

## Cap. 6 How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to an Arbitrary Directed Acyclic Graph
- From Global to Local Alignment [Matrici di Score...]
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- **Multiple Sequence Alignment** | OGGI “COMPLETIAMO”

1

**Multiple Alignment Problem:** *Find the highest-scoring alignment between multiple strings.*

- **Input:** A collection of  $t$  strings (and some way of scoring columns of a multiple alignment).
- **Output:** A multiple alignment of these strings having maximum score.

MSA is used for:

- Detection of conserved domains in a group of genes or proteins
- Construction of a phylogenetic tree
- Prediction of a protein 3D-structure starting with its amino acids sequence (e.g., AlphaFold, RoseTTAFold)
- Determination of a consensus sequence (e.g., transposons)

2

**Multiple Alignment Problem:** Find the highest-scoring alignment between multiple strings.

- **Input:** A collection of  $t$  strings (and some way of scoring columns of a multiple alignment).
- **Output:** A multiple alignment of these strings having maximum score.

The [sum-of-pairs score \(sop\)](#) of  $M$  is the sum of all pair-wise induced alignment scores

$$\begin{array}{c}
 \text{AAGAA\_A} \\
 \text{AT\_AATG} \\
 \text{CTG\_G\_G}
 \end{array}
 \quad \left\langle \quad \right\rangle \quad \left\{ \quad \right\} \quad \dots$$

$d/i =$   
 $r =$   
 $m =$

$$\begin{array}{c}
 \text{AAGAA\_A} \\
 \text{ATAATG} \\
 \text{C\_TGG\_G}
 \end{array}
 \quad \left\langle \quad \right\rangle \quad \left\{ \quad \right\}$$

The [optimal MSA for S](#) w.r.t. to [sop](#) is the MSA with the highest sop-score over all possible MSA for S

3

## Profile Representation of Multiple Alignment



4

## Profile Representation of Multiple Alignment

A *profile* is a probability for each letter to occur in each column, i.e., a description of the consensus of a multiple sequence alignment. It uses a position-specific scoring system to capture information about the degree of conservation at various positions in the multiple alignment.

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G

$p=1/\#$

A	0	1	0	0	0	0	1	0	0	.8	0	0	0
C	.6	0	0	0	1	0	0	.4	1	0	.6	.2	0
G	0	0	1	.2	0	0	0	0	0	.2	0	0	.4
T	.2	0	0	0	0	1	0	.6	0	0	0	.2	0
-	.2	0	0	.8	0	0	0	0	0	0	.4	.8	.4

5

### Is the “Profile representation” mentioned above suitable?

To predict an HIV phenotype, we need an *accurate* alignment: a single misalignment at a position influencing the SI phenotype leads to an error (e.g. **11-25 Rule**).

**STOP and Think:** Was it a good idea to use the *same* scoring matrix across different columns of an alignment?



We need a statistically solid **problem formulation** for alignment that uses a *different* scoring approach at different columns.

6

# Hidden Markov Model (HMM)



7

# Hidden Markov Model (HMM)

$\Sigma$ : an **alphabet** of emitted symbols H and T

**States** : a set of **hidden states** *F* and *B*

**Transition** =  $(transition_{i,k})$ : a  $|States| \times |States|$  matrix of **transition probabilities** changing from state  $i$  to state  $k$  *F*      *B*  
*F*    0.9    0.1  
*B*    0.1    0.9

**Emission**=  $(emission_k(b))$ : a  $|States| \times |\Sigma|$  matrix of **emission probabilities** emitting symbol  $b$  when the HMM is in state  $k$  *H*      *T*  
*F*    0.50    0.50  
*B*    0.75    0.25

Goal: Infer the most likely sequence of hidden states based on the sequence of emitted symbols

8

## Hidden Path

**Hidden path:** the sequence  $\pi = \pi_1 \dots \pi_n$  of states that the HMM passes through.

- $\Pr(x, \pi)$ : the probability that an HMM follows the hidden path  $\pi$  and emits the string  $x = x_1 x_2 \dots x_n$ .

$x:$	T	H	T	H	H	H	T	H	T	T	H
$\pi:$	F	F	F	B	B	B	B	B	F	F	F

$$\sum_{\text{all possible emitted strings } x} \sum_{\text{all possible hidden paths } \pi} \Pr(x, \pi) = 1$$

- $\Pr(x|\pi)$ : the **conditional probability** that an HMM emits the string  $x$  after following the hidden path  $\pi$ .

$$\sum_{\text{all possible emitted strings } x} \Pr(x|\pi) = 1$$

9

$$\Pr(x, \pi) = \Pr(x|\pi) * \Pr(\pi)$$

- $\Pr(x, \pi)$ : the probability that an HMM follows the hidden path  $\pi$  and emits the string  $x$ .
- $\Pr(x_i|\pi_i)$  – probability that  $x_i$  was emitted from the state  $\pi_i$  (equal to *emission* <sub>$\pi_i$</sub> ( $x_i$ )).
- $\Pr(\pi_{i-1} \rightarrow \pi_i)$  – probability that the HMM moved from  $\pi_{i-1} \rightarrow \pi_i$  (equal to *transition* <sub>$\pi_{i-1}, \pi_i$</sub> ).

$x:$	T	H	T	H	H	H	T	H	T	T	H
$\pi:$	F	F	F	B	B	B	B	B	F	F	F
$\Pr(\pi_{i-1} \rightarrow \pi_i)$	.5	.9	.9	.1	.9	.9	.9	.9	.1	.9	.9
$\Pr(x_i \pi_i)$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$

$$\Pr(\pi) = \prod_{i=1,n} \Pr(\pi_{i-1} \rightarrow \pi_i) = \prod_{i=1,n} \text{transition}_{\pi_{i-1}, \pi_i}$$

$$\Pr(x|\pi) = \prod_{i=1,n} \Pr(x_i|\pi_i) = \prod_{i=1,n} \text{emission}_{\pi_i}(x_i)$$

10

## Decoding Problem... Viterbi algorithm

**Decoding Problem:** Find an optimal hidden path in an HMM given its emitted string.

- **Input:** A string  $x = x_1 \dots x_n$  emitted by an HMM  $(\Sigma, States, Transition, Emission)$ .
- **Output:** A path  $\pi$  that maximizes the probability  $Pr(x, \pi)$  over all possible paths through this HMM.

$$\begin{aligned} Pr(x, \pi) &= Pr(x|\pi) * Pr(\pi) \\ &= \prod_{i=1,n} Pr(x_i|\pi_i) * Pr(\pi_{i-1} \rightarrow \pi_i) \\ &= \prod_{i=1,n} \text{emission}_{\pi_i}(x_i) * \text{transition}_{\pi_{i-1}, \pi_i} \end{aligned}$$

11

**HMM, PROFILE, MSA**

12

## From Alignment to Profile

	1	2	3	4	5	6	7	8	
<i>Alignment</i>	A	C	D	E	F	A C	A	D	F
	A	F	D	A	-	--	C	C	F
	A	-	-	E	F	D -	F	D	C
	A	C	A	E	F	--	A	-	C
	A	D	D	E	F	AA	A	D	F

Remove columns if the fraction of space symbols ("--") exceeds  $\theta$ , **the maximum fraction of insertions threshold.**

13

## From Alignment to Profile

	1	2	3	4	5	6	7	8	
<i>Alignment</i>	A	C	D	E	F	A C	A	D	F
	A	F	D	A	-	--	C	C	F
	A	-	-	E	F	D -	F	D	C
	A	C	A	E	F	--	A	-	C
	A	D	D	E	F	AA	A	D	F

	1	2	3	4	5	6	7	8
<i>Alignment*</i>	A	C	D	E	F	A	D	F
	A	F	D	A	-	C	C	F
	A	-	-	E	F	F	D	C
	A	C	A	E	F	A	-	C
	A	D	D	E	F	A	D	F

14

From Alignment to Profile									
	1	2	3	4	5	6	7	8	
Alignment	A	C	D	E	F	A C	A	D	F
	A	F	D	A	-	--	C	C	F
	A	-	-	E	F	D -	F	D	C
	A	C	A	E	F	--	A	-	C
	A	D	D	E	F	AA	A	D	F
Alignment*	A	C	D	E	F	A	D	F	
	A	F	D	A	-	C	C	F	
	A	-	-	E	F	F	D	C	
	A	C	A	E	F	A	-	C	
	A	D	D	E	F	A	D	F	
PROFILE(Alignment*)	A	1	0	1/4	1/5	0	3/5	0	0
	C	0	2/4	0	0	0	1/5	1/4	2/5
	D	0	1/4	3/4	0	0	0	3/4	0
	E	0	0	0	4/5	0	0	0	0
	F	0	1/4	0	0	1	1/5	0	3/5
Profile									

