAAATAGTG
↳ GGATGACTTTATAGTAGAGGTGAAAAATCAACCGGACCTAGTGATAGCTGGTTAGCCAAAAACGAATTTAAGTTCAACAATAAGTTTA
↳ TTT---CCTGATATACTAT-ATT-----AACCATAAACTTATAAGTTATTATAAGAGGTCAGCCCTTATGA-ATAAGG
↳ AAACAACCTCTAATAG-AGGGAATTTATTCT-----TATTTACATAGTAGGCCTAAAAGCAGCCATCAATTAAAAAA-G
↳ CGTTAAAGCTTAAA--CTTCATTCTTTCTTAATCCCTTAACCT--CTTAACTACCCCTAAA-ATAATATTGGCTTAATCTATGTCCCA
↳ TAGAAGAAATAATGCTAAAATAAGTAAACCAGAATT-----TATTCTCC-ATGCACTAGCTTAAATTAGAACGGAGC-ATCCACT
↳ AATAATTAACAGCTAAATAATCATAA-----ATATTAACCTAG-AAACATTATTTA-TAACTGTT-AACCCGACA
↳ CAGGAGTGCTATCTAAGGAAAGATTAAGGAGTAAAGGAACTCGGCAAA-CTAGGATTTCGCCTGTTTACCAAAAAACATCGCCTCTAG
↳ CA----TAACAAGTATTAGAGGTCTGCCTGCCAGTGATT---CTATTAACGGCCGCGGTAT-CCTGACCGTGCAAAGGTAGCATA
↳ ATCATTGTCTCCTAATTAGAGACTAGTATGAATGGCTAGACGAAAAATCCAAGTGTCTCTTACTCCCAATCAGTGAAATTGCCCTCCCC
↳ GTGCAGAGACGGGGATAACACCATAAGACGAGAAGACCCTGTGGAGCTTTAATCATAGATTTATTCTTT-----TTCATACAGACCC
↳ AAAAGCTCAACACCTTAAGAATT-----TTTTCAACTGATTTTGGTTGGGGTGACCTCGGAGTACAACAAAACCTCTGA--ATGCA
↳ TAG-----TAATAACTTACAAGTTTAAACGCAAACTGCCAG---TAATAGACCCAA-----ATTATTGATC
↳ AAAGGACCAAGTTACCCAGGGATAACAGCGCAATCCTATTCAAGAGTTTCATATCGACAATAGGGTTTACGACCTCGATGTTGGATCAG
↳ GACATCCAAATGTTGTCAGCCGCTATTAATGGTTTCGTTTGTTCACGATTAA-AGTCCTACGTGATCTGAGTTTCAGACCGGAGTAATCCA
↳ GGTGCGTTTCTATCTATGGGTAA--TTCTCCTAGTACGGAAAGGACCAGAGAAATCAGGCCAATCTTAAAAAGAAGCCTCC---AACT
↳ TAACAAATGATTTCTATCTCAATTTGCCACGTTACCCTTACTAATC-----CTAGACCAGGAATTC-T

>Anolis_carolinensis

```

```

GCAAAACAAC-----AACCTAATTAACAAACAAA----ATAGCACACACAACAA-----AACAAACCA
↳ TTTGA-----CAAAAAAGTAGAGGCGATCGAACCTT-AACCGTAGAGTTTTATT----GATAGTACCGCAAGGGAAAAA-TGAAATAA
↳ TAGTGAAA-----ACAAAAGCAAAACATAGCAAAGACTTCCCCTTGTACCTCTTGCATCACGATCTAGCAAGCATAAAACAAGCAAGA
↳ AGAAAAACACAGCCTGTTTCCCGAAC-TCAAGTGAGCTATTTTAAAGCAACTA----AAGAGTTAACCCGTCCCCTGTAGCAAAAGGGTG
↳ GGAAGACTTTAAATAGAGGTGAAAAGCCAACCGAAGCTTGCTGATAGCTGGTACCAGATAGACGAATTTTAGTTCAACTTAAACTTTT
↳ ATTAACCCGCCCTTAGTAA-AA-----ACTTAGTTTTTAAGATACTCAATGAGGGGACAGCCTTATTGA-GCCAGA
