## Slide 1

# BIOS 477/877

*Bioinformatics and Molecular Evolution*

## Lecture 15

**1**

## Slide 2

### TODAY'S TOPICS

➢ **Multiple Alignment**

- **T-Coffee**
- **MUSCLE and Mafft**
- **PRALINE, Clustal$\Omega$, etc.**

**2**

## Slide 3

### Progressive multiple alignment: Clustal W



**Pairwise alignment**

S1, S2 → D12
S1, S3 → D13
S1, S4 → D14

**Guide tree**

**Profile alignment**

| | | |
|---|---|---|
| – | | |
| D12 | – | |
| D13 | D23 | – |

**3**

## Slide 4

### Progressive multiple alignment: Clustal W

➢ **How sequence weighting works: Example 2**



```
        0.5  S1(1.1)     123
      1                  VAG
        0.5  S2(1.1)     VGA
  0.3                    IVG
          6  S3(6.1)      VG
          7  S4(7.0) VG
```

**[Simple average]**
(a) Alignment score = 1 + gap penalty
(b) Alignment score = 2 + gap penalty
(c) Alignment score = 1 + gap penalty

**[Weighted average]**
(a) Alignment score = -15.64 + gap penalty
(b) Alignment score = 27.06 + gap penalty
(c) Alignment score = 54.36 + gap penalty

| | (a) 123 | (b) 123 | (c) 123 |
|---|---|---|---|
| S1 | VAG | VAG | VAG |
| S2 | VGA | VGA | VGA |
| S3 | IVG | IVG | IVG |
| S4 | VG– | V–G | –VG |

*[Aligning S4 to the first 3 sequences (S1, S2, and S3) previously aligned]*

**4**

## Slide 5

### Progressive multiple alignment: Clustal W

➢ **How to choose scoring matrix:**

- **Choose only the scoring matrix series (BLOSUM, PAM, *etc.*)**

- ➔ **Specific matrix is determined based on distance between sequences**

| | |
|---|---|
| 80 - 100% identity | → Blosum80 |
| 60 - 80% identity | → Blosum62 |
| 30 - 60% identity | → Blosum45 |
| 0 - 30% identity | → Blosum30 |

**Thompson *et al.* (1994)**

**5**

## Slide 6

### Progressive multiple alignment: Clustal W

➢ **How gap penalties are determined:**

- ➔ **Initial gap penalties: GOP (gap opening) and GEP (gap extension) ➔ set by the user**
- ➔ **Weight (scoring) matrix dependent gap penalties**
- ➔ **Similarity level dependent gap penalties**
- ➔ **Sequence length dependent gap penalties**
- ➔ **Position specific gap penalties**
  - **if gaps already exist**
  - **residue specific (*e.g.*, hydrophilic stretches)**

**Thompson *et al.* (1994)**

**6**

## Progressive multiple alignment: Clustal W



**Pairwise alignment**

S1/S2 → D12
S1/S3 → D13
S1/S4 → D14

| | | |
|---|---|---|
| – | | |
| D12 | – | |
| D13 | D23 | – |

0.02 S1
0.15
0.09
0.08 S2
0.38
S3
0.46 S4

**Guide tree**

**Profile alignment**

- Progressive alignment
  - ➜ Greedy (finds local optima, but no guarantee for global optima)
  - ➜ Errors (incorrect gap positions) in the early alignments cannot be rectified later
- Global alignment only (local similarity may be missed)

**7**

---

## Clustal Web servers

**http://www.clustal.org/ (Clustal original website)**



**Clustal: Multiple Sequence Alignment**
Multiple alignment of nucleic acid and protein sequences

**Clustal Omega**
- Latest version of Clustal - fast and scalable (can align hundreds of thousands of sequences in hours), greater accuracy due to new HMM alignment engine
- Command line/web server only (GUI public beta available soon)

**ClustalW/ClustalX**
- "Classic Clustal"
- GUI (ClustalX), command line (ClustalW), web server versions available

**Use Clustal Ω instead of Clustal W**

**[Classic version]**
https://galaxy.pasteur.fr/ (part of Galaxy@Pasteur)
https://www.genome.jp/tools-bin/clustalw

**8**

---

## To solve the progressive-alignment problems

➢ **Incorporate more information to reduce early errors**
  - • Structural alignment (*e.g.*, Expresso, PROMALS3D, TM-Coffee, PRALINE, MAFFT-DASH)
  - • 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)

