





```
[15]: ! time t-coffee --infile=SUP35_10seqs.fa --outfile=SUP35_10seqs_tcoffee.fa --output=fastaaln > t_coffee
10..log

PROGRAM: T-COFFEE Version 11.00.8cbe486 (2014-08-12 22:05:29 - Revision 8cbe486 - Build 477)
- full_log S [0]
- genepred_score S [0] nsd
- run_name S [0]
- mem_mode S [0] mem
- extend D [1] 1
- quicktree FL [0] very_fast_triplet
- max_npair D [0] 10
- seq_name_for_quadruplet S [0] all
- compact S [0] default
- c_lib S [0] no
- do_self FL [0] 0
- do_normalise D [0] 1000
- template_r18 S [0]
- setenv S [0] 0
- template_mode S [0] 0
- flip D [0] 0
- remove_template_file D [0] 0
- profile_template_file S [0]
- in S [0]
- seq D [0]
- aln S [0]
- method_limits S [0]
- lib S [0]
- profile S [0]
- profile1 S [0]
- profile2 S [0]
- pdb S [0] 0
- relax_lib D [0] 1
- filter_lib D [0] 1
- shrink_lib D [0] 0
- out_lib W_F [0] no
- out_lib_mode S [0] primary
- lib_only D [0] 0
- outseqweight W_F [0] no
- dpa FL [0] 0
- run_name S [0] ANY
- cosmetic_penalty D [0] 0
- gapopen D [0] 0
- gapext D [0] 0
- fgapopen D [0] 0
- fgapext D [0] 0
- newtree D [0] 0
- newtree W_F [0] default
- tree W_F [0] NO
- use tree R_F [0] nj
- tree mode S [0] nj
- distance_matrix_sim_mode S [0] ktup
- distance_matrix_sim_mode S [0] 1dmat_sim1
- quicktree FL [0]
- outfile W_F [1] SUP35_10seqs_tcoffee.fa
- maximise FL [1] 1
- output S [1] fastaaln
- less D [0]
- infile R_F [1] SUP35_10seqs.fa
- matrix S [0] default
- t1 mode D [0] 1
- profile_mode S [0] cw_profile_profile
- profile_comparison S [0] profile
- dp_mode S [0] linked_pair_wise
- ktuple D [0] 1
- ndiag D [0] 0
- diag_threshold D [0] 0
- diag_mode D [0] 0
- sim_matrix S [0] vasiliky
- transform S [0]
- extend_seq FL [0] 0
- outorder S [0] input
- inorder S [0] aligned
- seqnos S [0] off
- case S [0] keep
- cpu D [0] 0
- ulimit D [0] -1
- maxnsq D [0] 1
- maxlen D [0] -1
- sample_dp D [0] 0
- weight D [0] default
- seq_weight S [0] no
- align FL [1] 1
- mcca FL [0] 0
- domain FL [0] 0
- start D [0] 0
- len D [0] 0
- scale D [0] 0
- mcca_interactive FL [0] 0
- method_evaluate_mode S [0] default
- evaluate_mode S [0] triplet
- get_type FL [0] 0
- clean_aln D [0] 0
- clean_threshold D [1] 1
- clean_iteration D [0] 1
- clean_evaluate_mode S [0] t_coffee_fast
- extend_matrix FL [0] 0
- prot_min_sim D [40] 40
- prot_max_sim D [90] 90
- prot_min_cov D [40] 40
- pdbe_type S [0] ANY
- pdbe_min_sim D [35] 35
- pdbe_max_sim D [100] 100
- pdbe_min_cov D [50] 50
- pdbe_max_cov D [50] 50
- t1 mode S [0] EBI
- blast W_F [0]
- blast_server W_F [0] EBI
- pdbe_db W_F [0] pdbe
- protein_db W_F [0] uniprot
- method_log W_F [0] no
- struct_use W_F [0]
- align_pdb_system_file W_F [0] use
- align_pdb_hash_mode W_F [0] no
- external_aligner S [0] hasch_ca_trace_bubble
- msa_mode S [0] tree
- master S [0] no
- blast_nseq S [0] default
- lalign_n_top D [0] 10
- iterate D [0] 0
- trim D [0] 0
- split D [0] 0
- trimfile S [0] default
- split D [0] 0
- split_seq_chres D [0] 0
- split_score_thres D [0] 0
- check_pdb_status D [0] 0
- clean_seq_name D [0] 0
- seq_to_keep S [0]
- dpa_master_aln S [0]
- dpa_maxnseq D [0] 0
- dpa_min_score1 D [0] 0
- dpa_min_score2 D [0] 0
- dpa_keep_tmfile FL [0] 0
- dpa_debug W_F [0]
- multi_core S [0] templates_seqs_relax_msa_evaluate
- n_core D [0] 0
- max_nproc D [0] 0
- tlib_list S [0]
- prune_lib_mode S [0] 5
- tip S [0] none
- rna_lib S [0]
- no_warning D [0] 0
- run_local_script D [0] 0
- plugins S [0] default
- proxy S [0] unset
- email S [0]
- clean_overaln D [0] 0
- overaln_parsa S [0]
- overaln_mode S [0]
- overaln_model S [0]
- overaln_threshold D [0] 0
- overaln_target D [0] 0
- overaln_P1 D [0] 0
- overaln_P2 D [0] 0
- overaln_P3 D [0] 0
- overaln_P4 D [0] 0
- rxon_boundaries S [0]
- dump D [0] no
- display D [0] 100

INPUT FILES
Input File (S) SUP35_10seqs.fa Format fasta_seq
Input File (M) proba_pair

Identify Master Sequences [no]:

Master Sequences Identified
INPUT SEQUENCES: 10 SEQUENCES [DNA]
Input File SUP35_10seqs.fa Seq SUP35_Agos_ATCC_10895_NM_211584 Length
2076 type DNA Struct unchecked Seq SUP35_Kla_AB039749 Length
2193 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Sarb_H_6_chrXIII_CM001575 Length
2043 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Sbou_unique28_CM003560 Length
2059 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Scer_74-D694_GCA_001578265.1 Length
2076 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Sker_beer078_CM005938 Length
2081 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Scer_CBS12357_chr_II_IV_DF968535 Length
2025 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Skud_IF01802T_36 Length
2047 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Smik_IF01815T_30 Length
2046 type DNA Struct unchecked Input File SUP35_10seqs.fa Seq SUP35_Spar_A12_Lit1 Length

Multi Core Mode: 8 processors:

--- Process Method/Library/Aln SSUP35_10seqs.fa
--- Process Method/Library/Aln Mproba_pair
xxx Retrieved SSUP35_10seqs.fa
xxx Retrieved Mproba_pair

All Methods Retrieved

MANUAL PENALTIES: gapopen=0 gapext=0

Library Total Size: [271542]

Library Relaxation: Multi_proc [4]

! [Relax Library][TOT= 1][100 %] [ELAPSED TIME: 0 sec.]
Relaxation Summary: [271542]-->[231898]

UN-WEIGHTED MODE: EVERY SEQUENCE WEIGHTS 1

MAKE GUIDE TREE
[MODE=nj][DONE]

PROGRESSIVE ALIGNMENT [Tree Based]
Group 1: SUP35_Agos_ATCC_10895_NM_211584
Group 2: SUP35_Kla_AB039749
Group 3: SUP35_Sarb_H_6_chrXIII_CM001575
Group 4: SUP35_Sbou_unique28_CM003560
Group 5: SUP35_Scer_74-D694_GCA_001578265.1
Group 6: SUP35_Sker_beer078_CM005938
Group 7: SUP35_Scer_CBS12357_chr_II_IV_DF968535
Group 8: SUP35_Skud_IF01802T_36
Group 9: SUP35_Smik_IF01815T_30
Group 10: SUP35_Spar_A12_Lit1

Group 11: [Group 7 ( 1 seq)] with [Group 3 ( 1 seq)]-->[Len= 2052][PID:23209]