↳ ATACAGCCTGAACTAG-AGAGAAACCTCAACA-----AAAACAACAGTAGGCCTTAAACAGCCATCTAATATAATA-A
↳ CGTCGCAGTCTTCA-----TAATCAAAATACCAACCGCA-CTTTAAACTCCTATA-CTACCACCTGGTAATTCAATAAATTTA
↳ TAGAAGATACTATGCTAGAACTAGT-AACAAGAAAT-----TTTCTCTCACGCACAACCATAAATCAGACATAGAAAACTACT
↳ GAAAATTAACAAAACAATAAAAAAC-----AACTTACATAG-----TATATTA-CTAGCTGTT-AACCCAACA
↳ CAGGTGTGCATAAA-AGAAAGATTAAGTTAAAAAGGAACTCGGCAACAAACAAATTCAACTGTTTACCAAAAAACATAGCCTTTAG
↳ CA----AACCAAGTATTAAGGTAAACGCCTGCCAGTGA-----AACTTAAACGGCCGCGGTAT-TCTAACCGTGCAAAGGTAGCGTA
↳ ATCACTGTCTCCTAAATAGGGACCTGTATGAACGGCTAAATGAATATTTAACTGTCTCTTTAACTAATCAGTGAAACTGATCTTCCA
↳ GTACAAAAGCTGGGATAATAATATAAGACGAGAAGACCCTGTGGAGCTTTAAATTTT-----
↳ AACCAACATATAATTATCACAGTTATTT---ATGCTTAACATTTTAAAGTTGGGGCGACTTCGGATAAAAACTAAACATCCGA--GCAAA
↳ TAA-----GGCACTGCCTA-AAAAAGGCCTACAAGCCAAAGCTA-----ATATTGACCCAG---TAATACTGATT
↳ AACGAACCAAGTTACCCAGGGATAACAGCGCCATCTTCTTTAAGAGTTTCATATCGACAAGAAGGTTTACGACCTCGATGTTGGATCAG
↳ GACACCCAAATGTTGACGCCGTATTAAGGTTTCGTTTGTTCACGATTAACAGTCTACGTGATCTGAGTTTACAGCCGGAGCAATCCA
↳ GGTCGGTTTCTATCTATACAGTTG-CTTTCTTTAGTAC-GAAAGGACCAAGAAAGCAGGACCAATATAACCTACGTCCTTTA-----C
↳ AATAGATTAAAAAAACTAAAAAT--TTAACAAACAACAAATTTTCCCC-----TAGACATAGGGGTTA--
>Gallus_gallus
CTTG-----CCCCCCTCTAGCCCGACAAACTCGTACCCTTAACATAAAAACTTACCTCCCCCTCTTA-----ACCAAAACA
↳ TTATA---AATTGTCCCAGTATAGGCGATAGAAAAGACTACCCCGCGCAATAGAGGCTAACTGTACCGCAAGGGAAAGA-TGAAATAG
↳ CAATGAAA-----ACCATAAGCAAAAAACAGCAAAGACCAACCCTTGTACCTTTTGCATCATGATTTAGCAAGAAC-AACCAAGCAAAG
↳ TG-AGCTAAAGTTTGCTTCCCGAAACCAAGCGAGCTACTTGCAGAGCAGTAAATTTGAGCGAACCCGTCTCTGTTGCAAAAGAGCG
↳ GGATGACTTGCCAGTAGAGGTGAAAAGCCTACCGAGCTGGGTGATAGCTGGTACCTGTCAAACGAATCTAAGTTCCCCCTTAACCCAC
↳ CCC---CTAAAGACACCCACCTT-TGTCAACCTTGAGAACGTTGGGGTTAAGAGCAATTCGATGGGGGTACAGCTCCATCGA-AAAAGA
↳ ACACAACCTCCTCCAG-CGGATAATAATACCCC-----TCCCGCACTGTGGGCCTTCAAGCAGCCACCAACAAAAGAGTG
↳ CGTCAAAGCTCCCT--CATTAATAATCTAAACC-----CTATTTGACTCCCTCA-ACCAAAGCAGGTTAACCTATG--ACAA
↳ TAGAAGAAATCAATGCTAAAATGAGT-AATCTGGAACCTA-----TCCTCCTACGGCGTAAACTTACATTA-----ATACATT