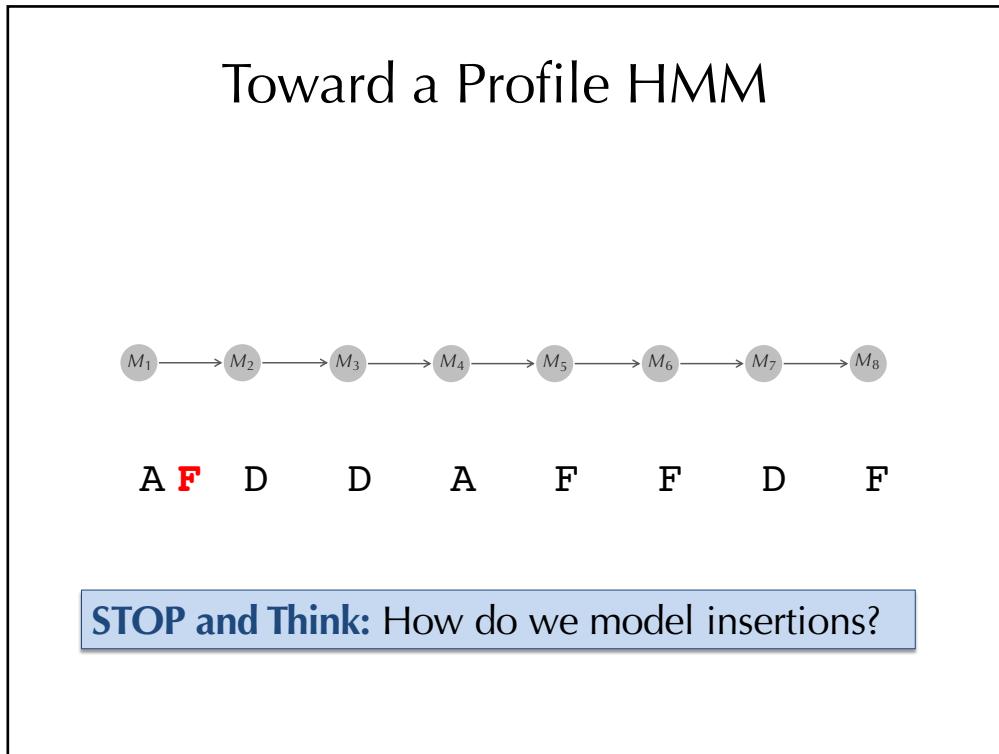
15

From Alignment to Profile									
	1	2	3	4	5	6	7	8	
Alignment	A	C	D	E	F	A C	A	D	F
	A	F	D	A	-	--	C	C	F
	A	-	-	E	F	D -	F	D	C
	A	C	A	E	F	--	A	-	C
	A	D	D	E	F	AA	A	D	F
Alignment*	A	C	D	E	F	A	D	F	
	A	F	D	A	-	C	C	F	
	A	-	-	E	F	F	D	C	
	A	C	A	E	F	A	-	C	
	A	D	D	E	F	A	D	F	
PROFILE(Alignment*)	A	1	0	1/4	1/5	0	3/5	0	0
	C	0	2/4	0	0	0	1/5	1/4	2/5
	D	0	1/4	3/4	0	0	0	3/4	0
	E	0	0	0	4/5	0	0	0	0
	F	0	1/4	0	0	1	1/5	0	3/5
Vice versa???									
GTCt/aGa/c									
									

16

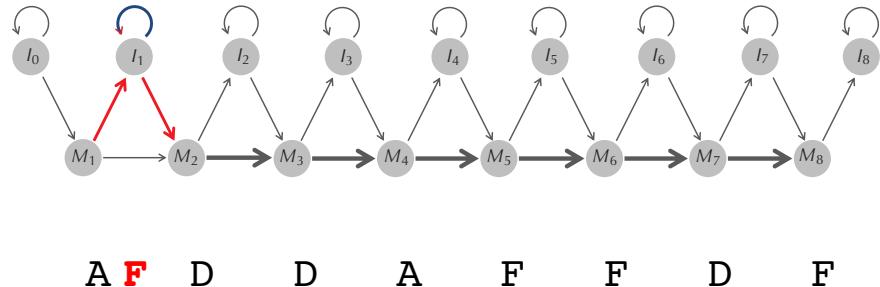
	1	2	3	4	5	6	7	8	
Alignment	A	C	D	E	F	A C	A	D	F
	A	F	D	A	-	--	C	C	F
	A	-	-	E	F	D -	F	D	C
	A	C	A	E	F	-	A	-	C
	A	D	D	E	F	A A	A	D	F
Alignment*	A	C	D	E	F	A	D	F	
	A	F	D	A	-	C	C	F	
	A	-	-	E	F	F	D	C	
	A	C	A	E	F	A	-	C	
	A	D	D	E	F	A	D	F	
PROFILE(Alignment*)	A 1	0	1/4	1/5	0	3/5	0	0	
	C 0	2/4	0	0	0	1/5	1/4	2/5	
	D 0	1/4	3/4	0	0	0	3/4	0	
	E 0	0	0	4/5	0	0	0	0	
	F 0	1/4	0	0	1	1/5	0	3/5	
<b>HMM diagram</b>									
1 * .25 * .75 * .20 * 1 * .20 * .75 * .60									
A D D A F F D F									

17



18

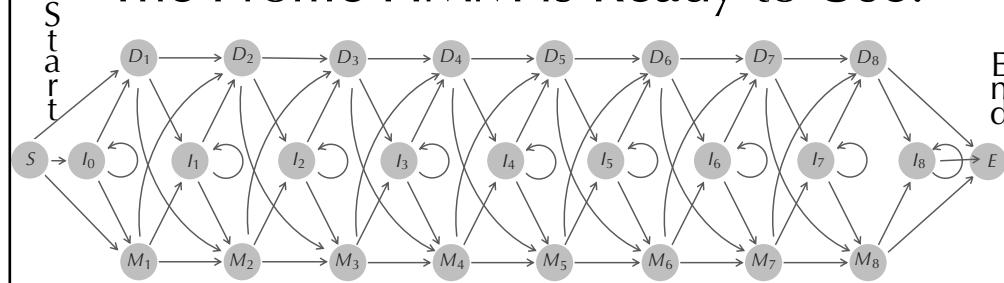
## Toward a Profile HMM: Insertions



**STOP and Think:** How do we model deletions?

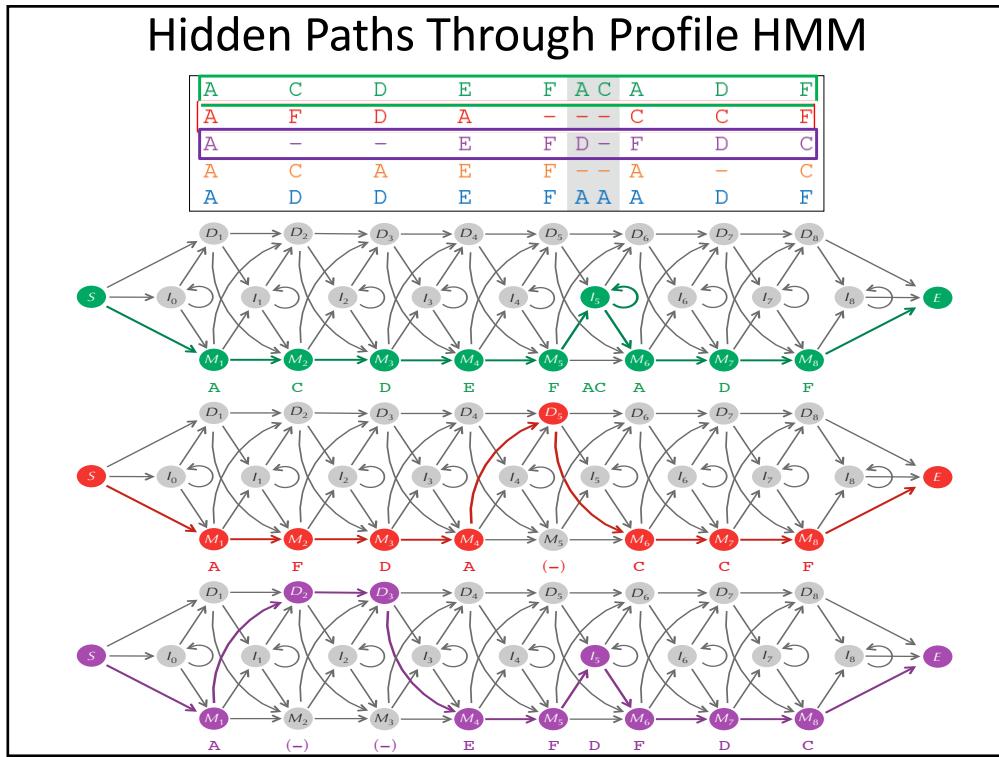
19

## The Profile HMM is Ready to Use!

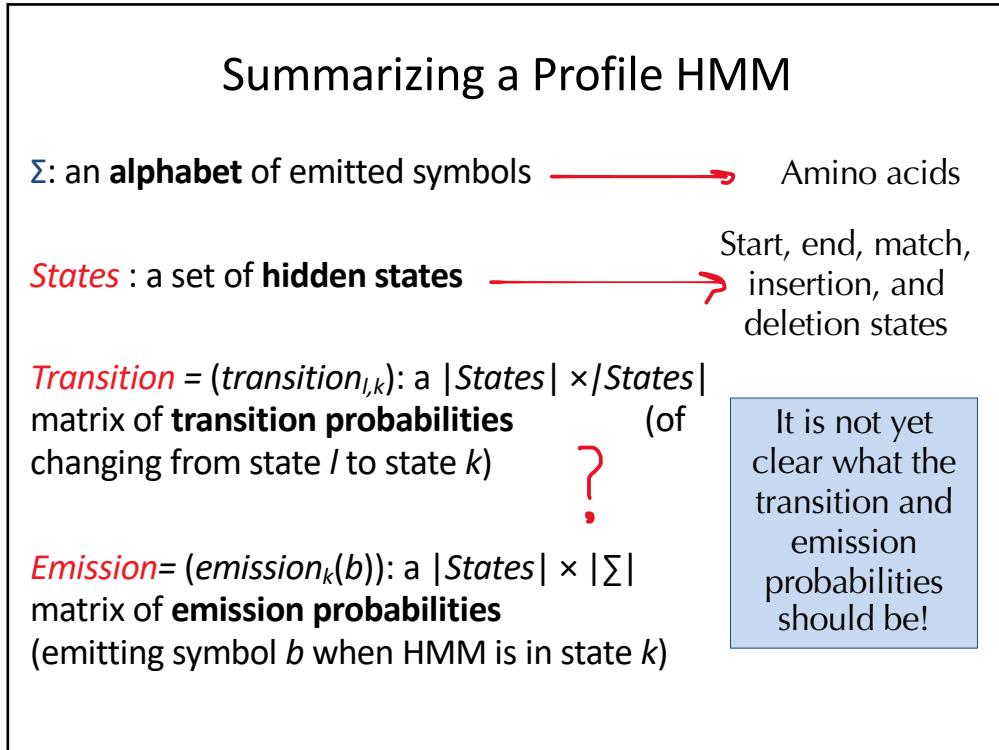


our «network»