Group 12: [Group 5 ( 1 seq)] with [Group 1 ( 1 seq)]-->[Len= 2058][PID:23210]
Group 13: [Group 6 ( 1 seq)] with [Group 12 ( 2 seq)]-->[Len= 2058][PID:23210]
Group 14: [Group 10 ( 1 seq)] with [Group 13 ( 3 seq)]-->[Len= 2058][PID:23210]
Group 15: [Group 9 ( 1 seq)] with [Group 14 ( 4 seq)]-->[Len= 2058][PID:23210]
Group 16: [Group 15 ( 5 seq)] with [Group 11 ( 2 seq)]-->[Len= 2064][PID:23178][For
ked]
Group 17: [Group 16 ( 7 seq)] with [Group 8 ( 1 seq)]-->[Len= 2064][PID:23178]
Group 18: [Group 2 ( 1 seq)] with [Group 17 ( 8 seq)]-->[Len= 2141][PID:23178]
Group 19: [Group 1 ( 1 seq)] with [Group 18 ( 9 seq)]-->[Len= 2210][PID:23178]

! [Final Evaluation][TOT= 276][100 %] [ELAPSED TIME: 0 sec.]

OUTPUT RESULTS
#### File Type= GUIDE_TREE Format= newick Name= SUP35_10seqs_dnd
#### File Type= MSA Format= fastaaln Name= SUP35_10seqs_tcoffee.fa

# Command Line: /home/isemenov/anaconda3/lib/t_coffee-11.0.0/bin/t_coffee -infile SUP35_10seqs.fa -ou
tfile SUP35_10seqs_tcoffee.fa -output fastaaln [PROGNAME: T-COFFEE]
# T-COFFEE Memory Usage: Current= 37.060 Mb, Max= 44.536 Mb
# Results Produced with T-COFFEE Version 11.00.8cbe486 (2014-08-12 22:05:29 - Revision 8cbe486 - Buil
d 477)
# T-COFFEE is available from http://www.tcoffee.org
# Register on: https://groups.google.com/group/tcoffee/
276.50user 5.67system 0.43.26elapsd 6506CPU (bavgtext+bavgdata 675740maxresident)k
137681inputs+4720outputs (131major+60438minor)pagefaults 0swaps

In [16]: record_dict = IO.to_dict(IO.parse("SUP35_10seqs_tcoffee.fa", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))

SUP35_Kla_AB039749
2193
SUP35_Agos_ATCC_10895_NM_211584
2210
SUP35_Scer_74-D694_GCA_001578265.1
2210
SUP35_Spar_A12_Lit1
2210
SUP35_Smik_IF01815T_30
2210
SUP35_Skud_IF01802T_36
2210
SUP35_Sbou_unique28_CM003560
2210
SUP35_Scer_beer078_CM005938
2210
SUP35_Sarb_H_6_chrXIII_CM001575
2210
SUP35_Scer_CBS12357_chr_II_IV_DF968535
2210

Prank

In [17]: ! time prank -d=SUP35_10seqs.fa -o=SUP35_10seqs_prank.fa

-----
PRANK v.170427:
-----

Input for the analysis
- aligning sequences in 'SUP35_10seqs.fa'
- using inferred alignment guide tree
- option '+F' is not used; it can be enabled with '+F'
- external tools available:
MAFFT for initial alignment
Exonerate for alignment anchoring

Correcting (arbitrarily) for multifurcating nodes.
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 1.

Alignment score: 4546
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 2.

Alignment score: 4981
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 3.

Alignment score: 4981
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 4.

Alignment score: 4981
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 5.

Alignment score: 4981

Writing
- alignment to 'SUP35_10seqs_prank.fa.best.fas'

Analysis done. Total time 7s
6.70user 0.15system 0.06.92elapsd 99%CPU (bavgtext+bavgdata 326360maxresident)k
23816inputs+4720outputs (131major+60438minor)pagefaults 0swaps

In [18]: record_dict = IO.to_dict(IO.parse("SUP35_10seqs_prank.fa.best.fas", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))

SUP35_Agos_ATCC_10895_NM_211584
2366
SUP35_Sarb_H_6_chrXIII_CM001575
2366
SUP35_Scer_CBS12357_chr_II_IV_DF968535
2366
SUP35_Skud_IF01802T_36
2366
SUP35_Scer_74-D694_GCA_001578265.1
2366
SUP35_Sbou_unique28_CM003560
2366
SUP35_Scer_beer078_CM005938
2366
SUP35_Spar_A12_Lit1
2366
SUP35_Smik_IF01815T_30
2366
SUP35_Kla_AB039749
2366

Compare algorithms on 10 sequences

algorithm alignment length (nt) time (seconds)
ClustalW 2148 3.30
MUSCLE 2275 2.02
MAFFT 2166 3.49
Kalign 2155 0.30
T-Coffee 2210 262.50
prank 2366 6.11

Conclusion: The best algorithm here is Prank. It has comparable speed with other algorithms and the highest length of alignment among them.

Strange file SUP35_10seqs_strange_aln.fa

As we can see there is one strange sequence in the file (highlighted) which does not match well with others.
It is connected with that fact that this sequence is a reverse-complement. So, we should obtain its reverse-complement and then realign
sequences.

In [19]: ! time clustalw -INFILE=SUP35_250seqs.fa -OUTPUT=FASTA -OUTFILE=SUP35_250seqs_clustalw.fa > clustalw_
250.log

1541.83user 0.18system 25.43.73elapsd 99%CPU (bavgtext+bavgdata 135960maxresident)k
4280inputs+4720outputs (20major+30423minor)pagefaults 0swaps

In [20]: record_dict = IO.to_dict(IO.parse("SUP35_250seqs_clustalw.fa", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))
break

SUP35_Scer_T_52_2N_CP007848
2179

Muscle

In [21]: ! time muscle -in SUP35_250seqs.fa -out SUP35_250seqs_muscle.fa

MUSCLE V3.8.1551 by Robert C. Edgar

This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

SUP35_250seqs 250 seqs, lengths min 1769, max 2183, avg 2054
00:00:00 19 MB(-2%) Iter 1 100.00% K-mer dist pass 1
00:00:00 19 MB(-2%) Iter 1 100.00% K-mer dist pass 2
00:00:14 332 MB(-43%) Iter 1 100.00% Align node
00:00:14 333 MB(-43%) Iter 1 100.00% Root alignment
00:00:19 333 MB(-43%) Iter 2 100.00% Refine tree
00:00:19 333 MB(-43%) Iter 2 100.00% Root alignment
00:00:19 333 MB(-43%) Iter 2 100.00% Refine biparts
00:01:02 333 MB(-43%) Iter 3 100.00% Refine biparts
00:01:46 333 MB(-43%) Iter 4 100.00% Refine biparts
00:01:46 333 MB(-43%) Iter 5 100.00% Refine biparts
00:01:46 333 MB(-43%) Iter 5 100.00% Refine biparts
1152inputs+1168outputs (9major+80479minor)pagefaults 0swaps

In [22]: record_dict = IO.to_dict(IO.parse("SUP35_250seqs_muscle.fa", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))
break

SUP35_Agos_ATCC_10895_NM_211584
2313

MAFFT

In [23]: ! time mafft --auto SUP35_250seqs.fa >SUP35_250seqs_mafft.fa

nthread = 0
nthreadpair = 0
nthreadb = 0
openalty ex = 0
stacksize: 8192 kb
generating a scoring matrix for nucleotide (dist=200) ... done
gap Penalty = -1.53, -0.00, -0.00

Making a distance matrix ..