↳ ATTAACAGAACTCAACTTATACCCC-----C-----ACACTAACAGCAATACGTATTCTT-CAATCTGTT-AAGCCAACC
↳ CAGGAGCGCCAC--AGGATGATTAACCTACAGAAGGAACTCGGCAAA-CCAAAGACCCGACTGTTTCCCAAAAACATAGCCTTCAG
↳ CTA--ACAACAAGTATTGAAGGTGATGCCGCGCCAGTGACCCCAAAGTTCAACGGCCGCGGTAT-CCTAACCGTGCGAAGGTAGCGCA
↳ ATCAATTGTCCCCTAAATTGAGACTTGTATGAATGGCTAAACGAGGTCTTAACTGTCTCTGTAGGTAATCTATGAAATTAGTATTCCC
↳ GTGCAAAAACGAGAATGTGAACATAAGACGAGAAGACCCTGTGGAACCTTTAAATCACGACCACCTTACAA----CCTTACAGACCCCC
↳ ACTGGGTCCACCCACACATAAACCCC-----TGTCGACATTTTTCGGTTGGGGCGACCTTGGAGAAAAAAAATCCTCAA--ACCCA
↳ CAG-----ACCACAACTCTT-ACTAAGACCAACTCTCAAAGTACCAACAGT-----AACCAGACCCAA----TATAATTGAGC
↳ AATGGACCAAGCTACCCAGGGATAACAGCGCAATCTCCTCCAAGAGCCCATATCGACAAGGAGGTTTACGACCTCGATGTTGGATCAG
↳ GACAACCTAATGTTGCAACCGTATTAAGGGTTCGTTTGTTCACGATTAACAGTCTACGTGATCTGAGTTTACAGCCGGAGCAATCCA
↳ GGTCGGTTTCTATCTATGGACA---CTCTCCTAGTAC-GAAAGGACCGGAGAAGTGGGGTCAATACCACTGAGCACACCCCAACCTTC
↳ TAAGCAATGAATACAACCTCAACTGCCAAGAACCCTCCCCACACCCGAACCTCTAGAAAAGGATCCA--
>Alligator_mississippiensis

```

```

AC-----CCATATCTAGCCCTACCTCCTTTCAACATGC-----TTAATACAAACCTCAA-----ACCAAAACA
↳ TTTATC--TAACACTTTAGTATTGGCGAAAGAAAATA-TAATCAGGCGCAACAGA----TAAAGTACCGCAAGGGAAAAAATGAAATAA
↳ AAATGAAA-----ACATAAGTAAACACAGCAAAGACCAACCCTTCTACCTTTCGCATCATGGTTTAGCAAACCA-CAACTGGCAAAA
↳ AG-AATTTAAGTCACACAACCCGAAA-CTAGGTGAGCTACTAAGTAGCAGCTAAA-TTTGAGCTAACCCCCCTCTGTGGCAAAAGAGCG
↳ GGAAGACTACTTAGTAGAAGTGAAAAGCCTACCGAACCTAGTGATAGCTGGCTGCTTGAAAAATGAATCTTAGTTCCACTGTAA-ATTC
↳ TCCTAACTACCACAACGAAGGAT-----AAGGAAGAATTTTCAAGCTATTAAATAGAGGTACAGCTCTATTAA-TGCAGG
↳ ACTCAACCTCCACTTA-AGGTAATCCATCCCAT-----TCCTCGGACTGTGGGCCCTTAAAGCAGCCACCA--CAAAAAA-G
↳ CGTCAAAGCTTAA-----CTCCAAAAAATACCAACAA--CAAATCGAACCCCTACA-TCATAACCAAGCCATCCTACAACCC--
↳ TAGAAAAGATTATGCTAAAAATGAGT-AATAAGAAAACGAC-----TTTCTCCTTTGCGTCAGCCTACATCCTACATGACA-CCCTATA
↳ GATTATTAACAGCCGCTTCCTACA-----CCACCCCAAA-AAAACAAAAAG-AAGACCT-----CACCTGTTGAACCCCAACA
↳ CAGGAACGCAACA--GGAAAGGCTAAACCTTGCAAAAGGAACCTCGGCAAA-CAAAGATTCCGACTGTTTACCAAAAAACAGCCCTAG
