

BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama
(School of Biological Sciences)

Spring 2026 Lecture 16

BIOS477/877 L16 - 1

1

Today's topics

- Multiple Sequence Alignment
 - ClustalΩ, MUSCLE, Mafft
 - Phylogeny-aware gap placement methods (PRANK, etc.)
 - and more ...
- Alignment Trimming/filtering

BIOS477/877 L16 - 2

2

How to solve progressive-alignment problems

- Incorporate more information to reduce **early errors**
 - **Structural alignment** (e.g., Expresso, PROMALS3D, TM-Coffee, PRALINE, MAFFT-DASH, MUSCLE-3D)
 - **Profile/profile-HMM alignment** (e.g., PRALINE, PSI-Coffee, PROMALS3D, ProbCons/CONTRAlign, ClustalΩ, MUSCLE5)
- Avoid the **greedy-algorithm** problem
 - **Iterative refinement** to search the global maxima
 - A good objective function is required (e.g., MUSCLE/MUSCLE5, MAFFT, ProbCons/CONTRAlign)
- **Global (or local) only** alignment problem
 - **Combine both methods** (e.g., T-Coffee)
- More accurate **insertion/deletion placement**
 - **Phylogeny aware gap-placement** (e.g., PRANK, ProPIP, Bali-Phy, SATé)

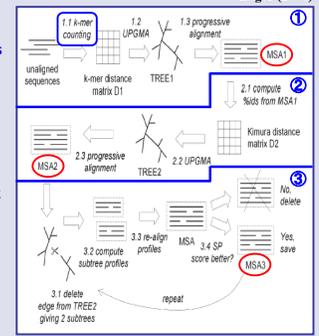
BIOS477/877 L16 - 3

3

MUSCLE v3

Available at [MUSCLE website](#) (in v3 legacy page) and [EBI Tools](#)

1. Draft progressive alignment:
 - Approximated distance based on **frequencies of shared k-mers (k-tuples) & UPGMA**
2. Improved progressive alignment:
 - **Kimura protein distance** (more accurate)
 - Alignment is done for subtrees where the branching patterns are changed between TREE1 and TREE2
3. Iterative refinement
 - **The tree is partitioned**
 - Profile is obtained from each subtree
 - **Profile alignment**
 - **Iteration using Sum-of-Pairs (SP) score**



4

MUSCLE v3 performance

BAliBASE Q scores (Sum-of-pairs: percentage of correctly aligned residue pairs)

Method	Equidistance Ref1	Family-orphan Ref2	< 25% identity Ref3	N/C-term extension Ref4	Internal insertion Ref5
MUSCLE	0.887	0.935	0.823	0.876	0.968
MUSCLE-p (w/o refinement)	0.871	0.928	0.813	0.857	0.974
T-Coffee	0.866	0.934	0.787	0.917	0.957
NWNSI (Mafft)	0.867	0.923	0.787	0.904	0.963
CLUSTALW	0.861	0.932	0.751	0.823	0.859
FFTNS1 (Mafft)	0.838	0.908	0.708	0.793	0.947

BAliBASE: Benchmark alignment database (includes many subsets representing various alignment problems)

BIOS477/877 L16 - 5

5

MAFFT

Available at [MAFFT website](#) or [EBI Tools](#)

1. First progressive alignment: **NS-1** (Kato et al. (2005, 2019))
 - 6-tuple distance from pairwise alignment & modified UPGMA for guide tree
 - Alignment strategy:

Alignment	FFT	NW
1st progressive	FFT-NS-1	NW-NS-1
2nd progressive	FFT-NS-2	NW-NS-2
Iterative refinement (i)	FFT-NS-i	NW-NS-i
2. Improved progressive alignment: **NS-2**
 - Better distance calculation from NS-1 alignment
3. Iterative refinement: **NS-i**
 - The tree-dependent restricted partitioning
 - Group-to-group (profile) alignment
 - Iteration based on the weighted SP score

Alignment	Consistency score
Global	G-INS-i
Local (SW)	L-INS-i
	E-INS-i
4. **COFFEE-like consistency score** from pairwise alignment information is used in **G-INS-i** (global), **L-INS-i** (local), and **E-INS-i** (local & more difficult alignment)

BIOS477/877 L16 - 6

6

MAFFT: add/merge alignment

There are also --seed and --addprofile options. Read the "tips" section on the MAFFT website.

