

T-COFFEE

Elena Montenegro, Nicolás Manosalva, Luis Cervera, Fernando Freire

January 4, 2019

Contents

1	Multiple Sequence Alignment: T-COFFEE	2
1.1	Sample files	2
1.2	Class AlignSequences	2
1.2.1	Tests	14
1.3	MSA	17
1.4	MSA generic methods	18
1.5	T-COFFEE methods	28
1.6	Main MSA method	31
1.7	T-COFFEE	35
1.8	CLUSTAL	40
1.8.1	Verification against CLUSTALW	46

1 Multiple Sequence Alignment: T-COFFEE

In this notebook we are going to explain a detailed description of T-Coffee MSA algorithm, or more precisely, the set of algorithms that converge on T-COFFEE strategy to multiple sequence alignment.

We do so over a code implementation, not the T-COFFEE standard but an implementation that covers the basic aspect of the T-COFFEE approach.

Where necessary, we compare the T-COFFEE approach with the CLUSTALW one, in order to make more understandable the T-COFFEE method.

1.1 Sample files

Script 1.1.1 (text)

```
1 %%writefile sample.fasta
2 >laboA
3 NLFVALYDFVASGDNTLSITKGEKLRVLGYNHNGEWCEAQTNGQGQWVPS
4 NYITPVN
5 >lycsB
6 KGVYIALWDYEPQNDELPMKEGDCMTIIHREDEDEIEWWARLNDKEGY
7 VPRNLLGLYP
8 >1pht
9 GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIG
10 WLNQYNETTGERGDFPGTYVEYIGRKKISP
11 >lvie
12 DRVRKKSGAAWQQQIVGWYCTNLTPEGYAVESEAHPGSVQIYPVAALERI
13 N
14 >1ihvA
15 NFRVYYRDSRDPVWKGPALLWKGEAVVIQDNSDIKVVPRRKAKIIRD
```

Output

```
Overwriting sample.fasta
```

1.2 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 1.2.1 (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 time
12  import sys
13
14  MIN = -sys.maxsize - 1
15  COMPAC = 100000
16  """int: Constant to compact max score."""
17  SCORE_MATCH = 2
18  """int: Default match score."""
19  SCORE_NO_MATCH = -3
20  """int: Default no match score."""
21  SCORE_GAP_INI = -10
22  """int: Default gap ini in affine gap penalty."""
23  SCORE_GAP_CONT = -2
24  """int: Default gap continuation in affine gap penalty"""
25  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}
26  """dict: Default substitution matrix (for "ACGT" common nucleotide alphabet)"""
27
28  sys.setrecursionlimit(5000)
29
30  class AlignSequences:
31      """Recursive implementation of global, local and long substring alignments methods.
32
33      Attributes:
34          sequences (list of str): Contains the two sequences to align. The first
35          one (index 0) is the query sequence (BLAST concept) or bottom sequence on
→   alignment prints
36          or vertical sequence in the common graphical representation of score matrix.
37          len_seq0 (int): Sequence 0 length.
38          len_seq1 (int): Sequence 1 length.
39          mode (str): Computation mode:
40              'GLOBAL'          Global Alignment
41              'LOCAL'          Local Alignment
42              'LONG_SUBSTRING' Obtain long common substring
43          score_match (int): Score of match characters.
44          score_no_match (int): Score of no match characters.
45          score_gap_ini (int): Score of gap init.
46          score_gap_cont (int): Score of gap continuation.
47          score (int): Score of last computed alignment.
48          gaps (int): Number of gaps of the last computed alignment.
49          matches (int): Number of matches of the last computed alignment.
50          unmatches (int): Number of unmatches of the last computed alignment.
51          align_seq0 (str): Sequence 0 with the gaps necessary for the alignment.
52          align_seq1 (str): Sequence 1 with the gaps necessary for the alignment.
53          matching (str): Printable line with the align relations ('|', '.', ' ') between

```

```

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

```

```

96     self.align_seq0 = ""
97     self.align_seq1 = ""
98     self.matching = ""
99     self.ini_time = 0
100    self.finish_time = 0
101    self.set_subst_matrix(subst_matrix)
102    self.set_matrix_mode()
103
104    def init_stores(self):
105        """Init dictionary that store temp data of the alignment"""
106        self.score_store = {}
107        self.matches_store = {}
108        self.gaps_store = {}
109        self.max_score_index = (0, 0, 0)
110        self.max_score = 0
111        self.forward_arrow = {}
112
113    def set_sequences(self, sequences):
114        """Update the target sequences of the alignment"""
115        self.sequences = sequences
116
117    def set_stacks(self, stack_0=[], stack_1=[],\
118                  stack_0_indexes=[], stack_1_indexes=[], stack_0_refs=[], stack_1_refs=[]):
119        """Update the stacks for msa"""
120        self.stacks = [stack_0, stack_1]
121        self.stacks_indexes = [stack_0_indexes, stack_1_indexes]
122        self.stacks_refs = [stack_0_refs, stack_1_refs]
123
124    def set_matrix_mode(self, mode="SUBST"):
125        """Update matrix mode"""
126        self.matrix_mode = mode
127
128    def set_subst_matrix(self, subst_matrix={}):
129        """Update the score matrix"""
130        self.subst_matrix = subst_matrix
131
132    def set_scores(self, score_match=SCORE_MATCH, score_no_match=SCORE_NO_MATCH,\
133                  score_gap_ini=SCORE_GAP_INI, score_gap_cont=SCORE_GAP_CONT):
134        """Update the weight scores of the alignment"""
135        self.score_match = score_match
136        self.score_no_match = score_no_match
137        self.score_gap_ini = score_gap_ini
138        self.score_gap_cont = score_gap_cont
139
140    def set_mode(self, mode="ALIGN"):
141        """Set computation mode"""
142        self.mode = mode
143
144    def forward_track(self, index):
145        """Calc alignments in forward direction.
146
147            The alignment strings are calculated from init cell (0,0) in global

```

```

148         alignments or maximum score cell in local alignments.
149
150         In local mode it's necessary to extend the alignments (local) to the total
151     → length of
152         the sequences to show the location of the alignment, and in order to compare with
153         BioPython outputs.
154
155     Args:
156         index (tuple of int): Cell coordinates of the starting cell
157
158     Returns:
159         string: align sequence 0 (bottom) for printing purposes
160         string: align sequence 1 (top) for printing purposes
161         tuple of int: Coordinates of the last cell
162
163     """
164     ret_align_seq0, ret_align_seq1 = "", ""
165     (i, j, gap_ini) = index
166     ret_final_pos = (self.len_seq0, self.len_seq1)
167     while i < self.len_seq0 or j < self.len_seq1:
168         if self.mode == "LOCAL" and self.score_store[(i, j, gap_ini)] == 0:
169             ret_final_pos = (i, j)
170             break
171         arrow = self.forward_arrow[(i, j, gap_ini)]
172         if arrow == "d":
173             ret_align_seq0 += self.sequences[0][i]
174             ret_align_seq1 += self.sequences[1][j]
175             i, j, gap_ini = i + 1, j + 1, 1
176         elif arrow == "h":
177             ret_align_seq0 += "-"
178             ret_align_seq1 += self.sequences[1][j]
179             i, j, gap_ini = i, j + 1, 0
180         elif arrow == "v":
181             ret_align_seq0 += self.sequences[0][i]
182             ret_align_seq1 += "-"
183             i, j, gap_ini = i + 1, j, 0
184     #compute the complete align in local mode
185     if self.mode == "LOCAL":
186         ret_align_seq0 = self.sequences[0][0:index[0]] + \
187             ret_align_seq0 + self.sequences[0][ret_final_pos[0]:]
188         ret_align_seq1 = self.sequences[1][0:index[1]] + \
189             ret_align_seq1 + self.sequences[1][ret_final_pos[1]:]
190         diff_pos_ini = index[1] - index[0]
191         if diff_pos_ini > 0:
192             ret_align_seq0 = '-' * diff_pos_ini + ret_align_seq0
193         else:
194             ret_align_seq1 = '-' * -diff_pos_ini + ret_align_seq1
195         diff_len = len(ret_align_seq1) - len(ret_align_seq0)
196         if diff_len > 0:
197             ret_align_seq0 += '-' * diff_len
198         else:
199             ret_align_seq1 += '-' * -diff_len
200     return ret_align_seq0, ret_align_seq1, ret_final_pos

```

```

199
200 def calc_matching(self, align_seq0, align_seq1, ini_pos=(), final_pos=()):
201     """Calc matching string
202
203     The matching string is the string line to print between the top and
204 bottom alignment strings. It contains the match (!), no match (.) and
205 gap ( ) indicators.
206
207     Args:
208         align_seq0 (string): Bottom sequence
209         align_seq1 (string): Top sequence
210         ini_pos (tuple of int): Initial cell coordinates
211         final_pos (tuple of int): Final cell coordinates
212
213     Returns:
214         string: Matching string
215
216     """
217     count = 0
218     ret_matching = ""
219     diff_pos_ini = ini_pos[1] - ini_pos[0]
220     if diff_pos_ini > 0:
221         delta_pos = diff_pos_ini
222     else:
223         delta_pos = 0
224     for n, (i, j) in enumerate(zip(align_seq0, align_seq1)):
225         if self.mode == "LOCAL" and not (n >= ini_pos[0] + delta_pos and n <
226             ↪ final_pos[0] + delta_pos):
227             ret_matching += ' '
228         else:
229             if i == j: ret_matching += '|'
230             elif i != j and i != '-' and j != '-': ret_matching += '.'
231             else: ret_matching += ' '
232         count += 1
233     return ret_matching
234
235 def store(self, key, score, matches, gaps):
236     """Store info related to a computed cell
237     The maximum score is computed having into account the number of matches, if there are
238     most than one solution. If the score are equal, the path with more matches is
239     ↪ selected.
240
241     Args:
242         key (tuple of int): Cell coordinates
243         score (int): Cell score
244         matches (int): Cell matches
245         gaps (int): Cell gaps
246
247     """
248     self.score_store[key] = score
249     super_score = score * COMPAC + 10 * matches
250     if super_score > self.max_score:
251         self.max_score_index = key
252         self.max_score = super_score

```

```

249     self.matches_store[key] = matches
250     self.gaps_store[key] = gaps
251
252 def calc_score_binary(self, seq_0, seq_1, i, j, seq_0_index=0, seq_1_index=1, pos_0=0,
↳ pos_1=0):
253     """Compute alignment scores for two sequences
254     If there are a substitution matrix (actually dictionary) defined,
255     the scores are computed from the dictionary.
256     Args:
257         seq_0 (int): Sequence 0
258         seq_1 (int): Sequence 1
259         i (int): Sequence 0 char index
260         j (int): Sequence 1 char index
261         seq_0_index (int): Sequence 0 index on original sequences (MSA)
262         seq_1_index (int): Sequence 1 index on original sequences (MSA)
263         pos_0 (int): Sequence 0 index on stack 0
264         pos_1 (int): Sequence 1 index on stack 0
265     """
266     if self.subst_matrix:
267         if self.matrix_mode == "SUBST":
268             #print("PAIR", i, j, seq_0[i], seq_1[j])
269             subst_matrix_index = (seq_0[i], seq_1[j])
270             subst_matrix_index_swap = (seq_1[j], seq_0[i])
271             if subst_matrix_index in self.subst_matrix:
272                 matrix_score = self.subst_matrix[subst_matrix_index]
273             elif subst_matrix_index_swap in self.subst_matrix:
274                 matrix_score = self.subst_matrix[subst_matrix_index_swap]
275             else: #weight matrix
276                 if pos_0 in self.stacks_refs and i in self.stacks_refs[pos_0]:
277                     i_orig = self.stacks_refs[pos_0][i]
278                 else:
279                     i_orig = i
280                 if pos_1 in self.stacks_refs and j in self.stacks_refs[pos_1]:
281                     j_orig = self.stacks_refs[pos_1][j]
282                 else:
283                     j_orig = j
284                 if i_orig in self.subst_matrix[seq_0_index][seq_1_index] and \
285                     j_orig in self.subst_matrix[seq_0_index][seq_1_index][i_orig]:
286                     matrix_score = self.subst_matrix[seq_0_index][seq_1_index][i_orig][j_ori
g]
287             else:
288                 matrix_score = 0
289             # Gaps in almost one of the sequences. This case only arises in MSA
290             # There is no matrix related entry. If matrix is a weight matrix we compute
291             # as zero (as defined in T-Coffee)
292             if seq_0[i] == "-" or seq_1[j] == '-':
293                 inc_matches = 0
294             if self.subst_matrix:
295                 if self.matrix_mode == "SUBST":
296                     inc_score = self.score_gap_cont
297                 else:
298                     inc_score = 0

```