There are 12 ambiguous characters.
200 / 250
done.

Constructing a UPGMA tree (effree=0) ...
240 / 250
done.

Progressive alignment 1/2...
STEP 247 / 249
Reallocating .done. *allocen = 5227
STEP 249 / 249
done.

Making a distance matrix from msa..
200 / 250
done.

Constructing a UPGMA tree (effree=1) ...
240 / 250
done.

Progressive alignment 2/2...
STEP 247 / 249
Reallocating .done. *allocen = 5227
STEP 249 / 249
done.

distbfast (nuc) Version 7.453
alg=A, model=DNA200 (2), 1.53 (4.59), -0.00 (-0.00), noshift, amax=0.0
0 thread(s)

generating a scoring matrix for nucleotide (dist=200) ... done
dndpre (nuc) Version 7.453
alg=X, model=DNA200 (2), 1.53 (4.59), 0.37 (1.11), noshift, amax=0.0
0 thread(s)

minimumpweight = 0.000018
autosubalignment = 0.000000
nthread = 0
randomseed = 0
blossum 62 / kimura 200
poffset = 0
niter = 2
suerr_global = 0.100000
nadd = 2
generating a scoring matrix for nucleotide (dist=200) ... done
240 / 250

Segment 1/ 20 1- 145
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 2/ 20 145- 248
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 3/ 20 248- 299
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 4/ 20 299- 367
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 5/ 20 361- 402
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 6/ 20 400- 462
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 7/ 20 462- 608
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 8/ 20 608- 714
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 9/ 20 714- 807
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 10/ 20 807- 911
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 11/ 20 911- 1045
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 12/ 20 1045-1196
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 13/ 20 1196-1347
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 14/ 20 1347-1498
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 15/ 20 1498-1649
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 16/ 20 1649-1800
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 17/ 20 1800-1951
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 18/ 20 1951-2192
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 19/ 20 2192-2245
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

Segment 20/ 20 2245-2323
STEP 002-248-1 identical. identical. identical. identical. identical. identical. 1
denical. identical. identical. identical. identical. identical. identical.
Coverged.

distbtr (nuc) Version 7.453
alg=A, model=DNA200 (2), 1.53 (4.59), -0.00 (-0.00), noshift, amax=0.0
0 thread(s)

FFT-NS-1 (Standard)
Iterative refinement method (max. 2 iterations)

If unsure which option to use, try 'mafft --auto input > output'.
For more information, see 'mafft --help', 'mafft --man' and the mafft page.

The default gap scoring scheme has been changed in version 7.110 (2013 Oct).
It tends to insert more gaps into gap-rich regions than previous versions.
To disable this change, add the -leavegappyregion option.

17.33user 0.15system 0:17.49elapsd 99%CPU (bavgtext+bavgdata 60600maxresident)k
115041inputs+50490outputs (52major+56496minor)pagefaults 0swaps

In [24]: record_dict = IO.to_dict(IO.parse("SUP35_250seqs_mafft.fa", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))
break

SUP35_Kla_AB039749
2322

Kalign
```



```
In [26]: ! time kalign <SUP35_250seqs.faa >SUP35_250seqs.kalign.faa
Kalign version 2.04q; Copyright (C) 2004, 2005, 2006 Timo Lassmann

Kalign is free software. You can redistribute it and/or modify
it under the terms of the GNU General Public License as
published by the Free Software Foundation.

reading from STDIN: Found 250 sequences
Aligning percent doneSequences with these parameters:
    217.00000000000000 gap open penalty
    39.40000153      gap extension
    282.0000000000000 terminal gap penalty
    283.0000000000000 bonus
Alignment will be written to stdout.

Distance Calculation:
100 percent done
Alignment:
    0 percent doneSaving mem...
    0 percent doneSaving mem...
    1 percent doneSaving mem...
    1 percent doneSaving mem...
    1 percent doneSaving mem...
    2 percent doneSaving mem...
    2 percent doneSaving mem...
    3 percent doneSaving mem...
    4 percent doneSaving mem...
    4 percent doneSaving mem...
    5 percent doneSaving mem...
    6 percent doneSaving mem...
    6 percent doneSaving mem...
    8 percent doneSaving mem...
    8 percent doneSaving mem...
    9 percent doneSaving mem...
    9 percent doneSaving mem...
    10 percent doneSaving mem...
    10 percent doneSaving mem...
    12 percent doneSaving mem...
    12 percent doneSaving mem...
    12 percent doneSaving mem...
    13 percent doneSaving mem...
    14 percent doneSaving mem...
    14 percent doneSaving mem...
    14 percent doneSaving mem...
    15 percent doneSaving mem...
    15 percent doneSaving mem...
    16 percent doneSaving mem...
    16 percent doneSaving mem...
    17 percent doneSaving mem...
    17 percent doneSaving mem...
    18 percent doneSaving mem...
    18 percent doneSaving mem...
    19 percent doneSaving mem...
    19 percent doneSaving mem...
    20 percent doneSaving mem...
    20 percent doneSaving mem...
    21 percent doneSaving mem...
    22 percent doneSaving mem...
    22 percent doneSaving mem...
    22 percent doneSaving mem...
    23 percent doneSaving mem...
    23 percent doneSaving mem...
    24 percent doneSaving mem...
    24 percent doneSaving mem...
    25 percent doneSaving mem...
    25 percent doneSaving mem...
    26 percent doneSaving mem...
    26 percent doneSaving mem...
    27 percent doneSaving mem...
    27 percent doneSaving mem...
    28 percent doneSaving mem...
    28 percent doneSaving mem...
    29 percent doneSaving mem...