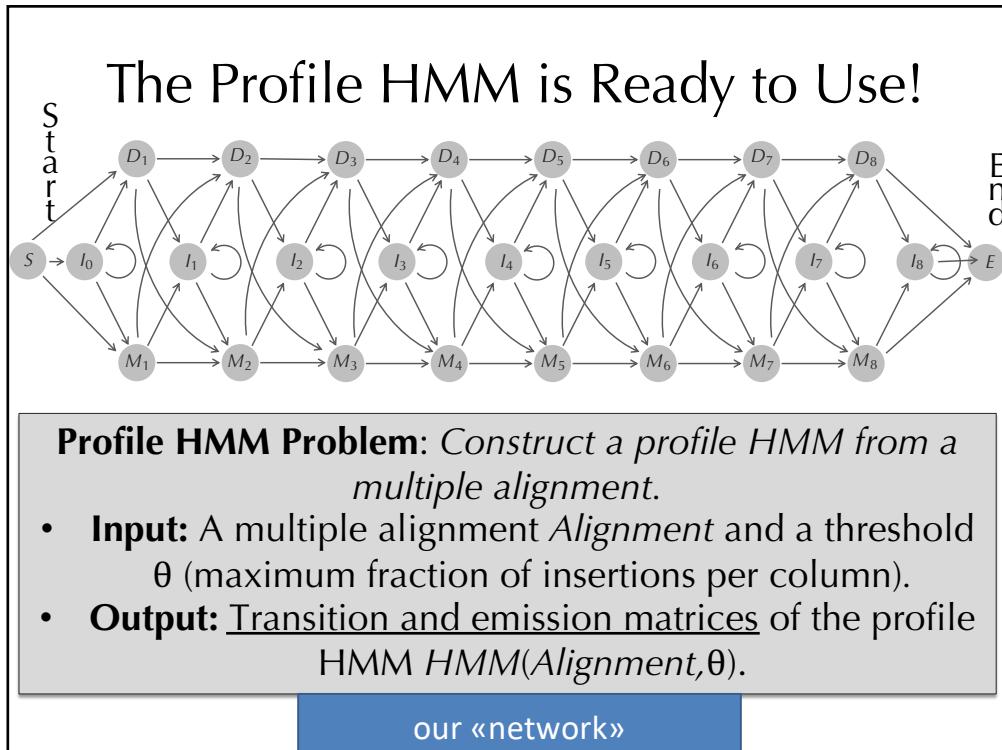
20



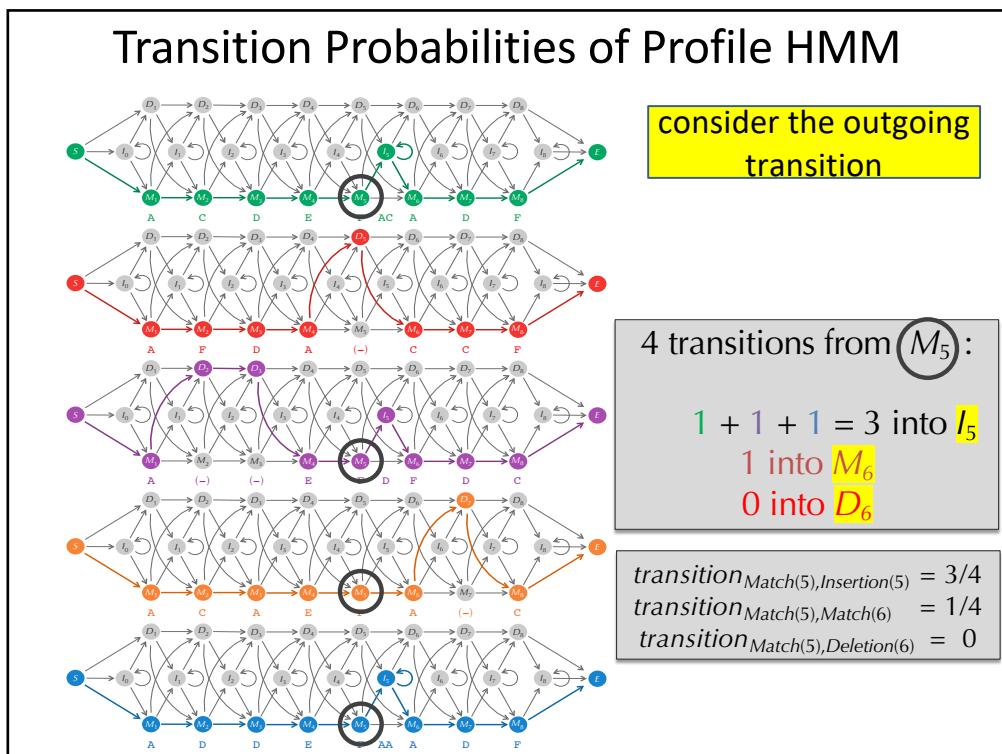
21



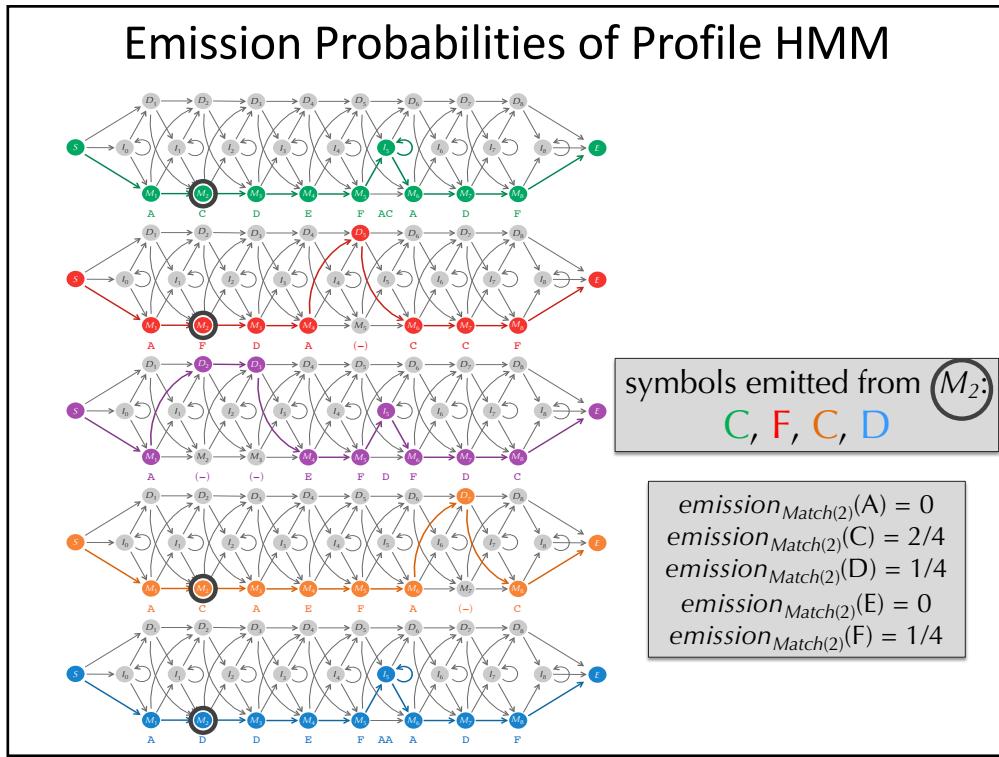
22



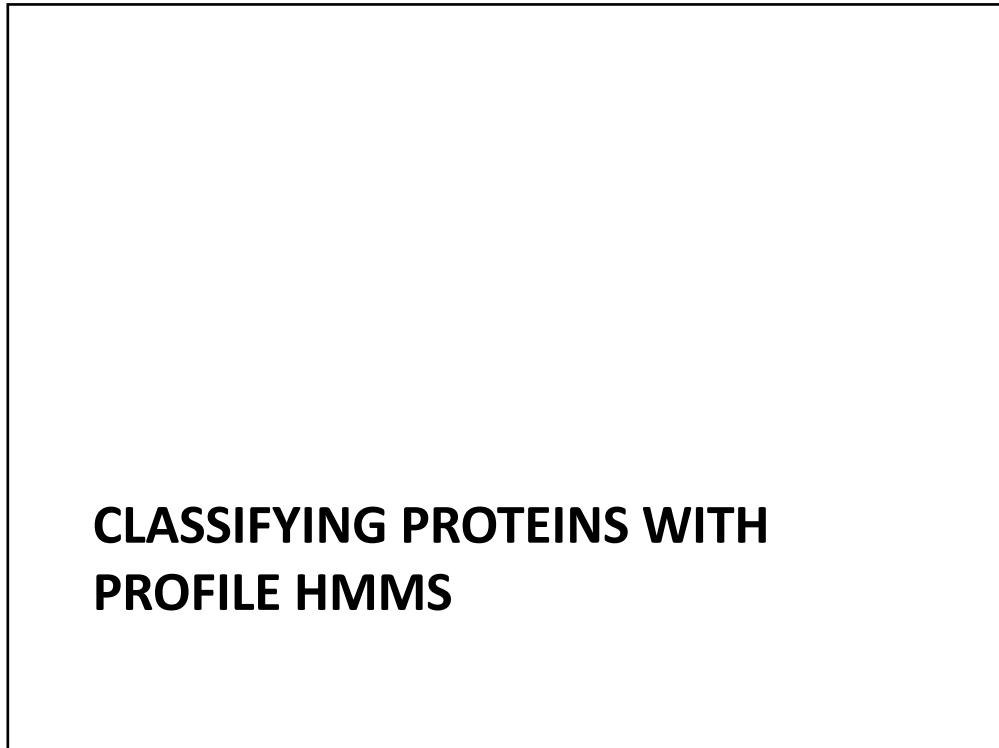
23



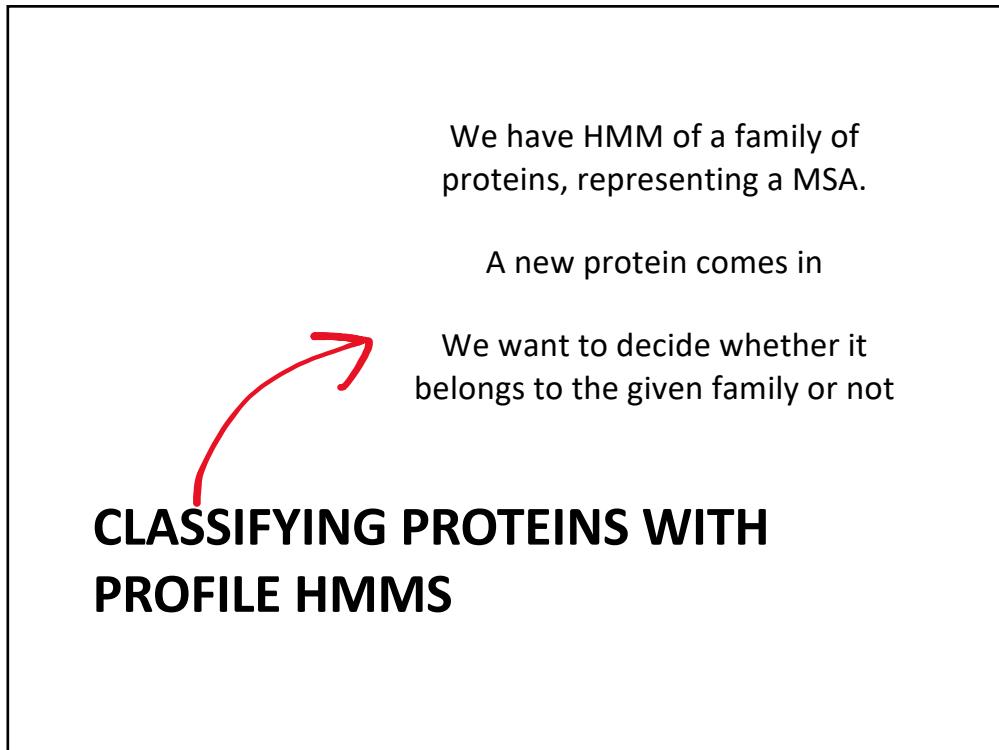
24



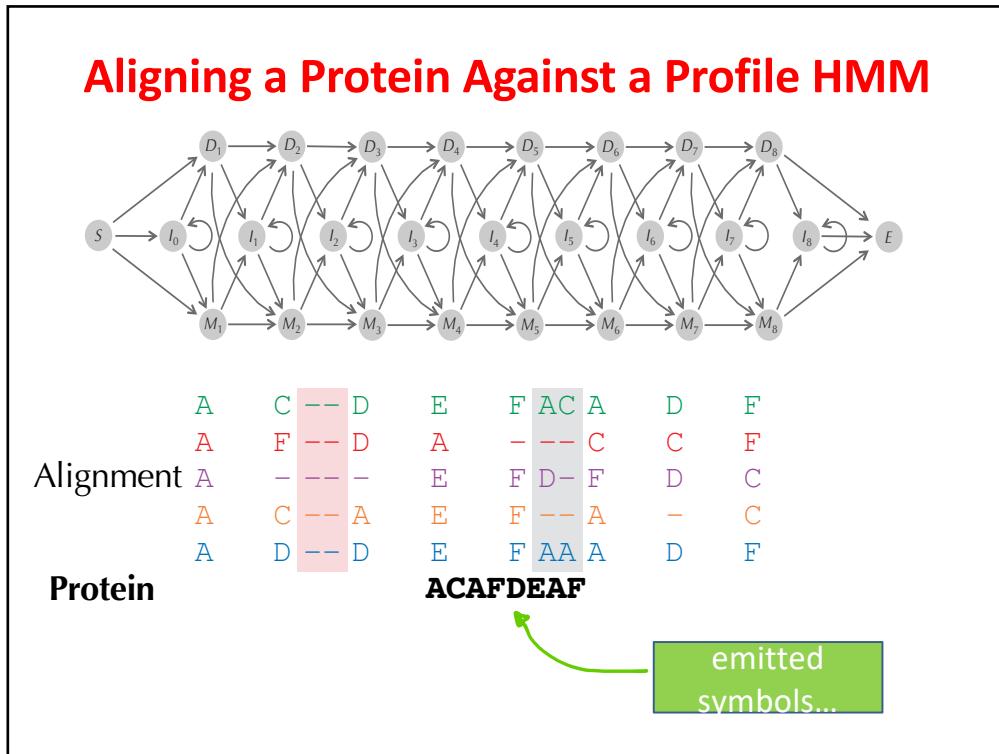
25



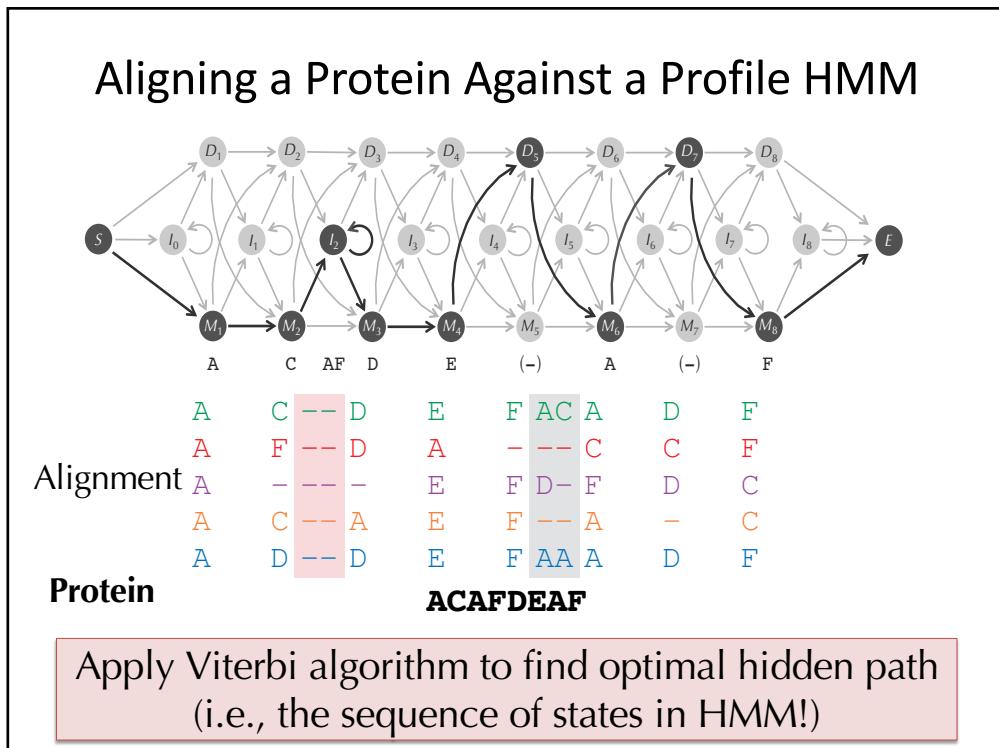
26



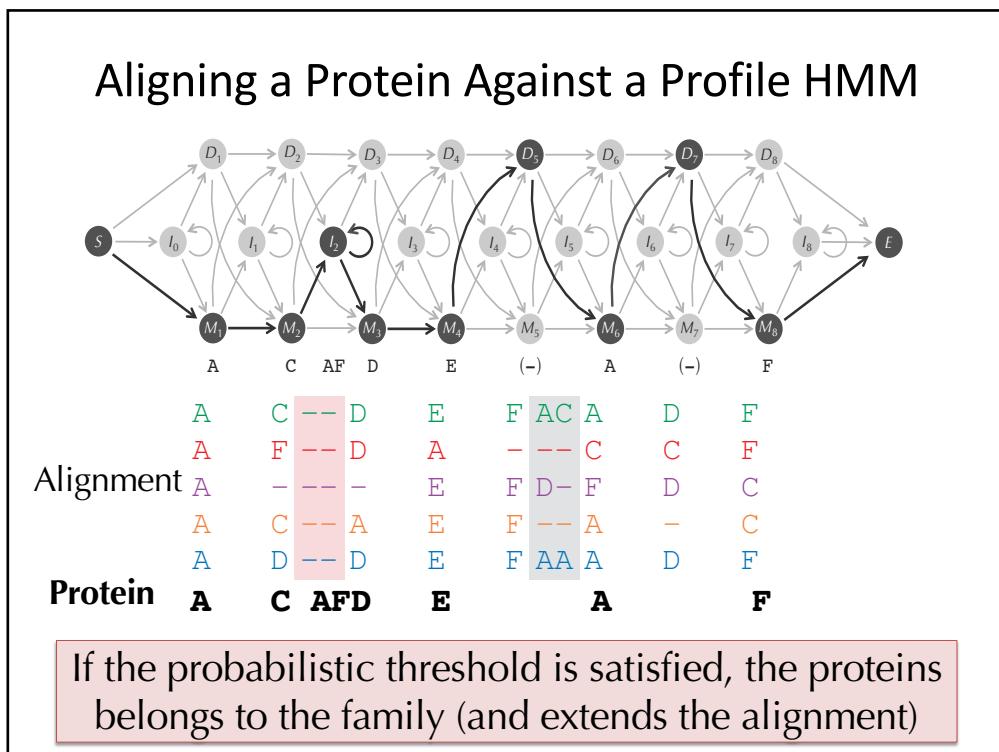
27



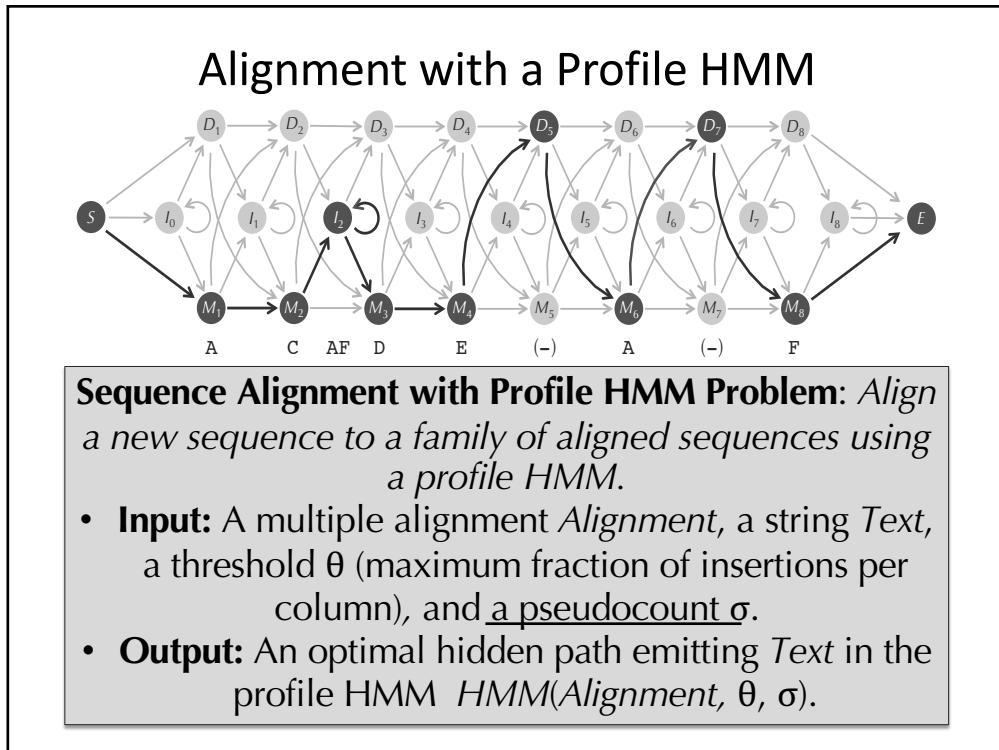
28



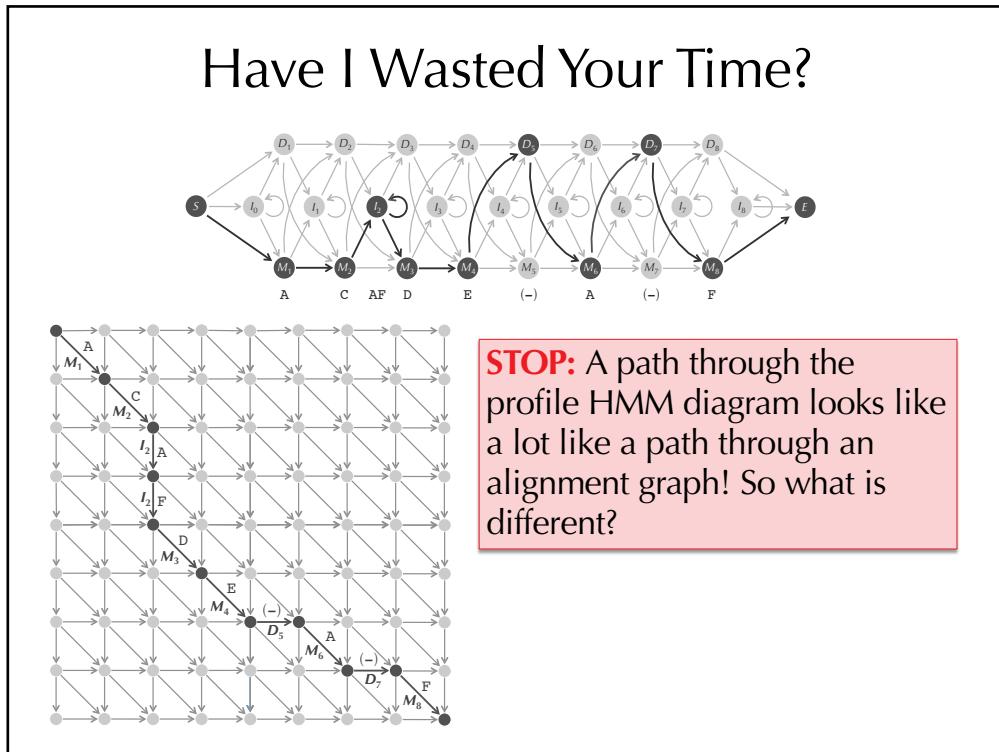
29



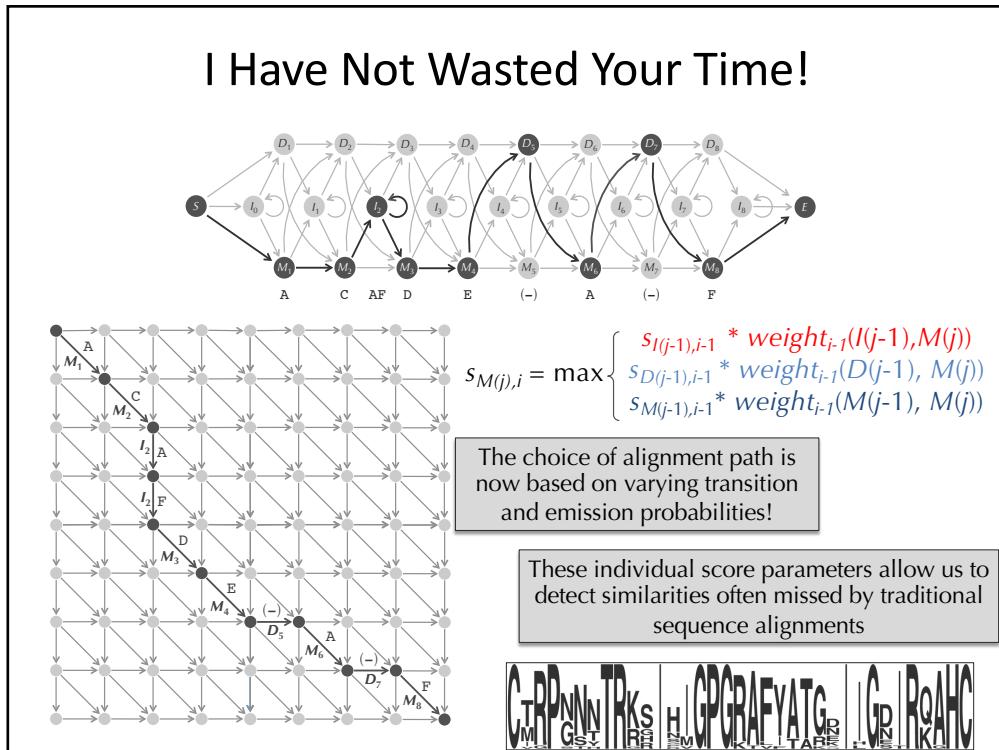
30



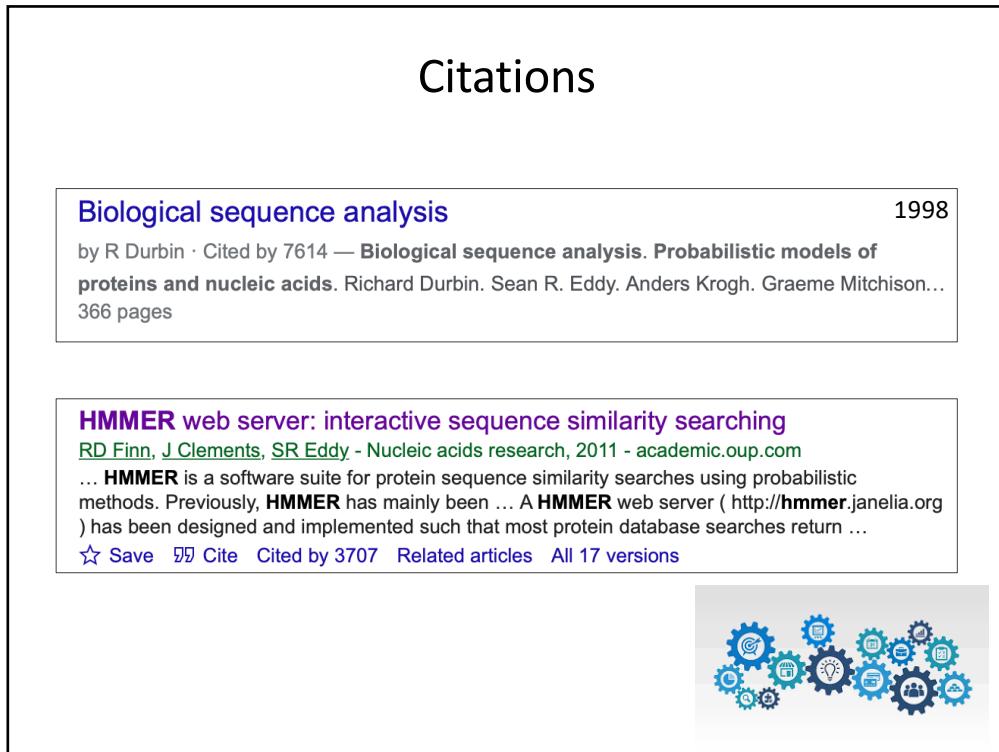
31



32



33



34

## Multiple Alignment: History

**1975 Sankoff**

Formulated multiple alignment problem and gave dynamic programming solution

**1988 Carrillo-Lipman**

Branch and Bound approach for MSA

**1990 Feng-Doolittle**

Progressive alignment

**1994 Thompson-Higgins-Gibson-ClustalW**

Most popular multiple alignment program

**1998 Morgenstern et al.-DIALIGN**

Segment-based multiple alignment

**2000 Notredame-Higgins-Heringa-T-coffee**

Using the library of pairwise alignments

**2004 MUSCLE**

**What's next?**

35

Inputs: N sequences

1. Using standard pairwise alignment, calculate a matrix of distances (alignment scores) between each pair of sequences. Consider this as an N-clique G, where edge {i,j} is labeled with the score of an optimal alignment of the i-th and j-th sequences.

2. Use [Kruskal's algorithm](#) to find a minimum spanning tree of G. Whenever a minimum spanning tree edge would connect two components, instead add a new root node with directed edges to the roots of the two components. This is the "guide tree".

3. Do [pairwise alignments according to the guide tree](#), working from the leaves to the root.