```

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

```

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.

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

388
 389 *Search algorithm of the long common substring:*

390
 391 *Modify the global algorithm in the following aspects:*
 392 *1. Only computes matches between characters or gaps in one or another*
 → *initial sequence.*

393
 394 *Args:*

395 *i (int): Sequence 0 index*
 396 *j (int): Sequence 1 index*
 397 *ini_gap (int): 1 if gap initiation, 0 if gap continuation*

398 *"""*
 399 score_diag, score_hor, score_ver = MIN, MIN, MIN
 400 matches_diag, matches_hor, matches_ver = MIN, MIN, MIN
 401 gaps_diag, gaps_hor, gaps_ver = MIN, MIN, MIN
 402 *#align and advance seq0 and seq1*
 403 *#in long_substring mode only matches are processed*
 404 if i < self.len_seq0 and j < self.len_seq1 and\
 405 (self.mode != "LONG_SUBSTRING" or self.sequences[0][i] == self.sequences[1][j]):
 406 inc_score, inc_matches = self.calc_score(i, j)
 407 key = (i + 1, j + 1, 1)
 408 if key in self.score_store:
 409 score_diag, matches_diag, gaps_diag = \
 410 self.score_store[key] + inc_score, self.matches_store[key] + inc_matches,
 → self.gaps_store[key]
 411 else:
 412 score, matches, gaps = self.align(i + 1, j + 1, 1)
 413 self.store(key, score, matches, gaps)
 414 score_diag, matches_diag, gaps_diag = score + inc_score, matches +
 → inc_matches, gaps
 415 *#don't align and gap in seq0 (advance seq1)*
 416 if j < self.len_seq1:
 417 gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
 418 key = (i, j + 1, 0)
 419 if key in self.score_store:
 420 score_hor, matches_hor, gaps_hor = self.score_store[key] + gap_score,\
 421 self.matches_store[key], self.gaps_store[key] + 1
 422 else:
 423 score, matches, gaps = self.align(i, j + 1, 0)
 424 self.store(key, score, matches, gaps)
 425 score_hor, matches_hor, gaps_hor = score + gap_score, matches, gaps + 1
 426 *#don't align and gap in seq1 (advance seq0)*
 427 if i < self.len_seq0:
 428 gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
 429 key = (i + 1, j, 0)
 430 if key in self.score_store:
 431 score_ver, matches_ver, gaps_ver = \
 self.gaps_store[key] + 1

```

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

```

x)

```

480         if not silent:
481             self.view()
482
483     def get_len_long_common_substring(self):
484         """Getter for the len of the common substring
485         That is equal to the number of matches of the alignment
486         """
487         return self.matches
488
489     def get_long_common_substring(self):
490         """Returns the longest common substring
491         whitout alignment (positional) information
492         """
493         long_common_substring = ""
494         for (char, match_char) in zip(self.align_seq1, self.matching):
495             if match_char == '|':
496                 long_common_substring += char
497         return long_common_substring
498
499     def view(self):
500         """Prints the alignment data"""
501         #unmatches = self.matching.count('.')
502         #gaps = self.matching.count(' ')
503         if self.matching:
504             gap_groups = self.matching.count('| ') + self.matching.count('. ') +
505                 ↪ self.matching[0].count(' ')
506         else:
507             gap_groups = 0
508         print(" ")
509         if self.mode == "LOCAL":
510             print("### AlignSequences. Local alignment (Smith-Waterman)")
511         elif self.mode == "LONG_SUBSTRING":
512             print("### AlignSequences. Long substring finder")
513         else:
514             print("### AlignSequences. Global alignment (Needleman-Wunsch)")
515         if self.subst_matrix:
516             print("\tUsing score matrix with matrix mode",self.matrix_mode)
517         print(self.align_seq1)
518         print(self.matching)
519         print(self.align_seq0)
520         print("\tScore:", self.score)
521         print("\tSimilarity (wo gaps):", self.matches / (self.matches + self.unmatches))
522         print("\tDistance (wo gaps):", self.unmatches / (self.matches + self.unmatches))
523         print("\tDistance:", self.unmatches / (self.matches + self.unmatches + self.gaps))
524         print("\tInit index:", self.max_score_index)
525         print("\tMatches:", self.matches, " Unmatches:", self.unmatches, " Gaps:",
526             ↪ self.gaps, " Gap groups:", gap_groups)
527         #simple scoring verification todo: apply to matrix
528         if not self.subst_matrix:
529             print("\tScore verified:", self.matches * self.score_match + self.unmatches *
530                 ↪ self.score_no_match \
531                 + self.gaps * self.score_gap_cont + gap_groups * self.score_gap_ini)

```

```

529     print("\tFinish. Execution milliseconds:", round((self.finish_time - self.ini_time)
530           ↪ * 1000))
531     print("\tScore Dictionary Size", len(list(self.score_store.keys())))
532
533     def edit_distance(self, score_match=0, score_no_match=-1, score_gap_ini=0,
534           ↪ score_gap_cont=-1):
535         """Calculates an edit distance as requested in questions 1 and 3
536           It's the same computation as a global alignment with -1 penalties applied to
537           score_gap_cont and score_nomatch and 0 in score_match and score_gap_ini
538
539           Args:
540             score_match (int): Score of match characters.
541             score_no_match (int): Score of no match characters.
542             score_gap_ini (int): Score of gap init.
543             score_gap_cont (int): Score of gap continuation.
544
545           """
546         self.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
547         self.compute("GLOBAL", True)
548         return abs(self.score)

```

1.2.1 Tests

Script 1.2.2 (python)

```

1  import re
2  from Bio import pairwise2
3  from Bio.pairwise2 import format_alignment
4  from Bio.SubsMat import MatrixInfo
5
6  failed = 0
7  passed = 0
8  launched = 0
9
10 def test_alignment(number, s1, s2, verbose=False, tipus="local", matrix={},\
11                   ↪ score_match=2, score_no_match=-3, score_gap_ini=-5,
12                   ↪ score_gap_cont=-2):
13     """Comparisons of global and local alignments between Biopython and AlignSequences
14       ↪ implementation.
15
16       Args:
17         number (int): The number(identifier) of the test.
18         s1 (str): Query string to align.
19         s2 (str): Subject string to align
20         verbose (bool): If True print outputs, default False
21         tipus (str) : If 'local' the alignment is local (Smith), if 'global' Waterman.
22         matrix (dict of int) : Substitution matrix
23         score_match (int): Score of character match
24         score_no_match (int): Score of character no match
25         score_gap_ini (int): Score of gap initiation
26         score_gap_cont (int): Score of gap continuation


```

```

26     """
27     global failed, passed, launched
28     try:
29         launched += 1
30         align = AlignSequences([s1, s2])
31         align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
32         if matrix != {}:
33             method = getattr(pairwise2.align, tipus + 'ds')
34             alignments = method(s2, s1, matrix,\
35                                 score_gap_ini + score_gap_cont, score_gap_cont)
36             align.set_subst_matrix(matrix)
37         else:
38             method = getattr(pairwise2.align, tipus + 'ms')
39             alignments = method(s2, s1, score_match, score_no_match,\
40                                 score_gap_ini + score_gap_cont, score_gap_cont)
41
42         align.compute(tipus.upper(), silent=not verbose)
43         m = re.match(r".*Score=([-1234567890]*)",
44                     → format_alignment(*alignments[0]).replace("\n", ""))
45         score = int(m.group(1))
46
47         #search AlignSequences alignment in all possibles alignments from Biopython
48         found = False
49         for a in alignments:
50             if verbose:
51                 print()
52                 print("BioPython alignment:")
53                 print(format_alignment(*a))
54             if align.align_seq0 == a[1] and align.align_seq1 == a[0]:
55                 if not verbose: print(format_alignment(*a))
56                 found = True
57                 break
58         assert(align.score == score)
59         print ("Passed test %s: scores are equal '%s'" % (number, align.score ))
60         assert(found)
61         print ("Passed test %s: alignments are equal '%s'" % (number, align.align_seq0 ))
62         passed += 1
63     except AssertionError:
64         print ("Failed test %s: alignments differ: \nBiopython:\n'%s'\nScore = %s \
65               \nAlignSequences\n'%s'\nScore = %s" \
66               % (number, alignments[0][1], score, align.align_seq0, align.score ))
67         failed += 1
68         exit(1)
69
70     prot1 = "GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNNGYNETTGERGDFPGTYVEYIGRKKISP"
71     prot2 = "NLFVALYDFVASGDNTLSITKGEKLRVLGYNHNGEWCEAQTNGQGWPSPNYITPVN"
72
73     test_alignment(1, prot1, prot2, False, "global", MatrixInfo.blosum62, 0, 0, 0, -8)
74
75     prot1 = "GARFIELD THE LAST FAT CAT"
76     prot2 = "GARFIELD THE FAST CAT"

```

```

77
78 test_alignment(2, prot1, prot2, True, "global", {}, 3, -2, 0, -8)
79 # test_alignment(2, prot1, prot2, False, "global", MatrixInfo.blosum62, 0, 0, 0, -8)
80 # test_alignment(3, prot1, prot2, True, "global", MatrixInfo.blosum40, 0, 0, 0, -8)
81
82 print(" ")
83 if launched == passed: print('Passed All Test')
84 else: print("ERROR: There are failed tests")

```

Output

```

Passed test 1: scores are equal '-88'
Failed test 1: alignments differ:
Biopython:
'GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNNGYNETTGERGDFPGTYVEYIGRKKISP'
Score = -88
AlignSequences
'GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNNGYNETTGERGDFPGTYVEYIGRKKISP'
Score = -88

### AlignSequences. Global alignment (Needleman-Wunsch)
GARFIELD THE FAS----T CAT
|||||.....||  ||||
GARFIELD THE LAST FAT CAT
      Score: 26
      Similarity (wo gaps): 0.9523809523809523
      Distance (wo gaps): 0.047619047619047616
      Distance: 0.04
      Init index: (0, 0, 1)
      Matches: 20  Unmatches: 1  Gaps: 4  Gap groups: 1
      Score verified: 26
      Finish. Execution milliseconds: 4
      Score Dictionary Size 1097

BioPython alignment:
GARFIELD THE FAST ----CAT
|||||.....|||  |||
GARFIELD THE LAST FAT CAT
      Score=26

BioPython alignment:
GARFIELD THE FAST---- CAT
|||||.....|||  |||
GARFIELD THE LAST FAT CAT
      Score=26

BioPython alignment:
GARFIELD THE FAS----T CAT
|||||.....||  ||||

```



```
GARFIELD THE LAST FAT CAT
```

```
Score=26
```

```
Passed test 2: scores are equal '26'
```

```
Passed test 2: alignments are equal 'GARFIELD THE LAST FAT CAT'
```

```
ERROR: There are failed tests
```

1.3 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.

1.4 MSA generic methods

Script 1.4.1 (python)