    29 percent doneSaving mem...
    30 percent doneSaving mem...
    30 percent doneSaving mem...
    31 percent doneSaving mem...
    31 percent doneSaving mem...
    32 percent doneSaving mem...
    32 percent doneSaving mem...
    33 percent doneSaving mem...
    33 percent doneSaving mem...
    34 percent doneSaving mem...
    34 percent doneSaving mem...
    35 percent doneSaving mem...
    35 percent doneSaving mem...
    36 percent doneSaving mem...
    36 percent doneSaving mem...
    37 percent doneSaving mem...
    37 percent doneSaving mem...
    38 percent doneSaving mem...
    38 percent doneSaving mem...
    39 percent doneSaving mem...
    39 percent doneSaving mem...
    40 percent doneSaving mem...
    40 percent doneSaving mem...
    41 percent doneSaving mem...
    41 percent doneSaving mem...
    42 percent doneSaving mem...
    42 percent doneSaving mem...
    43 percent doneSaving mem...
    43 percent doneSaving mem...
    44 percent doneSaving mem...
    44 percent doneSaving mem...
    45 percent doneSaving mem...
    45 percent doneSaving mem...
    46 percent doneSaving mem...
    46 percent doneSaving mem...
    47 percent doneSaving mem...
    47 percent doneSaving mem...
    48 percent doneSaving mem...
    48 percent doneSaving mem...
    49 percent doneSaving mem...
    49 percent doneSaving mem...
    50 percent doneSaving mem...
    50 percent doneSaving mem...
    51 percent doneSaving mem...
    51 percent doneSaving mem...
    52 percent doneSaving mem...
    52 percent doneSaving mem...
    53 percent doneSaving mem...
    53 percent doneSaving mem...
    54 percent doneSaving mem...
    54 percent doneSaving mem...
    55 percent doneSaving mem...
    55 percent doneSaving mem...
    56 percent doneSaving mem...
    56 percent doneSaving mem...
    57 percent doneSaving mem...
    57 percent doneSaving mem...
    58 percent doneSaving mem...
    58 percent doneSaving mem...
    59 percent doneSaving mem...
    59 percent doneSaving mem...
    60 percent doneSaving mem...
    60 percent doneSaving mem...
    61 percent doneSaving mem...
    61 percent doneSaving mem...
    62 percent doneSaving mem...
    62 percent doneSaving mem...
    63 percent doneSaving mem...
    63 percent doneSaving mem...
    64 percent doneSaving mem...
    64 percent doneSaving mem...
    65 percent doneSaving mem...
    65 percent doneSaving mem...
    66 percent doneSaving mem...
    66 percent doneSaving mem...
    67 percent doneSaving mem...
    67 percent doneSaving mem...
    68 percent doneSaving mem...
    68 percent doneSaving mem...
    69 percent doneSaving mem...
    69 percent doneSaving mem...
    70 percent doneSaving mem...
    70 percent doneSaving mem...
    71 percent doneSaving mem...
    71 percent doneSaving mem...
    72 percent doneSaving mem...
    72 percent doneSaving mem...
    73 percent doneSaving mem...
    73 percent doneSaving mem...
    74 percent doneSaving mem...
    74 percent doneSaving mem...
    75 percent doneSaving mem...
    75 percent doneSaving mem...
    76 percent doneSaving mem...
    76 percent doneSaving mem...
    77 percent doneSaving mem...
    77 percent doneSaving mem...
    78 percent doneSaving mem...
    78 percent doneSaving mem...
    79 percent doneSaving mem...
    79 percent doneSaving mem...
    80 percent doneSaving mem...
    80 percent doneSaving mem...
    81 percent doneSaving mem...
    81 percent doneSaving mem...
    82 percent doneSaving mem...
    82 percent doneSaving mem...
    83 percent doneSaving mem...
    83 percent doneSaving mem...
    84 percent doneSaving mem...
    84 percent doneSaving mem...
    85 percent doneSaving mem...
    85 percent doneSaving mem...
    86 percent doneSaving mem...
    86 percent doneSaving mem...
    87 percent doneSaving mem...
    87 percent doneSaving mem...
    88 percent doneSaving mem...
    88 percent doneSaving mem...
    89 percent doneSaving mem...
    89 percent doneSaving mem...
    90 percent doneSaving mem...
    90 percent doneSaving mem...
    91 percent doneSaving mem...
    91 percent doneSaving mem...
    92 percent doneSaving mem...
    92 percent doneSaving mem...
    93 percent doneSaving mem...
    93 percent doneSaving mem...
    94 percent doneSaving mem...
    94 percent doneSaving mem...
    95 percent doneSaving mem...
    95 percent doneSaving mem...
    96 percent doneSaving mem...
    96 percent doneSaving mem...
    97 percent doneSaving mem...
    97 percent doneSaving mem...
    98 percent doneSaving mem...
    98 percent doneSaving mem...
    99 percent doneSaving mem...
    100 percent done
MSQ0NDQDQGQQGQVHWKQYHQVNYNQNGQVGVHVGQQAQVQAGVQGPQGAQV
QGYNVQAAPAAQSQQMTLKDFQDNQGSITNAAXPKPKLLASSSGILVLGAKKPVAKT
EKDTSEKAFTTDONEAGESELFPKDOLKSIEAEPKTKEHTPSADOTSSEKTSAAKU
YTGMNSVDAALIKEDVEEVEEYVKWFSGDHVSIFTHGVNMGASGTNGNLVLVTGS
VPRKVKEYERAEAGRGQWLKSMVDNTKEERNODKTEIVGRAVFTEKRRYTLIDAP
GHMVYSMIGSAGADIGILVISARKEVETGFEGGGTGREHALAKTQCVNMIVMIVN
TSTGANSVDAALIKEDVEEVEEYVKWFSGDHVSIFTHGVNMGASGTNGNLVLVTGS
VTGSLLEYLDNMKTTDRHNIAFPLPIASKMKDMGTVEKIESGHIRKGQNTLLMPNR
TSVELITTYNEESVEDMAVCQVRILRIKGVEEIESAGFLVLTSPNMPKNVTRFVAQT
ALVELKSIHSAGFSQMDWHITAEIVTVLRLKLEKSDGNKSNKPPAFKGNKIIVAVI
ETNEPVCVETYDDYPQLGRFLTRODGTIIAIGKVIKLEN'
>SUP35_Agos_ATCC_10895_NM_211584_1
MSEEDIQSQNDQDGQQAQKQNDQNGQNDQGVYNPNFGNVQGVVPYGQVQAGVQGO
GNISYGVQGGQATAPVLMNFEGYVNPATAPPKPKTKLLASSSGILVLGAKKPVAKK
EKAKEPTKEEPSAEAGPKSESADATSSDKAVPSIEKLISEADTAQADAAGAATS
SUALIKEDVEEVEEYVKWFSGDHVSIFTHGVNMGASGTNGNLVLVTGSVDVEX
YEKAKEAGRGQWLKSMVDNTKEERNODKTEIVRSYFETEKRYTLLDAPGHMYMSE
Georf searches for ORFs in an input file and then outputs them.
-minsize: minimum size of ORF (otherwise it will output too many ORFs per sequence and most of them will be garbage)