BIOS477/877 L16 - 13

13

Refining alignment with MAFFT website

Refine the dataset based on sequence length, % identity, phylogenetic clustering, etc.

BIOS477/877 L16 - 14

14

Refining alignment with MAFFT website

BIOS477/877 L16 - 15

15

How to solve progressive-alignment problems

- Incorporate more information to reduce **early errors**
 - **Structural alignment** (e.g., Expresso, PROMALS3D, TM-Coffee, PRALINE, MAFFT-DASH, MUSCLE-3D)
 - **Profile/profile-HMM alignment** (e.g., PRALINE, PSI-Coffee, PROMALS3D, ProbCons/CONTRAlign, ClustalΩ, MUSCLE5)
- Avoid the **greedy-algorithm** problem
 - **Iterative refinement** to search the global maxima
 - A good objective function is required (e.g., MUSCLE/MUSCLE5, MAFFT, ProbCons/CONTRAlign)
- **Global (or local) only** alignment problem
 - **Combine both methods** (e.g., T-Coffee)
- More accurate **insertion/deletion placement**
 - **Phylogeny aware gap-placement** (e.g., PRANK, ProPIP, Bali-Phy, SATé)

BIOS477/877 L16 - 16

16

PRALINE

Simossis, Kleinjung & Heringa (2005)

PROGRESSIVE MULTIPLE ALIGNMENT

BIOS477/877 L16 - 17

17

PRALINE performance

Simossis, Kleinjung & Heringa (2005)

PSI-BLAST enhanced PRALINE performs better especially for alignment of highly divergent sequences

- From 624 HOMSTRAD pairwise alignments
- Q score = SP (sum-of-pairs) score: % correctly aligned residue pairs
- ΔQ: Q score difference from PRALINE without PSI-BLAST

BIOS477/877 L16 - 18

18

Clustal Ω

Available at [EBI Tools](#) Sievers *et al.* (2011, 2018)

- **Progressive alignment** following the guide tree
 - Features a fast method for making the “initial guide tree” (**mBed**)
 - Calculates distances to only a small number ($n = \log_2 N$) of seed sequences: if $N = 1000$, $n = \log_2 1000 \approx 10$
 - Sequences are represented as an n element vector, and clustered using Euclidean distances (UPGMA for subtrees)
 - Scalable for very large datasets
 - **Alignment is done using HAlign** (a profile hidden Markov model alignment)
 - Probabilistic, highly accurate alignment
 - **Simple iterative refinement**
 - Alignment is converted to **profile hidden Markov model (HMM)**
 - realigns input sequences against the profile HMM

BIOS477/877 L16 - 19

19

Clustal Ω

Sievers *et al.* (2011)

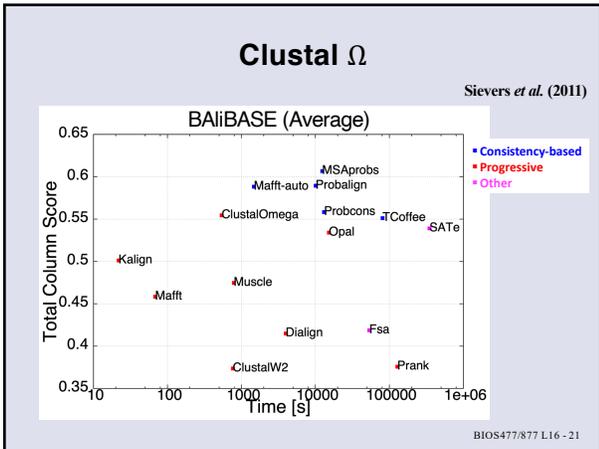
Table 1 BAiBASE results ➤ Top-rated methods use HMM or consistency function or both