```
1  """This methods shows alternative implementations of multiple sequence alignments, CLUSTAL
   → and 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,
→ score_gap_cont=-8):
96     """
97     Performs initial pairwise alignments against the class AlignSequences
98     returning the distance between 0 and 100.
99
100     Args:
101         s1 (str): First sequence to compare.
102         s2 (str): Second sequence to compare.
103         matrix (dict of tuples of int): Substitution matrix, Biopython format.
104         matrix_mode (str): Type of matrix
105             'SUBST'          Substitution matrix
106             'WEIGHT'         Weight matrix
107         mode (str): Computation mode:
108             'GLOBAL'         Global Alignment
109             'LOCAL'          Local Alignment
110             'LONG_SUBSTRING' Obtain long common substring
111         score_gap_ini (int): Score of gap init.
112         score_gap_cont (int): Score of gap continuation.
113         score_match (int): Score of match characters (used if no matrix informed)
114         score_no_match (int): Score of no match characters (used if no matrix informed)
115
116     Returns:
117         (int): distance between sequences
118     """
119     align = AlignSequences([s1, s2])
120     align.set_scores(0, 0, score_gap_ini, score_gap_cont)
121     align.set_subst_matrix(matrix)
122     align.set_matrix_mode(matrix_mode)
123     align.compute(mode.upper(), silent = True)
124     identity = round((align.matches + 1) * 100 / (align.matches + align.unmatches + 2))
125     return 100 - identity
126
127 def guide_tree_UPGMA(sequences, matrix, matrix_mode, mode,\
128                     score_gap_ini, score_gap_cont,\
129                     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
→ root
151         of the other two nodes.
152         (dict of int, boolean = True): contains all nodes
153
154     """
155     tree = {} #initial tree
156     guide_tree = [] #guided tree, pairs to align in sequence
157     max_score = MIN_SCORE
158     max_score_position = ()
159     for i in range(0, len(sequences)):
160         for j in range(0, i):
161             if (i,j) not in tree:
162                 score = pairwise_align(sequences[i], sequences[j], matrix, matrix_mode, mode,\
163                                         score_gap_ini, score_gap_cont, score_match, score_no_match)
164                 tree[(i,j)] = score
165                 if score >= max_score:
166                     max_score = score
167                     max_score_position = (i,j)
168
169     print(tree)
170     len_tree = len(sequences)
171     guide_tree_nodes = {}
172     # Generate guide tree. At every step we compute another row averaging the
173     # most closer rows and removing all their row coordinates from the tree
174     while len(tree) > 0:
175         (imax, jmax) = max_score_position
176         guide_tree.append((imax, jmax, len_tree))
177         guide_tree_nodes[imax] = True
178         guide_tree_nodes[jmax] = True
179         guide_tree_nodes[len_tree] = False
180
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 (lit 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 three, the third position of the tuple contains the
→      root
319         of the other two nodes.
320         (dict of int, boolean = True): contains all nodes
321
322     """
323     dmatrix = {} #initial distance matrix
324     n = len(sequences)
325     for i in range(0, n):
326         for j in range(0, i):
327             if (i,j) not in dmatrix:
328                 distance = pairwise_align_distance(sequences[i], sequences[j], matrix,
→                 matrix_mode,\
329                 mode, score_gap_ini, score_gap_cont)
330                 dmatrix[(i,j)] = distance
331                 dmatrix[(j,i)] = distance
332     nseq = n
333     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
438     """
439     align = AlignSequences([stack_0[0], stack_1[0]])
440     align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
441     align.set_subst_matrix(matrix)
442     align.set_matrix_mode(matrix_mode)
443     align.set_stacks(stack_0, stack_1, stack_0_indexes, stack_1_indexes, stack_0_refs,
444 ↪ stack_1_refs)
445     align.compute(mode.upper(), silent = True)
446     # align_seq0 align_seq1 are the seq0 and seq1 alignments
447     # we need to deduce the rest of alignments.
448     # what we do is perform the same gap insertions, if any, as the first sequence of the
449 ↪ stacks
450     # the gap insertions where performed taken into account the initial sequence

```

```

446     # to compute the references to initial sequence in order to employ a weight matrix if
447     ↳ informed
448     stack_0_gapped, stack_0_references = gapeator(stack_0[0], align.align_seq0, stack_0,
449     ↳ stack_0_refs)
450     stack_1_gapped, stack_1_references = gapeator(stack_1[0], align.align_seq1, stack_1,
451     ↳ stack_1_refs)
452     return stack_0_gapped, stack_1_gapped, stack_0_references, stack_1_references
453
454 def get_name(index, sequence_names):
455     """
456     Obtain sequence name from index
457     """
458     name = ""
459     if index < len(sequence_names):
460         name = sequence_names[index]
461     else:
462         name = str(index)
463     return name
464
465 def to_newick(tree, sequence_names):
466     """
467     Obtain guide tree in newick format
468     """
469     # Change format to intermediate roots
470     roots = {}
471     newick_tree = ""
472     for branch in tree:
473         (i, j, k) = branch
474         name_i = get_name(i, sequence_names)
475         name_j = get_name(j, sequence_names)
476         name_k = get_name(k, sequence_names)
477         if name_i in roots:
478             new_root_i = roots[name_i]
479         else:
480             new_root_i = name_i
481         if name_j in roots:
482             new_root_j = roots[name_j]
483         else:
484             new_root_j = name_j
485         roots[name_k] = [new_root_i, new_root_j]
486
487     for root in roots.values():
488         s_root = str(root)
489         if len(s_root) > len(newick_tree):
490             newick_tree = s_root.replace("[", "(").replace("]", ").").replace("'", "")
491
492     return newick_tree

```

1.5 T-COFFEE methods

Script 1.5.1 (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):
142                     if j != k and j != i:
143                         #print("Triplet:", i, k, j)
144                         extend_weights(weight_library, i, k, j)
145
146 def compute_libraries(sequences, matrix,\
147     score_gap_ini, score_gap_cont, score_match, score_no_match):
148     """
149     Compute initial library of identities based on scores of pairwise alignments
150     Args:
151         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

```

1.6 Main MSA method

Script 1.6.1 (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      """
7      Performs MSA alignments from fasta file
8      Args:
9          file (str): Name of the FASTA file.
10         main_alg (str): Main algorithm:
11             "CLUSTAL"      Clustal like
12             "T-COFFEE"     T-COFFEE like
13         method (str): NJ neighbor join / UPGMA
14         matrix (dict of tuples of int): Substitution matrix, Biopython format.
15         matrix_mode (str): Type of matrix
16             'SUBST'        Substitution matrix
17             'WEIGHT'        Weight matrix
18         mode (str): Computation mode:
19             'GLOBAL'        Global Alignment
20             '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     print(sequence_names)
38     return do_msa(sequences, sequence_names,\
39                   main_alg, method, matrix, matrix_mode,\
40                   mode, score_gap_ini, score_gap_cont, score_match,\
41                   score_no_match, verbose)
42
43 def do_msa(sequences, sequence_names, main_alg="CLUSTAL", method="NJ", matrix={},
    → matrix_mode="SUBST",\
44           mode="GLOBAL", score_gap_ini=-10, score_gap_cont=-5, score_match=3,\
45           score_no_match=-2, verbose=False):
46     """
47     Performs MSA alignments from sequences
48     Args:
49         sequences (list of str): Sequences
50         sequence_names (list of str): Names of sequences
51         main_alg (str): Main algorithm:
52             "CLUSTAL" Clustal like
53             "T-COFFEE" T-COFFEE like
54         method (str): NJ neighbor join / UPGMA
55         matrix (dict of tuples of int): Substitution matrix, Biopython format.
56         matrix_mode (str): Type of matrix
57             'SUBST' Substitution matrix
58             'WEIGHT' Weight matrix
59         mode (str): Computation mode:
60             'GLOBAL' Global Alignment
61             'LOCAL' Local Alignment
62             'LONG_SUBSTRING' Obtain long common substring
63         score_gap_ini (int): Score of gap init.
64         score_gap_cont (int): Score of gap continuation.
65         score_match (int): Score of match characters (used if no matrix informed)
66         score_no_match (int): Score of no match characters (used if no matrix informed)
67         verbose (bool): If True prints verbose info
68
69     Returns:

```



```

70         list of 3-tuples of int: guide three, the third position of the tuple contains the
    ⇨ root
71         of the other two nodes.
72     list of str: alignments
73     list of str: sequence_names
74     list of int: sequence indexes relating strings in alignment to original sequences
75     """
76     if main_alg == "T-COFFEE":
77         primary_library, weight_library = compute_libraries(sequences, matrix,\
78             score_gap_ini, score_gap_cont, score_match, score_no_match)
79         #print("Primary library", primary_library)
80         #print("Weight library", weight_library)
81         # Matrix mode and other MSA parameters
82         matrix = weight_library
83         matrix_mode = "WEIGHT"
84     #     score_gap_ini = 0
85     #     score_gap_cont = 0
86     # From here only we need is to compute a MSA with weight matrix as reference.
87     else:
88         matrix_mode = "SUBST"
89     if method == "NJ":
90         guide_tree, guide_tree_nodes = \
91             guide_tree_NJ(sequences, matrix, matrix_mode,\
92                 mode, score_gap_ini, score_gap_cont,\
93                 score_match, score_no_match)
94     else: #UPGMA
95         guide_tree, guide_tree_nodes = \
96             guide_tree_UPGMA(sequences, matrix, matrix_mode,\
97                 mode, score_gap_ini, score_gap_cont,\
98                 score_match, score_no_match)
99     sequences_store = {}
100     sequences_store_indexes = {}
101     sequences_store_refs = {}
102     print("Guide Tree",guide_tree, sequence_names)
103     #return guide_tree, "", sequence_names
104     # Create MSA
105     for i in guide_tree_nodes.keys():
106         if i < len(sequences):
107             sequences_store[i] = [sequences[i]]
108             sequences_store_indexes[i] = [i]
109             sequences_store_refs[i] = []
110             autorefs = {}
111             for k in range(0, len(sequences[i])):
112                 autorefs[k] = k
113             sequences_store_refs[i].append(autorefs)
114
115     if verbose: print(sequences_store)
116     for (i ,j ,k) in guide_tree:
117         stack_i = sequences_store[i]
118         stack_j = sequences_store[j]
119         stack_i_indexes = sequences_store_indexes[i]
120         stack_j_indexes = sequences_store_indexes[j]

```

```

121     stack_i_references = sequences_store_refs[i]
122     stack_j_references = sequences_store_refs[j]
123     if verbose: print("Stack i", i, stack_i)
124     if verbose: print("Stack j", j, stack_j)
125     if verbose: print("Stack i_indexes", i, stack_i_indexes)
126     if verbose: print("Stack j_indexes", j, stack_j_indexes)
127     stack_0, stack_1, stack_0_references, stack_1_references = \
128         pairwise_align_msa_step(stack_i, stack_j, sequences, matrix, matrix_mode, \
129                                 mode, stack_i_indexes, stack_j_indexes, \
130                                 stack_i_references, stack_j_references, \
131                                 score_match, score_no_match, score_gap_ini,
132                                 ↪ score_gap_cont)
133
132     sequences_store[k] = []
133     sequences_store_indexes[k] = []
134     sequences_store_refs[k] = []
135     if verbose: print("=====")
136     for s in stack_0:
137         if verbose: print(s)
138         sequences_store[k].append(s)
139     for s in stack_1:
140         if verbose: print(s)
141         sequences_store[k].append(s)
142     for seq_index in stack_i_indexes:
143         sequences_store_indexes[k].append(seq_index)
144     for seq_index in stack_j_indexes:
145         sequences_store_indexes[k].append(seq_index)
146     for seq_refs in stack_0_references:
147         sequences_store_refs[k].append(seq_refs)
148     for seq_refs in stack_1_references:
149         sequences_store_refs[k].append(seq_refs)
150     if verbose: print("=====")
151     if verbose: print("Sequences store indexes", sequences_store_indexes[k])
152     if verbose: print("Sequences store references", sequences_store_refs[k])
153     if verbose: print("New stack:", k, sequences_store[k])
154     alignment = sequences_store[k]
155     newick_tree = to_newick(guide_tree, sequence_names)
156     return newick_tree, alignment, sequence_names, sequences_store_indexes[k]
157
158 def score(alignment, matrix, score_gap_ini=0, score_gap_cont=0):
159     """
160     Score based on sum of pair scores (SOP) taking into account substitution matrix
161     Derived from the objective score of MUSCLE refinement stage
162     """
163     msa_score = 0
164     for k in range(0, len(alignment[0])): #columns of msa
165         score_column_k = 0
166         nvalues = 0
167         for i in range(0, len(alignment)):
168             for j in range(i + 1, len(alignment)):
169                 if alignment[i][k] == "-" and alignment[j][k] != "-":
170                     score_column_k += score_gap_cont
171                 if k == 0 or alignment[i][k-1] != "-":

```

