

T-COFFEE

Elena Montenegro, Nicolás Manosalva, Luis Cervera, Fernando Freire

January 5, 2019

Contents

1	Multiple Sequence Alignment: T-COFFEE	2
1.1	Pairwise alignment basis	2
1.1.1	Pairwise software	4
1.2	From pairs to sets	18
1.3	Which order? Guide trees.	19
1.4	Putting it all together: the CLUSTAL way.	19
1.5	From CLUSTAL to T-COFFEE	19
1.6	Beyond T-COFFEE	22
1.7	MSA software	22
1.7.1	MSA generic methods	22
1.7.2	T-COFFEE methods	32
1.7.3	Main MSA method	36
1.8	Some runs	40
1.8.1	T-COFFEE	40
1.8.2	CLUSTAL homemade	42
1.8.3	CLUSTALW official build.	44

1 Multiple Sequence Alignment: T-COFFEE

Authors: *Elena Montenegro, Nicolás Manosalva, Luis Cervera, Fernando Freire*

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 aspects 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.

We are going to cover the topics in this order:

1. Pairwise alignment basis.
2. From pairs to sets.
3. Which order? Guide trees.
4. Putting all together: CLUSTAL way.
5. From CLUSTAL to T-COFFEE.
6. Beyond T-COFFEE.

Sample file

This is a sample file of protein sequences that we use in the across the explanation.

Script 1.0.1 (text)

```
1 %%writefile sample.fasta
2 >1aboA
3 NLFVALYDFVASGDNTLSITKGEKLRVLGYNHNGEWCEAQTKNQGQWVPS
4 NYITPVN
5 >1ycsB
6 KGVIYALWDYEPQNDELPMKEGDCMTIIHREDEDEIEWWWARLNDKEGY
7 VPRNLLGLYP
8 >1pht
9 GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIG
10 WLNQYNETTGERGDFPGTYVEYIGRKKISP
11 >1vie
12 DRVRKKSGAAWQQQIVGWYCTNLTPEGYAVESEAHPGSVQIYPVAALERI
13 N
14 >1ihvA
15 NFRVYYRDSRDPVWKGPAKLLWKGEAVVIQDNSDIKVPRRKAKIIRD
```

Output

```
Overwriting sample.fasta
```

1.1 Pairwise alignment basis

These type of algorithms are in the kernel of all MSA methods and participating in one or more steps.

Two approaches

1. Precise methods (dynamic programming):

- Global alignment (Needleman-Wunsch)
- Local alignment (Smith-Waterman)

2. Heuristic methods:

- BLAST
- ...

Optimization criteria: All searching the highest alignment score or the lowest penalty score

Scoring

1. Amino acid or base independent:

- fixed match, no match, gap init, gap continuation, gap close
- implemented as parameter variables

2. Amino acid or base dependent:

- match, no match based on biological concepts:
- BLOSUM, PAM, ... for amino acid
- JC, F81, K80, ... for bases
- implemented as matrices or dictionaries (hashes)
- no biological basics to score gaps
 - affine gaps
 - gap initiation
 - gap continuation
 - other functions
 - penalties based on gap position
 - CLUSTAL scores to 0 final gap sequences
 - increase the complexity of algorithms (memory and CPU)

A common set of parameters of pairwise align software.

```
"""
    mode (str): Computation mode:
        'GLOBAL'          Global Alignment
        'LOCAL'           Local Alignment
        'LONG_SUBSTRING'  Obtain long common substring
    score_match (int): Score of match characters.
    score_no_match (int): Score of no match characters.
    score_gap_ini (int): Score of gap init.
    score_gap_cont (int): Score of gap continuation.
    subst_matrix (dict: tuples, float): Substitution matrix
"""
```

Some gap complexities

```
"""
```

```
For each cell, which coordinates are (i, j, ini_gap), calc the maximum score path from  
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, seting a gap in seq0 and advance seq1. Horizontal displacement.
- 3) To (i + 1, j, 0), that is, seting a gap in seq1 and advance seq0. Vertical displacement.

```
"""
```

1.1.1 Pairwise software

From the class below we'll take the global and local alignment methods that our implementation of t-coffee will use.

Script 1.1.1 (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 calculated 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 maximum
    → 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 (right),
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 as stack relatives to
71     char positions on original sequences
72     subst_matrix (dict: tuples, float): Substitution matrix
73     matrix_mode (str): If "SUBST" it's a substitution matrix, if not it's a weight matrix
74     with the keys
75         i = position of first sequence in stack
76         j = position of second sequence on stack

```

```

77         pos_i = coordinate of char on first sequence
78         pos_j = coordinate of char on second sequence
79         and the value is the weight to score this position
80         if not match, the score is 0.
81     """
82
83     def __init__(self, sequences, mode="ALIGN", score_match=SCORE_MATCH,
84         ↪ score_no_match=SCORE_NO_MATCH,\
85             score_gap_ini=SCORE_GAP_INI, score_gap_cont=SCORE_GAP_CONT, subst_matrix={}
86     ):
87         """Init parameters of alignment"""
88         self.set_sequences(sequences)
89         self.set_stacks()
90         self.len_seq0 = len(self.sequences[0])
91         self.len_seq1 = len(self.sequences[1])
92         self.init_stores()
93         self.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
94         self.set_mode(mode)
95         self.score = 0
96         self.matches = 0
97         self.unmatches = 0
98         self.gaps = 0
99         self.align_seq0 = ""
100        self.align_seq1 = ""
101        self.matching = ""
102        self.ini_time = 0
103        self.finish_time = 0
104        self.set_subst_matrix(subst_matrix)
105        self.set_matrix_mode()
106
107    def init_stores(self):
108        """Init dictionary that store temp data of the alignment"""
109        self.score_store = {}
110        self.matches_store = {}
111        self.gaps_store = {}
112        self.max_score_index = (0, 0, 0)
113        self.max_score = 0
114        self.forward_arrow = {}
115
116    def set_sequences(self, sequences):
117        """Update the target sequences of the alignment"""
118        self.sequences = sequences
119
120    def set_stacks(self, stack_0=[], stack_1=[],\
121        ↪ stack_0_indexes=[], stack_1_indexes=[], stack_0_refs=[], stack_1_refs=[]):
122        """Update the stacks for msa"""
123        self.stacks = [stack_0, stack_1]
124        self.stacks_indexes = [stack_0_indexes, stack_1_indexes]
125        self.stacks_refs = [stack_0_refs, stack_1_refs]
126
127    def set_matrix_mode(self, mode="SUBST"):
128        """Update matrix mode"""