Aligner	Av score (218 families)	BB11 (38 families)	BB12 (44 families)	BB2 (41 families)	BB3 (30 families)	BB4 (49 families)	BB5 (16 families)	Tot time (s)	Consistency
MSAprobs	0.607	0.441	0.865	0.464	0.607	0.622	0.608	12 382.00	Yes
Probalign	0.589	0.453	0.862	0.439	0.566	0.603	0.549	10 095.20	Yes
MAFFT (auto)	0.588	0.439	0.831	0.450	0.581	0.605	0.591	1475.40	Mostly (203/218)
Procons	0.558	0.417	0.855	0.406	0.544	0.532	0.573	13 086.30	Yes
Clustal Ω	0.554	0.358	0.789	0.450	0.375	0.578	0.533	539.91	No
T-Coffee	0.551	0.410	0.848	0.402	0.491	0.545	0.587	81 041.50	Yes
Kalign	0.501	0.365	0.790	0.360	0.476	0.504	0.435	21.88	No
MUSCLE	0.475	0.318	0.804	0.350	0.409	0.450	0.460	789.57	No
MAFFT (default)	0.458	0.258	0.749	0.316	0.425	0.480	0.496	68.24	No
FSA	0.419	0.270	0.818	0.187	0.259	0.474	0.398	53 648.10	No
Dialign	0.415	0.265	0.696	0.292	0.312	0.441	0.425	3977.44	No
PRANK	0.376	0.223	0.680	0.257	0.321	0.360	0.356	128 355.00	No
ClustalW	0.374	0.227	0.712	0.220	0.272	0.396	0.308	766.47	No

The figures are total column scores produced using ball score on core columns only. The average score over all families is given in the second column. The results for BAiBASE subgroupings are in columns 3–8. The total run time for all 218 families is given in the second last column. The last column indicates whether the method is consistency based.

BIOS477/877 L16 - 20

20



21

Muscle5

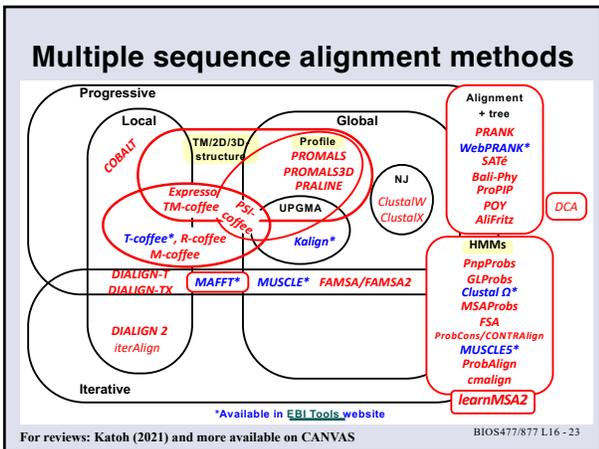
Available at [MUSCLE website](#) or [EBI Tools](#) Edgar (2022)

- **Improved re-implementation of ProbCons** (slow but highly accurate MSA method based on hidden Markov model)
 - More accurate than Clustal Ω, MAFFT, or ProbCons
 - Ensemble MSAs can be generated for reference-free estimate of MSA accuracy
- Parallelization for rapid alignment
- Can be used for large datasets
- For both protein and DNA alignments

Method	Balibase (protein)			Brallbase (RNA)		
	Time (mins.)	Max mem (Mb)	TC	Time (mins.)	Max mem (Mb)	TC
Clustal-Omega	5	0.98	0.56	5	0.5	0.73
MAFFT	35	0.5	0.61	7	0.4	0.73
MUSCLE v5	27	2.1	0.63	6	0.8	0.78
MUSCLE v3	8	0.08	0.53	2	0.2	0.75
PROBCONS	182	0.54	0.62	n/a	n/a	n/a

BIOS477/877 L16 - 22

22



23

PRANK, WebPRANK

Mind the gaps: Progress in progressive alignment

D. G. Higgins*, G. Blackshields, and L. M. Wallace
Conway Institute, University College Dublin, Belfield, Dublin 4, Ireland (2005 PNAS commentary)

"CLUSTALW attempts to compensate by using an elaborate scoring scheme to encourage gaps to end up on top of each other. ... results in alignments that are very "block-like"..."