36

# Multiple Alignment: History

## 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

## 1988 Carrillo-Lipman

Branch and Bound approach for MSA

## 1990 Feng-Doolittle

Progressive alignment

## 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

## 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

## 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

## 2004 MUSCLE

uses HMM profile-profile techniques and progressive alignment

Profile HMMs are probabilistic models that encapsulate the evolutionary changes that have occurred in a set of related sequences (i.e. a multiple sequence alignment).

linguaggio C

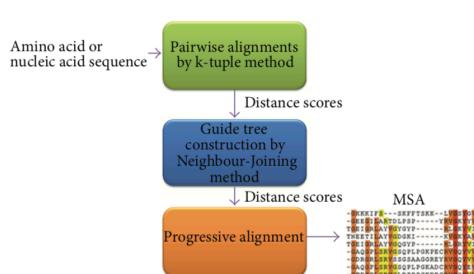
**What's next?**

37

# ClustalW

## (Cluster Alignment Weighted)

Iterative progressive alignment



I'idea è di iniziare  
a raggruppare  
sequenze simili

FIGURE 1: ClustalW algorithm, which works by taking an input of amino acid or nucleic acid sequences, completing a pairwise alignment using the k-tuple method, guide tree construction using the Neighbour-Joining method, followed by a progressive alignment to output a multiple sequence alignment.

(Clustal Omega: Iterative progressive alignment using hidden Markov models)

38

## Step 1 : Pairwise alignment of all sequences

Example : Alignment of 7 globins (Hbb\_human, Hbb\_horse, Hba\_human, Hba\_horse, Myg\_phyca, Glb5\_petma and Lgb2\_lupla)

Hbb_human	1	VHLTPEEKSAVTALWGVNVNDEVGGEALGRLLVVYPWTQRFESFGDLS...
Hbb_horse	2	VQLSGEEKAAVVLALWDKVNEEVGGEALGRLLVVYPWTQRFESFGDLSN...
Hbb_human	1	LTPEEKSAVTALWGVV.NVDEVGGEALGRLLVVYPWTQRFESFGDLS...
Hba_human	3	LSPADKTNVKAAGKVGAAHAGEYGAEALERMFLSFPTTKTYFPFH.DLS...