↳ CC----CCCTCAGTATTAGGGTGATGCCTGCCCAATGATTTC--CAATTGAATGGCCGCGGTATATACAACCGTGCGAAGGTAGCGTA
↳ ATCACTTGTCTCCTCAATAAGGACCAGTATGAACGGCTAAACGAGAATCTAACTGTCTCCTGCAAGCAACCAATGAAATTGATCTCCCT
↳ GTGCAAAAGCAGGAATGACCCACCAGACGAGAAGACCCTGTGAAACTTTAACCGACTAAGTCACACAC-----TAGGAACA
↳ ACACAACCCACAACCACCTTAAACCCA----TGACTTAGCGCTTTTGGTTGGGGTGACCCTAAACAAAAAAAACCTTTTAA--GACAA
↳ TCATAACAAAATTAGACTATTA-ACTAAGACCCACACCTCAAAGTACTTAACTGT-----AATTAGATCCGA----CAATGTCGATC
↳ CACGAATAAGCTACTCCAGGATAACAGCGCAATCCCTTCAAGAGCCCCTATCGACAAGGGGGTTTACGACCTCGATGTTGGATCAG
↳ GACACCCCATTTGGTGTAACCGTATTAATGGTTTCGTTTGTTCACGATTAA-AGTCCTACGTGATCTGAGTTTACAGCCGGAGTAATCCA
↳ GGTCGGTTTCTATCTATGACTATGTCTCTTCTAGTAC-GAAAGGACCGGGGAAACAAGGCCCATGCCACCAAGTACGCCTTACCTAT
↳ AATTTAATGCACGTAACATAAATAATTAGGATAAACACTACTTCT-----CAAAACAAGAGAAC---
>Alytes_obstetricans
GCT-----GTAAACCTAGCCC-----TCAAAAAACAA-----TAACCCTCAAATCCCA-----ACTAAATCA
↳ TTTAC-----ACTCTTTAGTATAGGTGATAAAAAAAG--AAATAAGCGCAACAGA----CAGAGTACCGTAAGGGAAAGA-TGAAATAG
↳ CTATGAAATAAA-CACCAAAGCACTAAAAAGCAGAGATACATCCTCGTACCTTTTGCATCATGGTCTAGCAAGAAA-CCTTAAGCAAAA
↳ CG-CATTAAAGTTTAACCCCCGAAA-CTAGACGAGCTACTCCGAGGCAGCTAT--AAGAGCAAACCATCTCTGTGGCAAAAGAGTG
↳ GGAAGACCCCGAGTAGAGGTGACAAGCCTACCGAGCCTAGTGATAGCTGGTTGCTTACGAAGAGAATATTAGTTCTGCCTCATAACTG
↳ TTTAAAGCACTATAAAGCAAAC-----AACAA-ACTTGAGAGTTAGTCAAAGGGGTACAGCCCATTTGA-CTCAGG
↳ AAACAACCTAACACCC-GGGGTAAAGATTATAATAAACAAGG---ATTTTGACCCAGTGGGCCTAAAAGCAGCCACCTGTTAAGAAA-G
↳ CGTCATAGCTCTAT--CAGCACCCGACCAACAATCCCAATAAC-CACTCTGATCCCTCCT-A--GTAGTAAGCCCTCCTATTACATA
↳ TAGAATTGTTTATGCTAGAAGTAGT-AACTAGAAT-----AATTCTCCTAATGCAAGTGCACATCAGCCCGAAAC-ACTCACT
↳ GATAACTAACGACTACTGAACCCA-----CTGTAGTA-ACACCCCCGAGAAAAACCTACACTA-ACCATCGTT-AACCTAACA
↳ CAAGAGCATCCCA--GAAAAGATTTAAGACGTGGAAGGAACCTCGGCAAA-CACGAACCTCGCCTGTTTACCAAAAAACATCGCTCTTG
↳ TC-----AACTATAAGAGGTCCCGCTGCCCAGTGACTAC--ATGTTCAACGGCCGCGGTAT-TATGACCGTGCAAAGGTAGCGTA
↳ ATCACTTGTCTTTTAAATAAAGTCTAGTATGAAAGGACACGAAAGTTCAACTGTCTCCACATCTAATCAGTGAAACTGATCCCCC
