# BIOS 477/877
# Bioinformatics and Molecular Evolution

**Instructor: Etsuko Moriyama**
**(School of Biological Sciences)**

| Spring 2026 | Lecture 15 |

**1**

---

# Today's topics

➢ **Multiple Sequence Alignment**

- **Progressive alignment (ClustalW)**
- **T-Coffee**

**2**

---

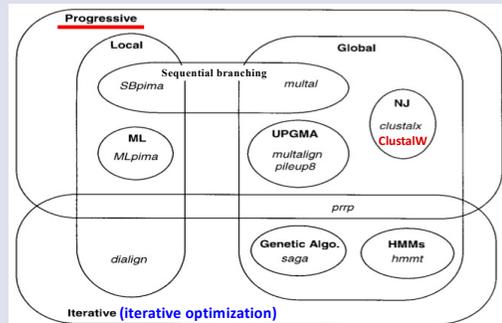## Multiple alignment as an extension of pairwise alignment

➢ **Dynamic programming algorithm**
→ Guarantees to find the optimal alignment
→ Optimal alignments are searched based on alignment score
- How can we score multiple alignment?
  - Sum of pairs score: $S(A) = \Sigma_{i,j} S(A_{ij})$
    → **No statistical justification**
→ **Not very efficient for many sequences**

> **We need good/efficient heuristic methods!**

**3**

---

# Multiple alignment methods
## (heuristic approach)



Thompson *et al.* (1999)
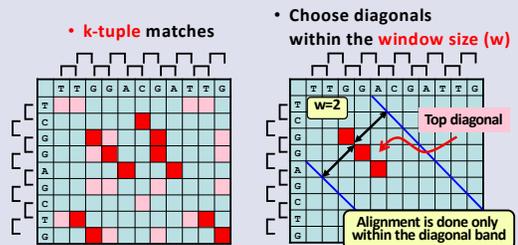
**4**

---

# ClustalW

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

**1. Pairwise alignment**
(fast approximation or **full dynamic programming)**

**2. Generate a distance matrix**
(% identities converted to **distances**)

**3. Construct a guide tree**
(**neighbor-joining** phylogenetic tree)

**4. Progressive alignment following the guide tree**
(scoring matrix, sequence weight, gap penalties, etc.)

Updated to **ClustalΩ** and should not be used; to try, see Course Web 'Links' page

**5**

---

# ClustalW: 1. Pairwise alignment

➢ **Fast approximation using a diagonal band**

- **k-tuple** matches
- Choose diagonals within the **window size (w)**



**The default is to use full dynamic programming: better**

**6**

## Slide 7

### ClustalW: 2 & 3. Guide tree generation

➤ **A guide tree is generated from a distance matrix**

Pairwise alignment | % identities | Distances (100 - % identifies)

Seq1 xxxxx
Seq2 yyyyy → S12 → D12

Seq1 xxxxx
Seq3 zzzzz → S13 → D13

Seq2 yyyyy
Seq3 zzzzz → S23 → D23

Distance matrix

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

Neighbor-joining phylogenetic reconstruction

**Guide tree**

BIOS477/877 L15- 7

**7**

## Slide 8

### ClustalW: steps

Thompson *et al.* (1994)

**1. Pairwise alignment**
(fast approximation or **full dynamic programming**)

**2. Generate a distance matrix**
(% identities converted to **distances**)

**3. Construct a guide tree**
(**neighbor-joining** phylogenetic tree)

**4. Progressive alignment following the guide tree**
(scoring matrix, sequence weight, gap penalties, etc.)

BIOS477/877 L15 - 8

**8**

## Slide 9

### ClustalW: 4. progressive alignment

➤ **Alignment is done following the guide tree**

0.02 S1 peeksav
0.15   peeksav geekaav
0.08 S2 geekaav
0.09
0.38 S3 egewglv
0.46 S4 aaektki

**Guide tree**

- Closest sequences are aligned first
- Pairwise alignment can be done by using simply the dynamic programming algorithm
- Next closest sequence is aligned against the existing alignment
  ➔ Profile alignment (can be alignment between alignments)

BIOS477/877 L15 - 9

**9**

## Slide 10

### ClustalW: profile alignment

➤ **Alignment against existing alignment**

Profile

S1 peeksav
S2 geekaav

| | | P | E | E | K | S | A | V | S1 |
| | | G | E | E | K | A | A | V | S2 |

S3 egewglv

E G E W G L V (S3)

