

# Multiple Sequence Alignment (MSA)

# Multiple Sequence Alignment

Most important contribution of MB to **evolutionary analysis** is the discovery that DNA sequences of different organisms are often **related**.

i.e., genes are **conserved** across widely divergent species, often performing a **similar** or even identical function, and at other times, mutating or rearranging to perform an **altered** function.

Through **simultaneous alignment** of gene sequences, sequence patterns that have been subject to alteration may be analyzed.

# Multiple Sequence Alignment

Aligning more than two sequences

```
A C - - B C D B
- C A D B - D -
A C A - B C D -
```

In an MSA, homologous residues among a set of sequences are aligned together in columns.

'Homologous' is meant in both the structural and evolutionary sense.

# Motivation for MSA

- MSA helps identify conserved regions and those allowed to vary – regions resistant to change are functionally most important to the molecule.
- Carries more information than mere pair-wise alignment
- Multiple sequence similarity suggests common structure for the protein, a common function or evolutionary origin
- MSA requirements are different in the various applications

# Motivation for MSA

Multiple alignments can improve pairwise alignments:

(A) p110 $\alpha$  TFILGIGDRHNSNIMVKDDG-QLFHI DFGHFLDHKKKKFGYKRERVPFVLT--QDFLIVI 142

cAMP-kinase QIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVT DFGFAKRVKGRTWXLCTPEYLAPE 179

(B) p110 $\beta$  SYVLGIG-----DRHSDNINVKKTGQLFHI DFGHILGNFKSKFGIKRERVPFILT 136

p110 $\delta$  TYVLGIG-----DRHSDNIMIRESGQLFHI DFGHFLGNFKTKFGINRERVPFILT 136

p110 $\alpha$  TFILGIG-----DRHNSNIMVKDDGQLFHI DFGHFLDHKKKKFGYKRERVPFVLT 135

p110 $\gamma$  TFVLGIG-----DRHNDNIMITETGNLFHI DFGHILGNYKSFLGINKERVPFVLT 135

p110<sub>dicti</sub> TYVLGIG-----DRHNDNLMVTKGGRLFHI DFGHFLGNYKKKFGFKRERAPFVFT 135

cAMP-kinase QIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVT DFGFAKRVKGRTWXLCTPEYLA 177

Catalytic domains of 5 P13-kinases and cAMP-dependent protein kinase

# MSA for DNA Sequences

In DNA sequences MSA is used in

- **Genome sequence assembly - shotgun sequencing**
- Discovering new regulatory elements
- **Inferring evolutionary relationships**
- DNA barcoding
- **SNP identification**
- Develop primers & probes - use conserved regions to develop
  - Primers for PCR
  - Probes for DNA microarrays

**In which of these applications do we look for similarity/differences?**

# MSA for Protein Sequences

In protein sequences, MSA is used in

- Homology modeling of proteins
- Building phylogenetic tree
- Constructing scoring matrices – PAM, BLOSUM
- Predicting secondary & tertiary structures of new sequences
- Identifying conserved patterns, motifs, blocks in protein sequences – to characterize protein families
- Identify related proteins in database searches, e.g., Profiles, PSI-BLAST, HMMs

# MSA

**Visual alignment of MSA tables use different colours for displaying AAs of different physico-chemical type - aids in identifying conserved patterns:**

Colour	Residue Type	Amino acids
Yellow	Small nonpolar	Gly, Ala, Ser, Thr
Green	Hydrophobic	Cys, Val, Ile, Leu, Pro, Phe, Tyr, Met, Trp
Magenta	Polar	Asn, Gln, His
Red	Negatively charged	Asp, Glu
Blue	Positively charged	Lys, Arg