In [28]: ! time prank -d SUP35_250seqs.faa -o SUP35_250seqs.prank.faa
-----
PRANK v.17B427:
-----
Input for the analysis
- aligning sequences in 'SUP35_250seqs.faa'
- using inferred alignment guide tree
- option '+F' is not used; it can be enabled with '+f'
- external tools available:
  Mafft for initial alignment
  Exonerate for alignment anchoring
Correcting (arbitrarily) for multifurcating nodes.
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 1.
Alignment score: 9063
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 2.
Alignment score: 7787
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 3.
Alignment score: 8220
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 4.
Alignment score: 7867
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 5.
Alignment score: 7869
Writing
- alignment to 'SUP35_250seqs.prank.faa.best.fas'
Analysis done. Total time 119s
116.9User 1.16Sysmem 1.59:50elapseds 98%CPU (avgtext+avgdata 66236maxresident)k
0Inputs+1674Outputs (0major+629152minor)pagefaults 0swaps

In [29]: record_dict = IO.to_dict(IO.parse(("SUP35_250seqs.prank.best.fas", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))
    break
SUP35_Kla_AB039749_1
2214

t-coffee
It is too slow to launch it on 250 sequences.

In [27]: # ! time t-coffee -infile=SUP35_250seqs.faa -outfile=SUP35_250seqs.tcoffee.faa -output=fasta_aln

Prank
In [28]: ! time prank -d SUP35_250seqs.faa -o SUP35_250seqs.prank.faa
-----
PRANK v.17B427:
-----
Input for the analysis
- aligning sequences in 'SUP35_250seqs.faa'
- using inferred alignment guide tree
- option '+F' is not used; it can be enabled with '+f'
- external tools available:
  Mafft for initial alignment
  Exonerate for alignment anchoring
Correcting (arbitrarily) for multifurcating nodes.
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 1.
Alignment score: 9063
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 2.
Alignment score: 7787
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 3.
Alignment score: 8220
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 4.
Alignment score: 7867
Correcting (arbitrarily) for multifurcating nodes.
Generating multiple alignment: iteration 5.
Alignment score: 7869
Writing
- alignment to 'SUP35_250seqs.prank.faa.best.fas'
Analysis done. Total time 119s
116.9User 1.16Sysmem 1.59:50elapseds 98%CPU (avgtext+avgdata 66236maxresident)k
0Inputs+1674Outputs (0major+629152minor)pagefaults 0swaps

In [30]: record_dict = IO.to_dict(IO.parse(("SUP35_250seqs.prank.best.fas", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))
    break
SUP35_Agos_ATCC_10895_NM_211584_1
2472

Compare algorithms on 250 sequences

algorithm          alignment length (nt)   time (seconds)
-----
ClustalW           2179             1552.93
MUSCLE              2332              93.31
MAFFT               2231              15.33
Kalign              2214              15.07
t-coffee            2179             102.53
prank                2472              102.53

Conclusion: here we can highlight two algorithms - MAFFT (it is the fastest one end has high alignment length) and Prank - it is a bit slower but has the highest length of alignment

Translation into amino acids

In [30]: # sudo apt-get update -y
# sudo apt-get install -y emboss

TransE translates nucleotide sequence into amino acids just directly reading codons

In [31]: ! transeq -sequence SUP35_10seqs.faa -outseq SUP35_10seqs.t.faa
Translate nucleic acid sequences

In [32]: ! head -20 SUP35
```



```
[In [45]:] record_dict = IO.to_dict(IO.parse("SUP3S_10seqs.kallign.faa", "fasta"))
for key in record_dict.items():
    print(key[0], "\n", len(key[1].seq))

SUP3S_Kla_AB039749_1
752
SUP3S_Agos_ATCC_10895_NM_211584_1
720
SUP3S_Scer_74-D694_GCA_001578265_1_1
752
SUP3S_Spar_A12_Lit1_1
720
SUP3S_Smk_IFO18157_30_1
752
SUP3S_Skud_IFO18027_36_1
720
SUP3S_Shou_uniq28_CM093560_1
752
SUP3S_Scer_beer078_CM095938_1
720
SUP3S_Sarb_H-6_chrxIII_CM091575_1
720
SUP3S_Seub_CB512357_chr_II_IV_DF968535_1
720

t-coffee

In [46]: ! time t-coffee -infile=SUP3S_10seqs.g.faa -outfile=SUP3S_10seqs.tcoffee.faa -output=fastaaln

Command Line Arguments: [/home/isemenov/anondaca/lib/t_coffee/, '-infile=SUP3S_10seqs.g.faa', '-out
file=SUP3S_10seqs.tcoffee.faa', '-output=fastaaln']
Install folder: /home/isemenov/anondaca/lib/t_coffee-11.0.8

PROGRAM: T-COFFEE Version 11.00.8cbe486 (2014-08-12 22:05:29 - Revision 8cbe486 - Build 477)
- full_log S [0] nsd
- nemmed_score S [0] nsd
- mem_name S [0] mem
- extend 0 [1] 1
- extend_mode S [0] very_fast_triplet
- max_n_dial S [0] 0
- seq_name_for_quadruplet S [0] all
- compact S [0] default
- clean_seq S [0] ANY
- do_self FL [0] 0
- do_normalise D [0] 1000
- delete_file FL [0] 0
- setenv S [0] 0
- template_mode S [0] 0
- rlib FL [0] 0
- remove_template_file S [0] 0
- profile_template_file S [0] 0
- in S [0] 0
- out S [0] 0
- aln S [0] 0
- method_limits S [0] 0
- method S [0] 0
- lib_min S [0] 0
- lib_max S [0] 0
- profile S [0] 0
- profile1 S [0] 0
- profile2 S [0] 0
- pdb S [0] 0
- relax_lib 0 [0] 1
- rlib FL [0] 0
- shrink_lib 0 [0] 0
- out_lib W_F [0] no
- lib_only D [0] primary
- outsequence W_F [0] no
- dpa FL [0] 0
- seq_source S [0] ANY
- cosmetic_penalty D [0] 0
- gapopen 0 [0] 0
- gapext 0 [0] 0
- fgapopen 0 [0] 0
- hominch 0 [0] 0
- newtree W_F [0] default
- tree W_F [0] NO
- userdef R_P [0] 0
- tree_order S [0] nj
- distance_matrix_mode S [0] ktup
- distance_matrix_sim_mode S [0] 1dmat_sml
- quicktree FL [0] 0
- outfile W_F [1] SUP3S_10seqs.tcoffee.faa
- maximise FL [1] 1
- output_pos S [1] fasta_aln
- len D [0] 0
- infile R_F [1] SUP3S_10seqs.g.faa
- matrix S [0] default
- igmode D [0] 0
- update_mode S [0] cw_profile_profile
- profile_comparison S [0] profile
- dp_mode S [0] linked_prank_wise
- ktuple D [0] 1
- ndiag D [0] 0