"... there may be a price for this prettiness and detachment from phylogenetic reality. CLUSTALW (and other programs) may be guilty of "overalignment", that is where sequences that should not go together are forced into neat-looking blocks. These overligned regions may be neat looking but misleading."

"There is an understandable tendency for users of multiple alignment software to want their residues neatly aligned in blocks and columns. This is fine when such blocks are biologically accurate as will happen in parts of protein alignments. In cases where insertions or deletions have happened in a less organized manner, as will happen in many noncoding DNA sequences and in less organized parts of protein sequences, such block-like alignments may be biologically meaningless. Perhaps we need to reeducate our eyes to see beauty in what actually happened rather than what looks nice on paper."

BIOS477/877 L16 - 24

24

PRANK, WebPRANK

Löytynoja & Goldman (2008)

Insertions are more penalized than deletions in progressive sequence alignment.

	Insertion	Deletion
Evolution		
Progressive alignment	<p>Cost: -6 (↘ ↘ ↘)</p>	<p>Cost: -2 (↖)</p>

BIOS477/877 L16 - 25

25

PRANK, WebPRANK

Available at EBI website Löytynoja & Goldman (2005, 2008, 2010)

- PRANK: **Probabilistic Alignment Kit**
 - A probabilistic multiple alignment program for DNA, codon, and amino-acid sequences.
 - Treats insertions correctly.
 - Avoids over-estimation of the number of deletion events.
 - **Not meant for the alignment of very diverged protein sequences.**

BIOS477/877 L16 - 26

26

PRANK, WebPRANK

Löytynoja & Goldman (2008)

Different sequence alignment approaches can give contradicting pictures of evolutionary mechanisms behind functional sequence changes.

	ClustalW	PRANK
A		
B		

MSAs generated by traditional methods show excess substitutions
→ Can be erroneously thought to be under positive selection

BIOS477/877 L16 - 27

27

ProPIP: Progressive MSA with Poisson Indel Process

Maiolo *et al.* (2018, 2021)

Method	Alignment	Phylogenetic Tree
PRANK		
ProPIP		
MAFFT		

Rigorous mathematical indel model is incorporated
 ▪ Gaps can be inferred in a phylogenetically consistent way
 ▪ Over-alignment can be avoided

BIOS477/877 L16 - 28

28

Bali-Phy: Statistical coestimation method

Nute *et al.* (2019), Redelings (2021), Gupta *et al.* (2021)

[1,192 datasets]

SP score (recall) = How much of correct alignment is recovered
Modeller score (precision) = How much of the reconstructed alignment is correct

Bali-Phy website.

▪ Co-estimates MSA and phylogeny iteratively
 ▪ Bali-Phy v3/v4 is much faster for large datasets

BIOS477/877 L16 - 29

29

Large-scale MSA

(large-scale MSA methods)

- MAFFT DPPartTree
- ClustalΩ
- FAMS/FAMS2
- Kalign3 ...

- PASTA
- UPP/UPP2
- MAGUS
- MAFFT/Sparsecore
- Regressive/T-Coffee
- MUSCLE5 ...