```

```

127         self.matrix_mode = mode
128
129     def set_subst_matrix(self, subst_matrix={}):
130         """Update the score matrix"""
131         self.subst_matrix = subst_matrix
132
133     def set_scores(self, score_match=SCORE_MATCH, score_no_match=SCORE_NO_MATCH,\
134                   score_gap_ini=SCORE_GAP_INI, score_gap_cont=SCORE_GAP_CONT):
135         """Update the weight scores of the alignment"""
136         self.score_match = score_match
137         self.score_no_match = score_no_match
138         self.score_gap_ini = score_gap_ini
139         self.score_gap_cont = score_gap_cont
140
141     def set_mode(self, mode="ALIGN"):
142         """Set computation mode"""
143         self.mode = mode
144
145     def forward_track(self, index):
146         """Calc alignments in forward direction.
147
148         The alignment strings are calculated from init cell (0,0) in global
149         alignments or maximum score cell in local alignments.
150
151         In local mode it's necessary to extend the alignments (local) to the total
152     → length of
153         the sequences to show the location of the alignment, and in order to compare with
154         BioPython outputs.
155
156         Args:
157         index (tuple of int): Cell coordinates of the starting cell
158
159         Returns:
160         string: align sequence 0 (bottom) for printing purposes
161         string: align sequence 1 (top) for printing purposes
162         tuple of int: Coordinates of the last cell
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]

```

```

178         i, j, gap_ini = i , j + 1, 0
179     elif arrow == "v":
180         ret_align_seq0 += self.sequences[0][i]
181         ret_align_seq1 += "-"
182         i, j, gap_ini = i + 1, j, 0
183     #compute the complete align in local mode
184     if self.mode == "LOCAL":
185         ret_align_seq0 = self.sequences[0][0:index[0]] + \
186             ret_align_seq0 + self.sequences[0][ret_final_pos[0]:]
187         ret_align_seq1 = self.sequences[1][0:index[1]] + \
188             ret_align_seq1 + self.sequences[1][ret_final_pos[1]:]
189         diff_pos_ini = index[1] - index[0]
190         if diff_pos_ini > 0:
191             ret_align_seq0 = '-' * diff_pos_ini + ret_align_seq0
192         else:
193             ret_align_seq1 = '-' * -diff_pos_ini + ret_align_seq1
194         diff_len = len(ret_align_seq1) - len(ret_align_seq0)
195         if diff_len > 0:
196             ret_align_seq0 += '-' * diff_len
197         else:
198             ret_align_seq1 += '-' * -diff_len
199     return ret_align_seq0, ret_align_seq1, ret_final_pos
200
201 def calc_matching(self, align_seq0, align_seq1, ini_pos=(), final_pos=()):
202     """Calc matching string
203
204     The matching string is the string line to print between the top and
205     bottom alignment strings. It contains the match (/), no match (.) and
206     gap ( ) indicators.
207
208     Args:
209         align_seq0 (string): Bottom sequence
210         align_seq1 (string): Top sequence
211         ini_pos (tuple of int): Initial cell coordinates
212         final_pos (tuple of int): Final cell coordinates
213
214     Returns:
215         string: Matching string
216
217     """
218     count = 0
219     ret_matching = ""
220     diff_pos_ini = ini_pos[1] - ini_pos[0]
221     if diff_pos_ini > 0:
222         delta_pos = diff_pos_ini
223     else:
224         delta_pos = 0
225     for n, (i, j) in enumerate(zip(align_seq0, align_seq1)):
226         if self.mode == "LOCAL" and not (n >= ini_pos[0] + delta_pos and n <
227             ↪ final_pos[0] + delta_pos):
228             ret_matching += ' '
229         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     Args:
241         key (tuple of int): Cell coordinates
242         score (int): Cell score
243         matches (int): Cell matches
244         gaps (int): Cell gaps
245     """
246     self.score_store[key] = score
247     super_score = score * COMPAC + 10 * matches
248     if super_score > self.max_score:
249         self.max_score_index = key
250         self.max_score = super_score
251     self.matches_store[key] = matches
252     self.gaps_store[key] = gaps
253
254 def calc_score_binary(self, seq_0, seq_1, i, j, seq_0_index=0, seq_1_index=1, pos_0=0,
255     ↪ pos_1=0):
256     """Compute alignment scores for two sequences
257     If there are a substitution matrix (actually dictionary) defined,
258     the scores are computed from the dictionary.
259     Args:
260         seq_0 (int): Sequence 0
261         seq_1 (int): Sequence 1
262         i (int): Sequence 0 char index
263         j (int): Sequence 1 char index
264         seq_0_index (int): Sequence 0 index on original sequences (MSA)
265         seq_1_index (int): Sequence 1 index on original sequences (MSA)
266         pos_0 (int): Sequence 0 index on stack 0
267         pos_1 (int): Sequence 1 index on stack 0
268     """
269     if self.subst_matrix:
270         if self.matrix_mode == "SUBST":
271             #print("PAIR", i, j, seq_0[i], seq_1[j])
272             subst_matrix_index = (seq_0[i], seq_1[j])
273             subst_matrix_index_swap = (seq_1[j], seq_0[i])
274             if subst_matrix_index in self.subst_matrix:
275                 matrix_score = self.subst_matrix[subst_matrix_index]
276             elif subst_matrix_index_swap in self.subst_matrix:
277                 matrix_score = self.subst_matrix[subst_matrix_index_swap]
278             else: #weight matrix
279                 if pos_0 in self.stacks_refs and i in self.stacks_refs[pos_0]:
280                     i_orig = self.stacks_refs[pos_0][i]

```

```

279         else:
280             i_orig = i
281             if pos_1 in self.stacks_refs and i in self.stacks_refs[pos_0]:
282                 j_orig = self.stacks_refs[pos_1][j]
283             else:
284                 j_orig = j
285             if i_orig in self.subst_matrix[seq_0_index][seq_1_index] and \
286                 j_orig in self.subst_matrix[seq_0_index][seq_1_index][i_orig]:
287                 matrix_score = self.subst_matrix[seq_0_index][seq_1_index][i_orig][j_ori
g]
288             else:
289                 matrix_score = 0
290             # Gaps in almost one of the sequences. This case only arises in MSA
291             # There is no matrix related entry. If matrix is a weight matrix we compute
292             # as zero (as defined in T-Coffee)
293             if seq_0[i] == "-" or seq_1[j] == '-':
294                 inc_matches = 0
295                 if self.subst_matrix:
296                     if self.matrix_mode == "SUBST":
297                         inc_score = self.score_gap_cont
298                     else:
299                         inc_score = 0
300                 else:
301                     inc_score = self.score_gap_cont
302             else:
303                 if seq_0[i] == seq_1[j]:
304                     if self.subst_matrix:
305                         inc_score = matrix_score
306                     else:
307                         inc_score = self.score_match
308                     inc_matches = 1
309                 else:
310                     if self.subst_matrix:
311                         inc_score = matrix_score
312                     else:
313                         inc_score = self.score_no_match
314                     inc_matches = 0
315
316             return inc_score, inc_matches
317
318     def calc_score(self, i, j):
319         """Compute alignment scores.
320         If there are stacks associated with the sequence, we compute the score weighing the
321         scores of the stacks (SOP: Score of Pairs). Stacks contains also the guiding
322         ↪ sequences.
323             Args:
324                 i (int): Sequence 0 index
325                 j (int): Sequence 1 index
326         """
327         if self.stacks == [[], []]:
328             return self.calc_score_binary(self.sequences[0], self.sequences[1], i, j, 0, 1,
329             ↪ 0, 0)