```

172         score_column_k += score_gap_ini
173     if alignment[j][k] == "-" and alignment[i][k] != "-":
174         score_column_k += score_gap_cont
175         if k == 0 or alignment[j][k-1] != "-":
176             score_column_k += score_gap_ini
177     elif (alignment[i][k], alignment[j][k]) in matrix:
178         score_column_k += matrix[(alignment[i][k], alignment[j][k])]
179         nvalues += 1
180     elif (alignment[j][k], alignment[i][k]) in matrix:
181         score_column_k += matrix[(alignment[j][k], alignment[i][k])]
182         nvalues += 1
183     if nvalues > 0:
184         #score += score_column_k / nvalues
185         msa_score += score_column_k
186     return msa_score
187
188 def score_from_fasta(file, matrix, score_gap_ini=0, score_gap_cont=0):
189     seq_fasta = readFasta(file)
190     sequences = list(seq_fasta.values())
191     return score(sequences, matrix, score_gap_ini, score_gap_cont)

```

1.7 T-COFFEE

Script 1.7.1 (python)

```

1 matrix = MatrixInfo.blosum80
2 file = "sample.fasta"
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,
20     ↪ quality_score_gap_cont))
21 print()
22
23 guide_tree_nj, align, sequence_names, indexes = do_msa_from_fasta(file,\
24     main_alg = "T-COFFEE", method = "NJ", \
25     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(aligned):
31     print(">" + sequence_names[indexes[i]])
32     print(s)
33 print()
34 print("Score", score(aligned, matrix, quality_score_gap_ini, quality_score_gap_cont))

```

Output

```

['Oncorhynchus_mykiss', 'Carcharhinus_leucas', 'Latimeria_menadoensis', 'Protopterus_dolloi',
→ 'Alytes_obstetricans', 'Anolis_carolinensis', 'Gallus_gallus',
→ 'Alligator_mississippiensis', 'Ornithorhynchus_anatinus', 'Homo_sapiens']
{(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

```
RTVKAVTGRQIFQPLHALRTAEKALLPGYHSFEWKPLKNVSANTEVGIIIDGLSGLPHTVDDYPIDTIAKFRFYDAALVSALMDMEEDILEGM
↳ KAHDLDYYLNG-PFTVVVKESCDGMGDVSEKHGCGPAVPEKAVRFSFTVMTIAI--THGNTNVRIFEEVKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERETMKNSVLLLEMGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVLHSITRSHTE
↳ LERYEVWRSNPYHESVDELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERKRWQSALDKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
```

>Gallus_gallus

```
RTVKAVTGRQIFQPLHALRTAEKALLPGYHPFEWKPLKNVSTNTEVGIIIDGLSGLPLSIDDYPIDTIAKFRFYDTALVSALKDMEEEILEGM
↳ KAKNLDDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNIAI--DHENERIRIFEEVKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKNSSELLLEIGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVHFSITRSHAEN
↳ LERYEIWRSNPYHESVDELDRVKGVSAPFIIETVPSIDALHCDIGNATEFYRIFQMEIGEYVKNPDATKEERKRWQLTLDKHLRKKMK
↳ LKPMRMMSGNFARKLMSKETVEAVCELIKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSYNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
```

>Homo_sapiens

```
RTVKAITGRQIFQPLHALRNAEKVLLPGYHHFEWQPPLKNVSSSTDVGIIDGLSGLSSSVDDYPVDTIAKFRFYDSALVSALMDMEEDILEGM
↳ RSQDLDDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPVPEKAVRFSFTIMKITI--AHSSQNVKVFEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKSELMLLEGGILRTFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVHFSITRSHAEN
↳ LERYEVWRSNPYHESVEELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEYVKNPNASKEERKRWQATLDKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVDAVCELIPEERHEALRELMMDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
```

>Ornithorhynchus_anatinus

```
RTVKAITGRQIFQPLHSLRTAEKALLPGYHPFEWDPLKNVSANTEVGIMDGLSGLPVSVDDYPVDTIAKFRFYDAALVSALMDMEEDILEGM
↳ KSQDLDDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNITL--AYEQENVKIFEEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKDSELKLEMGILRSFRFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVLHSITRSHAEN
↳ LERYEVWRSNPFHESVEELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEAYKNPNASKEERKRWQATLDKHLRKKMK
↳ LKPIMRMNGNFARKLMTKETVEAVCELVHCEERHEALRELMMDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
```

>Anolis_carolinensis

```
RTVKAVTGRQIFQPLHALRTAEKALLPGYHQFEWKPLKNVSSNTEVGIIIDGLSGIQLHVDYPVDTIAKFRFYDAALASALMDMEEDILEGL
↳ KRQDLDDYFYG-PFTVVVKESCDGMGDVSEKHGCGPAVPEKAVRFSFTLMTISV--THGNASIRVFEECKPNSELCKPLCLMLADESD
↳ HETLTAILSPLVAERAMKDSVLILDMAIPRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLILHSITRSHAEN
↳ LERYELWRTNPYHETVDELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERRRQSTLDKHLRKKMN
↳ LKPMTRMNGNFARKLMTKETVEAVCELIKSEERHEALRELMMDLYLKMKPVWRSSCPTKECPPELVQYSFNSQRFAELLATKFRYRYAGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
```

>Alytes_obstetricans

```
RTVKATTGRQIFQPLHALRNAEKALLPGYHPFEWKPLKNVSTCTDTGIIIDGLSGLNRSIDEPVEAISKFRFYDTALVSALKDMEEDILEGL
↳ RSHDMDYYLNG-PFTVVVKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMTISV--PNNNECVRIFDETVPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKTSELMLLEMGILRNFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDSTRLEASQNLVHFSITRCHTEN
↳ LQRYETWRANPHESVDELDRVKGVSAPFIIETLPSIDALHCDIGNAAEFYRLFQLEIGEYVKNPNATKEERKRWQSTLDKHLRKKMN
↳ LKPIMRMNGNFARKLMSKETVEAVCELVHSEERQEILRELMMDLYLKMKPVWRSSCPAKECEPELLYQYSFHSQRFAELLSTKFKYRYAGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
```

>Latimeria_menadoensis

```
RTVKATTGKQIFQPLHSLRNAEKALLPGFHPFEWQPPLKNVSSSTTEVGIIIDGMSGMTQFVDEYPLDTISKFRFYDAALVSALKDLEEEELKGL
↳ IEEDLEDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRYSFTIMTISV--ANSHNENVTFEEGKPNSELCKPLCLMLADESD
↳ HETLTAILGPVIAERAMKNSSELFLEMGILRSFKFIFRGTYDEKLIRDVEGLEASGSYICTLCDSTRSEASQNFILHSITRSHKEN
↳ LERYEIWRSNPYQEPVEELDRVKGVSAPFIIETLPSIDALHCDIGNATEFYKIFQDEIGEYVKNPNPSREEKKRWHSVLDKHLRKNMN
↳ LKPVMRMNGNYARKLMTKETVNAVCELIPEERQEALKELVDLYLKMKPVWRSTCPAKECEPELLCQYSFHSQRFAELLSTMYRYRYEGK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKYYELEDVL
```

>Protopterus_dolloi

```

RTVKAATGRQIFQPLHALRSAEKALLPGYHPFEWQPPLVGVSSTDVGIINGLSGLTSSVDEYPVEALAKFRFYDAALVSALKDIEENILEGM
↳ KQNGLDEYLSG-PFTVVIKESCDGMGDVSEKHGSGPPVPEKAVRFSFTIMSISV-AMSDSENVQIFEEFKPNSELCKPLCLMIADESD
↳ HETLTIVILGPVIAEREAAMKTSELMLELGGIILRTFKFFFRGTGYDEKLVREVEGLEASGSHYICTLCDATRQEASRNLVLHSITRSHAEN
↳ LERYEVWRSNPYNESVDELDRDVKGVSAPFIETRPCIDALQCDIGNATEFYKIFQDEVGEVYKRPNPSKEDRKRWHMTLDKHLRKKLS
↳ LKPVMRMNGNFARKLITKEAVDAVCELIPSEERRAAIRDLVHLYMLMKPVWRSTYPAKECPELLCQYSFNSQRFAELLSTKFQYRYEAK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYELEDVL
>Carcharhinus_leucas
KTVKAITGKQIFQPLHALRNAEKTLLPGYYSFEWQPPLANISTNTRVGIIDGLSGWVQCVDYPMETISRRLSYDVALASAVKEMEDDILEGL
↳ RSQNVDEFVSG-PFTVVIKESCDGMGDVSEKHGCGPTVPEKAVRFSFTIMSISV-MNENNEKVVFEEFKPNSELCCRPLCLMLADESD
↳ RETLTAILGPVIAERQSMKTSDLIVEIGDLYRSFQIFRGTGYDEKLVREVEGLEASGSIYICTLCDSTRSEASKNMVLHSITRNHAEN
↳ LERYEIWRSNPYHETADELDRDVKGVSAPFIETQPSIDALHCDIGNATEFYRIFQDEIGEYVYKNSNSKEERKRWQSMCLKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVEAVCELIPSEERREILRELMHLYLLMKPVWRSTFPTTECPDLLCQYSFNSQRFAELLHTEFSHRYEGK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKSYELEDIL
>Oncorhynchus_mykiss
RTVKATSGRQIFQPLHLTLRTAEKELLPGYHPFEWQPALKSVSTSCHVGIIDGLSGWIASVDDSPADTVTRRFRYDVALVSALKDLEEDIMEGL
↳ RERGLEDSACTSGFSVMIKESCDGMGDVSEKHGGPPVPEKPVRFSTIMSISVIAEGEDEAITIFREPKNSEMSCKPLSLMFVDES
↳ HETLTGVLGPVVAERNAMKHSRLILSVGGLSRFRHFHFRGTGYDEKVMREMEGLEASGSTYICTLCDSTRAEASQNTLHVSITRSHDEN
↳ LERYELWRTNPHSESAEELDRDVKGVSAPFMETQPTLDALHCDIGNATEFYKIFQDEIGEYVYKANPSREQRSSWAALDKQLRKKMK
↳ LKPVMRMNGNYARKLMTREAVEAVCELCSEERQEALRELMGLYIQMKPVWRSTCPAKECPDELCRYFSNSQRFAELLSTVFKYRYDGK
↳ ITNYLHKTLAHVPEIVERDGSIGAWASEGNESGNKLFRRFRKMNRQSKTFELEDVL