**9**

---

## Multiple alignment: T-Coffee

**T-Coffee Home page**
**https://tcoffee.crg.eu/**
**https://tcoffee.org**

**https://tcoffee.crg.eu/apps/tcoffee/index.html**
**(mirror site)**

**https://www.ebi.ac.uk/jdispatcher/msa/tcoffee**
**(only T-Coffee, without other associated programs)**

**Notredame, Higgins & Heringa (2000); Taly *et al*. (2011)**

**10**

---

## Multiple alignment: T-Coffee

- **T-Coffee: Tree-based Consistency Objective Function for alignment Evaluation**

- **Based on the progressive alignment algorithm**
  - ➜ Uses a guide tree, fast

- **Tries to avoid the greedy nature of the progressive algorithm**
  - ➜ Using alignment libraries derived from a mixture of alignment programs (global, local, *etc.*)

**11**

---

## Multiple alignment: T-Coffee
(Tree-based Consistency Objective Function for alignment Evaluation)

➢ **Primary libraries of alignments**



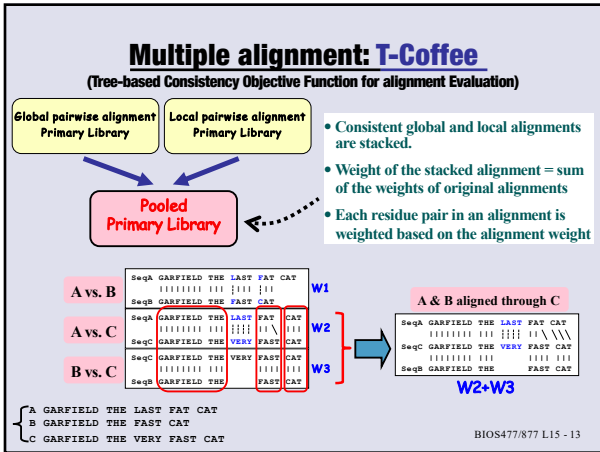**ClustalW Primary Library**
(Global Pairwise Alignment)

**Lalign Primary Library**
(Local Pairwise Alignment)

- ➜ Global and local, or any combination of pairwise alignment methods
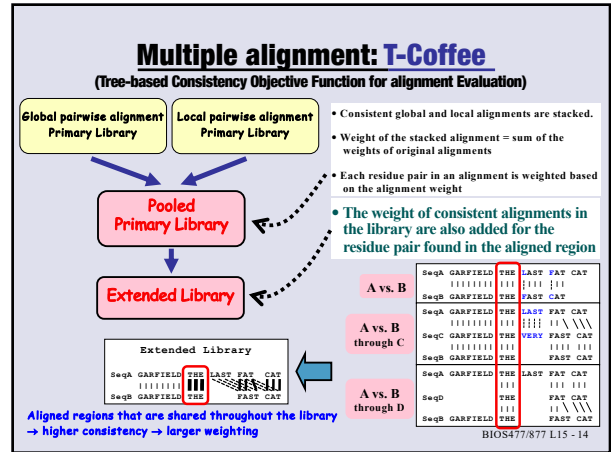- ➜ Each pairwise alignment is given a weight based on % identity ignoring gap sites   ATTCGG / ATAGCG ⟹ 3/6 = 50%

**12**

**2**

## Slide 13

**Multiple alignment: T-Coffee**
(Tree-based Consistency Objective Function for alignment Evaluation)

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

- **Consistent global and local alignments are stacked.**
- **Weight of the stacked alignment = sum of the weights of original alignments**
- **Each residue pair in an alignment is weighted based on the alignment weight**

Pooled Primary Library

```
A vs. B
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||| ||||| |||
SeqB GARFIELD THE FAST CAT            W1

A vs. C
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||| || \ |||           W2
SeqC GARFIELD THE VERY FAST CAT

B vs. C
SeqC GARFIELD THE VERY FAST CAT
     |||||||| ||| |||| |||           W3
SeqB GARFIELD THE FAST CAT
```