↳ GTGCAGAAACGGGATACTCTCATAAGACGAGAAGACCCCATGGAGCTTCAAACACGCAAGCAGCCGCTTTAAACCTCATCTAACCACA
↳ AGCAAGCGGACA-----TGACCAGCTGTTTTCGGTTGGGGCGACCACGGAGAACAAAAAATCCTCCGA--GACGA
↳ AAG---GGCCAACCCCTTA-CCCAAGAATCACCCTCTAAGGAATTGAATATCTAAC--ATCATGATCCAA-----TATTTGATC
↳ AACGAACCAAGTTACCCTGGGATAACAGCGCAATCCATTTTATAGAGCCCCTATCGACAAATGGGTTTACGACCTCGATGTTGGATCAG
↳ GGTACCCAGTGGTGACGCCGCTACTAAAGGTTTCGTTTGTTCACGATTAA-AACCCTACGTGATCTGAGTTTACAGCCGGAGCAATCCA
↳ GGTCAGTTTCTATCTATGCAATG--CTCTCTTAGTAC-GAAAGGACCGAAAGAACCTGGCCAATGCTAATGCAAGCCCTAA---CTTA
↳ GCCTTAGTGAAATCAGCTCAACTAAGACAAAAATTGCAACTAAGC-----CCGAGATAAGGGCT--
>Carcharhinus_leucas

```

```

ACCTTAAAGCTAGCCTATATAT-TAATTCAACTAGACCTTATAAATCCAATC-----TATATTACAGTTTTTA-----ATCAAAACA
↳ TTCTC----AACCTTTTAGTATGGGTGACAGAACAAT-AACCCAGCGCGATAGC----TTATGTACCGCAAGGGAAGC-TGAAAAAG
↳ AAATGAAATAAATCATTAAAGTACTAAAAAGCAGAGATGATACCTCGTACCTTTTGCATCATGATTTAGCTAGAAA-AACTAGGCAAAA
↳ AG-ACCTTAAGTCTATCTCCCGAAA-CTAAACGAGCTACTCCGAAGCAGCATT--TAGGGCTAACCCGTCTCTGTGGCAAAAGAGTG
↳ GGAAGACTTCCGAGTAGCGGTGAAAAGCCTACCGAGTTTAGTGATAGCTGGTTACCCAAGAAAAGAACTTTAGTTCTGCATTA-ATTTT
↳ TCACTATCTAGACAAGACTTACTCGTCAAA---GAATTCCATAAGAATTAATAGTTATTTAGAAGAGGTACAACCCTTCTAA-ATCAAG
↳ ACACAACCTTTTTAAGG-TGGGAAATGATCATAATTATTAAGG---TTACCATCCTAGTGGGCCTAAAAGCAGCCACCTGTTAAGTAA-G
↳ CGTCGAGCTCTAATCTAATAATTAACCTTTAATTAGATATT-CATTCAACCCCCCTTATTCTATTGGGTATTTTATATTTATA
↳ TAAAAGAACTTATGCTAAAAAGAGT-AATAAAGAGAATAA-----ATCTCTCCCGACACAAGTGTATGTGAGAAAGAAATTAATCACT
↳ GATAATTAACGACCCCAAACTGAG---GTCATTATATCAATCAATCAACTAGAAAACCTATTTTC-CTACTCGTT-ACCCCTACA
↳ CAGGAGTGTCACTA-GGAAAGATTAAGAAAATAAAGGAACCTCGGCAAA-CATAAACTCCGCCTGTTTACCAAAAAACATCGCCTCTTG
↳ AT-----AACTATAAGAGGTCCCGCCTGCCCCTGTGATA----ATGTTCAACGGCCGCGGTAT-TTTGACCGTGCAAAGGTAGCGTA
↳ ATCACTTGTCTTTTAAATGAAGACCTGTATGAAAGGCATCACGAGAGTTTAACTGTCTCTATTTTCTAATCAATGAAATTGATCTATTCT
↳ GTGCAGAAGCGAATATAATAACATTAGACGAGAAGACCCTATGGAGCTTTAAACACTTAAATTAATTATGTAACCATCCACTTCTCAGG
↳ ATATAAACAAAATATACAATACTTCT-----AATTTAACTGTTTTTGGTTGGGGTGACCAAGGGGAAAAACAAATCCCCCTT--ATCGA