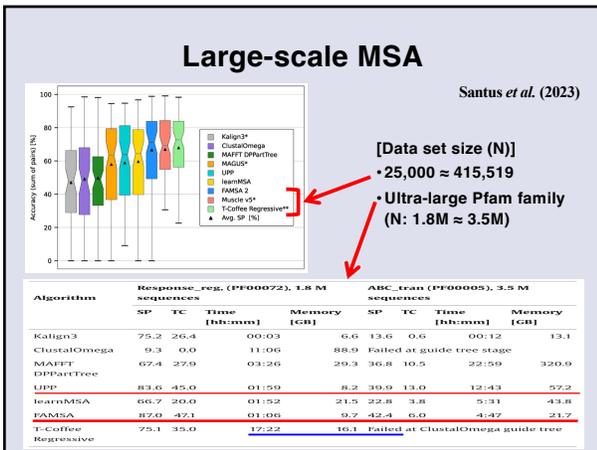
- MMseq2
- learnMSA/learnMSA2

Using pHMM, deep learning protein language modeling, etc.

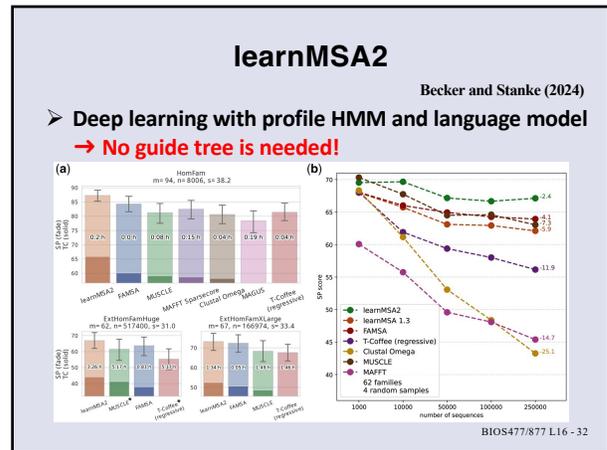
Santus *et al.* (2023)