```

Score 99825

```

['Oncorhynchus_mykiss', 'Carcharhinus_leucas', 'Latimeria_menadoensis', 'Protopterus_dolloi',
↳ 'Alytes_obstetricans', 'Anolis_carolinensis', 'Gallus_gallus',
↳ 'Alligator_mississippiensis', 'Ornithorhynchus_anatinus', 'Homo_sapiens']
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
RTVKAVTGRQIFQPLHALRTAEKALLPGYHSFEWKPLKNVSANTEVGIIDGLSGLPHTVDDYPIDTIAKFRFYDAALVSALMDMEEDILEGM
↪ KAHDLDYDLNG-PFTVVVKESCDGMGDVSEKHGCGPAVPEKAVRFSFTVMTIAI--THGNTNVRIFEEVKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERETMKNVLLLEMGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVLHSITRSHTEN
↪ LERYEVRNSNPYHESVDELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERKRWQSALDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSFNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Gallus_gallus
RTVKAVTGRQIFQPLHALRTAEKALLPGYHPFEWKPLKNVSTNTEVGIIDGLSGLPLSIDDYPIDTIAKFRFYDTALVSALKDMEEEILEGM
↪ KAKNLDDYDLNG-PFTVVVKECCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNIAI--DHENERIRIFEEVKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKNSSELLLEIGILRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLVHFSITRSHAEN
↪ LERYEIWRNSNPYHESVDELDRVKGVSAPFIIETVPSIDALHCDIGNATEFYRIFQMEIGEYVKNPDATKEERKRWQLTLDKHLRKKMK
↪ LKPMRMMSGNFARKLMSKETVEAVCELKCEERHEALKELMDLYLKMKPVWRSSCPAKECEPELLCQYSYNSQRFAELLSTKFKYRYEGK
↪ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Anolis_carolinensis
RTVKAVTGRQIFQPLHALRTAEKALLPGYHQFEWKPLKNVSSNTEVGIIDGLSGIQLHVDDYPVDTIAKFRFYDAALASALMDMEEDILEGL
↪ KRQDLDDYDFKG-PFTVVIKESCDGMGDVSEKHGCGPAVPEKAVRFSFTLMTISV--THGNASIRVFEECKPNSELCKPLCLMLADESD
↪ HETLTAILSPLVAERAMKDSVLILDMAGIPRTFKFIFRGTYDEKLVREVEGLEASGSTYICTLCDATRLEASQNLILHSITRSHAEN
↪ LERYELWRNPNYHETVDELDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQFEIGEYVKNPDASKEERRRQSTLDKHLRKKMN
↪ LKPMTRMNGNFARKLMTKETVEAVCELKSEERHEALRELMMDLYLKMKPVWRSSCPTKECPVELVCQYSFNSQRFAELLATKFRYRYAGK
↪ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Alytes_obstetricans
RTVKATTGRQIFQPLHALRNAEKALLPGYHPFEWKPLKNVSTCTDTGIIDGLSGLNRSIDEPVEAISKFRFYDTALVSALKDMEEDILEGL
↪ RSHDMDYDLNG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMTISV--PNNNECVRIFDETKPNSELCKPLCLMLADESD
↪ HETLTAILSPLIAERAMKTSSELMLLEMGILRNFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDSTRLEASQNLVHFSITRCHTEN
↪ LQRYETWRANPHHESVDELDRVKGVSAPFIIETLPSIDALHCDIGNAAEFYRLFQLEIGEYVKNPNATKEERKRWQSTLDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELHSEERQEILRELMMDLYLKMKPVWRSSCPAKECEPELLYQYSFHSQRFAELLSTKFKYRYAGK
↪ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKFYEMEDVL
>Latimeria_menadoensis
RTVKATTGKQIFQPLHSLRNAEKALLPGFHPFEWQPPLKNVSSSTTEVGIIDGMSGMTQFVDEYPLDTISKFRFYDAALVSALKDLEEELLKGL
↪ IEEDLEDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRYSFTIMTISV-ANSHNENVTIFFEEGKNSELCKPLCLMLADESD
↪ HETLTAILGPVIAERAMKNSSELFLEMGILRSFKFIFRGTYDEKLIRDVEGLEASGSSYICTLCDSTRSEASQNFILHSITRSHKEN
↪ LERYEIWRNSNPYQEPVEELDRVKGVSAPFIIETLPSIDALHCDIGNATEFYKIFQDEIGEYVKNPNPSREEKKRWHSVLDKHLRKNMN
↪ LKPVMRMNGNYARKLMTKETVNAVCELIPSEERQEALKELVDLYLKMKPVWRSTCPAKECEPELLCQYSFHSQRFAELLSTMYRYRYEGK
↪ ITNYLHKTALAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKYEELEDVL
>Carcharhinus_leucas
KTVKAITGKQIFQPLHALRNAEKTLLPGYYSFEWQPLANISTNTRVGIIDGLSGVQCVDDYPMETISRRLSYDVALASAVKEMEDDILEGL
↪ RSQNVDEFVSG-PFTVVIKESCDGMGDVSEKHGCGPTVPEKAVRYSFTIMTISV-MNENNEKVVFEEEMKNSELCCRPLCLMLADESD
↪ RETLTAILGPVIAERQSMKTSDLIVEIGDLRSFQFIFRGTYDEKLVREVEGLEASGSIYICTLCDSTRSEASKNMVLHSITRNHAEN
↪ LERYEIWRNSNPYHETADELDRVKGVSAPFIIETQPSIDALHCDIGNATEFYRIFQDEIGEYVKNNSKEERKRWQSMMLDKHLRKKMN
↪ LKPIMRMNGNFARKLMTKETVEAVCELIPSEERREILRELMHLYLLMKPVWRSTFPTTECPDLLCQYSFNSQRFAELLHTEFSHRYEGK
↪ ITNYLHKTALAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKSYEELEDIL
>Homo_sapiens

```



```

RTVKAITGRQIFQPLHALRNAEKVLLPGYHHFEWQPPLKNVSSSTDVGIIDGLSGLSSSVDDYPVDITIAKFRYDSALVSALMDMEEDILEGM
↳ RSQDLDDYLSG-PFTVVVKESCDGMGDVSEKHGSGPVPEKAVRFSFTIMKITI--AHSSQNVKVFEEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKSSSELMLELGGILRTFKFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVFSITRSHAEN
↳ LERYEVWRSNPYHESVEELRDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEYKPNPNASKEERKRWQATLDKHLRKKMN
↳ LKPIMRMNGNFARKLMTKETVDVAVCELIPSEERHEALREMLDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Ornithorhynchus_anatinus
RTVKAITGRQIFQPLHSLRTAEKVLLPGYHPFEWDPPLKNVSANTEVGIMDGLSGLPVSVDDYPVDITIAKFRYDAALVSALMDMEEDILEGM
↳ KSQDLDDYLSG-PFTVVIKESCDGMGDVSEKHGSGPAVPEKAVRFSFTVMNITL--AYEQENVKIFEEAKPNSELCKPLCLMLADESD
↳ HETLTAILSPLIAERAMKDSELKLEMGILRSFRFIFRGTYDEKLVREVEGLEASGSVYICTLCDATRLEASQNLVLHSITRSHAEN
↳ LERYEVWRSNPFHESVEELRDRVKGVSAPFIIETVPSIDALHCDIGNAAEFYKIFQLEIGEAYKPNPNASKEERKRWQATLDKHLRKKMK
↳ LKPIMRMNGNFARKLMTKETVEAVCELVHCEERHEALREMLDLYLKMKPVWRSSCPAKECPESLCQYSFNSQRFAELLSTKFKYRYEGK
↳ ITNYFHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYEMEDVL
>Protopterus_dolloi
RTVKAATGRQIFQPLHALRSAEKALLPGYHPFEWQPPLVGVSSSTDVGIINGLSGLTSSVDEYPVEALAKFRYDAALVSALKDIEENILEGM
↳ KQNGLDEYLSG-PFTVVIKESCDGMGDVSEKHGSGPPVPEKAVRFSFTIMSISV-AMSDSENVQIFEEFKPNSELCKPLCLMLADESD
↳ HETLTIVILGPVIAERAMKTSELMLELGGILRTFKFFFRGTGYDEKLVREVEGLEASGSHYICTLCDATRQEASRNLVLHSITRSHAEN
↳ LERYEVWRSNPYNESVDELDRVKGVSAPFIIETRPCIDALQCDIGNATEFYKIFQDEVGEVYKPNPSKEDRKRWHMTLDKHLRKKLS
↳ LKPVMRMNGNFARKLITKEAVDAVCELIPSEERRAAIRDLVHLYMLMKPVWRSTYPAKECPPELLCQYSFNSQRFAELLSTKFQYRYEAK
↳ ITNYLHKTLAHVPEIIERDGSIGAWASEGNESGNKLFRRFRKMNRQSKCYELEDVL
>Oncorhynchus_mykiss
RTVKATSGRQIFQPLHLTLRTAEKELLPGYHPFEWQPALKSVSTSVCHVGIIDGLSGWIASVDDSPADTVTRFRYDVALVSALKDLEEDIMEGL
↳ RERGLEDSACTSGFSVMIKESCDGMGDVSEKHGGGPPVPEKPVRFSTIMSVSIIAQEGEDEAITIFREPKNSEMSCKPLSLMFVDES
↳ HETLTGVLGPVVAERNAMKHSRLILSVGLSRSFRFHFRGTGYDEKVMREMGLEASGSTYICTLCDSTRAEASQNMTLHSVTRSHDEN
↳ LERYELWRTNPHSESAAEELRDRVKGVSAPFIMETQPTLDALHCDIGNATEFYKIFQDEIGEYVHKANPSREQRRSWRAALDKQLRKKMK
↳ LKPVMRMNGNYARKLMTREAVEAVCELCSEERQEALREMLGLYIQMKPVWRSTCPAKECPDELCRYFNSQRFAELLSTVFKYRYDYGK
↳ ITNYLHKTLAHVPEIVERDGSIGAWASEGNESGNKLFRRFRKMNRQSKTFELEDVL

```

Score 107036

1.8 CLUSTAL

Script 1.8.1 (python)

```

1 matrix = {}
2 file = "sample.fasta"
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()