**Structure prediction tools also give more reliable results when based on MSAs than on single sequences.**



Synechocystis\_sp.\_PCC\_6803/1-40

*Halotheca\_sp.\_PCC\_7418/1-42*

*Cyanotheca\_sp.\_PCC\_8802/1-42*

*Microcystis\_aeruginosa\_DIANCHI905/1-42*

*Microcystis\_aeruginosa\_NIES-843/1-56*

*Cyanotheca\_sp.\_ATCC\_51472/1-42*

*cyanobacterium\_UCYN-A/1-42*

*Crocospaera\_watsonii\_WH\_8501/1-42*

*Microcoleus\_sp.\_PCC\_7113/1-42*

*Cyanotheca\_sp.\_PCC\_7822/1-39*

*Chamaesiphon\_minutus\_PCC\_6605/1-41*

*Thermosynechococcus\_elongatus\_BP-1/1-4*

*Moorea\_producens\_3L/1-41*

*Synechococcus\_sp.\_PCC\_6312/1-43*

*Oscillatoria\_nigro-viridis\_PCC\_7112/1-42*

*Synechococcus\_sp.\_PCC\_7002/1-37*

*Leptolyngbya\_sp.\_PCC\_7376/1-39*

*Trichodesmium\_erythraeum\_IMS101/1-41*

*Leptolyngbya\_sp.\_PCC\_7375/1-41*

*Synechococcus\_sp.\_PCC\_7335/1-41*

*Pseudanabaena\_sp.\_PCC\_7367/1-42*

*Cyanobacterium\_stanieri\_PCC\_7202/1-42*

*Cyanobacterium\_aponinum\_PCC\_10605/1-42*

*Dactylococcopsis\_salina\_PCC\_830/1-42*

*Gloeocapsa\_sp.\_PCC\_73106/1-42*

*Nodularia\_spumigena\_CCY9414/1-50*

*Oscillatoriales\_cyanobacterium\_JSC-12/1-43*

*Anabaena\_sp.\_90/1-48*

*Coleofasciculus\_chthonoplastes\_PCC\_7420/1-42*

*Calothrix\_sp.\_PCC\_7507/1-49*

*Synechococcus\_sp.\_CB0101/1-38*

*Calothrix\_sp.\_PCC\_6303/1-52*

*Synechococcus\_elongatus\_PCC\_7942/1-44*

*Leptolyngbya\_sp.\_PCC\_6406/1-40*

*Synechococcus\_sp.\_CB0205/1-38*

*Synechococcus\_sp.\_CC9311/1-39*

*Nostoc\_sp.\_PCC\_7107/1-49*

*Nostoc\_sp.\_PCC\_7120/1-49*

*'Nostoc\_azollae'\_0708/1-49*

*Nostoc\_sp.\_PCC\_7524/1-49*

*Nostoc\_punctiforme\_PCC\_73102/1-49*

*Cylindrospermum\_stagnale\_PCC\_7417/1-49*

*Synechococcus\_sp.\_BL107/1-38*

*Lyngbya\_sp.\_PCC\_8106/1-41*

*Synechococcus\_sp.\_RCC307/1-38*

*Synechococcus\_sp.\_WH\_7805/1-39*

*Synechococcus\_sp.\_RS9917/1-39*

*Synechococcus\_sp.\_WH\_7803/1-37*

*Synechococcus\_sp.\_RS9916/1-38*

-----MIGLKSFLSTAPVMIALLTFTAGILIEFNRFYPLLFFHP--  
-----MDGLKTFSSAPVLIALLTFTAGILIEFNRFYPLLFFHMP  
-----MEGLTKFLSTAPVLIALLTFTAGLLIEFNRFYPLLFFHPLG  
-----MEGLTKFLSSAPVLIALLTFTAGILIEFNRFYPLLFFHPLG  
MCQTLDFLFPRRTMEGLTKFLSSAPVLIALLTFTAGILIEFNRFYPLLFFHPLG  
-----MEGLTKFLSTAPVLIALLTVTAGILIEFNRFYPLLFFHPLG  
-----MENLTKFLSTAPILIMYLLTFTAGLLIEFNRFYPLLFFHPLG  
-----MAFTKFLSTAPVLIALLTFTAGILIEFNRFYPLLFFHPLG  
-----MIGLPRFLSSAPVLIALLSVTAGILIEFNRFYPLLFFHMS  