A & B aligned through C
```
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||| |||| |||
SeqC GARFIELD THE VERY FAST CAT
     |||||||| ||| ||||
SeqB GARFIELD THE FAST CAT
```
**W2+W3**

A GARFIELD THE LAST FAT CAT
B GARFIELD THE FAST CAT
C GARFIELD THE VERY FAST CAT

BIOS477/877 L15 - 13

**13**

## Slide 14

**Multiple alignment: T-Coffee**
(Tree-based Consistency Objective Function for alignment Evaluation)

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

- **Consistent global and local alignments are stacked.**
- **Weight of the stacked alignment = sum of the weights of original alignments**
- **Each residue pair in an alignment is weighted based on the alignment weight**
- **The weight of consistent alignments in the library are also added for the residue pair found in the aligned region**

Pooled Primary Library

Extended Library

```
Extended Library
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||||
SeqB GARFIELD THE         FAST CAT
```

A vs. B
```
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||| |||
SeqB GARFIELD THE FAST CAT
```

A vs. B through C
```
SeqA GARFIELD THE LAST FAT CAT
     |||||||| ||| |||||  || \ \\\
SeqC GARFIELD THE VERY FAST CAT
     |||||||| |||          |||| |||
SeqB GARFIELD THE         FAST CAT
```

A vs. B through D
```
SeqA GARFIELD THE LAST FAT CAT
     |||||||| |||            ||| |||
SeqD GARFIELD THE              || \ \\\
SeqB GARFIELD THE         FAST CAT
```

**Aligned regions that are shared throughout the library → higher consistency → larger weighting**

BIOS477/877 L15 - 14

**14**

## Slide 15

**Multiple alignment: T-Coffee**
(Tree-based Consistency Objective Function for alignment Evaluation)

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

- **Consistent global and local alignments are stacked.**
- **Weight of the stacked alignment = sum of the weights of original alignments**
- **Each residue pair in an alignment is weighted based on the alignment weight**
- **The weight of consistent alignments in the library are also added for the residue pair found in the aligned region**

Pooled Primary Library

Extended Library

```
Extended Library
SeqA GARFIELD THE LAST FAT CAT
     |||||||| |||
SeqB GARFIELD THE         FAST CAT
```
**For each pair of sequences**

**Each residue pair is weighted based on the consistency found in the primary library**

Generate a sequence-pair specific scoring matrix

➜ **Pre-made scoring matrices (BLOSUM etc.) are not used**
➜ **No gap penalty is used**

BIOS477/877 L15 - 15

**15**

## Slide 16

**Multiple alignment: T-Coffee**
(Tree-based Consistency Objective Function for alignment Evaluation)

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

- **Consistent global and local alignments are stacked.**
- **Weight of the stacked alignment = sum of the weights of original alignments**
- **Each residue pair in an alignment is weighted based on the alignment weight**
- **The weight of consistent alignments in the library are also added for the residue pair found in the aligned region**

Pooled Primary Library

Extended Library

```
Extended Library
SeqA GARFIELD THE LAST FAT CAT
     |||||||| |||
SeqB GARFIELD THE         FAST CAT
```

**Each residue pair is weighted based on the consistency found in the primary library**

**Progressive Alignment with NJ guide tree using the generated scoring matrix for each pair of sequences**

**T-Coffee does not use pre-made scoring matrices (BLOSUM etc.) nor gap penalties**

BIOS477/877 L15 - 16

**16**

## Slide 17

**Multiple alignment: T-Coffee**

https://tcoffee.crg.eu/

**For protein sequence alignment:**
- **Structural alignments (Expresso)**
- **Combine popular aligners (M-Coffee)**
- **Transmembrane proteins (PSI/TM-Coffee)**
- **Homology extension (PSI-Coffee)**

PROTEINS | RNA

T-COFFEE SIMPLE MSA

**For RNA sequence alignment:**
- **Secondary structure (R-Coffee)**
- **Tertiary structure (SARA-Coffee)**
- **Combine popular aligners (M-Coffee)**