Hba_horse	2	LSPADKTNVKAAGKVGAAHAGEYGAEALERMFLSFPTTKTYFPFH.DLS...

The alignment can be obtained with:

- global or local methods
- heuristic methods
- ktuple method

Adapted from Julie Thompson, IGBMC

39

## Step 2 : Distance matrix construction

$$\text{Distance between two sequences} = 1 - \frac{\text{No. identical residues}}{\text{No. aligned residues}}$$

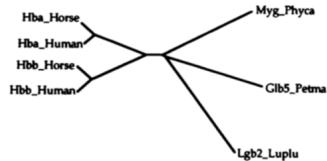
Hbb_human	1	-					
Hbb_horse	2	.17	-				
Hba_human	3	.59	.60	-			
Hba_horse	4	.59	.59	.13	-		
Myg_phyca	5	.77	.77	.75	.75	-	
Glb5_petma	6	.81	.82	.73	.74	.80	-
Lgb2_lupla	7	.87	.86	.86	.88	.93	.90
	1	2	3	4	5	6	7

Adapted from Julie Thompson, IGBMC

40

## Step 3 : Guide tree construction

### Unrooted NJ tree



Hbb_human	1	-					
Hbb_horse	2	.17	-				
Hba_human	3	.59	.60	-			
Hba_horse	4	.59	.59	.13	-		
Myg_phyc	5	.77	.77	.75	.75	-	
Glb5_petma	6	.81	.82	.73	.74	.80	-
Lgb2_luplu	7	.87	.86	.86	.88	.93	.90
	1	2	3	4	5	6	7

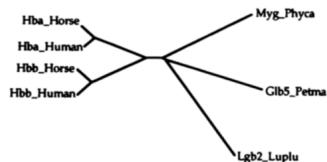
### **Neighbor-Joining** clustering method:

The tree is then built by linking the least distant pair of nodes. When two nodes are linked, their common ancestral node is added to the tree and the terminal nodes with their respective branches are removed from the tree.

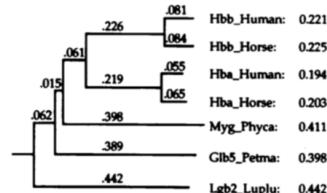
41

## Step 3 : Guide tree construction

### Unrooted NJ tree



### Guide tree



Hbb_human	1	-					
Hbb_horse	2	.17	-				
Hba_human	3	.59	.60	-			
Hba_horse	4	.59	.59	.13	-		
Myg_phyc	5	.77	.77	.75	.75	-	
Glb5_petma	6	.81	.82	.73	.74	.80	-
Lgb2_luplu	7	.87	.86	.86	.88	.93	.90
	1	2	3	4	5	6	7

### **Neighbor-Joining** clustering method:

The tree is then built by linking the least distant pair of nodes. When two nodes are linked, their common ancestral node is added to the tree and the terminal nodes with their respective branches are removed from the tree.