↳ TTG--AGTACTAAGTACTTAAA-AATTAGAATGACAATTCTAATTAATAAAATATTTATCGAAAAATGACCCAGG---ATTTCTGATC
↳ AATGAACCAAGTTACCCTAGGATAACAGCGCAATCCTTTCTCAGAGTCCCTATCGAAGAAAGGGTTACGACCTCGATGTTGGATCAG
↳ GACATCCTAATGTGCAACCGCTATTAAAGGGTTCGTTTGTTCACGATTAATAGTCCTACGTGATCTGAGTTTCAGACCGGAGAAATCCA
↳ GGTCAGTTTCTATCTATGAATACA-CTTTTCTAGTAC-GAAAGGACCGGAAAAGTAGGGCCAATGCCACTAGTACGCCCTA---TTTT
↳ CATCTATTGAATAAACTAAAAATAGATAAGAAAAGATCACCTAGTGCC-----CAAAAAAGGGTT-----

>Latimeria_menadoensis
GCT-----ATATAGCTAGCCAACAAAACCACAACCACTATACCA-----TCAACCAAAACCCAAACAAGCAAAACCTAAACCA
↳ TTTA-----CCATCCAAGTATAGGCGATAGAAAAGA--CACCAGGCACAATAAT---GAAAGTACCGCAAGGGAACGC-TGAAAAAG
↳ AAATGAAACAACCTCGTAAAAGCAACAAAACCAAGACTAACCCTTGTACCTTTTGCATTATGATCTAGTTAGACCCACCGGGCAAAA
↳ TG-AATTTAAGTCCAACCCCCGAAA-CTAAGTGAGCTACTTCGAAACAGCCTAT--AAGGGCAAACCTTCTCTGTAGCAAAAGAGTG
↳ GGAAGATTTCCAAGTAGAGCGGATAAACCTAACGAACCTCAGTGATAGCTGGTTATTCGAGAAATGATTTTTAGTTTCAGCCCCTCAGCCT
↳ CCA-AACCCTAACACACAAAACCAATG-----TGAACCAGAGAGAGTTATTCAAAGGGGGTACAGCCCCTTTGA-AAAAGG
↳ ATACAACCTTCTAAGACCGGATAAAGACCATATTAACCTCAAGGATCAGCTGCCTAAGTGGGCCTAAAAGCAGCCACC---TAAAAAA-G
↳ CGTTAAAGCTTAAG--CAGCACTCAACCAACCTATAATGATAAATTAACCTCATTCCACCA-CCACTATCGAATTATTCTATA--TATA
↳ TAGAAGAATAAATGCTAGAATTAGT-AACAAGAAGGCAATAAACCTTCTCTAAGTGCATAAGTGTACATCAGATTAGATT-AACCACT
↳ GATAATTAACGACTTCAAAGAGAA-----TACTATGA-CATAAAACAAGAAAAACACACACCCACATCGTT-AATCCAACA
↳ CAGGAATGCAACCA-CGAAAGATTAAGAAAGAAAAGGAACTCGGCAAAACCATAAGCCCCGCCTGTTTACCAAAAAACATCGCCTCCTG
↳ CCAACAACAGAAGTATTGGAGGTCCCGCCTGCCAGTGACA---AGATTTAACGGCCGCGGTAT-CCTGACCGTGCAAAGGTAGCGAA
↳ ATCACTTGTCTTTTAAATGAAGACCTGTATGAATGGCACCACGAGGGCTTAACTGTCTCTCTTTCCAATCAGTAAATTTGATCTGTCC
↳ GTGCAGAAGCGGACTTAACCACATTAGACGAGAAGACCCTGTGGAGCTTCAGACACAAAGCCAACCTACAAATAA-ACACTTAGATAACA
↳ AGTTGATAAACCAAGTAGCCTAACACT-----GGCCCTATTGTCTTTGGTTGGGGCGACCACGGAGAAAAAGTAATCCTCCAA--GCCGA
↳ TTG--GTACCACCTGTACTAAA-ACAAAGGGTGACACCCCAAGTAATAAAAAATTTATC-GGACATGACCCAGGAACCTGAACCTGATC
↳ AACGAACCAAGTTACCCAGGATAACAGCGCAATCCTTTCCAAGAGTCCGAATCGACGAAAGGGTTTACGACCTCGATGTTGGATCAG