```



```

# Alignment:
>Oncorhynchus_mykiss
CGCACCGTCAAGGCCACCAGTGGGCGTCAGATCTTCCAGCCCCTACACACCTTACGCACTGCAGAGAAGGAGCTCCTCCAGGCTACCACCCC
└ TTTGAGTGGCAGCCGGCCCTCAAGAGTGTGTCCACATCCTGCCATGTGGGGATCATTGATGGGCTATCAGGGTGGATCGCTTCGGTAGA
└ CGACTCCCCAGCAGATACAGTCACGCGACGGTTTTCGCTACGACGTGGCCCTGGTGTACGCCCTGAAGGACCTGGAGGAGGACATCATGG
└ AGGGGCTGAGAGAGCGAGGCCCTGGAGGACAGTGCTTGACCTCGGGCTTCAGCGTTATGATCAAGGAGTCTTGCGATGGTATGGGGGAC
└ GTCAGTGAGAAGCATGGCGGAGGGCCCGCTCCCGGAAAAGCCTGTGCGTTTCTCCTTACCATCATGTCCGTCTCTATTCAAGCTGA
└ GGGAGAAGATGAGGCGATCACCATTTTCCGGGAGCCCAAGCCCAACTCAGAGATGTCTGCAAGCCGCTAAGCCTGATGTTTGTGGACG
└ AGTCGGACCACGAGACTCTCACAGGCGTCTTGGGGCTGTGGTGGCCGAAAGGAATGCTATGAAGCACAGCCGTCTCATCCTGTCTGTG
└ GCGGCGCTTTCTCGCTCCTTCCGCTTCCACTTCCGGGGCAGGGCTATGATGAGAAGATGGTGCAGAGATGGAGGGTTTGGAGGCCTC
└ TGGCTCCACTTACATCTGCACGCTGTGTGACTCCACTCGGGCAGAGGCCTCCAAAACATGACTCTCCACTCTGTACCCGCGAGCCATG
└ ACGAGAACCTGGAGCGCTACGAACTTTGGAGGACCAACCCTCATTCTGAGTCAGCTGAAGAGCTGCGAGAGCCGAGTCAAAGCGTCTCT
└ GCCAAGCCCTTATGAGAGCCAGCCACACTGGAGCCCTGCACTGTGATATCGGCAATGCCACTGAGTTCTACAAGATCTTCCAGGA
└ TGAGATAGGGGAGGTCTATCACAAGGCAAAACCCAGCCGGGAGCAGCGTCGGAGCTGGCGGGCCCGCTGGACAAGCAGCTGAGGAAGA
└ AGATGAAGCTGAAGCCTGTGATGAGGATGAATGGGAACTATGCACGGAAGCTGATGACCCGGGAGGCAGTGGAGGCAGTGTGTGAGCTG
└ GTGTGCTCAGAGGAGCGTCAGGAAGCTCTGAGGGAGCTGATGGGGCTCTACATCCAGATGAAGCCTGTGTGGCGCTCCACCTGCCCGGC
└ CAAGGAGTGCCCGAGACGAGCTCTGCCGTATAGCTTCAACTCCCAACGCTTTGCAGAGCTGCTCTCCACCGTCTTCAAGTACAGGTATG
└ ACGGAAAGATCACCAACTACCTGCACAAGACCTTGGCCCATGTGCCAGAGATTGTGGAGAGGGATGGTCCATCGGGGCCTGGGCCAGC
└ GAGGGGAATGAGTCTGGGAACAAGCTGTTTCAGACGGTTCAGGAAGATGAATGCCCGCCAGTCCAAGACCTTTGAGCTGGAGGACGTGCTG
>Alligator_mississippiensis
AGAAGTGTAAAGCTGTACCCGGGAGGCAGATCTTCCAGCCATTGCATGCTCTTCGCACTGCTGAGAAAGCCCTGCTCCAGGTTATCATTCA
└ TTTGAGTGGAAAGCCTCCTTTGAAAAATGTGTCTGCTAATACAGAAAGTGGGAATTATAGATGGGCTTTCAGGCCTGCCACACACAGTTGA
└ TGACTACCCAATTGACACAATTGCAAAGAGGTTTTGATATGATGCAGCCTTGGTTTCTGCCTTAATGGATATGGAAGAAGACATCTTGG
└ AAGGCATGAAAGCACATGACCTGGATGACTATTTGAATGGT---CCTTTCACTGTGGTGGTGAAGAGTCTTGTGATGGGATGGGAGAT
└ GTCAGTGAGAAGCATGGATGTGGACCAGCAGTCCCAGAAAAGGCAGTTTCGCTTTTCTTTCACTGTCATGACCATCGCTATA-----AC
└ TCATGGCAATACAAATGTAAGGATCTTTGAAGAAGTCAAGCCCAATTGAGAGCTATGTTGCAAGCCATTGTGTCTTATGTTGGCTGATG
└ AGTCAGACCATGAGACTCTGACAGCCATCCTGAGCCCTCTCATAGCAGAAAGAGAGACCATGAAAAACAGTGTACTACTTCTTGAAATG
└ GGAGGCATCCTCAGAACATTCAAATTCATCTTTAGGGGTACAGGGTATGATGAGAACTTGTACGTGAGGTGGAAGGCCTTGAAAGCCTC
└ TGGCTCCACTTACATTTGTACCCTCTGTGATGCAACCCGCGCTGGAAGCCTCTCAGAATTTGGTCTCCTCCACTCCATAACAAGGAGTCACA
└ CTGAAAATCTAGAGCGATATGAGGTGTGGAGGTCCAACCCATACCATGAATCAGTTGATGAGCTCCGTGATAGAGTGAAGGGTGTCTTCT
└ GCAAAACCATTTATTGAGACTGTTCCCTCTATAGATGCGTTGCACTGTGACATTGGCAATGCAGCTGAGTTTTACAAGATATTCCAATT
└ TGAGATAGGTGAGGTCTACAAGAACCCTGATGCATCTAAAGAAGAGAGGAAGAGGTGGCAATCAGCTCTTGACAAGCACCTGAGGAAGA
└ AGATGAACCTGAAACCAATAATGAGGATGAATGGAACTTTGCTAGAAAGCTGATGACTAAAGAGACAGTGAAGCAGTTTGTGAATTA
└ ATAAAGTGCAGGAAAGGCATGAAGCCCTCAAAGAAGCTGATGGATCTCTACCTGAAGATGAAACCAAGTGTGGAGATCTTCATGTCTGCG
└ CAAGGAGTGCCCGAAGTGTGTGCCAGTACAGTTTCAACTCACACGTTTTGCTGAGCTCCTATCCACAAAGTTCAAGTATAGATATG
└ AGGGCAAAATTACAAATTACTTTCATAAACTCTTGCTCATGTTCTGAAATCATAGAAAGAGATGGGTCCATTGGTGCCTGGGCAAGT
└ GAAGGAAATGAATCTGGGAACAAATTGTTTAGACGTTTCCGAAAAATGAATGCCAGGCAGTCAAAATGCTATGAAATGGAGGATGTCCTG
>Gallus_gallus

```

```

CGAACTGTAAAAGCTGTCACTGGGAGGCAGATTTTCCAGCCACTGCATGCTCTTCGCACCGCCGAGAAAGCCCTCTTACCAGGTTATCATCCT
↳ TTTGAATGGAAACCTCCCTTAAAAATGTATCCACTAACACAGAAAGTGGGAATTATAGATGGACTATCAGGACTACCCCTCTCGATTGA
↳ TGACTACCCAATAGACACAATTGCAAAGAGATTTCCGATATGATACAGCCTTGTTTCTGCTTTAAAGGACATGGAGGAAGAAATCTTGG
↳ AGGGCATGAAGGCTAAAAACCTGGATGACTATCTGAATGGT---CCCTTCACTGTGGTAGTAAAGAGTGCTGTGACGGAATGGGAGAT
↳ GTCAGCGAGAAGCATGGAAGTGGTCCCTGCTGTCCCAGAGAAGGCTGTACGCTTTTCATTACAGTCATGAACATTGCTATA-----GA
↳ CCACGAGAACGAAAGGATAAGGATCTTTGAAGAAGTAAAGCCCAATTAGAGTTGTGTTGCAAACCCTTGTGTCTTATGCTGGCTGATG
↳ AATCAGACCATGAAACTCTGACAGCAATCCTGAGCCCTCTCATAGCAGAGAGAGAGGCTATGAAGAACAGTGAACCTGCTTCTTGAAATA
↳ GGAGGCATCCTGAGAACATTCAAATTCATCTTTGAGGTACAGGCTATGATGAGAACTTGTAAGGGAAGTGAAGGGCTGGAGGCCTC
↳ AGGTTCCACTTATATCTGTACCCTCTGTGATGCAACCCGCTAGAGGCCTCCAGAAATTTGGTCTTCCACTCCATCACCAGGAGCCACG
↳ CTGAAAATCTGGAGCGATATGAAATATGGAGGTCCAACCCATATCAGCAATCTGTTGATGAGCTCCGTGACAGAGTGAAGGGTGTTCG
↳ GCCAAACCTTTTATAGAGACTGTTCCCTCCATAGACGCGTTGCACTGTGACATTGGCAACGCAACAGAAATTTACAGGATTTTCCAGAT
↳ GGAGATTGGTGAAGTCTATAAGAATCCTGATGCGACTAAAGAGGAGAGGAAGAGGTGGCAGTTGACTCTTGACAAACATCTCAGGAAGA
↳ AGATGAAATTAACCTATGATGAGGATGAGTGGGAACTTTGCTAGAAAGCTCATGTCCAAAGAGACAGTAGAGGCAGTATGTGAATTA
↳ ATAAAGTGTGAGGAAAGGCATGAAGCCCTAAAAGAACTAATGGACCTTTATCTGAAAATGAAACCAGTATGGCGATCTTCGTGCCCTGC
↳ CAAAGAGTGTCCAGAATTGCTGTGCCAGTACAGCTATAATTACAGCGTTTTTGCGGAGCTCCTGTCTACCAAGTTTAAATACAGATATG
↳ AGGGCAAGATTACCAATTATTTCCACAAAACCTTGCTCATGTACCTGAAATCATTGAAAGAGATGGGTCCATTGGTGCTTGGGCAAGT
↳ GAAGGAAATGAGTCCGGAACAAACTATTTAGGAGGTTCCGGAATGAATGCGAGGCAGTCCAAATCTATGAAATGGAGGATGTCTTA
>Anolis_carolinensis
AGGACTGTAAAAGCTGTCACTGGGAGGCAGATATTCCAGCCACTCCATGCACTCCGAACTGCTGAAAAGGCCCTCTTGCCCGGTTACCATCAA
↳ TTTGAGTGGAAACCACCCTTAAAAATGTCTCCAGTAACACAGAAAGTAGGCATTATTGATGGACTGTCAGGCATACAACATTTGGTTGA
↳ TGACTACCCAGTGGACACAATTGCAAAGAGATTCGGATATGATGCAGCTTTGGCTTCTGCCTTGATGGATATGGAAGAAGACATCCTAG
↳ AAGGCCTGAAAAGACAGGACTTGGACGACTACTTCAAAGGC---CCTTCACTGTGGTGATTAAGAGTCCTGTGATGGGATGGGAGAT
↳ GTTAGTAAAAACATGGCTGTGGCCCTGCTGTTCCCTGAGAAAGCAGTTTCGATTCTCTTTCACACTCATGACCATCTCTGTC-----AC
↳ TCGTGGCAATGCAAGCATAAGGTTTTTGAAGAATGAAGCCCAATTAGAGCTGTGTTGTAAGCCCTTGTGCCTTATGCTGGCTGATG
↳ AATCAGACCATGAGACACTCACAGCCATCCTGAGTCTCTTGTGGCAGAAAGAGAGGCCATGAAAGACAGTGTACTGATACTTGATATG
↳ GCTGGAATCCCAAGAACATTCAAATTCATCTTTAGGGGCACTGGATATGATGAAAACTTGTCCTGTAAGTAGAGGTCTTGAAAGCTTC
↳ AGGCTCCACTTATATTTGCACTCTTTGTGATGCAACACGCTTAGAAGCCTCTCAGAACTTGATCCTTCATTCCATCACAAGGAGTCATG
↳ CTGAAAACCTAGAACGATATGAATTGTGGAGGACCAACCCCTACCATGAGACCGTTGATGAACTGCGTGACAGAGTCAAAGGGTTTCT
↳ GCAAAACCTTTTATGAGACCGTTCTTCAATAGATGCCTTGCACTGTGACATTGGCAATGCAGCTGAATTTTACAAGATATTCCAGTT
↳ TGAGATTGGTGAAGTCTACAAAAACACTGATGCATCAAAAGAAGAGAGAAGGAGATGGCAGTCAACCCCTTGACAAACACCTCAGAAAGA
↳ AGATGAACTTGAAAGCCTATGACAAGGATGAATGGGAATTTTGCCAGAAAGCTCATGACCAAGGAGACTGTGGAAGCAGTCTGTGAATTA
↳ ATAAAGAGTGAGGAGAGACATGAAGCCCTCAGAGAACTCATGGACCTTTACCTGAAGATGAAACCAGTGTGGCGGTCTTCATGTCCAC
↳ CAAAGAATGCCCAGAATTAGTATGCCAGTATAGCTTCAATTCTCAACGGTTTGCAAGAGCTGCTGGCTACAAAATTCGGTTACAGATATG
↳ CAGGCAAGATTACTAATTACTTTCATAAAACCTTGCTCATGTTCAGAAATATTGAAACGAGATGGATCTATTGGTGCTTGGGCAAGT
↳ GAAGGAAATGAGTCTGGGAACAACTTTTCAGACGTTTTCGAAAAATGAATGCCAGGCAATCAAAATCTATGAAATGGAAGATGTCTTA
>Ornithorhynchus_anatinus