42

## Step 3 : Guide tree construction

At each step, the pair (X,Y) that diverges least among all is selected. In the matrix, the entry X and Y is replaced with X/Y. Each distance involving X and Y is computed as the arithmetic mean of the distances.

Alignment Weighting  
if you align two alignments with 2 and 4 sequences respectively,  
the score at each position is the **average** of 8 (2x4) comparisons.

Without sequence Weights:	
1	peeksav <u>hal</u>
2	geekav <u>hal</u>
3	padktnv <u>kaa</u>
4	aadktnv <u>kaa</u>
5	egewg <u>lhlhv</u>
6	aaekt <u>krsaa</u>

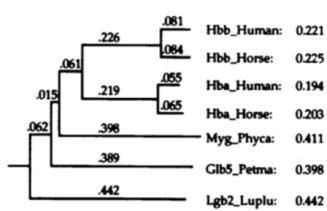
With sequence Weights $W_C$ :	
Score =	$M(t, v)$
	$+ M(t, l)$
	$+ M(l, v)$
	$+ M(l, k)$
	$+ M(k, v)$
	$+ M(k, l)$
	$+ M(v, l)$
	$+ M(v, k)$
	$+ M(l, k)/8$

43

## Step 3 : Guide tree construction

At each step, the pair (X,Y) that diverges least among all is selected. In the matrix, the entry X and Y is replaced with X/Y. Each distance involving X and Y is computed as the arithmetic mean of the distances.

Alignment Weighting  
if you align two alignments with 2 and 4 sequences respectively,  
the score at each position is the **average** of 8 (2x4) comparisons.



The guide tree is then used to calculate WEIGHT for each sequence, which depends on the distance from branch to the root. If a sequence shares a common branch with another sequence, then the two or more sequences will share the weight calculated from the shared branch, and the SEQUENCE LENGTHS will be added together and divided by the number of sequences sharing the same branch.