- dia3_threshold D [0] 0
- dia4_threshold D [0] 0
- sim_matrix S [0] vasiliky
- transform S [0] 0
- reorderer S [0] input
- in_order S [0] aligned
- sequence S [0] 0
- case S [0] keep
- cpu 0 [0] 0
- lutil D [0] -1
- maxseq D [0] 7
- maxlen D [0] -1
- sample_dp D [0] 0
- weight D [0] default
- seq_weight S [0] no
- align FL [1] 1
- prot_min_sim D [0] 40
- domain FL [0] 0
- start D [0] 0
- len D [0] 0
- scale D [0] 0
- mcca_interactive FL [0] 0
- method_evaluate FL [0] default
- evaluate_mode S [0] triplet
- get_type FL [0] 0
- clean_aln D [0] 0
- criteria_threshold D [0] [1] 1
- clean_evaluation 0 [1] 1
- clean_evaluate_mode S [0] t_coffee.fast
- extend_matrix FL [0] 0
- fast_mode FL [0] 0
- prot_max_sim D [0] 90
- prot_min_cov 0 [40] 40
- seq_id D [0] 0
- pdb_min_sim D [35] 35
- pdb_max_sim D [100] 100
- pdb_blast_server W_F [0] EBI
- blast W_F [0] 0
- blast_server W_F [0] EBI
- pdb_db D [0] 0
- protein_db W_F [0] uniprot
- method_log W_F [0] no
- tool_to_use W_F [0] 0
- cache W_F [0] use
- align_pdb_param_file W_F [0] no
- align_pdb_hatch_mode W_F [0] hatch_ca_trace_bubble
- external_aligner S [0] NO
- nsa_mode S [0] tree
- master S [0] no
- blast_seq S [0] 0
- laalign_top 0 [0] 10
- iterate 0 [0] 0
- trie D [0] 0
- split D [0] 0
- trisafe S [0] default
- split D [0] 0
- split_inuse_thres D [0] 0
- split_score_status 0 [0] 0
- check_pdb_thresh 0 [0] 0
- clean_seq_name S [0] 0
- seq_to_keep S [0] 0
- dpa_master_aln S [0] 0
- dpa_maxseq D [0] 0
- dpa_min_score1 D [0] 0
- dpa_min_score2 D [0] 0
- dpa_keep_tmfile D [0] 0
- multi_core S [0] templates_jobs_rels_nsa_evaluate
- n_core 0 [0] 0
- overaln_proc 0 [0] 0
- lib_list S [0] none
- prune_lib_mode S [0] 0
- mode D [0] none
- rna_lib S [0] 0
- no_warning 0 [0] 0
- run_local_script 0 [0] 0
- glues D [0] default
- proxy S [0] unset
- email S [0] 0
- clean_overaln S [0] 0
- overaln_param S [0] 0
- overaln_mode S [0] 0
- overaln_model S [0] 0
- overaln_target 0 [0] 0
- overaln_P1 D [0] 0
- overaln_P2 D [0] 0
- overaln_P4 D [0] 0
- exon_boundaries S [0] 0
- dump S [0] 0
- display 0 [0] 100

INPUT FILES
Input File (S) SUP3S_10seqs.g.faa Format fasta.seq
Input File (M) prmobr.pair

Identify Master Sequences [no]:

Master Sequences Identified
INPUT SEQUENCES: 10 SEQUENCES [PROTEIN]
Input File SUP3S_Agos_ATCC_10895_NM_211584_1 Seq SUP3S_Agos_ATCC_10895_NM_211584_1 Le
nth 681 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Kla_AB039749_1 Le
nth 700 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Sarb_H-6_chrxIII_CM091575_1 Le
nth 680 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Shou_uniq28_CM093560_1 Le
nth 685 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Scer_74-D694_GCA_001578265_1_1 Le
nth 666 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Scer_beer078_CM095938_1 Le
nth 674 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Skud_IFO18027_36_1 Le
nth 678 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Smk_IFO18157_30_1 Le
nth 681 type PROTEIN Struct Unchecked
Input File SUP3S_10seqs.g.faa Seq SUP3S_Spar_A12_Lit1_1 Le
nth 681 type PROTEIN Struct Unchecked

Multi Core Mode: 0 processors:

--- Process Method/Library/Aln SSUP3S_10seqs.g.faa
--- Process Method/Library/Aln Mproba_pair
xxx Retrieved Mproba.pair

All Methods Retrieved

MANUAL PENALTIES: gapopen=0 gapext=0

Library Total Size: [87364]

Library Relaxation: Multi_proc [8]

! [Relax Library][Total= 1][100 %] ELAPSED TIME: 0 sec.]
Relaxation Summary: [87364]->[77444]

UN-WEIGHTED MODE: EVERY SEQUENCE WEIGHTS 1

MAKE GUIDE TREE
[MODE=no][DONE]

PROGRESSIVE_ALIGNMENT [Tree Based]
Group 1: SUP3S_Agos_ATCC_10895_NM_211584_1
Group 2: SUP3S_Kla_AB039749_1
Group 3: SUP3S_Sarb_H-6_chrxIII_CM091575_1
Group 4: SUP3S_Shou_uniq28_CM093560_1
Group 5: SUP3S_Scer_74-D694_GCA_001578265_1_1
Group 6: SUP3S_Scer_beer078_CM095938_1
Group 7: SUP3S_Skud_IFO18027_36_1
Group 8: SUP3S_Smk_IFO18157_30_1
Group 9: SUP3S_Spar_A12_Lit1_1
Group 10: SUP3S_Spar_A12_Lit1_1

Group 11: [Group 5 ( 1 seq)] with [Group 4 ( 1 seq)]->[Len= 685][PID:28734]
Group 12: [Group 6 ( 1 seq)] with [Group 11 ( 2 seq)]->[Len= 685][PID:28734]
Group 13: [Group 12 ( 3 seq)] with [Group 10 ( 1 seq)]->[Len= 685][PID:28734]
Group 14: [Group 13 ( 4 seq)] with [Group 12 ( 4 seq)]->[Len= 685][PID:28734]
Group 15: [Group 8 ( 1 seq)] with [Group 14 ( 5 seq)]->[Len= 685][PID:28734]
Group 16: [Group 3 ( 1 seq)] with [Group 15 ( 6 seq)]->[Len= 685][PID:28734]
Group 17: [Group 16 ( 7 seq)] with [Group 17 ( 7 seq)]->[Len= 687][PID:28734]
Group 18: [Group 18 ( 2 seq)] with [Group 17 ( 8 seq)]->[Len= 752][PID:28789][For
ked]

! [Final Evaluation][Tot= 94][Libach [http://www.tcoffee.org] [MODE: ], CPU=0.00 sec, S
CORE=958, Nseq=10, Len=752

CLUSTAL FORMAT for T-COFFEE version 11.00.8cbe486 [http://www.tcoffee.org] [MODE: ], CPU=0.00 sec, S
CORE=958, Nseq=10, Len=752

SUP3S_Kla_AB039749_1 MSDNQNDQDQGGGQGVYNYNGVYNGVYNYQQ-QGVYGWQGQ-GAP-QGVQ
SUP3S_Agos_ATCC_10895_NM_211584_1 MSEED-QIQSQ-GNDQ--GSGAQDQGDQGGQGNFYQY-NPS-NFQ
SUP3S_Scer_74-D694_GCA_001578265_1_1 MSDSN-QGNQ-QGVYQSGNQGQGNRYGYQAQAQA-QPAGGYQ
SUP3S_Spar_A12_Lit1_1 HSDN-QGNQ-QSYQYQNPDQGNQGNRYGYQAQAQA-QPAGGYQ
SUP3S_Smk_IFO18157_30_1 MSDSN-QGNQ-QGVYQNPNQDGQGNRYGYQAQAQA-QPAGHYQ
SUP3S_Shou_uniq28_CM093560_1 MSDSN-QGNQ-QGVYQSGNQGQGNRYGYQAQAQA-QPAGGYQ
SUP3S_Scer_beer078_CM095938_1 MSDSN-QGNQ-QGVYQSGNQGQGNRYGYQAQAQA-QPAGGYQ
SUP3S_Sarb_H-6_chrxIII_CM091575_1 MSDFP-NNQN-QGVYQGNQGNQGNRYGYQAQAQA-QPAGGYQ
SUP3S_Seub_CB512357_chr_II_IV_DF968535_1 MSDSN-QGNQ-QGVYQGNQGNQGNRYGYQAQAQA-QPAGGYQ

SUP3S_Kla_AB039749_1 AYQAVQGPQGGAY-QGVNPQQA-Q-Q-QVPP-Q-----QY-----
SUP3S_Agos_ATCC_10895_NM_211584_1 GYGE-QSATATPLTNFNKGGVTPMTAKPKMKTKLASSSSGLVKGK
SUP3S_Scer_74-D694_GCA_001578265_1_1 NVGVYSYQGGY-QGVNPDAQVQ--QYNRP-
SUP3S_Spar_A12_Lit1_1 NVGVYSYQGGY-QGVNP
```