**Simple T-Coffee**

DNA

**For DNA sequence alignment:**
- **Combine popular aligners (M-Coffee)**
- **Homologous promoter regions (Pro-coffee)**

BIOS477/877 L15 - 17

**17**

## Slide 18

**Multiple alignment: T-Coffee**

https://tcoffee.crg.eu/apps/tcoffee/all.html

**T COFFEE**
Home | History | Tutorial | References | Contacts | Projects | Download

**T-Coffee**
*A collection of tools for Computing, Evaluating and Manipulating Multiple Alignments of DNA, RNA, Protein Sequences and Structures*

**Alignment**

T-Coffee — Aligns DNA, RNA or Proteins using the default T-coffee  >> Cite
M-Coffee — Aligns DNA, RNA or Proteins by combining the output of popular aligners  >> Cite
R-Coffee — Aligns RNA sequences using predicted secondary structures  >> Cite
SARA-Coffee — Aligns RNA sequences using tertiary structure  >> Cite
Expresso — Aligns protein sequences using structural information  >> Cite
PSI-Coffee — Aligns distantly related proteins using homology extension (slow and accurate)  >> Cite
PSI/TM-Coffee — Align Proteins using Homology Extension against Reduced Databases  >> Cite
Pro-Coffee — Aligns homologous promoter regions  >> Cite
Accurate — Automatically combine the most accurate modes for DNA, RNA and Proteins (experimental)  >> Cite
Combine — Combines two (or more) multiple sequence alignments into a single one  >> Cite

**Evaluation**

TCS — Evaluates your Alignment and outputs a Colored version indicating the local reliability  >> Cite
iRMSD-APDB — Evaluates Multiple Sequence Alignment using structural information with APDB and iRMSD  >> Cite
T-RMSD — Allows fine-grained structural clustering of a given group of related protein domains  >> Cite
Strike — Evaluation of protein MSAs using a single 3D structure  >> Cite

**Other**

Advanced — Run your alignment using full featured T-Coffee options  >> Cite

BIOS477/877 L15 - 18

**18**

**19**



**20**



**21**



**22**



**23**



**24**

## Slide 25

# Evaluating multiple alignments using TCS

- TCS (transitive consistency score) Chang *et al.* (2014)

**T-Coffee**

```
SCORE=395
BAD AVG GOOD
1aboA : 48
1ycsB : 47
1pht  : 45
1vie  : 32
1ihvA : 32
cons  : 39

1aboA  NL-FVAL----YDFVASGDNTLSITKGEKLR-------VLGYNH------NGEWCE--AQTKNGQGWVPSNYITPV-N-----
1ycsB  RGVIYAL---WDYEPQNDDELPMKEGDCMT-------IIHREDED-----EIEWWW--ARLNDKEGYVPRNLLGLY------P
1pht   GYQYRAL---YDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNGYNETTGERGDFPGTVYEYIGRKKISP
1vie   DR----------VRKK--SGAAWQGQIVGWYCTNLTPEGYAVESEAHPGSV-----------------QIYPVAALERI----N
1ihvA  NF-RVYYRDSRDPVWKGPAK-LLWKGEGAV-------VIQDNS-------DIK------------VVPRRKAKIIRD------
cons   
```

**ClustalW2**

```
SCORE=305
BAD AVG GOOD
1aboA : 42
1ycsB : 41
1pht  : 39
1vie  : 25
1ihvA : 20
cons  : 30

1aboA  -NLFV-ALYDFVASGDNTLSITKGEKLRV-------LGYNHNG-------EWCEAQ--TK42NGQGWVPSNYITPVN------57
1ycsB  RGVIY-ALWDYEPQNDDELPMKEGDCMYI-------IHREDEDEI-----EWWWAR--LN45DKEGYVPRNLLGLYP------60
1pht   -GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNGYNETT59GERGDFPGTYVEYIGRKKISP80
1vie   --------DRVRKKSG--AAWQGGIVGW---------YCTNLTP----EGYAVESEAHP36GSVOIYPVAALERIN------51
1ihvA  --------NFRVYYRDSRD--PVWKGPAKLL-------WKGEG-------AVVIQD---N33SDIKVVPRRKAKIIRD------49
cons   
```