A sequence (S3) is aligned against the alignment (S1, S2)

**Existing alignment = profile**
➔ **Each position of alignment shows the configuration of possible amino acids**

BIOS477/877 L15 - 10

**10**

## Slide 11

### ClustalW: profile alignment

➤ **How cell scores can be calculated**

S1 peeksav
S2 geekaav

S3 egewglv

From a scoring matrix

S1 vs. S3 = S(P,E)
S2 vs. S3 = S(G,E)

S1/S2 vs. S3
= {S(P,E)+S(G,E)}/2

| | P G |
|---|---|
| 0 | –g |
| E –g | ? –2g |

(gap penalty = –g)

Simple average from all pairwise scores

**Existing alignment = profile**
➔ **Each position of alignment shows the configuration of possible amino acids**

BIOS477/877 L15 - 11

**11**

## Slide 12

### ClustalW: profile alignment

➤ **Cell scores using simple vs. weighted average**

S1 peeksav (W1)
S2 geekaav (W2)

S3 egewglv (W3)

**Sequence weights**
from Thompson *et al.* (1994)

[Simple average (without weighting)]
Score(P/G,E) = {S(P,E) + S(G,E)} / 2

or

[Weighted average]
Score(P/G,E) = {W1xW3xS(P,E) + W2xW3xS(G,E)} / 2

If W1=W2=W3=1, the same as simple average

| | P G |
|---|---|
| 0 | –g |
| E –g | ? |

**Closely related sequences share information**
➔ **The sequences with duplicated information should receive smaller weights**

BIOS477/877 L15 - 12

**12**

2

**13**



**14**



**15**



**16**



**17**



**18**

## Slide 19

# ClustalW: profile alignment

➤ **How sequence weighting works: example 2**

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

**[Simple average] (match=2, mismatch=-1)**

**Alignment (a)**   [Score = 1 + gap penalty]
1: Score = {S(V,V)+S(V,V)+S(I,V)}/3 = (2+2-1)/3 = 1
2: Score = {S(A,G)+S(G,G)+S(V,G)}/3 = (-1+2-1)/3 = 0
3: Score = (gap penalty x 3)/3

**Alignment (b)**   [Score = 2 + gap penalty]
1: Score = {S(V,V)+S(V,V)+S(I,V)}/3 = (2+2-1)/3 = 1
2: Score = (gap penalty x 3)/3
3: Score = {S(G,G)+S(A,G)+S(G,G)}/3 = (2-1+2)/3 = 1

**Alignment (c)**   [Score = 1 + gap penalty]
1: Score = (gap penalty x 3)/3
2: Score = {S(A,V)+S(V,V)+S(V,V)}/3 = (-1-1+2)/3 = 0
3: Score = {S(G,G)+S(A,G)+S(G,G)}/3 = (2-1+2)/3 = 1

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

BIOS477/877 L15 - 19

---

## Slide 20

# ClustalW: profile alignment

➤ **How sequence weighting works: example 2**

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

**[Weighted average] (match=2, mismatch=-1)**

**Alignment (a)**

1: Score =
  {1.1x7xS(V,V)+1.1x7xS(V,V)+6.1x7xS(I,V)}/3
  = {7.7x2+7.7x2+42.7x(-1)}/3 = -3.97
2: Score =
  {1.1x7xS(A,G)+1.1x7xS(G,G)+6.1x7xS(V,G)}/3
  = {7.7x(-1)+7.7x2+42.7x(-1)}/3 = -11.67
3: Score = (gap penalty x 3)/3

**Alignment Score = -15.64 + gap penalty**

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

BIOS477/877 L15 - 20

---

## Slide 21

# ClustalW: profile alignment

➤ **How sequence weighting works: example 2**

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

**[Weighted average] (match=2, mismatch=-1)**

**Alignment (b)**

1: Score =
  {1.1x7xS(V,V)+1.1x7xS(V,V)+6.1x7xS(I,V)}/3
  = {7.7x2+7.7x2+42.7x(-1)}/3 = -3.97
2: Score = (gap penalty x 3)/3
3: Score =
  {1.1x7xS(G,G)+1.1x7xS(A,G)+6.1x7xS(G,G)}/3
  = {7.7x2+7.7x(-1)+42.7x(2)}/3 = 31.03