```

```

AGGACCGTGAAAGCTATCACTGGCAGGCAGATCTTTCAGCCCTTACACTCCCTCCGCACTGCCGAGAAGGTCTCTTGCCAGGTTATCATCCA
↳ TTCGAGTGGGATCCACCCTTAAAAAATGTCTCGGCCAACACCGAGGTGGGGATCATGGACGGGCTCTCTGGGCTGCCAGTCTCGGTGGA
↳ CGACTACCCTGTAGACACCATTGCCAAAAGGTTCCGCTACGATGCGGCCCTAGTGTCCGCTTGATGGATATGGAGGAAGATATTCTGG
↳ AGGGCATGAAATCCCAAGACCTCGATGACTACCTAAGTGGC---CCCTTTACCGTAGTGATCAAAGAGTCTTGCGACGGGATGGGAGAT
↳ GTGAGTGAAAAGCACGGGAGCGGCCCGGCAGTCCCTGAAAAGCGGTTTCGCTTTTCTTTCACCGTCATGAATATCACTCTG-----GC
↳ TTATGAGCAAGAGAACGTAAAGATCTTCGAAGAGGCCAAGCCAACTCAGAGCTGTGCTGCAAACCACTCTGTCTGATGCTGGCTGATG
↳ AATCGGATCATGAGACCCTGACAGCCATCTTGAGCCCACTTATAGCAGAGAGAGAGGCCATGAAGGACAGCGAATTGAAGCTGGAGATG
↳ GGAGGCATCCTGAGGTCTTCAGATTCTCTTTAGGGGCACTGGCTATGATGAGAACTCGTTCGGGAAGTGGAAGGTCTCGAGGCTTC
↳ TGGCTCGGTCTACATTTGCACTCTCTGTGATGCCACCCGCTGGAAGCATCTCAAAATCTCGTCTTCACTCCATAACCCGGAGTCACG
↳ CTGAGAACCTGGAGCGCTATGAAGTTTGGAGATCCAATCCCTTCCACGAGTCTGTGGAAGAACTGCGGGACAGGGTGAAAGGGTTTCG
↳ GCAAAGCCATTTCGAGACCGTTCTTCCATAGACGCACTCCACTGTGATATTGGCAACGCGGCCGAGTTTTATAAGATTTTCCAGCT
↳ GGAGATCGGGGAGGCATATAAGAACCCCAACGCGTCCAAAGAGGAAAAGAAAAGATGGCAAGCCACCCTAGACAAACACCTCCGGAAGA
↳ AGATGAAGCTGAAGCCAATCATGAGGATGAATGGCAATTTTGCTAGGAAGCTGATGACCAAAGAGACAGTGAAGCAGTCTGTGAATTG
↳ GTTCACTGTGAAGAGAGGCATGAAGCCCTGAGGGAGCTGATGGACCTTTATCTGAAGATGAAACCAGTGTGGCGCTCATCTGTCCCGC
↳ CAAAGAGTGCCCGGAGTCTCTGTGCCAGTACAGTTTCAACTCACAGCGCTTTGCCGAGCTCCTGTCTACCAAGTTCAAATACAGATATG
↳ AGGGCAAAATCACTAATTACTTCCACAAAACCTCTGGCTCATGTTCTGAGATCATAGAGCGGGATGTTCCATTGGGGCCTGGGCTAGT
↳ GAGGGGAATGAGTCTGGGAACAACTGTTAGGCGTTTCCGAAAAATGAATGCCAGGCAGTCCAAGTGTTATGAGATGGAGGATGTTTTG
>Homo_sapiens
AGGACTGTGAAAGCCATCACAGGGAGACAGATTTTTCAGCCTTTGCATGCCCTTCGGAATGCTGAGAAGGTACTTCTGCCAGGCTACCACCAC
↳ TTTGAGTGGCAGCCACCTCTGAAGAATGTGTCTTCCAGCACTGATGTTGGCATTATTGATGGGCTGTCTGGACTATCATCTCTGTGGA
↳ TGATTACCCAGTGGACACCATTGCAAAGAGGTTCCGCTATGATTACGCTTTGGTGTCTGCTTTGATGGACATGGAAGAAGACATCTTGG
↳ AAGGCATGAGATCCCAAGACCTTGATGATTACCTGAATGGC---CCCTTCACTGTGGTGGTGAAGGAGTCTTGTGATGGAATGGGAGAC
↳ GTGAGTGAGAAGCATGGGAGTGGGCCTGTAGTTCCAGAAAAGGCAGTCCGTTTTTCATTACAATCATGAAAATTACTATT-----GC
↳ CCACAGCTCTCAGAATGTGAAAGTATTTGAAGAAGCCAAACCTAACTCTGAACTGTGTTGCAAGCCATTGTGCCTTATGCTGGCAGATG
↳ AGTCTGACCACGAGACGCTGACTGCCATCCTGAGTCTCTCATTGCTGAGAGGGAGGCCATGAAGAGCAGTGAATTAATGCTTGAGCTG
↳ GGAGGCATTCTCCGACTTTCAAGTTCATCTTCAGGGGCACCGGCTATGATGAAAACTTGTGCGGAAGTGGAAGGCCTCGAGGCTTC
↳ TGGCTCAGTCTACATTTGTACTCTTTGTGATGCCACCCGTCTGGAAGCCTCTCAAAATCTTGTCTTCACTCTATAACCAGAAGCCATG
↳ CTGAGAACCTGGAACGTTATGAGGTCTGGCGTTCCAACCTTACCATGAGTCTGTGGAAGAACTGCGGGATCGGGTGAAAGGGTCTCA
↳ GCTAAACCTTTTCATTGAGACAGTCCCTTCCATAGATGCACTCCACTGTGACATTGGCAATGCAGCTGAGTTCTACAAGATCTTCCAGCT
↳ AGAGATAGGGGAAGTGTATAAGAATCCCAATGCTTCCAAAGAGGAAAAGGAAAGGTGGCAGGCCACACTGGACAAGCATCTCCGGAAGA
↳ AGATGAACCTCAAACCAATCATGAGGATGAATGGCAACTTTGCCAGGAAGCTCATGACCAAAGAGACTGTGGATGCAGTTTGTGAGTTA
↳ ATTCTTCCGAGGAGAGGCACGAGGCTCTGAGGGAGCTGATGGATCTTTACCTGAAGATGAAACCAGTATGGCGATCATCATGCCCTGC
↳ TAAAGAGTGCCCGAAGTCCCTCTGCCAGTACAGTTTCAATTACAGCGTTTTGCTGAGCTCCTTTCTACGAAGTTCAAGTATAGGTATG
↳ AGGGAAAAATCACCAATTATTTTCAAAAACCTGGCCCATGTTCTGAAATTATTGAGAGGGATGGCTCCATTGGGGCATGGGCAAGT
↳ GAGGGAAATGAGTCTGGTAACAACTGTTTAGGCGCTTCCGAAAAATGAATGCCAGGCAGTCCAATGCTATGAGATGGAAGATGTCCTG
>Alytes_obstetricans

```

```

AGAACTGTTAAAGCTACAACCTGGCAGACAGATCTTCCAGCCATTGCATGCCTTGAGGAATGCAGAGAAGGCCTTATTGCCAGGGTATCATCCT
↳ TTTGAATGGAAGCCACCTCTGAAAAATGTATCCACTTGTACAGACACTGGGATTATTGATGGACTTTCTGGACTGAACCGGTCCATAGA
↳ TGAATACCCCTGTGGAAGCCATTTCAAAAAGGTTTAGGTATGATACCGCTCTAGTGTCAGCTTTAAAAGACATGGAGGAAGATATTCTAG
↳ AAGGCTTGAGGTCTCATGACATGGATGATTACTTAAATGGT---CCCTTCACAGTGGTCATTAAAGAGTCTTGTGATGGGATGGGAGAT
↳ GTGAGTGAAAAACATGGCAGCGGACCAGCTGTGCCTGAAAAGGCTGTCCGATTTTCTTTCACAGTAATGTACATCAGTGTG-----CC
↳ AAACAACAATGAATGTGTGAGGATCTTCGATGAGACCAAACAACTCTGAGTTGTGCTGCAAGCCCCCTGTGCTTAATGCTAGCAGATG
↳ AATCTGATCATGAGACCCTGACTGCTATACTAAGTCTCTCATAGCAGAAAGAGAGGCCATGAAAACCACTGAGCTAATGCTGGAATG
↳ GGAGGTATTCTCAGGAACCTTCAAATTTATATTCCGTGGCACAGGATATGATGAGAAGCTGGTGCAGGGAAGTGGAGGGACTGGAAGCATC
↳ TGGCTCGGTCTACATCTGTACATTATGTGATTCCACCCGCTGGAAGCTTCTCAGAACTTGGTTTTCCACTCTATAACCAGGTGCCACA
↳ CTGAGAACCTGCAGCGCTATGAGACGTGGAGAGCCAATCCGCACCATGAATCTGTGATGAGCTGCGAGACCGAGTTAAAGGGGTTTCT
↳ GCTAAGCCATTTATTGAAACCTCCCATCAATTGATGCATTGCACTGTGATATTGGAATGCGGCAGAGTTTTACAGACTTTTTTCAGCT
↳ GGAGATAGGGGAGGTCTACAAAAACCCCAATGCCACCAAAGAGGAGAGAAAGAGATGGCAGTCAACCCCTTGACAAGCACCTAAGAAAAA
↳ AGATGAACCTTAAGCCCATTATGAGAATGAATGGGAACTTTGCCAGGAAGCTTATGAGTAAAGAGACTGTTGAAGCTGTCTGTGAATTA
↳ GTGCACAGTGAGGAGAGGCAAGAAATCTTAGAGAGCTGATGGACTTGTACCTTAAATGAAACCCGTATGGCGCTCGTCTGCCCTGC
↳ CAAAGAGTGTCCAGAACTGCTGTATCAATACAGCTTCCACTCCCAGCGCTTTGCTGAGTTACTCTCCACCAAGTTCAAGTACAGATACG
↳ CAGGGAAGATCACAACTACTTCCACAAAACCTCTAGCCCATGTGCCAGAAATCATTGAGCGGACGGCTCCATTGGTGCTTGGGCCAGC
↳ GAAGGCAATGAGTCCGGTAACAACTCTTCAGGCGCTTCCGTAATAATGAATGCCAGGCAGTCCAAATTTTATGAGATGGAAGATGTATTG
>Latimeria_menadoensis
AGGACAGTGAAAGCTACAACCTGGAAGCAGATCTTCCAACCATTGCATTCTCTTCGAAATGCAGAGAAGGCTCTTCTTCCAGGTTTCCATCCT
↳ TTTGAGTGGCAACCTCCTTTAAAAAATGTGTCCAGTACAACCTGAAGTTGGCATCATTGATGGAATGTCTGGAATGACACAGTTTGTGTA
↳ TGAATACCCACTAGACACAATTTCAAAAAGATTCAAGTATGATGCAGCTTTGGTTTCAGCCTTAAAGGACCTTGAGGAAGAGTTACTTA
↳ AGGGACTGATAGAGGAAGACCTGGAAGACTATCTAAGTGGC---CCGTTACAGTCATAATTAAGAGTCTTGTGATGGTATGGGAGAT
↳ GTCAGTGAGAAGCATGGAAGTGGGCCAGCAGTCCCTGAAAAGCAGTAAGGTATTCTTTTACAATCATGACCATCAGTGTG---GCAAA
↳ CAGTCATAATGAGAATGTAACCATTTTGAAGAAGGCAAGCCAACTCAGAGCTGTGTTGCAAGCCCCCTATGTCTCATGCTTGCTGACG
↳ AATCTGATCATGAGACACTGACTGCCATCTTGGGTCTGTCTATTGCTGAAAGAGAGGCAATGAAAAACAGTGAAGTCTTCTTGAAATG
↳ GGGGGAATCCTCAGGTCCTTCAAATTTATCTTCCGAGGCACTGGATATGATGAAAAGCTCATAAGAGATGTTGAAGGTCTCGAGGCTTC
↳ AGGTTCAAGTTATATTTGCACCCTCTGTGATTCCACACGCAGTGAAGCTTCACAAAACCTCATTCTTCATTCTATAACGAGGAGTCATA
↳ AAGAGAATCTGGAAGGTATGAAATTTGGAGGTCCAACCCATACCAAGAGCCAGTAGAAGAACTGCGTGACAGAGTTAAGGGGGTTTCA
↳ GCCAAACCATTCATTGAAACTCTGCCTTCAATAGATGCCTTACATTGCGATATTGGGAATGCAACTGAGTTCTACAAAATATTCCAAGA
↳ TGAGATAGGGGAAATTTACAAAAACCCCTAACCTTCCAGAGAAGAAAAAAGAGATGGCATTTCAGTTCTTGATAAACATCTCAGGAAAA
↳ ATATGAATTTAAAGCCGGTCATGAGAATGAATGGGAATTATGCAAGGAAATTAATGACAAAGGAGACAGTGAATGCAGTATGTGAGTTG
↳ ATTCCTTCTGAAGAGAGACAGGAAGCCCTCAAGGAACTGGTGGACCTCTATTGAAAATGAAACCACTATGGCGTTCCACCTGCCAGC
↳ CAAAGAATGCCAGAGTTGTTATGCCAATACAGCTTCCACTCTCAACGATTTGCTGAGCTTCTGTCTACAATGTATAGGTATAGATATG
↳ AAGGGAAAAATCACAACTATCTTCACAAAACCTCTGGCTCATGTTCCAGAAATATTGAAAGAGATGGTCCATTGGGGCCTGGGCAAGC
↳ GAGGTAATGAGTCAGGGAACAACTATTTAGACGCTTCAGAAAAATGAATGCCCGGCAGTCCAAGTATTATGAACTGGAGGATGTCCTA
>Protopterus_dolloi