---

## Slide 26

# Evaluating multiple alignments using TCS



Ref

**T-Coffee**

**M-Coffee** (kalign+mafft +muscle)

**PSI-Coffee**

**Expresso**

Taken from Taly *et al.* (2011)

---

## Slide 27

# To solve the progressive-alignment problems

- ➤ **Incorporate more information to reduce early errors**
  - **Structural alignment** (*e.g.*, Expresso, PROMALS3D, TM-Coffee, PRALINE, MAFFT-DASH)
  - **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 L15 - 27

---

## Slide 28

# MUSCLE (Edgar 2004)

http://www.drive5.com/muscle/  https://www.ebi.ac.uk/jdispatcher/msa/muscle

1. **Draft progressive alignment:**
   - ➜ **K-mer distance & UPGMA** (Word or k-tuple)
2. **Improved progressive alignment:**
   - ➜ **Kimura protein distance**
   - ➜ **Tree comparison (branching orders are changed or not)**
   - ➜ **Iteration until the tree stays the same**
3. **Iterative refinement**
   - ➜ **The tree is partitioned**
   - ➜ **Profiles are obtained from each subtree**
   - ➜ **Profile alignment**
   - ➜ **Iteration based on SP score**



Count the number of shared k-mers

SP score: sum-of-pairs score

---

## Slide 29

# MUSCLE (Edgar 2004)

**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 L15 - 29

---

## Slide 30

# MAFFT (Katoh *et al.* 2005, 2019)

https://mafft.cbrc.jp/alignment/software/index.html
https://www.ebi.ac.uk/jdispatcher/msa/mafft

1. **First progressive alignment: FFT-NS-1**
   - ➜ **6-tuples distance & UPGMA**
   - ➜ **Fast Fourier transform (FFT) is used to detect highly similar segments**
   - ➜ **Segment to segment dynamic programming**
2. **Improved progressive alignment: FFT-NS-2**
   - ➜ **A better distance matrix from FFT-NS-1 tree**
3. **Iterative refinement: FFT-NS-i**
   - ➜ **The tree-dependent restricted partitioning**
   - ➜ **Group-to-group alignment**
   - ➜ **Iteration based on the weighted SP score**
4. **Instead of FFT, full dynamic programming can be used: NW-NS-[12i] [after version 5.0]**
   - ➜ **COFFEE-like consistency score with pairwise alignment information is included for Global (G-INS-i) or Local (L-INS-i, E-INS-i*) alignments**
   - ***E-INS-i is for more difficult alignment**

| Alignment | FFT | NW | Consistency score |
|---|---|---|---|
| 1st progressive | NS-1 | NS-1 | |
| 2nd progressive | NS-2 | NS-2 | |
| Iterative refinement (i) | NS-i | NS-i | |
| Global | | | G-INS-i |
| Local (SW) | | | L-INS-i E-INS-i |