**Alignment Score = 27.06 + gap penalty**

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

BIOS477/877 L15 - 21

---

## Slide 22

# ClustalW: profile alignment

➤ **How sequence weighting works: example 2**

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

**[Weighted average] (match=2, mismatch=-1)**

**Alignment (c)**

1: Score = (gap penalty x 3)/3
2: Score =
  {1.1x7xS(A,V)+1.1x7xS(G,V)+6.1x7xS(V,V)}/3
  = {7.7x(-1)+7.7x(-1)+42.7x(2)}/3 = 23.33
3: Score =
  {1.1x7xS(G,G)+1.1x7xS(A,G)+6.1x7xS(G,G)}/3
  = {7.7x2+7.7x(-1)+42.7x(2)}/3 = 31.03

**Alignment Score = 54.36 + gap penalty**

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

BIOS477/877 L15 - 22

---

## Slide 23

# ClustalW: profile alignment

➤ **How sequence weighting works: example 2**

```
      0.5
   1      S1(1.1)        123
      0.5  S2(1.1)       VAG
0.3                      VGA
         6  S3(6.1)      IVG
         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   or     VGA   or     VGA
S3  IVG          IVG          IVG
S4  VG-          V-G          -VG
```

BIOS477/877 L15 - 23

---

## Slide 24

# ClustalW: parameter selection

➤ **How scoring matrix is chosen**

- **Users choose only a scoring matrix series (BLOSUM, PAM, *etc.*)**
- **Specific matrix (BLOSUM80, *etc.*) is determined based on distance between sequences**
  - 80 - 100% identity  → BLOSUM80
  - 60 - 80% identity   → BLOSUM62
  - 30 - 60% identity   → BLOSUM45
  - 0 - 30% identity    → BLOSUM30

BIOS477/877 L15 - 24
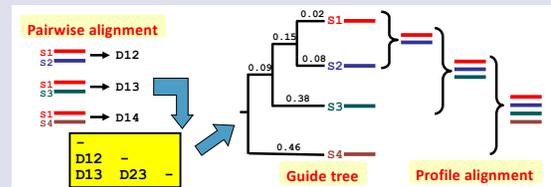
---

4

## ClustalW: parameter selection

➢ **How gap penalties are determined**
  ▪ **Initial gap penalties: GOP** (gap opening) and **GEP** (gap extension) are set by the user
  ▪ **Position- or residue-specific** gap penalties
    o Weight (scoring) matrix dependent
    o Similarity level dependent
    o Sequence length dependent
    o Position-specific
      ➔ if gaps already exist
      ➔ Residue-specific (*e.g.*, hydrophilic stretches)
  **See Thompson *et al.* (1994) for more details**

BIOS477/877 L15 - 25

**25**

---

## ClustalW: summary & limitation



■ **Progressive alignment**
  ➔ **Greedy** (finds local optima, but no guarantee for global optima)
  ➔ **Guide tree** is built only once at the beginning (cannot be fixed)
  ➔ **Errors** in the early alignments **(incorrect gap positions)** cannot be rectified later
  ➔ ClustalW does not rely on one set of alignment parameters
    o **Position- and various feature-specific scoring matrix and gap penalties**
■ **Global alignment** only (local similarity may be missed)

BIOS477/877 L15 - 26

**26**

---

## How to solve 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

**27**

---

## T-Coffee

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

➢ **T-Coffee: Tree-based Consistency Objective Function for alignment Evaluation**
➢ **Progressive alignment**
  ➔ **Uses a guide tree, fast**
➢ **Tries to avoid the greedy nature of the progressive algorithm**
  ➔ **Using alignment libraries** derived from a mixture of alignment methods (global, local, *etc.*)

**[T-Coffee on the Web]**
  o **T-Coffee home page**
  o **T-Coffee site from Notredame Lab**
  o **T-Coffee @ EBI Tools** (only T-Coffee, without other associated programs)

BIOS477/877 L15 - 28

**28**

---

## T-Coffee: primary libraries

➢ **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% ⇨ W=50

BIOS477/877 L15 - 29

**29**

---

## T-Coffee: pooled primary library



**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**

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

**A & B aligned through C**

Weight of the stacked alignment = W2+W3

BIOS477/877 L15 - 30

**30**

5

## Slide 31

### T-Coffee: extended library

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

→ Pooled Primary Library

→ Extended 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 to the weight of the residue pair found in the aligned region