```

```

CGAACAGTGAAAGCTGCAACTGGAAGGCAGATCTTTCAACCATTACATGCTCTTAGAAGTGCTGAAAAGGCCCTCCTTCAGGGTACCATCCA
↳ TTTGAATGGCAGCCTCCTCTAGTAGGTGTTTCCTCAAGTACAGATGTTGGAATCATAAATGGGTTGTCTGGTCTGACTTCTTCTGTTGA
↳ TGAATACCCAGTAGAAGCCCTTGCAAAGAGATTTAGATATGATGCAGCTTTGGTTTCTGCTTTAAAGGACATTGAGGAAAACATCCTAG
↳ AAGGTATGAAACAAAATGGCTTGGATGAGTATCTGAGTGGC---CCTTTCACTGTAGTGATAAAAGAATCTTGTGATGGTATGGGAGAT
↳ GTCAGTGAAAAACATGGAAGTGGCCACCAGTTCCAGAAAAAGCTGTGAGGTTTTCTTTTACAATTATGTCTATTCTGTA---GCAAT
↳ GAGTGACAGTGAAAAATGTACAGATTTTTGAAGAATCAAGCCTAATTCTGAACTCAGCTGTAAACCACTATGCCTCATGATTGCTGATG
↳ AATCTGACCATGAAACATTGACTGTCTCTTGGGCCCTGTGATTGCAGAACGTGAAGCCATGAAAACAAGTGAAGTATGCTTGAATTA
↳ GGAGGCATTCTGAGAACTTTCAAATTTTTCTTCAGAGGTACAGGTATGATGAAAAGCTTGTGAGAGAGGTTGAAGGATTAGAGGCTTC
↳ AGGTTACATTATATATGTACCCTCTGCGATGCGACTCGACAAGAAGCCTCTCGGAACTTAGTTCTTCATTCAATAACAAGGAGCCATG
↳ CTGAAAACCTGGAGAGGTATGAAGTATGGAGGTCAAACCCATATAACGAATCAGTTGATGAACTGCGTGACAGAGTAAAGGGTGTTC
↳ GCCAAGCCTTTTCATTGAAACACGGCCCTGCATTGATGCATTACAGTGTGATATTGGAATGCAACTGAATTTTATAAAATATTTCAAGA
↳ TGAGGTTGGTGAAGTATACAAACGTCCAAATCCATCAAAAGAAGACAGAAAAAGATGGCATATGACCCTTGACAAACACCTCAGGAAGA
↳ AGCTAAGCCTTAAGCCAGTTATGAGAATGAATGGAATTTTGCAAGAAAACCTATTACAAAAGAGGCAGTAGATGCTGTGTGTGAACTG
↳ ATTCCTTCAGAAGAAAGCGTGCAGCTATCAGGGATTTGGTACACCTTTATATGCTTATGAAACCAGTATGGCGATCTACTTATCCAGC
↳ CAAAGAGTGCCAGAGTTGCTCTGCCAGTACAGCTTTAATTACAAAAGGTTTGCAAGAACTGTTATCTACTAAATTCCAGTACAGATATG
↳ AAGCCAAAATTACAACTACCTTCACAAGACTTTGGCTCATGTTCTGAAATAATTGAAAGGGATGGGTCTATAGGTGCTTGGGCAAGT
↳ GAAGGCAATGAGTCGGGTAAATAAACTGTTCCGCAGATTCAAAAAATGAATGCTAGACAGTCAAAATGCTATGAATTGGAAGACGTATTG
>Carcharhinus_leucas
AAGACAGTTAAGGCTATTACTGGGAAGCAGATCTTCAGCCATTACACGCACTCCGAAATGCAGAGAAAACATTATTGCCAGGTTATTATTCA
↳ TTCGAATGGCAGCCTCCTCTCGAAACATCTCCACAAACACACGAGTGGGGATAATCGATGGTTTGAGTGGTTGGTTTCAGTGTGTTGA
↳ TGATTACCCAATGAAACCATATCAAGGCGGTTGTCTACGACGTAGCGCTTGCTTCAGCTGTGAAGGAGATGGAGGACGATATCTTGG
↳ AAGGGTTAAGAAGCCAAAATGTAGATGAGTTTGTGAGTGA---CCGTTTACGTTGTAATCAAGGAATCTTGTGATGGGATGGGGAT
↳ GTGAGTGAAAAACATGGATGCGGTCCACAGTTCCAGAAAAGGCAGTGAGATATTCCTTACAATCATGAGTATAAGTGTG---ATGAA
↳ TGAAAAAATGAAAAAGTGAAGGTGTTTGAAGAGATGAAACCAAAATTCWGAGCTGTGCTGTAGGCCACTTTGCCTGATGTTAGCTGATG
↳ AATCTGACCGTGAACATTGACCGCAATACTGGGGCCAGTAATTGCTGAAAGACAATCAATGAAGACCAGTGATTAAATCGTTGAAATT
↳ GGAGATTGTGACCGATCTTTCCAGTTTATCTTTGAGGCACGGGGTACGATGAGAAGCTTGTTTCGGGAAGTGAAGGACTTGAAGCCTC
↳ GGGTTCTATCTATATCTGCACATTGTGTGACTCCACACGAAGTGAAGCTTCAAAGAACATGTTTCTTCATTCCATAACAAGGAACCATG
↳ CAGAGAACTTGGAACGCTATGAAATCTGGCGTTCAAACCCATATACGAGACAGCTGATGAGTTACGTGATAGAGTGAAAGGAGTTTCA
↳ GCGAAACCATTTATAGAACTCAACCCTCTATAGATGCCTTGCAATTGTGACATTGGAATGCGACAGAATTCTACAGAATCTCCAGGA
↳ TGAGATTGGGGAAGTATACAAAAATTCAAATTCATCAAAGGAGGAAAGGAAAAGGTGGCAATCGATGCTGGATAAAACACTTGAGGAAAA
↳ AGATGAATTTGAAACCTATAATGAGGATGAATGGAACCTTTGCAAGAAAACCTAATGACAAAAGAGACTGTGGAAGCCGTGTGTGAACTG
↳ ATTCCATCTGAAGAAAGGAGAGAGATTCTCAGAGAGTTAATGCACCTTTATCTACTTATGAAACCAGTATGGCGTCCACATTTCACAC
↳ TACGGAATGTCCAGACCTTCTGTGTGAGTACAGTTTAAATCCCAGAGATTTGCAGAACTGCTTCATACAGAATTCAGTCACAGATATG
↳ AAGGGAATAACCAACTATCTGCACAAGACCCTCGCTCATGTCCCTGAGATCATCGAGAGAGACGGTTCCATAGCGCGTGGGCAAGT
↳ GAAGGCAATGAATCGGGCAACAACTTTTCCGTCGTTCAAGAAAATGAATGCAAGGCAGTCCAAAAGCTATGAGTTGGAAGACATTTTG

```

1.8.1 Verification against CLUSTALW

As I can see, one of the sequences S3 : PPDGKSDS-- has the two gaps at the end of the sequence. With our algorithm the gaps are inside favoring the match of the last S to the corresponding final S of the rest of sequences.

The scores printed on the output also are different, very similar but not at all. Nor the guide trees. Further insights into CLUSTAL code could provide the clues to understand where the discrepancies came from.

This discrepancies are due to subtle differences in the algorithms, for sure.

Script 1.8.2 (text)

```
1 %%bash
```

```

2 ~/clustalw2 -OUTPUTTREE=phylip -NEGATIVE -INFILE=sample.fasta -OUTORDER=ALIGN
  ↳ -STATS=align.log -TREE -ALIGN -CLUSTERING=NJ -OUTFILE=align.fasta -OUTPUT=CLUSTAL
  ↳ -MATRIX=BLOSUM -TYPE=PROTEIN -PWGAPOPEN=10 -PWGAPEXT=5
3 cat align.fasta
4 cat sample.dnd

```

Output

CLUSTAL 2.1 Multiple Sequence Alignments

Sequence type explicitly set to Protein

Sequence format is Pearson

Sequence 1: laboA 57 aa

Sequence 2: 1ycsB 60 aa

Sequence 3: 1pht 80 aa

Sequence 4: 1vie 51 aa

Sequence 5: 1ihvA 49 aa

Start of Pairwise alignments

Aligning...

Sequences (1:2) Aligned. Score: 22

Sequences (1:3) Aligned. Score: 12

Sequences (1:4) Aligned. Score: 5

Sequences (1:5) Aligned. Score: 6

Sequences (2:3) Aligned. Score: 11

Sequences (2:4) Aligned. Score: 9

Sequences (2:5) Aligned. Score: 4

Sequences (3:4) Aligned. Score: 15

Sequences (3:5) Aligned. Score: 12

Sequences (4:5) Aligned. Score: 6

Guide tree file created: [sample.dnd]

There are 4 groups

Start of Multiple Alignment

Aligning...

Group 1: Delayed

Group 2: Delayed

Group 3: Delayed

Group 4: Delayed

Alignment Score -155

CLUSTAL-Alignment file created [align.fasta]

CLUSTAL 2.1 multiple sequence alignment

```

1aboA      -NLFVALYDFVASGDNTLSITKGEKLRVLGY-----NHNG-----EWCEAQ--TKN
1ycsB      KGVIIYALWDYEPQNDDELPMKEGDCMTII-----HREDEDEIEWWWAR--LND
1pht       GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNQYNETTG

```

```

1vie      -----DRVRKKSGAAWQG-----QIVGWYCTN---LT
1ihvA     -----NFRVYYRDSRDPVWKGPakLLWK
                                     *

```

```

1aboA     GQGW-----VPSNYI--TPVN-----
1ycsB     KEGY-----VPRNLLGLYP-----
1pht      ERGD-----FPGTYVEYIGRKKISP---
1vie      PEGYAVESEAHPGSVQIYPVAALERIN-----
1ihvA     GEGAVVIQDNSD-----IKVVP RRKAKIIRD
          .*                :

```

```

(
(
1aboA:0.38808,
1ycsB:0.38385)
:0.08490,
(
1pht:0.39594,
1vie:0.44720)
:0.00848,
1ihvA:0.47811);

```