BIOS477/877 L15 - 30

# Slide 31

## MAFFT

| Method | Ref11 | Ref12 | Ref2 | Ref3 | Ref4 | Ref5 | Overall average | CPU time (s) |
|---|---|---|---|---|---|---|---|---|
| **Consistency based methods** | | | | | | | | |
| MAFFT 5.662 L-INS-i | **67.11 / 44.61** | 93.62 / 83.73 | 92.57 / 45.17 | 85.58 / 56.83 | **91.91 / 59.47** | 90.15 / 58.44 | **87.05 / 58.64** | 5,500 |
| MAFFT 5.662 E-INS-i | 66.13 / 44.53 | 93.54 / 83.18 | **92.64 / 44.32** | **85.08 / 58.53** | 91.42 / 59.02 | 89.93 / 59.13 | 86.91 / 58.55 | 6,000 |
| ProbCons 1.10 (default) | 66.99 / 41.66 | 94.12 / **85.52** | 91.67 / 40.54 | 84.60 / 54.30 | 90.52 / 54.37 | 89.28 / 56.50 | 86.46 / 55.99 | 43,000 |
| ProbCons 1.10 (trained) | 66.73 / 41.47 | **94.13** / 85.38 | 91.85 / 42.00 | 84.47 / 54.03 | 89.79 / 51.94 | 89.34 / 57.69 | 86.27 / 55.71 | (44,000) |
| MAFFT 5.662 G-INS-i | 60.46 / 34.53 | 92.42 / 81.32 | 90.34 / 38.71 | 85.27 / 52.73 | 88.37 / 52.51 | 87.87 / 52.75 | 84.23 / 52.64 | 6,900 |
| TCofee 2.46 | 61.48 / 33.63 | 93.04 / 82.36 | 91.71 / 39.68 | 81.61 / 48.87 | 89.22 / 52.90 | 89.03 / 57.13 | 84.56 / 52.76 | (210,000) |
| **Iterative refinement methods** | | | | | | | | |
| MAFFT 5.662 FFT-NS-i | 58.87 / 33.47 | 91.64 / 80.11 | 89.54 / 40.37 | 83.27 / 49.97 | 87.11 / 47.37 | 86.27 / 52.44 | 82.95* / 50.97** | 2,800 |
| Muscle 3.52 (most accurate option) | 56.62 / 30.87 | 90.96 / 79.59 | 88.90 / 35.17 | 81.07 / 37.87 | 85.90 / 45.06 | 85.17 / 46.19 | 81.67** / 46.79* | 3,400 |
| PRRN 3.11 | 58.21 / 34.74 | 92.16 / 79.20 | 90.46 / 41.66 | 82.68 / 47.63 | 85.83 / 47.98 | 83.83 / 47.56 | 82.61* / 50.73* | 250,000 |
| MAFFT 3.89 † FFT-NS-i | 54.56 / 30.26 | 90.78 / 78.61 | 90.12 / 37.46 | 82.65 / 49.33 | 87.83 / 50.76 | 85.65 / 49.31 | 82.18* / 50.27* | 3,600 |
| ClustalW 2.0 (Iteration=tree; Sep, 2007) | 49.94 / 25.08 | 86.91 / 75.32 | 85.80 / 21.61 | 72.78 / 30.43 | 81.20 / 40.84 | 76.49 / 35.06 | 76.67** / 39.58** | (58,000) |
| **Progressive methods** | | | | | | | | |
| Kalign 1.0 | 54.51 / 27.79 | 91.17 / 78.59 | 87.79 / 29.56 | 79.69 / 35.47 | 83.02 / 42.57 | 84.59 / 44.75 | 80.25** / 44.00** | 480 |
| MAFFT 5.662 FFT-NS-2 | 51.80 / 26.92 | 88.79 / 71.55 | 88.61 / 36.76 | 80.78 / 40.17 | 84.57 / 40.06 | 83.59 / 46.81 | 79.88** / 44.01** | 250 |
| MAFFT 5.662 FFT-NS-1 | 50.15 / 22.76 | 88.16 / 72.32 | 88.03 / 32.98 | 79.47 / 34.37 | 82.96 / 41.92 | 81.18 / 42.06 | 78.63** / 42.00** | 140 |
| Muscle 3.52 (fastest option) | 53.36 / 26.97 | 88.79 / 72.32 | 86.39 / 29.37 | 77.74 / 32.93 | 79.38 / 34.47 | 76.59 / 35.56 | 77.63** / 39.71** | 160 |
| ClustalW 1.83 | 50.06 / 22.74 | 86.43 / 71.14 | 85.20 / 21.98 | 72.50 / 27.23 | 78.82 / 39.55 | 74.244 / 30.75 | 75.34** / 37.35** | 2,000 |

The SP and TC scores are shown for each method. † Obsolete version of MAFFT. For iterative options of MAFFT and Muscle, the maximum numbers of iteration were set at 1,000. The significance of differenc[...] most accurate method is indicated by * (p<0.05) or ** (p<0.01) (Wilcoxon test) only for overall average.