**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        THE      FAT CAT
     |||      ||  \\\
SeqB GARFIELD THE      FAST CAT
```

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

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

BIOS477/877 L15 - 31

**31**

---

## Slide 32

### T-Coffee: residue-pair scoring

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

→ Pooled Primary Library

→ Extended Library

→ Generate a sequence-pair specific scoring matrix

- 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 to the weight of the residue pair found in the aligned region

```
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

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

BIOS477/877 L15 - 32

**32**

---

## Slide 33

### T-Coffee: progressive alignment

Global pairwise alignment Primary Library | Local pairwise alignment Primary Library

→ Pooled Primary Library

→ Extended Library

→ **Progressive alignment with the neighbor-joining guide tree**

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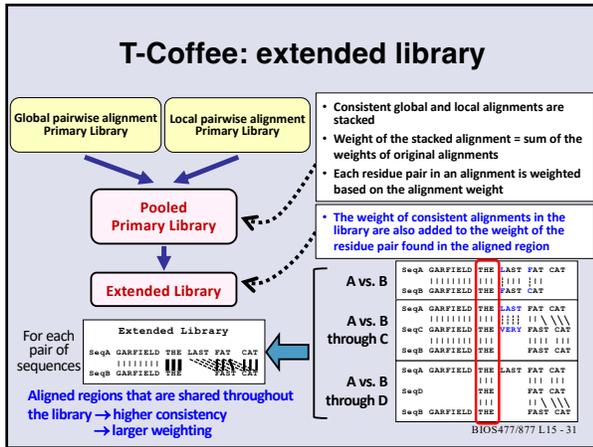- Weight of the stacked alignment = sum of the weights of original alignments
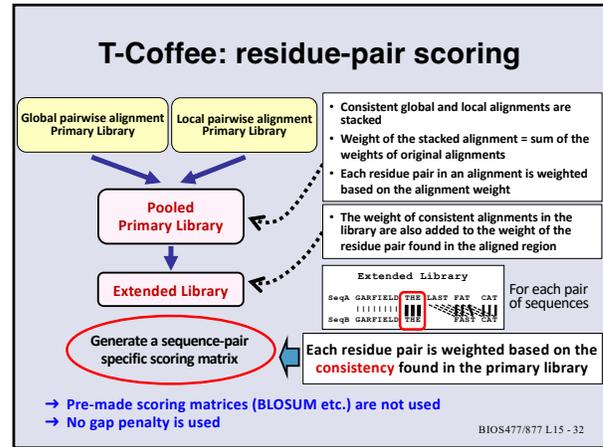- Each residue pair in an alignment is weighted based on the alignment weight
- The weight of consistent alignments in the library are also added to the weight of the residue pair found in the aligned region
  → Each residue pair is weighted based on the **consistency** found in the primary library