SUP35\_Smk1\_IF018157\_38\_1[1\_--2843]  
 791  
 SUP35\_Smk1\_IF018157\_38\_1[1\_--2843]  
 791  
 SUP35\_Sarb\_H-6\_chrXIII\_CM081575\_1[1\_--2848]  
 791  
 SUP35\_Skud1\_IF018827\_36\_1[1\_--2834]  
 791  
 SUP35\_Saub\_CBS12357\_chr\_II\_IV\_DF968535\_1[1\_--2822]  
 791

### Compare algorithms on 10 protein sequences

algorithm	alignment length (nt)	time (seconds)
ClustalW	719	0.38
CLustalO	757	1.22
MUSCLE	743	0.16
MAFFT	759	0.37
Kalign	720	0.04
t-coffee	752	31.96
prank	791	9.91

**Confusion:** the best algorithm here is **MAFFT**. It has rather good length of alignment and speed (balance between these two metrics)

## Adding sequences to alignments

MAFFT and Muscle allows to add new sequences to existing alignments without their recalculations

```
In [50]: | cat SUP35_2addseqs.fa
```

```
>SUP35_scer_bcer090_C0M05746
ATTGCGGATTCAACACGAGCAGCAATCARGAACCAATCAGCAATCAGCAGCAGACAGACCAACACAG
AAGTATCAACAGATTCACCAAGGTTATCAAGTCAATCTTGAAACCGATTCAGATTCGAGCGGTGATACCA
CAACAGTATATATCTTCAGAGGAGGATATCAACAGATATCAATCTTCAGCGGGTGTATCAGCAGCAATTCATC
CAAGATGTTGGCTGGAAHATACAAAGATTCAGCTATCAATCTTCAGAGATTCAGAGATGTTGAGCTGGTTC
CCACAGCATCTCTTCAGGATGATATCTTTGAGACGCTTTGACAGATTCAGACAGAGGCGCTCTCCACACCA
AAGAGATCTTGAGAGCTTGATCTGATCTTCGGATGATCTGATCTGAGCAATCTGCGCAGATAGATCTGACACA
CTGCTGCGCATCTCTTAAAGAGAGAGAGAGATCTGATCTGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAG
AAGTCTGGAAGGACCAAGTTTAAAGAGAGAGAGAGACCAATCAGCTAGGAAAGAGAGAGAGAGAGAAATG
GAATCTCCAAAGTAGAGAGACCTTAAATCTCTGTAATCAACATATAATCAACAGATCCGATATCTACCA
AGCTCTGCTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCT
TAAAGATGATCTCTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTT
CTATCATCTGACCTGGCTCTTGAGATAGAGAGCAATTTGAGAAATATAAGAAAGAGAGACAGGATGACGCA
GATCTGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAG
AGTTGGTAGAGGAGGCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTT
TACGTTCTCCGAGATGATCCGGGGTGGCTCTCTCAGGCACTGTTGGTGGTTTGGTCACTTCCGCCAGGAGG
GATCTGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAG
TCTTAAATAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAGAG
GATGACCATCTGAGATGATTCAGCAATTTCTTGAGAGCAATCTGTTTACACATTAAGAGACAGATCTGTGAT
TATCTGAGATGATCTGCGCTACAGTTGGCTGGAAHATTTGAGATGATGATGATGATGATGATGATGATGAT
CACCGGCCCAATCTGTTTAAAGATTTGTGATCAATGAGACAGCTGCGATCGCATCACTGCTTCATCT
```