↳ GACACCCCAATGTGGAAGCCGTATTAAAGGTTTGTGTTGTTCAACAATTAAGTCTACGTGATCTGAGTTTCAGACCGGAGAGATCCA
↳ GGTCAGTTTCTATCTATGACGTTA-ATTCTCCAGTAC-GAAAGGACCGGAGAACTTGAGCCAATACCACAAGCACGCTCAT---CTTC
↳ AACCTGATGAAAACAACCTAAAAAGATAAAGAAGAACGTAACCACCC-----CGTAAACAACGGGTAAT

>Oncorhynchus_mykiss

```

```

GCTGACTAGCTAGCCAACATATTTGGTCCAACACCACAACATACATACCCCAA-----TAAAACTTAGAATTAAGTC-----AACAAACCA
↳ TTTTT-----CCACCTTAGTAGGGGCGACCGAAAAGGAGATAATTGAGCAACAGA----AAAAGTACCGCAAGGGGAAGC-TGAAAGAG
↳ AATTGAAATAACCCATTTAAGCCTAGAGAAGCAGAGATTAATCTCGTACCTTTTGCATCATGATTTAGCCAGCAC-ACCTGAGCAAAG
↳ AG-AACTTTAGTTTAGGCCCCCGAAA-CTAGACGAGCTACTCCGGGACAGCCTATTGTAGGGCCAACCCGTCTCTGTGGCAAAAGAGTG
↳ GGACGAGCCCCGAGTAGAGGTGATAAACCTATCGAGCCTAGTTATAGCTGGTTGCTTAGGAAATGAATAGAAGTTCAGCCCCCGGCTT
↳ TCT-TAGGACCTTAAGGTAAAACTAATATT--GTCCCAAAGAAACCAGGAGAGTTAGTCAAAGGAGGTACAGCTCCTTTGA-ACAAGG
↳ ACACAACCTTAACAGG-CGGCTAAGGATCATAGTTCGAAGGT---AACCTGTTACAGTGGGCCTAAGAGCAGCCACCTGCACAGAAA-G
↳ CGTTAAAGCTCAGA--CAG-ATACAAACCTCTTATCCTGATAAG-AAATCCCACCCCCCTA-ACCGTACTAAGCCGTTCCATGCCCCCA
↳ TGGAAGAGATTATGCTAGAATGAGT-AATAAGAGAGT-----ACAACCTCTCTCCAGCACATGTGTAAGTCGGACCGGACC-CCCCACC
↳ GACAAATAACGAACCCAAACCAAGAGGGAACCTGTAGGCCAGAA-CAAACACCAAGAAAAACCTACACCAACAAATCGTT-ACCCCCACA
↳ CAGGAGTGCCCCAAGGGAAAGACCCAAAGGAAGAGAAGGAACCTCGGCAAA-CACAAGCCTCGCCTGTTTACCAAAAAACATCGCCTCTTG
↳ CA----AATCAAAACATAGAGGT-CCGCCTGCCCTGTGACTAT--GGGTTTAAACGGCCGCGGTAT-TTTGACCGTGCGAAGGTAGCGCA
↳ ATCACTTGTCTTTTAAATGAAGACCTGTATGAATGGCATCACGAGGGCTTAGCTGTCTCTCTTCCAAAGTCAATGAAATTGATCTGCCC
↳ GTGCAGAAGCGGACATAAGCACATAAGACGAGAAGACCCTATGGAGCTTTAGACACCAGGCAGATCACGTCAAGCAACCTT---GAATT
↳ AACAAGTAAAAACGCAGTAGACCCCT-----AGCCCATATGTCTTTGGTTGGGGCGACCGCGGGGGGAAAAATTAAGCCCCACTGTGGA
↳ CTG-GGGGGCACTGCCCCACCA-GCCGAGAGCTACAGCTCTAAGCACCAGAATATCTGACCAAAATATGATCCGG----CGAACGCATTTC
↳ AACGGACCGAGTTACCCTAGGGATAACAGCGCAATCCTCTCCAGAGTCCCTATCGACGAGGGGGTTTACGACCTCGATGTTGGATCAG
↳ GACATCCTAATGTGTCAGCCGTATTAAGGGTTCGTTTGTTCACGATTAA-AGTCCTACGTGATCTGAGTTTCAGACCGGAGTAATCCA