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

BIOS477/877 L15 - 33

**33**

---

## Slide 34

### T-Coffee programs

**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 | DNA

T-COFFEE SIMPLE MSA

Regular T-Coffee

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

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

BIOS477/877 L15 - 34

**34**

---

## Slide 35

### T-Coffee programs

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



BIOS477/877 L15 - 35

**35**

---

## Slide 36

### T-Coffee programs

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

T-Coffee: aligns DNA, RNA, or protein sequences

R-Coffee and SARA-Coffee: align RNA sequences considering their 2D/3D structures
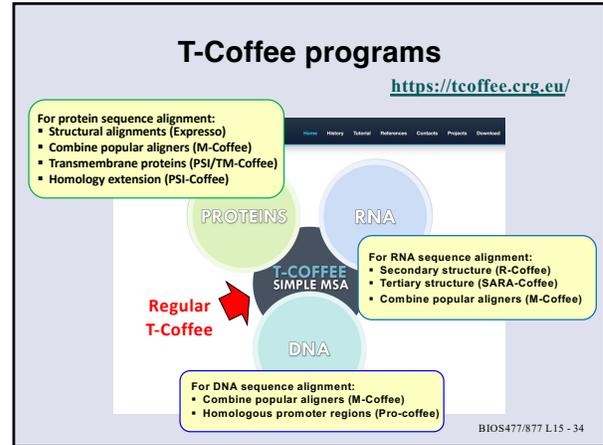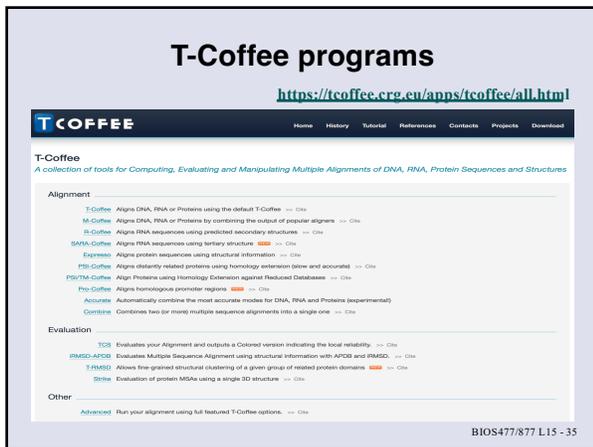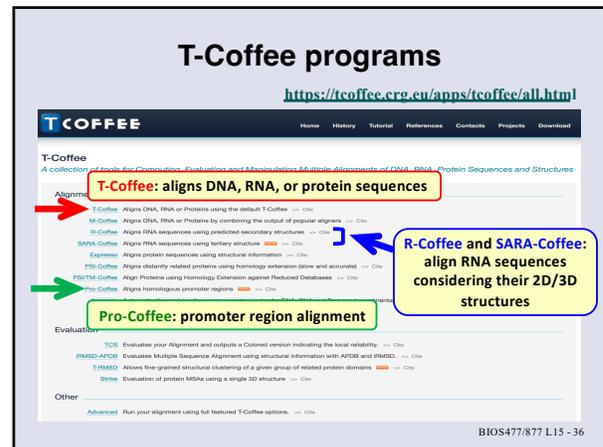
Pro-Coffee: promoter region alignment



BIOS477/877 L15 - 36

**36**

## T-Coffee programs

**Expresso (3D-Coffee): aligns protein sequences with structural information incorporated as pairwise structural alignments (library)**

**PSI-Coffee: extends each sequence with PSI-BLAST**
**→ Structural information can be incorporated**

**PSI/TM-Coffee: PSI-BLAST against a reduced database (for fast search)**
**→ Only TM proteins can be targeted (optional)**

BIOS477/877 L15 - 37

**37**

## T-Coffee programs to improve MSA

**M-Coffee: combines multiple alignments**
**→ shows consistency between multiple alignments**

**Combine: combines your own multiple alignments**

BIOS477/877 L15 - 38

**38**

## T-Coffee programs to evaluate MSA

**TCS: evaluates alignments based on the consistency between pairwise alignment library and multiple alignment**

**Structural information based evaluation and clustering**

BIOS477/877 L15 - 39

**39**

## Evaluating multiple alignments using TCS

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

SCORE=395

BAD AVG GOOD

**T-Coffee**

```
1aboA : 48
1ycsB : 47
1pht  : 45
1vie  : 32
1ihvA : 32
cons  : 39
```

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

cons

SCORE=305

BAD AVG GOOD

**ClustalW2**

```
1aboA : 42
1ycsB : 41
1pht  : 39
1vie  : 25
1ihvA : 20
cons  : 30
```

```
1aboA   -NLFV-ALYDFVASGDNTLSITKGEKLRV-------LGYNHNG-------EWCEAQ--TK42NGQGWVPSNYITPVN------57
1ycsB   KGVIY-ALWDYEPQNDDELPMKEGDCMTI-----IHREDEDEI-----EWWWAR--LN45DKEGYVPRNLLGLYP------60
1pht    -GYGYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQEARPEEIGWLNGYNETT59GERGDFPGTYVEYIGRKKISP80
1vie    ---------DRVRKKSG--AAWQGQIVGW----------YCTNLTP-----EGYAVESEAHP36GSVQIYPVAALERIN------51
1ihvA   -----NFRVYYRDSRD--PVWKGPAKLL-------WKGEG------AVVIQD--N33SDIKVVPRRKAKIIRD------49
```

cons

**40**

## Evaluating multiple alignments using TCS

Reference alignment

a

**T-Coffee**

b

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

c

**PSI-Coffee**

d

**Expresso**

from Taly *et al.* (2011)

**41**

7