Without sequence Weights:	
1	peeksav <u>hal</u>
2	geekav <u>hal</u>
3	padktnv <u>kaa</u>
4	aadktnv <u>kaa</u>
5	egewg <u>lhlhv</u>
6	aaekt <u>krsaa</u>

With sequence Weights $W_C$ :	
Score =	$M(t, v)$
	$+ M(t, l)$
	$+ M(l, v)$
	$+ M(l, k)$
	$+ M(k, v)$
	$+ M(k, l)$
	$+ M(v, l)$
	$+ M(v, k)$
	$+ M(l, k)/8$

44

## Step 3 : Guide tree construction

At each step, the pair (X,Y) that diverges least among all is selected. In the matrix, the entry X and Y is replaced with X/Y. Each distance involving X and Y is computed as the arithmetic mean of the distances.

Construct the MSA by seeing the guided tree

**Alignment Weighting**  
if you align two alignments with 2 and 4 sequences respectively,  
the score at each position is the **average** of 8 (2x4) comparisons.

**Without sequence Weights:**

1	peeksav <u>hal</u>	Score = $M(t, v)$
2	geekav <u>hal</u>	$\star M(t, l)$
3	padktv <u>nya</u> a	$\star M(l, v)$
4	aadktv <u>nya</u> a	$\star M(l, l)$
		$\star M(t, l)$
		$\star M(l, t)$
		$\star M(t, t)$
		$M(t, l)/8$

**With sequence Weights  $W_C$ :**

5	egewq <u>l</u> lhv	Score = $M(t, v) * W_C$
6	aaekt <u>k</u> rsa	$\star M(t, l) * W_C$
		$\star M(l, v) * W_C$
		$\star M(l, l) * W_C$
		$\star M(t, l) * W_C$
		$\star M(l, t) * W_C$
		$\star M(t, t) * W_C$
		$M(t, l) * W_C/8$

The guide tree is then used to calculate **WEIGHT** for each sequence, which depends on the distance from branch to the root. If a sequence shares a common branch with another sequence, then the two or more sequences will share the weight calculated from the shared branch, and the **SEQUENCE LENGTHS** will be added together and divided by the number of sequences sharing the same branch.

45

## Step 3 : Guide tree construction

Guide tree

The time complexity for traditional approach is  $O(N^2)$ . CLUSTER OMEGA uses a different method (mBed) to generate the guide tree, which has a time complexity of  $O(N \log N)$  — much faster than ClustalW.

Hbb_human	1	-
Hbb_horse	2	.17 -
Hba_human	3	.59 .60 -
Hba_horse	4	.59 .59 .13 -
Myg_phycia	5	.77 .77 .75 .75 -
Glb5_petma	6	.81 .82 .73 .74 .80 -
Lgb2_luplu	7	.87 .86 .86 .88 .93 .90 -

**Neighbor-Joining** clustering method:  
The tree is then built by linking the least distant pair of nodes. When two nodes are linked, their common ancestral node is added to the tree and the terminal nodes with their respective branches are removed from the tree.

46

## Step 4 : Progressive alignment

The sequences are aligned progressively (global or local algorithm). I can find these cases, seeing the guided tree:

- alignment of 2 sequences, create profile (consensus)
- alignment of 1 sequence and a profile (group of sequences)
- alignment of 2 profiles (groups of sequences)

Adapted from Julie Thompson, IGBMC

Cluster Omega uses HHalign to increase accuracy. Sequences alignment is converted to HMM, which then helps the alignment.

47

## Iterative alignment

Iterative alignment refines an initial progressive multiple alignment by iteratively dividing the alignment into two profiles and realigning them.

Adapted from Julie Thompson, IGBMC

Cluster Omega does iteration on both the guide tree and HMM alignment.

48

## Multiple Alignment: History

### 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

### 1988 Carrillo-Lipman

Branch and Bound approach for MSA

### 1990 Feng-Doolittle

Progressive alignment

### 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

### 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

### 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

### 2004 MUSCLE

**What's next?**

comparison of whole segments of sequences instead of comparison of single nucleic/amino acids.

The program DiAlign constructs alignments from gapfree pairs of similar segments of the sequences.

49

## Multiple Alignment: History

### 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

### 1988 Carrillo-Lipman

Branch and Bound approach for MSA

### 1990 Feng-Doolittle

Progressive alignment

### 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

### 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

### 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

### 2004 MUSCLE

**What's next?**

pre-process a data set of all pairwise alignments between the sequences. This provides us with a library of alignment information that can be used to guide the progressive alignment.

50

## Multiple Alignment: History

### 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

### 1988 Carrillo-Lipman

Branch and Bound approach for MSA

### 1990 Feng-Doolittle

Progressive alignment

### 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

### 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

### 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

### 2004 MUSCLE

**What's next?**

### linguaggio C++

1. Kmer distance between two sequences is defined by first collecting the set of k-mers (subsequences of length k) occurring in the two sequences, then measuring how different the two sets are.
2. re-estimates the tree using the Kimura distance, which is more accurate but requires an alignment.
3. The tree is broken into subtrees, and the sub-alignments refined.

51

## Multiple Alignment: History

### 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

### 1988 Carrillo-Lipman

Branch and Bound approach for MSA

### 1990 Feng-Doolittle

Progressive alignment

### 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

### 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

### 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

### 2004 MUSCLE

**What's next?**

**Big alignments, do they make sense?**

52

## Multiple Alignment: History

### 1975 Sankoff

Formulated multiple alignment problem and gave dynamic programming solution

### 1988 Carrillo-Lipman

Branch and Bound approach for MSA

### 1990 Feng-Doolittle

Progressive alignment

### 1994 Thompson-Higgins-Gibson-ClustalW

Most popular multiple alignment program

### 1998 Morgenstern et al.-DIALIGN

Segment-based multiple alignment

### 2000 Notredame-Higgins-Heringa-T-coffee

Using the library of pairwise alignments

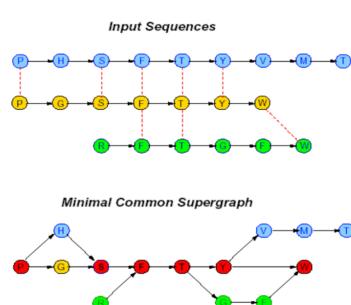
### 2004 MUSCLE

## What's next?

If k is unknown, the problem is NP-complete.

53

## Solution: Representing Sequences as Paths in a Graph



Each protein sequence is represented by a path. Dashed edges connect “equivalent” positions; vertices with identical labels are fused.

54

## Partial Order Multiple Alignment

### **Two objectives:**

- Find a graph that represents domain structure
- Find mapping of each sequence to this graph

### **Solution**

- **PO-MSA (Partial Order Multiple Sequence Alignment)** – for a set of sequences  $S$  is a graph such that every sequence in  $S$  is a path in  $G$ .

55

## MAFFT

### Introduction

MAFFT (Multiple Alignment using Fast Fourier Transform) is a high speed multiple sequence alignment program which implements the Fast Fourier Transform (FFT) to optimise protein alignments based on the physical properties of the amino acids. The program uses progressive alignment and iterative alignment. MAFFT is useful for hard-to-align sequences such as those containing large gaps (e.g., rRNA sequences containing variable loop regions).

56

To do... To see..

BWT-based aligner...

Alignment-free techniques

57

Tandy Warnow | Tutorial | Multiple Sequence  
Alignment | CGSI 2022

<https://www.youtube.com/watch?v=MgkTZgR1ErQ>

58

29

*Review Article*

## An Overview of Multiple Sequence Alignments and Cloud Computing in Bioinformatics

Jurate Daugelaitė,<sup>1</sup> Aisling O' Driscoll,<sup>2</sup> and Roy D. Sleator<sup>1</sup>

<sup>1</sup> Department of Biological Sciences, Cork Institute of Technology, Rossa Avenue, Bishopstown, Cork, Ireland  
<sup>2</sup> Department of Computing, Cork Institute of Technology, Rossa Avenue, Bishopstown, Cork, Ireland

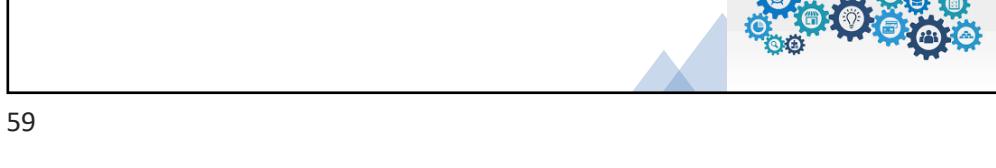
Correspondence should be addressed to Roy D. Sleator; roy.sleator@cit.ie

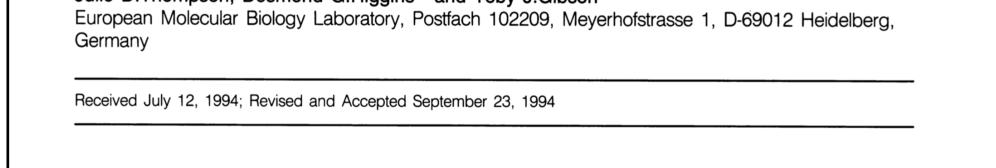
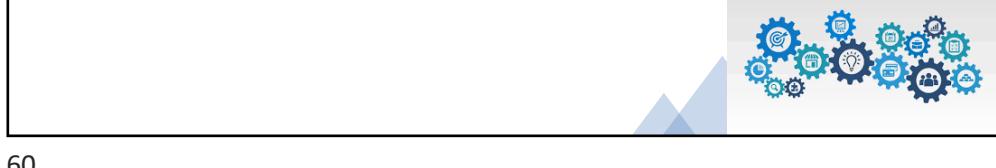
Received 24 May 2013; Accepted 23 June 2013

Academic Editors: M. Glavinovic and X.-Y. Lou

Copyright © 2013 Jurate Daugelaitė et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Multiple sequence alignment (MSA) of DNA, RNA, and protein sequences is one of the most essential techniques in the fields of molecular biology, computational biology, and bioinformatics. Next-generation sequencing technologies are changing the biology landscape, flooding the databases with massive amounts of raw sequence data. MSA of ever-increasing sequence data sets is becoming a significant bottleneck. In order to realise the promise of MSA for large-scale sequence data sets, it is necessary for existing MSA algorithms to be run in a parallelised fashion with the sequence data distributed over a computing cluster or server farm. Combining MSA algorithms with cloud computing technologies is therefore likely to improve the speed, quality, and capability for MSA to handle large numbers of sequences. In this review, multiple sequence alignments are discussed, with a specific focus on the ClustalW and Clustal Omega algorithms. Cloud computing technologies and concepts are outlined, and the next generation of cloud base MSA algorithms is introduced.