```

```

328     else:
329         computed_score = 0
330         computed_matches = 0
331         nvalues = 0
332         for pos_0, (seq_0, index_0) in enumerate(zip(self.stacks[0],
333             ↪ self.stacks_indexes[0])):
334             for pos_1, (seq_1, index_1) in enumerate(zip(self.stacks[1],
335                 ↪ self.stacks_indexes[1])):
336                 score, matches = self.calc_score_binary(seq_0, seq_1, i, j, index_0,
337                     ↪ index_1, pos_0, pos_1)
338                 computed_score += score
339                 computed_matches += matches
340                 nvalues += 1
341         ret_score = computed_score / nvalues
342         ret_matches = computed_matches / nvalues
343         return ret_score, ret_matches
344
345     def align(self, i=0, j=0, ini_gap=1):
346         """Recursive align of sequences
347         For each cell, which coordinates are (i, j, ini_gap), calc the maximum score path
348         ↪ from
349         three alternative displacements:
350
351         1) To (i + 1, j + 1, 1), that is, matching or no matching the seq0(i) and seq1(i)
352         ↪ characters.
353         This is a diagonal displacement.
354         2) To (i, j + 1, 0), that is, seting a gap in seq0 and advance seq1. Horizontal
355         ↪ displacement.
356         3) To (i + 1, j, 0), that is, seting a gap in seq1 and advance seq0. Vertical
357         ↪ displacement.
358
359         The scores of these displacements are calculated adding the score ot the target cells
360         (that are computed recursively) and the matrix, default of gap scores in each case.
361
362         The score, matches, gaps and forward_arrow are stored at related dictionary entry
363         ↪ based on
364         coordinates (i, j, ini_gap), all of them asociated to the maximum score of
365         the three possible paths starting from the cell, avoiding recomputation
366         of the cell if it's called from another recuersive path.
367
368         Each cell has a third score coordinate, because a cell could be called from a cell
369         ↪ with yet
370         has a gap (only from horizontal or vertical prior displacement) or from a cell with
371         ↪ has a match/no match.
372         Then we need to store two scores, matches, gaps and forward_arrows related to the
373         ↪ two possible
374         cell incarnations at coordinates (i, j, 0) and (i, j, 1).
375
376         We store matches and gaps in order to have one additional criterion to tiebreaker
377         if some of the scores are equal. We are using this approach in local alignment
378         ↪ computation. If two
379         scores are equal we choose the solucion with the greatest number of matches.

```

368 We store the displacement directions in forward_arrow dict to compute the alignment.
 369 It's possible to avoid this, using only the score information, but we have let this
 370 approach
 371 as proof of concept and for clarity in the algorithm.
 372
 373 In this scenario we observe that the differences between the global,
 374 local and long substring algorithms are minimal.
 375
 376 Local algorithm:
 377
 378 Starting from the global algorithm, which would be the most general,
 379 the local algorithm only changes two aspects:
 380 1. Rejection of the roads with negative values of the score, equaling these
 381 values to 0, that is, not letting previous alignments of poor quality affect the
 382 final result.
 383 2. Use as cell of beginning of the alignment the one with the highest scores.
 384 In our implementation we also take into account the number of matches,
 385 as we have already mentioned.
 386
 387 Finally, but outside the algorithm of alignment itself (at forward_track and
 388 matching methods)
 389 it only remains to extend the alignment obtained to show its location within the
 390 chains to be aligned.
 391
 392 Search algorithm of the long common substring:
 393
 394 Modify the global algorithm in the following aspects:
 395 1. Only computes matches between characters or gaps in one or another
 396 initial sequence.
 397
 398 Args:
 399 i (int): Sequence 0 index
 400 j (int): Sequence 1 index
 401 ini_gap (int): 1 if gap initiation, 0 if gap continuation
 402
 403 """
 404 score_diag, score_hor, score_ver = MIN, MIN, MIN
 405 matches_diag, matches_hor, matches_ver = MIN, MIN, MIN
 406 gaps_diag, gaps_hor, gaps_ver = MIN, MIN, MIN
 407 #align and advance seq0 and seq1
 408 #in long_substring mode only matches are processed
 409 if i < self.len_seq0 and j < self.len_seq1 and\
 410 (self.mode != "LONG_SUBSTRING" or self.sequences[0][i] == self.sequences[1][j]):
 411 inc_score, inc_matches = self.calc_score(i, j)
 412 key = (i + 1, j + 1, 1)
 413 if key in self.score_store:
 414 score_diag, matches_diag, gaps_diag = \
 415 self.score_store[key] + inc_score, self.matches_store[key] + inc_matches,
 416 self.gaps_store[key]
 417 else:

```

413         score, matches, gaps = self.align(i + 1, j + 1, 1)
414         self.store(key, score, matches, gaps)
415         score_diag, matches_diag, gaps_diag = score + inc_score, matches +
            ↪ inc_matches, gaps
416         #don't align and gap in seq0 (advance seq1)
417         if j < self.len_seq1:
418             gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
419             key = (i, j + 1, 0)
420             if key in self.score_store:
421                 score_hor, matches_hor, gaps_hor = self.score_store[key] + gap_score,\
422                     self.matches_store[key], self.gaps_store[key] + 1
423             else:
424                 score, matches, gaps = self.align(i, j + 1, 0)
425                 self.store(key, score, matches, gaps)
426                 score_hor, matches_hor, gaps_hor = score + gap_score, matches, gaps + 1
427         #don't align and gap in seq1 (advance seq0)
428         if i < self.len_seq0:
429             gap_score = self.score_gap_cont + ini_gap * self.score_gap_ini
430             key = (i + 1, j, 0)
431             if key in self.score_store:
432                 score_ver, matches_ver, gaps_ver =\
433                     self.score_store[key] + gap_score, self.matches_store[key],
434                     ↪ self.gaps_store[key] + 1
435             else:
436                 score, matches, gaps = self.align(i + 1, j, 0)
437                 self.store(key, score, matches, gaps)
438                 score_ver, matches_ver, gaps_ver = score + gap_score, matches, gaps + 1
439         #choose the high score path
440         matcher_diag, matcher_hor, matcher_ver = score_diag, score_hor, score_ver
441         if i < self.len_seq0 or j < self.len_seq1:
442             if self.mode == "LOCAL" and matcher_diag < 0 and matcher_hor < 0 and matcher_ver
443                 ↪ < 0:
444                 score_diag, score_hor, score_ver = 0, 0, 0
445                 #matcher_diag, matcher_hor, matcher_ver = 0, 0, 0
446             if matcher_diag > matcher_hor and matcher_diag > matcher_ver:
447                 ret_score, ret_matches, ret_gaps, ret_arrow =\
448                     score_diag, matches_diag, gaps_diag, "d"
449             elif matcher_hor > matcher_ver:
450                 ret_score, ret_matches, ret_gaps, ret_arrow =\
451                     score_hor, matches_hor, gaps_hor, "h"
452             else:
453                 ret_score, ret_matches, ret_gaps, ret_arrow =\
454                     score_ver, matches_ver, gaps_ver, "v"
455         else:
456             ret_score, ret_matches, ret_gaps, ret_arrow =\
457                 0, 0, 0, ""
458         self.forward_arrow[(i, j, ini_gap)] = ret_arrow
459         if i == 0 and j == 0:
460             self.store((0, 0, 1), ret_score, ret_matches, ret_gaps)
461             if self.mode in ["GLOBAL", "LONG_SUBSTRING"]: self.max_score_index = (0, 0, 1)
462             else: ret_score = self.max_score // COMPAC
463             ret_matches = self.matches_store[self.max_score_index]