**Scores: SP (sum of pairs)/TC (total column)**
**tested using BAliBASE benchmarking alignment datasets**

BIOS477/877 L15 - 31

**31**

---

# Slide 32

## MAFFT website

BIOS477/877 L15 - 32

**32**

---

# Slide 33

## Mafft-homologs

**Accuracy of an alignment of distantly related sequences can be improved when aligned with their close homologs**

**MAFFT options**



**PSI-BLAST is used to collect similar sequences**

**Similar approach is used in PSI-Coffee, PRALINE, etc.**

BIOS477/877 L15 - 33

**33**

---

# Slide 34

## MAFFT-DASH (Rozewicki *et al.* 2019)

### Integrated protein sequence and structural alignment



**Conserved Aspartic Acid**

**DASH (Database of Aligned Structural Homologs) is used to incorporate the structural information to improve the alignment**

BIOS477/877 L15 - 34

**34**

---

# Slide 35

## MAFFT-DASH (Rozewicki *et al.* 2019)

### Integrated protein sequence and structural alignment

Table 1. Benchmarks using reference MSAs

| Methods \ Data | HMFM | MBSF | MBTL | OXFM | BB11 | BB12 | BB20 | BB30 | BB40 | BB50 | SY |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **[SP scores]** | | | | | | SP | | | | | |
| MAFFT | 0.916** | 0.571** | 0.203** | 0.894** | 0.649** | 0.937** | 0.927** | 0.862 | 0.917 | 0.899* | 0.751** |
| Promals | 0.947** | 0.726** | 0.475** | 0.947** | 0.791 | 0.936 | 0.933* | 0.883 | 0.898 | 0.903 | 0.848** |
| T-Coffee | 0.922** | 0.585** | 0.224** | 0.909** | 0.657** | 0.945 | 0.916** | 0.837** | 0.897 | 0.895* | 0.778** |
| Expresso | 0.950** | 0.708** | 0.330** | 0.954** | 0.734** | 0.903** | 0.878** | 0.827** | 0.867** | 0.874** | 0.805** |
| MAFFT-DASH | **0.971** | 0.770** | 0.436** | **0.974** | 0.764* | **0.943** | **0.937** | **0.880** | **0.909** | 0.918 | 0.838* |
| MAFFT-DASH Homologs | **0.976** | 0.787 | 0.530* | **0.975** | 0.793 | 0.946 | 0.938 | 0.885 | 0.889 | 0.919 | 0.851 |
| Promals3D | 0.965** | 0.780** | 0.598 | 0.972** | 0.807 | 0.897** | 0.926** | 0.881 | 0.899 | 0.899* | 0.873 |
| T-Coffee DASH† | 0.966** | 0.740** | 0.396** | 0.970** | 0.756** | 0.941* | 0.934* | 0.868 | 0.899 | 0.917 | 0.830** |
| **[TC scores]** | | | | | | TC | | | | | |
| MAFFT | 0.798** | 0.254** | 0.075** | 0.852** | 0.407** | 0.838* | 0.456** | 0.586 | 0.598 | 0.591** | 0.554** |
| Promals | 0.851** | 0.393** | 0.298** | 0.919** | 0.582* | 0.817 | 0.496** | 0.516** | 0.508* | 0.572* | 0.663** |
| T-Coffee | 0.808** | 0.262** | 0.098** | 0.871** | 0.411** | 0.855 | 0.403** | 0.474** | 0.550 | 0.587 | 0.591** |
| Expresso | 0.845** | 0.372** | 0.173** | 0.919** | 0.518** | 0.752** | 0.369** | 0.391** | 0.440** | 0.514** | 0.579** |
| MAFFT-DASH | **0.909** | 0.440** | 0.259** | **0.961** | **0.550** | 0.853 | 0.557 | 0.610 | 0.533 | 0.643* | 0.666 |
| MAFFT-DASH Homologs | 0.922 | 0.464 | 0.335 | 0.957 | 0.588 | 0.855 | 0.576 | 0.603 | 0.490 | 0.652 | 0.684 |
| Promals3D | 0.892** | 0.451** | 0.407 | 0.952** | 0.630 | 0.755** | 0.502** | 0.580** | 0.495** | 0.555* | 0.690 |
| T-Coffee DASH† | 0.890** | 0.410** | 0.217** | 0.950** | 0.528** | 0.852 | 0.466** | 0.533* | 0.519 | 0.646 | 0.642** |
| Number of cases | 87 | 225 | 34 | 165 | 38 | 44 | 41 | 30 | 49 | 16 | 149 |

HMFM, HomFam; MBSF, Mattbench-Superfamily; MBTL, Mattbench-Twilight; OXFM, OxFam; BB11-BB50, BAliBASE subsets 11-50; SY, SISYPHUS. Scores that are significantly worse than the best are marked with * (*P* < 0.05) and ** (*P* < 0.01) as calculated with Wilcoxon signed-rank test. Others are in bold.