[illegible]

```

ATGTTGCCATTCCGGCATGAAGGATGAGGATACCATCTGTGAAGGATAAATTTAGCCGCGGTACATA
CTCAAAAGGAGGTCATCCACDCCATCTACTGTTGATCAACAAACCGCTGTGGAAATTTCAAAATTTATACAGCA
ACTCTGAATTAATTCATTTTATATGCTGATCTGTGTGTGGACAGTTTAAATTAAGATCAAAAGTGTGTGAAGA
GAAGCAATTTCTCCAGAGTTTGTGTACTAATACCTCGCAAGAACCCCTCAAGAAGTTTATCCAAAGTTTGTAG
CTCAAAATTCGCTATGATTAATAATCTATGATGACGGCGGCTTTTGTATGCTTGTATCGATCTATCTCAT
AGCCTACGAAGAGAGGATCTTTGTTTGAATTTATGTCACAAATTTAGAGAAAGGTATCCACAGCTGAGTCAAG
AAACCCDCCGTTTTCGTAAAGAGAGGATGAAGGATCTGCGCTGTTTGTAGAAATCGAAGCTCCAGTTTGTG
TGGAAGATCTACCAAGATACCCCTCAATTAGGTAGATCTTTGAGAGACCTCAAGGTACACATAGCAAT
TGGTAAATTTGTTAAATTTCCCGAGTAA

```

Firstly, we align these two sequences to each other

```
In [51]: ! muscle -in SUP35_2ddseqs.fa -out SUP35_2ddseqs_muscle.fa
```

MUSCLE v3.8.1551 by Robert C. Edgar

<http://www.drive5.com/muscle>  
This software is donated to the public domain.  
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

```

SUP35_2ddseqs 2 secs, lengths min 2088, max 2088, avg 2088
00:00:00   15 MB(-) Citer 1 100.00% k-mer dist pass 2
00:00:00   15 MB(-) Citer 1 100.00% k-mer dist pass 2
00:00:00   27 MB(-) Citer 1 100.00% Align node
00:00:00   27 MB(-) Citer 1 100.00% Root alignment

```

And then we start to existing alignments of our 250 sequences

```
In [52]: ! muscle -profile -i1n SUP35_250seqs_muscle.fa -i2n SUP35_2addseqs.fa -out SUP35_252seqs_muscle.fa
MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R. C. Nucleic Acids Res 32(5), 1792-97.

00:00:00 15 MB(-2%) Reading SUP35_250seqs_muscle.fa
00:00:00 17 MB(-2%) 259 seqs 2313 cols
00:00:00 17 MB(-2%) Reading SUP35_2addseqs.fa
00:00:00 17 MB(-2%) 2 seqs 2058 cols
00:00:01 18 MB(-2%) Aligning profiles
00:00:01 29 MB(-4%) Building output
00:00:01 30 MB(-4%) Writing output

Similar procedure can be done with MAFFT

In [53]: ! mafft --auto SUP35_2addseqs.fa > SUP35_2addseqs_mafft.fa

outputthat23=16
treein = 0
compacttree = 0
stacksize= 8192 kb
generating a scoring matrix for nucleotide (dist=200) ... done
All-to-all alignment.
thfast-pair (mcs) version 7.453
align, model=DNA200 (2), 2.00 (6.00), -0.10 (-0.30), nohsift, amax=0.0
0 thread(s)
```

```

outputhat23=16
Loading 'hat3.seed' ...
done.
Writing hat3 for iterative refinement
generating a scoring matrix for nucleotide (dist=200) ... done
Gap Penalty = -1.50, +0.00, +0.00
tubtree = 1, compacttree = 0
Constructing a UPGMA tree ...
    0 / 2
done.

Progressive alignment ...
STEP      1 /1
done.
tbfast (nuc) Version 7.453
alg=A, model=DNA200 (2), 1.53 (4.50), -0.00 (-0.00), noshift, amax=0.0
1 thread(s)

minimumweight = 0.000010
autosubalignment = 0.000000
nthread = 0
randomseed = 0
blossum 62 / kimura 200
poffset = 0
niter = 16
sueff_global = 0.100000
nadd = 16
Loading 'hat3' ... done.
generating a scoring matrix for nucleotide (dist=200) ... done

```

```
Segment 1/ 1 1-2059
done 001-001-1 identical.
divditr (nuc) Version 7.453
alpha, model=IAC90 (2), 1.53 (4.50), -0.00 (-0.00), noshift, amax=0.0
0 thread(s)
```

Strategy:

- I-INS-i (Probably most accurate, very slow)
- Iterative refinement method (<i>i</i>) with LOCAL pairwise alignment information

If unsure which option to use, try 'mafft --auto input > output'.

For more information, see 'mafft -help', 'mafft --man' and the mafft page.

The default gap scoring scheme has been changed in version 7.11b (2013 Oct).

It tends to insert more gaps into gap-rich regions than previous versions.

To disable this change, add the --leavegappyregion option.



```
In [54]: | mafft --add SUP35_2addseqs_mafft.fa SUP35_250seqs_mafft.fa > SUP35_252seqs_mafft.fa

nadd = 2
nthread = 0
nthreadb = 0
nthreadtb = 0
penalty_ex = 0
stacksize: 8192 kb
generating a scoring matrix for nucleotide (dist=200) ... done
Gap Penalty = -1.53, +0.00, +0.00

Making a distance matrix ..

There are 12 ambiguous characters.
 201 / 252
done.

Constructing a UPGMA tree (efffree=0) ...
 250 / 252
done.

Progressive alignment 1/2...
STEP  208 / 251
done.

Making a distance matrix from msa..
 200 / 252
done.

Constructing a UPGMA tree (efffree=1) ...
 250 / 252
done.

Progressive alignment 2/2...
STEP  209 / 251
done.

disttbfast (nuc) Version 7.453
align4, model=DNA200 (2), 1.53 (4.59), -0.00 (-0.00), noshift, amax=0.0
0 thread(s)

Strategy:
FFT-NS-2 (Fast but rough)
Progressive method (guide trees were built 2 times.)

If unsure which option to use, try 'mafft --auto input > output'.
For more information, see 'mafft --help', 'mafft --man' and the mafft page.

The default gap scoring scheme has been changed in version 7.110 (2013 Oct).
It tends to insert more gaps into gap-rich regions than previous versions.
To disable this change, add the --leavegappyregion option.
```