59

© 1994 Oxford University Press      *Nucleic Acids Research*, 1994, Vol. 22, No. 22 4673–4680

---

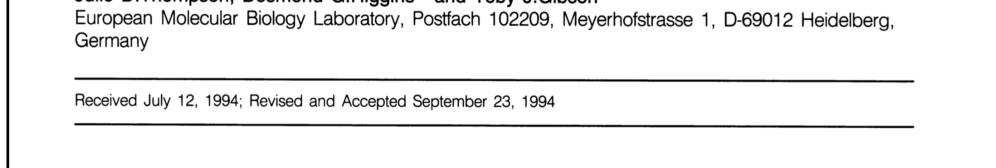
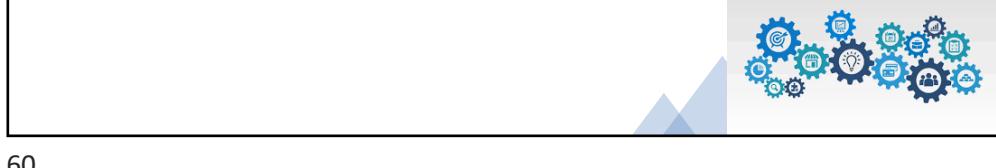
## CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice

---

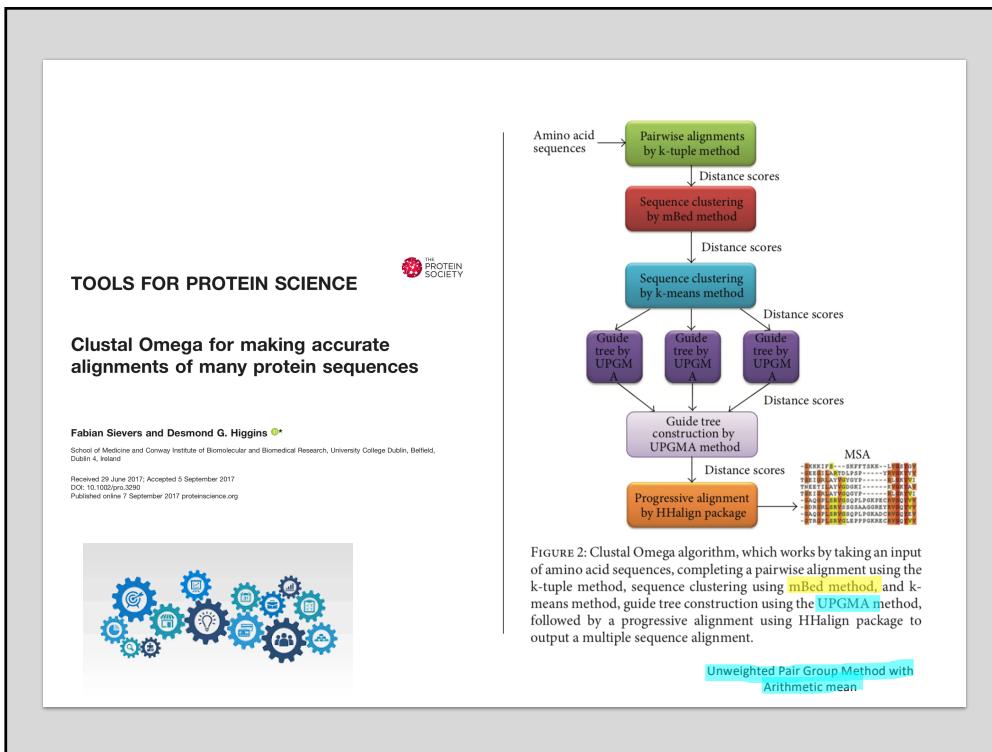
Julie D.Thompson, Desmond G.Higgins<sup>+</sup> and Toby J.Gibson<sup>\*</sup>  
European Molecular Biology Laboratory, Postfach 102209, Meyerhofstrasse 1, D-69012 Heidelberg, Germany

---

Received July 12, 1994; Revised and Accepted September 23, 1994


60



61



62

**THE EFFECT OF THE GUIDE TREE ON MULTIPLE  
SEQUENCE ALIGNMENTS AND SUBSEQUENT  
PHYLOGENETIC ANALYSES**

S. NELESEN, K. LIU, D. ZHAO, C. R. LINDER, AND T. WARNOW \*

*The University of Texas at Austin  
Austin, TX 78712*

*E-mail: {serita,kliu,wzhao,tandy}@cs.utexas.edu, rlinder@mail.utexas.edu*

analizza quanto l'albero guidato  
iniziale possa influenzare l'allineamento  
finale

63

Online Magazine of the European Molecular Biology Laboratory

---

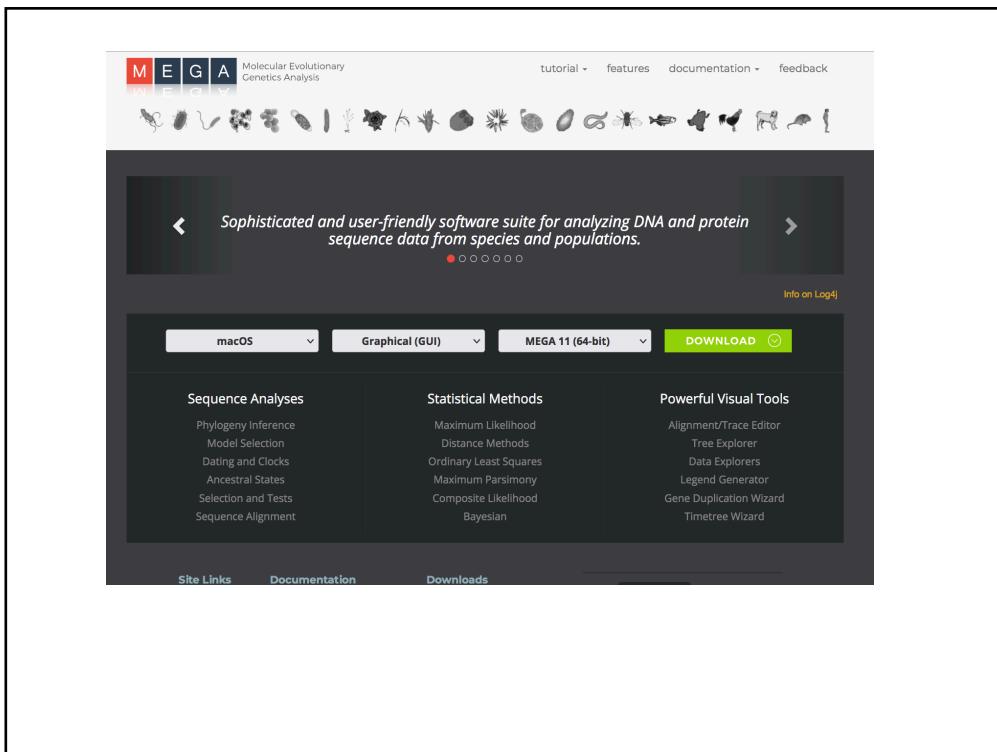
## The story of Clustal: democratising sequence alignments

"I figured out how to make multiple alignments work on these tiny little computers," said Higgins.  
"It meant that now anyone could make their own multiple alignments in their offices."

[https://www.embl.org/  
news/embletc/issue-  
100/the-story-of-  
clustal-democratising-  
sequence-alignments/](https://www.embl.org/news/embletc/issue-100/the-story-of-clustal-democratising-sequence-alignments/)

"Omega is the last letter of the Greek alphabet," said Higgins, adding that many new programs for multiple alignment have come up in recent years, including MAFFT, also hosted by EMBL-EBI. "Life goes on," he added philosophically.

64



65

# EMBL-EBI

[https://www.ebi.ac.uk/seqdb/confluence/display/JDSAT/  
Multiple+Sequence+Alignment](https://www.ebi.ac.uk/seqdb/confluence/display/JDSAT/Multiple+Sequence+Alignment)



**Clustal Omega**

Multiple sequence alignment of DNA or protein sequences. Clustal Omega replaces the older ClustalW alignment tools.



**HMMER**

Fast sensitive protein homology searches using profile hidden Markov models (HMMs) for querying against both sequence and HMM target databases.

abbiamo usato il tool con il file contenente 4 proteine della ferritina (vedi allegato)

66