**Incorporates 3D information**

BIOS477/877 L15 - 35

**35**

---

# Slide 36

## MAFFT website

BIOS477/877 L15 - 36

**36**

## Slide 37

# MAFFT website

**—add** (Updated to use more resources, 2020/Apr/11)

Align **full length sequences** to **an MSA**

Input — Output (default) — Output (—keeplength)

**—addfragments** (Updated to use more resources, 2020/Apr/11 *New!*)

Align **fragment sequences** to **an MSA**

Input — Output (default) — Output (—keeplength)

**—addlong** (Experimental; Updated to use more resources, 2020/Apr/11 *New!*)

Align **long sequences** to **a short MSA**

Input — Output (default) — Output (—keeplength)

Merge two or more sub MSAs into a single MSA *In alpha testing* (2015/Jun) Help

Two or more sub MSAs are merged into a single MSA. Sub MSAs are assumed to be **phylogenetically separated** from each other. If it cannot be assumed, try **—add** or **—addfragments**.

https://mafft.cbrc.jp/alignment/server/merge.html

BIOS477/877 L15 - 37

**37**

## Slide 38

# Alignment Refinement with MAFFT



S477/877 L15 - 38

**38**

## Slide 39

# To solve the progressive-alignment problems

➢ **Incorporate more information to reduce early errors**
- **Structural alignment** (*e.g.*, Expresso, PROMALS3D, TM-Coffee, PRALINE, MAFFT-DASH)
- **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
  �safter 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 L15 - 39

**39**

## Slide 40

# PRALINE (Simossis, Kleinjung & Heringa 2005)

https://www.ibi.vu.nl/programs/pralinewww/



BIOS477/877 L16 - 40

**40**

## Slide 41

# PRALINE (Simossis, Kleinjung & Heringa 2005)



- From 624 HOMSTRAD pairwise alignments
- Q score: Sum-of-pairs (percentage of correctly aligned residue pairs)
ΔQ: Q score difference from PRALINE without PSI-BLAST

BIOS477/877 L16 - 41

**41**

## Slide 42

# Motif-Aware PRALINE (Dijkstra *et al.* 2021)

**Copper-binding site (PS00196):**
**[PGA]-x(0,2)-[YSA]-x(0,1)-[VFYL]-{SEDT}-C-x(1,2)-[PGA]-x(0,1)-H-x(2,4)-[MQ]**



https://github.com/ibivu/MA-PRALINE

BIOS477/877 L16 - 42

**42**

## Slide 43

# Clustal Ω http://www.clustal.org/omega/
### https://www.ebi.ac.uk/jdispatcher/msa/clustalo

- **Progressive alignment** following the guide tree
- **Features a fast method for making "guide trees"**
  - ➔ calculates only distances to *n* references (**mBed method**)
  - ➔ scalable for very large datasets
- **Alignment is done using HHalign (a profile hidden Markov model alignment)**
  - ➔ highly accurate alignment
- **Simple iterative refinement**
  - ➔ Alignment is converted to hidden Markov model (HMM)
  - ➔ Realign input sequences against the HMM

**Sievers *et al.* (2011, 2018)**

BIOS477/877 L16 - 43

**43**

## Slide 44

# Clustal Ω http://www.clustal.org/omega/
### https://www.ebi.ac.uk/jdispatcher/msa/clustalo

**Table I** BAliBASE 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) |
| Probcons | 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 | 0.575 | 0.579 | 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 bali score on core columns only. The average score over all families is given in the second column. The results for BAliBASE 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.

**Sievers *et al.* (2011)**

BIOS477/877 L16 - 44

**44**

## Slide 45

# Clustal Ω http://www.clustal.org/omega/
### https://www.ebi.ac.uk/jdispatcher/msa/clustalo



**Sievers *et al.* (2018)**

BIOS477/877 L16 - 45

**45**