```

```

462         ret_gaps = self.gaps_store[self.max_score_index]
463
464     return ret_score, ret_matches, ret_gaps
465
466 def compute(self, mode="LOCAL", silent=False):
467     """Calc alignment
468         Args:
469             mode (str): Type of algorithm (local, global or long substring)
470             silent (bool): If true don't show alignment output
471     """
472     self.ini_time = time.time()
473     self.init_stores()
474     self.set_mode(mode)
475     self.score, self.matches, self.gaps = self.align()
476     self.align_seq0, self.align_seq1, final_pos = self.forward_track(self.max_score_index,
x)
477     self.matching = self.calc_matching(self.align_seq0, self.align_seq1,
        ↪ self.max_score_index, final_pos)
478     self.unmatches = self.matching.count('.')
479     self.gaps = self.matching.count(' ')
480     self.finish_time = time.time()
481     if not silent:
482         self.view()
483
484 def get_len_long_common_substring(self):
485     """Getter for the len of the common substring
486     That is equal to the number of matches of the alignment
487     """
488     return self.matches
489
490 def get_long_common_substring(self):
491     """Returns the longest common substring
492     without alignment (positional) information
493     """
494     long_common_substring = ""
495     for (char, match_char) in zip(self.align_seq1, self.matching):
496         if match_char == '|':
497             long_common_substring += char
498     return long_common_substring
499
500 def view(self):
501     """Prints the alignment data"""
502     #unmatches = self.matching.count('.')
503     #gaps = self.matching.count(' ')
504     if self.matching:
505         gap_groups = self.matching.count('| ') + self.matching.count('. ') +
        ↪ 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:",
    ↪ self.gaps, " Gap groups:", gap_groups)
526 #simple scoring verification todo: apply to matrix
527 if not self.subst_matrix:
528     print("\tScore verified:", self.matches * self.score_match + self.unmatches *
    ↪ self.score_no_match \
529         + self.gaps * self.score_gap_cont + gap_groups * self.score_gap_ini)
530 print("\tFinish. Execution milliseconds:", round((self.finish_time - self.ini_time)
    ↪ * 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,
    ↪ score_gap_cont=-1):
534     """Calculates an edit distance as requested in questions 1 and 3
535     It's the same computation as a global alignment with -1 penalties applied to
536     score_gap_cont and score_nomatch and 0 in score_match and score_gap_ini
537
538     Args:
539         score_match (int): Score of match characters.
540         score_no_match (int): Score of no match characters.
541         score_gap_ini (int): Score of gap init.
542         score_gap_cont (int): Score of gap continuation.
543     """
544     self.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
545     self.compute("GLOBAL", True)
546     return abs(self.score)

```

Script 1.1.2 (python)

1

Tests

Script 1.1.3 (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
27
28     """
29     global failed, passed, launched
30     try:
31         launched += 1
32         align = AlignSequences([s1, s2])
33         align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
34         if matrix != {}:
35             method = getattr(pairwise2.align, tipus + 'ds')
36             alignments = method(s2, s1, matrix,\
37                               ↪ score_gap_ini + score_gap_cont, score_gap_cont)
38             align.set_subst_matrix(matrix)
39         else:
40             method = getattr(pairwise2.align, tipus + 'ms')
41             alignments = method(s2, s1, score_match, score_no_match,\
42                               ↪ score_gap_ini + score_gap_cont, score_gap_cont)
43
44         align.compute(tipus.upper(), silent=not verbose)
45         m = re.match(r".*Score=([-1234567890]*)",
46                     ↪ format_alignment(*alignments[0]).replace("\n", ""))
47         score = int(m.group(1))
48
49         #search AlignSequences alignment in all possibles alignments from Biopython
50         found = False
51         for a in alignments:
52             if verbose:
53                 print()
54                 print("BioPython alignment:")
55                 print(format_alignment(*a))

```



```

53         if align.align_seq0 == a[1] and align.align_seq1 == a[0]:
54             if not verbose: print(format_alignment(*a))
55             found = True
56             break
57     assert(align.score == score)
58     print ("Passed test %s: scores are equal '%s'" % (number, align.score ))
59     assert(found)
60     print ("Passed test %s: alignments are equal '%s'" % (number, align.align_seq0 ))
61     passed += 1
62
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 = "NLFVALYDFVASGDNTLSITKGEKLRVLGYNHNGEWCEAQTKNGQGWVPSNYITPVN"
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: 5
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

```

1.2 From pairs to sets

In each step of multiple alignment we need to **align two sets of sequences previously intra-aligned**, 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.

The precise algorithm would need to align all the sequences at the same time, but this is a tremendous CPU intensive problem and a huge memory consumer one. A NP-complete problem.

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 our case, it would be easy to extend it to address MSA?. The answer is affirmative. With slight modifications in the pairwise methods, 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 (method *gapeator*).

1.3 Which order? Guide trees.

But which order to follow 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.

To compute the guide tree two steps are necessary:

1. Perform a pairwise alignments between all pairs of 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 we 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. There are other approaches.

1.4 Putting it all together: the CLUSTAL way.

The alignment has three phases:

1. Perform pairwise alignments between all the sequences involved and assign them a score. We use the global alignment Needleman-Wunsch.
2. Build the guide tree. 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. There are other strategies.
3. Multiple alignment. Progressive alignment following the order indicated by the guide tree.

NOTE: the substitution matrices used in pairwise could not be the same matrices used in step three.

1.5 From CLUSTAL to T-COFFEE

1. Substitution matrices to weight matrices
 - But some versions of CLUSTAL uses weight matrices

2. Weight matrices calculated combining information of several alignments:

- Global one (Waterman)
- Local one (Smith-Waterman and others)
- In our implementation we combine the scoring information of a local alignment with a global one. In T-COFFEE terminology this object is the primary library, that we have compute as a dictionary of four keys: index of first sequence, index of second sequence, position of the first sequence, position of the second sequence. The positions are the position computed by a pairwise (global) alignment for every pair of sequences. These are the most important methods to do so:

```
def compute_libraries(sequences, matrix,\
                      score_gap_ini, score_gap_cont, score_match, score_no_match):
    """
    Compute initial library of identities based on scores of pairwise alignments
    Args:
        sequences (list of str): Sequences to compare.
        matrix (dict of tuples of int): Substitution matrix, Biopython format.
        score_gap_ini (int): Score of gap init.
        score_gap_cont (int): Score of gap continuation.
        score_match (int): Score of match characters (used if no matrix informed)
        score_no_match (int): Score of no match characters (used if no matrix informed)

    Returns:
        dict : primary library of alignments
        dict : weight matrix
    """
    weight_library = {}
    primary_library = compute_library(sequences, matrix,\
                                      weight_library, "GLOBAL", score_gap_ini, score_gap_cont,\
                                      score_match, score_no_match)

    _ = compute_library(sequences, matrix,\
                        weight_library, "LOCAL", score_gap_ini, score_gap_cont,\
                        score_match, score_no_match)

    extend_library_weights(sequences, weight_library)

    return primary_library, weight_library

def compute_library(sequences, matrix={}, weight_library={}, mode="GLOBAL",\
                    score_gap_ini=0, score_gap_cont=-8,\
                    score_match=3, score_no_match=-2):
    """
    Compute initial library of identities based on scores of PA
    Args:
        sequences (list of str): Sequences to compare.
        matrix (dict of tuples of int): Substitution matrix, Biopython format.
        mode (str): Computation mode:
```

```

        'GLOBAL'           Global Alignment
        'LOCAL'           Local Alignment
        'LONG_SUBSTRING'   Long substring alignment
    score_gap_ini (int): Score of gap init.
    score_gap_cont (int): Score of gap continuation.
    score_match (int): Score of match characters (used if no matrix informed)
    score_no_match (int): Score of no match characters (used if no matrix informed)

```

Returns:

```

    list of str: primary library of alignments
    """

```

```

primary_library = {}
for i in range(0, len(sequences)):
    for j in range(0, i):
        if (i,j) not in primary_library:
            identity, align_i, align_j = pairwise_align_coffee(sequences[i], sequences[j],\
                matrix, mode, score_gap_ini, score_gap_cont,\
                score_match, score_no_match)
            update_weight_library(weight_library, i, j, identity,\
                sequences[i], sequences[j], align_i, align_j)
            primary_library[(i,j)] = (align_i, align_j, identity)
return primary_library

```

```

def update_weight_library(weight_library, i, j, identity,\
    seq_i, seq_j, align_i, align_j):
    """
    Update weights library from alignments and %identity
    """
    for k, (c_i, c_j) in enumerate(zip(align_i, align_j)):
        if c_i != "-" and c_j != "-":
            pos_i = get_pos(k, seq_i, align_i)
            pos_j = get_pos(k, seq_j, align_j)
            update_weight_at_pos(weight_library, i, j, pos_i, pos_j, identity)
            update_weight_at_pos(weight_library, j, i, pos_j, pos_i, identity)

```

2. Extend libraries taking into account all the intermediate alignments using every k sequence between two other pair of sequences. It's a way to correct the weights computed in the primary library trying to circumvent the CLUSTAL problem that derived for poor initial alignments, i.e. all the sequences contribute to the weights of the primary library.

```

def extend_library_weights(sequences, weight_library):
    """
    Extend library for all triplets of sequences
    Taken into account the simetry i -> k -> j
    """
    len_sequences = len(sequences)
    for i in range(0, len_sequences):
        for k in range(0, len_sequences):

```

```

        if k != i:
            for j in range(0, len_sequences ):
                if j != k and j != i:
                    #print("Triplet:", i, k, j)
                    extend_weigths(weight_library, i, k, j)

def extend_weigths(weight_library, i, k, j):
    """
    Extend weigths for pair of sequences (i,j) at pos (pos_i_posj)
    taken into account the routes using k as
    intermediate, by means of the alignments (i, k) and (k, j).
    """
    for pos_i, pos_i_j in weight_library[i][j].items():
        for pos_j in pos_i_j.keys():
            if pos_j in weight_library[j][k].keys():
                for pos_k in weight_library[j][k][pos_j].keys():
                    if pos_k in weight_library[k][i].keys():
                        for pos_i_new in weight_library[k][i][pos_k].keys():
                            if pos_i_new == pos_i:
                                m = min(\
                                    weight_library[i][k][pos_i][pos_k], \
                                    weight_library[j][k][pos_j][pos_k])
                                weight_library[i][j][pos_i][pos_j] += m

```

3. Guide trees by NJ or UPGMA as CLUSTAL. T-COFFEE prefers CLUSTAL one.
4. Progressive alignment based on weight matrices, not on substitution matrices as CLUSTAL, but other than that, the same strategy.

1.6 Beyond T-COFFEE

TODO: general comments about MUSCLE or other approaches

1.7 MSA software

1.7.1 MSA generic methods

Script 1.7.1 (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,
96                             ↪ score_gap_cont=-8):
97     """
98     Performs initial pairwise alignments against the class AlignSequences
99     returning the distance between 0 and 100.
100
101     Args:
102         s1 (str): First sequence to compare.
103         s2 (str): Second sequence to compare.
104         matrix (dict of tuples of int): Substitution matrix, Biopython format.
105         matrix_mode (str): Type of matrix
106             'SUBST'          Substitution matrix
107             'WEIGHT'         Weight matrix
108         mode (str): Computation mode:
109             'GLOBAL'         Global Alignment
110             'LOCAL'          Local Alignment
111             'LONG_SUBSTRING' Obtain long common substring
112         score_gap_ini (int): Score of gap init.

```



```

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
151         ↪ root
152         of the other two nodes.
153         (dict of int, boolean = True): contains all nodes
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 Args:
410     stack_0 (list of str): First stack of sequences to align.
411     stack_1 (list of str): Second stack of sequences to align.
412     sequences (list of str): Sequences to align.
413     matrix (dict of tuples of int): Substitution matrix, Biopython format.
414     matrix_mode (str): Type of matrix
415         'SUBST' Substitution matrix
416         'WEIGHT' Weight matrix
417     mode (str): Computation mode:
418         'GLOBAL' Global Alignment