↳ GGTCAGTTTCTATCTATGAAGTGA-TGTTTCTAGTAC-GAAAGGACCGGAAAGAAGGGGCCCATGCTTGAGGCACGCCCCA---CCCC
↳ CACCTGATGAAGGCAACTAAAAACAGACAAGGGGGCACACCAAGATTGC-----CTAAAAGAACGGCGC--
>Protopterus_doloi
ACTAAATGG-----TAAGCCTGTCAATTTTTTAAACATGTAAATTTT-ATATAACACACACTCTGTA-----AGTAAACCA
↳ TTTA-----TACTTCTAGTATTGGAGAAAGAAAGAA--ATTCAAGCGCAT-----AAAAGTACCGCAAGGGGAAAC-TGAAAAAC
↳ TAGTGAAG-----AATTAAGTTTAAAAAGCAAGACTAACACTTGTACCTTTTGCATCATGGTCTAGCTAGTCA-TGAAGGGCAGAA
↳ AG-AATTTTAGTCC-CACCCCGAAA-CTAGGCGAGCTACTCCGAGACAGCCAA---ACGGCCAACCTCGTCCATGTGGCAAAATGGTG
↳ AGAAGAGCTCCGAGTAGCGGTGAAAAGCCAAACGAGCCTAGAGATAGCTGGTTGCGCGAGAAACGAATCTTAGTTCTACCCTAAATTTT
↳ TCACGGCTACAGCCTTAAACCCCGTA-----ATTAAATTTAGAG-GCTACTCAAAGGGGGGACAGCCCTTTGT-GAAAGG
↳ ACACAACCTCAACAAA-CGGATAATGATTAATCTACA---AG---GTAAATTTAAGTGGGCCTAAAAGCAGCCACCAACAAAGAAA-G
↳ CGT-ATAGCTCCCC--CATAACTTTATCCATCAATTCAATTTAA-CAATCATAATCCATCTTAC-GTATCGAGCTGACTTATATCATT
↳ TAAGTGCAAAATGCTAGAATGAGT-AACAAGAAAAT-----ACAATTTTCTCTCACATAGGTGTTAGTCAGAACAGATT-ACCTACT
↳ GACAATTACTGATAATGAACACAA-----TGTATTAACCTG-TCATAAGCAAGAAAAACCTTACACTGATATATCATA-AATTCTACA
↳ CTGAAGTGTGCTT--GGAAAAATTAAGGGGGGAGAAAGGAATTCGGCAA--CTATGGCCTCGCCTGTTTACCAAAAAACATCGCCTCCTG
↳ CC-----AAATATAGGAGGTACTGCCTGCCCTGTGACTC--TGTTTAACTGGCCGCGGTAT-TTTGACCGTGCGAAGGTAGCGTA
↳ ATCACTTGTCTCTTAATTAGGGACCTGTATGAATGGCAACACGAGGGTCCAAGTGTCTCTCTCCCCAGATTAGTGAAATTGATCTATCC
↳ GTTCAAAAGCGGATATTTTTTCATAAGACGAGAAGACCCTGTGGAGCTTAAAGTTCTAATATTAAACAT-----
↳ -GAGCGTAAATATATACTTTAAGTTTTG-TAATATTAATACTTTTCGGTTGGGGCGACCGGAGTATAAAAAACCTCCGC-----
↳ -----AATACAATTCTATTTTAAAGAGGACACTTCTCA-AACTAGAATATCTAGC-ACAATTGACCCAG----TCTAACTGAGC
↳ AATGAACCAAGTTACCCAGGGATAACAGCGCAATCCCTTTAAGAGTCCCATCGACGAGGGGGTTTACGACCTCGATGTTGGATCAG
↳ GGTATCCTGGTGTGTCAGCCGTACCAAGGGTTTGTGTTGTTCAACAATTAA-TATCCTACGTGATCTGAGTTTCAGACCGGAGCAATCCA
↳ GGTCAGTTTCTATCTATGACTTC--TTTTTCTAGTAC-GAAAGGACTGAAAAGGGGGGCCTATAT-AAAAATATGCCCCA---CCCC
↳ CTACTACTGAATTTATATAAAGTAGCCAAGTGGGAAACCCCCACACG-----GGAGAAAACCACT--

```