BIOS477/877 L16 - 30

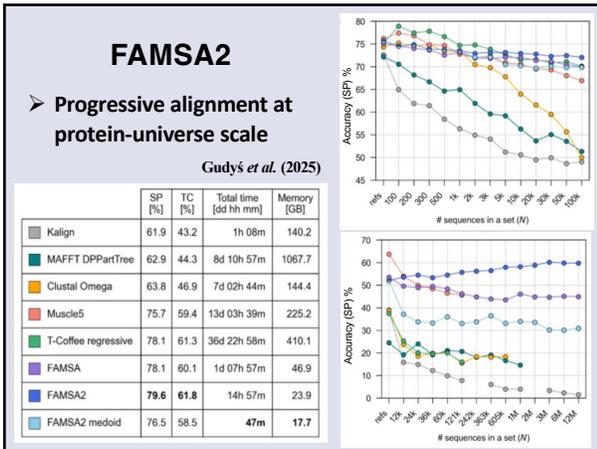
30



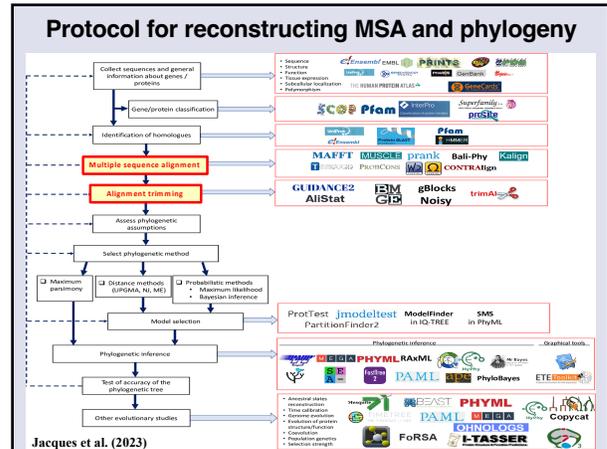
31



32



33



34

Alignment trimming

TABLE 1. Overview of filtering methods considered in this study

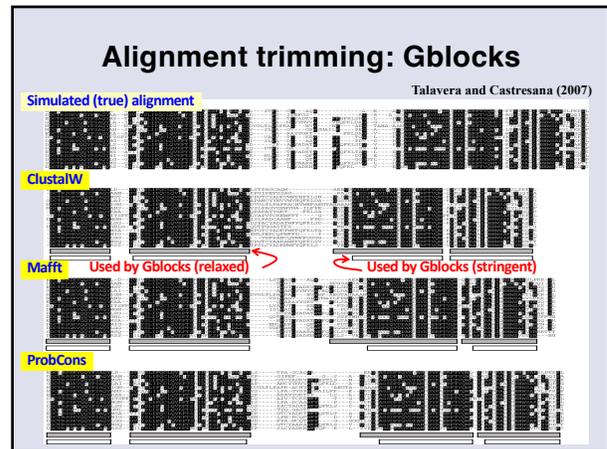
Filtering methods	Type of "undesirable" sites filtered out by the method	Accounts for tree structure?	Uses a substitution matrix or model of evolution?	Adapts parameters for particular data sets?	References
Gblocks	Gap-rich and variable sites	No	No	No	Talavera and Castresana (2007)
TrimAl	Gap-rich and variable sites	No	Yes	Yes	Capella-García <i>et al.</i> (2009)
Noisy	Homoplastic sites	In part	No	No	Dress <i>et al.</i> (2008)
Aliscore	Random-like sites	No	No	No	Kück <i>et al.</i> (2010)
BMGE	High entropy sites	No	Indirectly	No	Crisuolo and Gibaldini (2010)
Zorro	Sites with low posterior probability	Yes	Yes	No	Wu <i>et al.</i> (2012)
Guidance	Sites sensitive to the alignment guide tree	Yes	Indirectly	No	Penn <i>et al.</i> (2010)

- **Gblocks**: selection of conserved blocks from multiple alignments (included in Phylogeny.fr)
- **trimAl**: a tool for automated alignment trimming in large-scale phylogenetic analysis (included in Phylogen2)

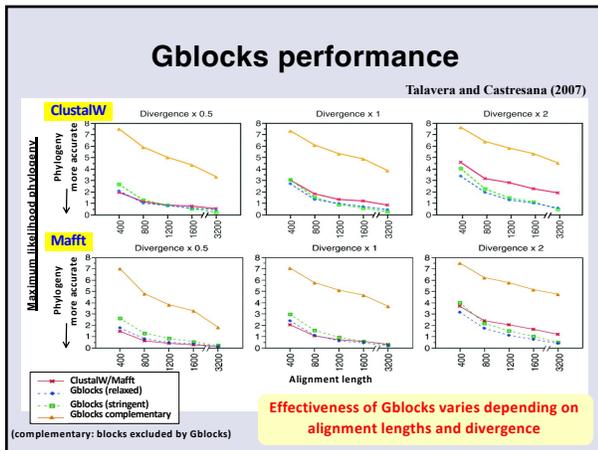
Tan *et al.* (2015); See also TCS paper by Chang *et al.* (2014)

BIOS477/877 L16 - 35

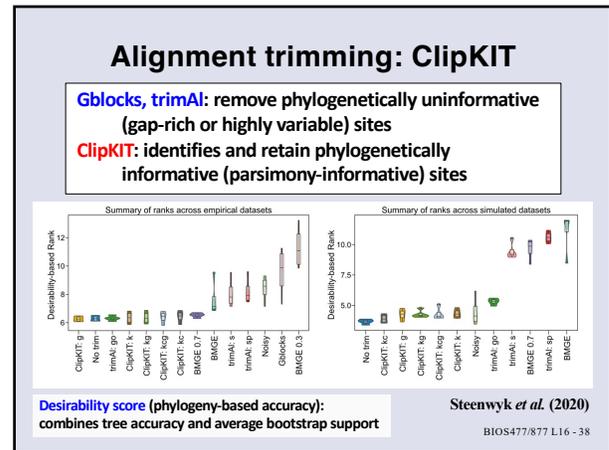
35



36



37



38