```

```

419         'LOCAL'           Local Alignment
420         'LONG_SUBSTRING'  Obtain long common substring
421         stack_0_indexes(list of int) : Indexes of initial sequences related to stack
→ sequences 0
422         stack_1_indexes(list of int) : Indexes of initial sequences related to stack
→ sequences 1
423         stack_0_refs(list of dict) : stack_0 references to original sequences
424         stack_1_refs(list of dict) : stack_1 references to original sequences
425         score_match (int): Score of match characters (used if no matrix informed)
426         score_no_match (int): Score of no match characters (used if no matrix informed)
427         score_gap_ini (int): Score of gap init.
428         score_gap_cont (int): Score of gap continuation.
429
430     Returns:
431         list of str: stack_0 gapped (with the gaps necessary for the alignment)
432         list of str: stack_0 gapped (with the gaps necessary for the alignment)
433         list of dict: stack_0 references to original sequences
434         list of dict: stack_1 references to original sequences
435     """
436     align = AlignSequences([stack_0[0], stack_1[0]])
437     align.set_scores(score_match, score_no_match, score_gap_ini, score_gap_cont)
438     align.set_subst_matrix(matrix)
439     align.set_matrix_mode(matrix_mode)
440     align.set_stacks(stack_0, stack_1, stack_0_indexes, stack_1_indexes, stack_0_refs,
→ stack_1_refs)
441     align.compute(mode.upper(), silent = True)
442     # align_seq0 align_seq1 are the seq0 and seq1 alignments
443     # we need to deduce the rest of alignments.
444     # what we do is perform the same gap insertions, if any, as the first sequence of the
→ stacks
445     # 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
→ informed
447     stack_0_gapped, stack_0_references = gapeator(stack_0[0], align.align_seq0, stack_0,
→ stack_0_refs)
448     stack_1_gapped, stack_1_references = gapeator(stack_1[0], align.align_seq1, stack_1,
→ stack_1_refs)
449     return stack_0_gapped, stack_1_gapped, stack_0_references, stack_1_references
450
451 def get_name(index, sequence_names):
452     """
453     Obtain sequence name from index
454     """
455     name = ""
456     if index < len(sequence_names):
457         name = sequence_names[index]
458     else:
459         name = str(index)
460     return name
461
462 def to_newick(tree, sequence_names):
463     """

```

```

464 Obtain guide tree in newick format
465 """
466 # Change format to intermediate roots
467 roots = {}
468 newick_tree = ""
469 for branch in tree:
470     (i, j, k) = branch
471     name_i = get_name(i, sequence_names)
472     name_j = get_name(j, sequence_names)
473     name_k = get_name(k, sequence_names)
474     if name_i in roots:
475         new_root_i = roots[name_i]
476     else:
477         new_root_i = name_i
478     if name_j in roots:
479         new_root_j = roots[name_j]
480     else:
481         new_root_j = name_j
482     roots[name_k] = [new_root_i, new_root_j]
483
484 for root in roots.values():
485     s_root = str(root)
486     if len(s_root) > len(newick_tree):
487         newick_tree = s_root.replace("[", "(").replace("]", ")").replace("'", "")
488
489 return newick_tree

```

1.7.2 T-COFFEE methods

Script 1.7.2 (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_weigths(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:

```

```

123         #print("Extension", pos_i, pos_j, pos_k, weight_library[i][k]
124         ↪ ] [pos_i][pos_k], weight_library[j][k][pos_j][pos_k])
125         m = min(\
126             weight_library[i][k][pos_i][pos_k], \
127             weight_library[j][k][pos_j][pos_k])
128         #print("++", m, weight_library[i][j][pos_i][pos_j])
129         weight_library[i][j][pos_i][pos_j] += m
130
131 def extend_library_weights(sequences, weight_library):
132     """
133     Extend library for all triplets of sequences
134     Taken into account the simetry i -> k -> j
135     """
136     len_sequences = len(sequences)
137     for i in range(0, len_sequences):
138         for k in range(0, len_sequences):
139             if k != i:
140                 for j in range(0, len_sequences):
141                     if j != k and j != i:
142                         #print("Triplet:", i, k, j)
143                         extend_weights(weight_library, i, k, j)
144
145 def compute_libraries(sequences, matrix, \
146                     score_gap_ini, score_gap_cont, score_match, score_no_match):
147     """
148     Compute initial library of identities based on scores of pairwise alignments
149     Args:
150         sequences (list of str): Sequences to compare.
151         matrix (dict of tuples of int): Substitution matrix, Biopython format.
152         score_gap_ini (int): Score of gap init.
153         score_gap_cont (int): Score of gap continuation.
154         score_match (int): Score of match characters (used if no matrix informed)
155         score_no_match (int): Score of no match characters (used if no matrix informed)
156
157     Returns:
158         dict : primary library of alignments
159         dict : weight matrix
160     """
161     weight_library = {}
162     primary_library = compute_library(sequences, matrix, \
163                                     weight_library, "GLOBAL", score_gap_ini, score_gap_cont, \
164                                     score_match, score_no_match)
165
166     _ = compute_library(sequences, matrix, \
167                       weight_library, "LOCAL", score_gap_ini, score_gap_cont, \
168                       score_match, score_no_match)
169
170     # _ = compute_library(sequences, matrix, \
171     #                   weight_library, "LONG_SUBSTRING", score_gap_ini, score_gap_cont, \
172     #                   score_match, score_no_match)
173

```

```

174     extend_library_weights(sequences, weight_library)
175
176     return primary_library, weight_library

```

1.7.3 Main MSA method

Script 1.7.3 (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
21             'LONG_SUBSTRING' Obtain long common substring
22         score_gap_ini (int): Score of gap init.
23         score_gap_cont (int): Score of gap continuation.
24         score_match (int): Score of match characters (used if no matrix informed)
25         score_no_match (int): Score of no match characters (used if no matrix informed)
26         verbose (bool): If True prints verbose info
27
28     Returns:
29         list of 3-tuples of int: guide tree, the third position of the tuple contains the
30         → root
31         of the other two nodes.
32         list of str: alignments
33         list of str: sequence_names
34         list of int: sequence indexes relating strings in alignment to original sequences
35     """
36     seq_fasta = readFasta(file)
37     sequences = list(seq_fasta.values())
38     sequence_names = list(seq_fasta.keys())
39     print(sequence_names)
40     return do_msa(sequences, sequence_names,\
41                  main_alg, method, matrix, matrix_mode,\
42                  mode, score_gap_ini, score_gap_cont, score_match,\
43                  score_no_match, verbose)

```

```

43 def do_msa(sequences, sequence_names, main_alg="CLUSTAL", method="NJ", matrix={},
44   ↪ matrix_mode="SUBST",\
45       mode="GLOBAL", score_gap_ini=-10, score_gap_cont=-5, score_match=3,\
46       score_no_match=-2, verbose=False):
47     """
48     Performs MSA alignments from sequences
49     Args:
50         sequences (list of str): Sequences
51         sequence_names (list of str): Names of sequences
52         main_alg (str): Main algorithm:
53             "CLUSTAL"      Clustal like
54             "T-COFFEE"     T-COFFEE like
55         method (str): NJ neighbor join / UPGMA
56         matrix (dict of tuples of int): Substitution matrix, Biopython format.
57         matrix_mode (str): Type of matrix
58             'SUBST'       Substitution matrix
59             'WEIGHT'       Weight matrix
60         mode (str): Computation mode:
61             'GLOBAL'      Global Alignment
62             'LOCAL'       Local Alignment
63             'LONG_SUBSTRING' Obtain long common substring
64         score_gap_ini (int): Score of gap init.
65         score_gap_cont (int): Score of gap continuation.
66         score_match (int): Score of match characters (used if no matrix informed)
67         score_no_match (int): Score of no match characters (used if no matrix informed)
68         verbose (bool): If True prints verbose info
69
70     Returns:
71     ↪ list of 3-tuples of int: guide three, the third position of the tuple contains the
72         of the other two nodes.
73         list of str: alignments
74         list of str: sequence_names
75         list of int: sequence indexes relating strings in alignment to original sequences
76     """
77     if main_alg == "T-COFFEE":
78         primary_library, weight_library = compute_libraries(sequences, matrix,\
79             score_gap_ini, score_gap_cont, score_match, score_no_match)
80         #print("Primary library", primary_library)
81         #print("Weight library", weight_library)
82         # Matrix mode and other MSA parameters
83         matrix = weight_library
84         matrix_mode = "WEIGHT"
85         # From here only we need is to compute a MSA with weight matrix as reference.
86     else:
87         matrix_mode = "SUBST"
88     if method == "NJ":
89         guide_tree, guide_tree_nodes = \
90             guide_tree_NJ(sequences, matrix, matrix_mode,\
91                 mode, score_gap_ini, score_gap_cont,\
92                 score_match, score_no_match)
93     else: #UPGMA

```

```

93     guide_tree, guide_tree_nodes = \
94         guide_tree_UPGMA(sequences, matrix, matrix_mode, \
95             mode, score_gap_ini, score_gap_cont, \
96             score_match, score_no_match)
97 sequences_store = {}
98 sequences_store_indexes = {}
99 sequences_store_refs = {}
100 print("Guide Tree", guide_tree, sequence_names)
101 #return guide_tree, "", sequence_names
102 # Create MSA
103 for i in guide_tree_nodes.keys():
104     if i < len(sequences):
105         sequences_store[i] = [sequences[i]]
106         sequences_store_indexes[i] = [i]
107         sequences_store_refs[i] = []
108         autorefs = {}
109         for k in range(0, len(sequences[i])):
110             autorefs[k] = k
111         sequences_store_refs[i].append(autorefs)
112
113 if verbose: print(sequences_store)
114 for (i, j, k) in guide_tree:
115     stack_i = sequences_store[i]
116     stack_j = sequences_store[j]
117     stack_i_indexes = sequences_store_indexes[i]
118     stack_j_indexes = sequences_store_indexes[j]
119     stack_i_references = sequences_store_refs[i]
120     stack_j_references = sequences_store_refs[j]
121     if verbose: print("Stack i", i, stack_i)
122     if verbose: print("Stack j", j, stack_j)
123     if verbose: print("Stack i_indexes", i, stack_i_indexes)
124     if verbose: print("Stack j_indexes", j, stack_j_indexes)
125     stack_0, stack_1, stack_0_references, stack_1_references = \
126         pairwise_align_msa_step(stack_i, stack_j, sequences, matrix, matrix_mode, \
127             mode, stack_i_indexes, stack_j_indexes, \
128             stack_i_references, stack_j_references, \
129             score_match, score_no_match, score_gap_ini,
130             ↪ score_gap_cont)
131
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)

```

```

144     for seq_refs in stack_0_references:
145         sequences_store_refs[k].append(seq_refs)
146     for seq_refs in stack_1_references:
147         sequences_store_refs[k].append(seq_refs)
148     if verbose: print("=====")
149     if verbose: print("Sequences store indexes", sequences_store_indexes[k])
150     if verbose: print("Sequences store references", sequences_store_refs[k])
151     if verbose: print("New stack:", k, sequences_store[k])
152 alignment = sequences_store[k]
153 newick_tree = to_newick(guide_tree, sequence_names)
154 return newick_tree, alignment, sequence_names, sequences_store_indexes[k]
155
156 def score(alignment, matrix, score_gap_ini=0, score_gap_cont=0):
157     """
158     Score based on sum of pair scores (SDP) taking into account substitution matrix
159     Derived from the objective score of MUSCLE refinement stage
160     """
161     msa_score = 0
162     for k in range(0, len(alignment[0])): #columns of msa
163         score_column_k = 0
164         nvalues = 0
165         for i in range(0, len(alignment)):
166             for j in range(i + 1, len(alignment)):
167                 if alignment[i][k] == "-" and alignment[j][k] != "-":
168                     score_column_k += score_gap_cont
169                     if k == 0 or alignment[i][k-1] != "-":
170                         score_column_k += score_gap_ini
171                 if alignment[j][k] == "-" and alignment[i][k] != "-":
172                     score_column_k += score_gap_cont
173                     if k == 0 or alignment[j][k-1] != "-":
174                         score_column_k += score_gap_ini
175                 elif (alignment[i][k], alignment[j][k]) in matrix:
176                     score_column_k += matrix[(alignment[i][k], alignment[j][k])]
177                     nvalues += 1
178                 elif (alignment[j][k], alignment[i][k]) in matrix:
179                     score_column_k += matrix[(alignment[j][k], alignment[i][k])]
180                     nvalues += 1
181             if nvalues > 0:
182                 #score += score_column_k / nvalues
183                 msa_score += score_column_k
184     return msa_score
185
186 def score_from_fasta(file, matrix, score_gap_ini=0, score_gap_cont=0):
187     seq_fasta = readFasta(file)
188     sequences = list(seq_fasta.values())
189     return score(sequences, matrix, score_gap_ini, score_gap_cont)

```

1.8 Some runs

1.8.1 T-COFFEE

Script 1.8.1 (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 guide_tree_nj, align, sequence_names, indexes = do_msa_from_fasta(file,\
23     main_alg = "T-COFFEE", method = "NJ", \
24     matrix = matrix, matrix_mode = "SUBST",\
25     mode = "GLOBAL", score_gap_ini = -10,\
26     score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
27 print("# Guide Tree:", guide_tree_nj)
28 t_nj, ts_nj = draw_guide_tree(guide_tree_nj)
29 print("# Alignment:")
30 for i, s in enumerate(align):
31     print(">" + sequence_names[indexes[i]])
32     print(s)
33 print()
34 print("Score", score(align, matrix, quality_score_gap_ini, quality_score_gap_cont))
```

Output

```
['1aboA', '1ycsB', '1pht', '1vie', '1ihvA']
{(1, 0): 12, (2, 0): 9, (2, 1): 8, (3, 0): 8, (3, 1): 9, (3, 2): 8, (4, 0): 10, (4, 1): 6,
 ↪ (4, 2): 10, (4, 3): 16}
Guide Tree [(4, 3, 5), (1, 0, 6), (5, 2, 7), (7, 6, 8)] ['1aboA', '1ycsB', '1pht', '1vie',
 ↪ '1ihvA']
# Guide Tree: (((1ihvA, 1vie), 1pht), (1ycsB, 1aboA))

    /-1ihvA
```



```

      /-|
    /-|  \-1vie
    |  |
--|  \-1pht
    |
    |  /-1ycsB
    \-|
      \-1aboA
# Alignment:
>1ihvA
NFR---VY-YRDSRD-----PVWKGPAKLLWKGEAVVIQDNSDIKVPPRR-----KAKIIRD-----
>1vie
D-----RVRKKSGAAWQQQIVGWYCTNLTPEGYAVE-----SEAHPGSVQIYPVAALERIN
>1pht
GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSQGQEARPEEIGWLNQYNETTGERGDFPGTYVEYIG-----RKKISP
>1ycsB
KGV-IYALWDYEPQNDDELPMKEGDCMTIIHREDE-DEIEWWWARLNDK-----EGYVPRNLLG---LYP-----
>1aboA
N-L-FVALYDFVASGDNTLSITKGEKLRVLGYNHN-G--EWCEAQTKNG-----QGWVPSNYIT---PVN-----

Score -1648

['1aboA', '1ycsB', '1pht', '1vie', '1ihvA']
Guide Tree [(4, 3, 6), (6, 2, 7), (7, 0, 8), (1, 8, 9)] ['1aboA', '1ycsB', '1pht', '1vie',
↳ '1ihvA']
# Guide Tree: (1ycsB, (((1ihvA, 1vie), 1pht), 1aboA))

    /-1ycsB
    |
    |      /-1ihvA
--|      /-|
    |  /-|  \-1vie
    |  |  |
    \-|  \-1pht
        |
        \-1aboA
# Alignment:
>1ycsB
KGVYIYALWDYEPQNDDELPMKEGDCMTIIHREDEDEIEWWWARLNDKEGYVP-----RNLL-----GLYP-----
>1ihvA
NFR---VY-YRDSRD-----PVWKGPAKLLWKG--EGAVVIQDNSDIKVPPRR-----KAKIIRD-----
>1vie
D-----RVRKKSGAAWQG--QIVGWYCTNLTPEGYAVE-----SEAHPGSVQIYPVAALERIN
>1pht
GYQYRA-LYDYKKEREEDIDLHLGDILTVNKGSLV--ALGFSQGQEARPEEIGWLNQYNETTGERGDFPGTYVEYIG-----RKKISP
>1aboA
N-LFVA-LYDFVASGDNTLSITKGEKLRVLGYNHN--G-EWCEAQTKNGQG--VPSNYITPVN-----

Score -1885

```

1.8.2 CLUSTAL homemade

Script 1.8.2 (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()
17 print("Score", score(align, MatrixInfo.blosum62, quality_score_gap_ini,
18     ↪ quality_score_gap_cont))
19 print()
20 guide_tree_nj, align, sequence_names, indexes = do_msa_from_fasta(file,\
21     main_alg = "CLUSTAL", method = "NJ", \
22     matrix = matrix, matrix_mode = "SUBST",\
23     mode = "GLOBAL", score_gap_ini = -10,\
24     score_gap_cont = -5, score_match = 3, score_no_match = -2, verbose = False)
25 print("# Guide Tree:", guide_tree_nj)
26 t_nj, ts_nj = draw_guide_tree(guide_tree_nj)
27 print("# Alignment:")
28 for i, s in enumerate(align):
29     print(">" + sequence_names[indexes[i]])
30     print(s)
31 print()
32 # print("Score", score(align, matrix, quality_score_gap_ini, quality_score_gap_cont))
```

Output

```
['1aboA', '1ycsB', '1pht', '1vie', '1ihvA']
{(1, 0): 24, (2, 0): 19, (2, 1): 23, (3, 0): 19, (3, 1): 13, (3, 2): 25, (4, 0): 16, (4, 1):
↪ 18, (4, 2): 20, (4, 3): 18}
Guide Tree [(3, 2, 5), (1, 0, 6), (5, 4, 7), (7, 6, 8)] ['1aboA', '1ycsB', '1pht', '1vie',
↪ '1ihvA']
# Guide Tree: (((1vie, 1pht), 1ihvA), (1ycsB, 1aboA))

      /-1vie
     /-|
    /-| \-1pht
   |  |
  --|  \-1ihvA
```

```

      |
      | /-1ycsB
      \-|
        \-1aboA
# Alignment:
>1vie
-----DRVRKKSGAAWQQQIVGWYCTNLTPEGYAVESEAH-----GSVQIYPVAALERI-----N
>1pht
GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNNGYNETTGERGDFPGTYVEYIGRKKISP
>1ihvA
-----NFRVYYRDSRDPVWKGPALLWKGEAVVIQ-----DNSDIKVVPRRKAKIIRD
>1ycsB
KGVIIYALWDYEPQNDDELPMKEGDCMTIIHREDEDEIEWWWA-----RLNDKEGYVPRNLLGLYP
>1aboA
-NLFVALYDFVASGDNTLSITKG--EKLRVLGYNHNGEWCEA-----QTKNGQGWPVPSNYITPVN

Score -1105

['1aboA', '1ycsB', '1pht', '1vie', '1ihvA']
Guide Tree [(1, 0, 6), (6, 2, 7), (7, 4, 8), (3, 8, 9)] ['1aboA', '1ycsB', '1pht', '1vie',
↪ '1ihvA']
# Guide Tree: (1vie, (((1ycsB, 1aboA), 1pht), 1ihvA))

      /-1vie
      |
      | /-1ycsB
--|    /-|
    | /-| \-1aboA
    | | |
    \-| \-1pht
      |
      \-1ihvA
# Alignment:
>1vie
DR-----VRKKSGAAWQQQIVGWYCTNLTPEGYAVESEAH-----GSVQIYPVAALE-----RIN
>1ycsB
KGVIIYALWDY--EPQNDDELPMKEGDCMTIIHREDEDEIEWWWARLNDKEGYVPRNLLG-----LYP
>1aboA
-NLFVALYDF--VASGDNTLSITKG--EKLRVLGYNHNGEWCEAQTKNGQGWPVPSNYIT-----PVN
>1pht
GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNNGYNETTGERGDFPGTYVEYIGRKKISP
>1ihvA
NFRVYYRDSRDPVWKGPALLWKGEAVVIQDNSDIKVVPRRKAKI-----IRD

```

1.8.3 CLUSTALW official build.

Script 1.8.3 (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: lyCSB 60 aa

Sequence 3: lpht 80 aa

Sequence 4: lvie 51 aa

Sequence 5: lihvA 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       GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGGEARPEEIGWLNNGYNETTG
1vie       -----DRVRKKSGAAWQG-----QIVGWYCTN---LT
1ihvA      -----NFRVYYRDSRDPVWKGPAPKLLWK

```

*

```

1aboA      GQGW-----VPSNYI--TPVN-----
1ycsB      KEGY-----VPRNLLGLYP-----
1pht       ERGD-----FPGTYVEYIGRKKISP---
1vie       PEGYAVESEAHPGSVQIYPVAALERIN-----
1ihvA      GEGAVVIQDNSD-----IKVVPRRKAKIIRD

```

.*

:

```

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

```