-----MAKFLSSAPVLIALLTFTAGLLIEFNRFYPLLFFYPLG  
-----MSHLLRFLSTAPVLAAVVMTFTAGILIEFNRFYPLLFFHPL  
-----MKHFLTLYLSTAPVLAIVMTITAGILIEFNRFYPLLFFHPL  
-----MQLLKYLSTAPVLLAVVMTITAGILIEFNRFYPLLFFHPL  
MQVKYLLTYLSTAPVLAAVVMAFTAGLLIEFNRFYPLLFFHPL  
-----MQYFLKYLSTAPVLAIVMTITAGILIEFNRFYPLLFFHMP  
-----MKFLSSAPVLLTANMVFTAGLLIEFNRFYPLLFFHPL  
-----MGKFLSSAPVLLTANMVFTAGLLIEFNRFYPLLFFHPLG  
-----MQLLKYLSTAPVLATVVMITAGILIEFNRFYPLLFFHMP  
-----MPEGLVKYLSTAPVLATVVMITAGILIEFNRFYPLLFFHPL  
-----MSSNLLKYLSTAPVIATVVMITAGILIEFNRFYPLLFFHPL  
-----MNLKYLSTAPVLATVVMITAGILIEFNRFYPLLFFHPL  
-----MKGLAFLLSTAPVLIALLVFTAGLLIEFNRFYPLLFFHMPG  
-----MKGLTFLSTAPVLIALLVFTAGLLIEFNRFYPLLFFHMPG  
-----MDNFKTFLSSAPVLLTALLTFTAGLLIEFNRFYPLLFFHMP  
-----MKLFTAFLLSTAPVLIALLTFTAGMLIEFNRFYPLLFFHPLQ  
MAEEKGAQSSYFMTFLSTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MQYFMKYLSTAPVIAAIWLTITAGILIEFNRFYPLLFFHPLV  
MAEKGNETNYLITFISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MQYFLKYLSTAPVLIALLVFTAGLLIEFNRFYPLLFFHMP  
MADKGDSSKSYFVFTLTTAPVITITWLTITAGILIEFNRFYPLLFFHPL  
-----MKKFLTAPVFAAIWFTVFTAGILIEFNRFYPLLFFHMP  
MNIILGDLNMIANFLRFLSTAPVMIALLSFTAGLLIEFNRFYPLLFFHPL  
-----MLAMGLKRYLSSAPILATIWFAITAGILIEFNRFYPLLFFHPL  
-----MNLKYLSTAPVIATVVFVITAGILIEFNRFYPLLFFHPL  
-----MKKFLTAPVFAAIWFTVFTAGILIEFNRFYPLLFFHMP  
-----MKKFLTAPVFAAIWFTLTAGILIEFNRFYPLLFFHMP  
MAEKSDQSSYLKFIISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MADKADQSSYLKFIISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MADKSDQSSYLKFIISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MADKTDQSSYLKFIISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MADKGDQSSYLKFIISTAPVAATIWLTITAGILIEFNRFYPLLFFHPL  
-----MADKSDQSSYLKFIISTAPVAATLWLTITAGILIEFNRFYPLLFFHPL  
-----MKKFLTAPVFAAIWFTATAGILIEFNRFYPLLFFHMP  
MQYFLKYLSTAPVLAIAVFTITAGILIEFNRFYPLLFFHPL  
-----MKKFLTAPVFAAIWFTVFTAGILIEFNRFYPLLFFHPL  
-----MQKFLTAPVFAAIWFTLTAGILIEFNRFYPLLFFHPLA  
-----MQKFLTAPVFAAIWFTLTAGILIEFNRFYPLLFFHMPG  
-----MQKFLTAPVFAAIWFTLTAGILIEFNRFYPLLFFHPL  
-----MQKFLTAPVFAAIWFTLTAGILIEFNRFYPLLFFHMP

# MSA

To be informative a MSA should

- contain a distribution of **closely-** and **distantly-** related sequences.

If all closely-related - information contained is largely redundant

⇒ **few inferences can be drawn.**

If all very distantly-related - difficult to construct an accurate alignment

⇒ **quality of results & inferences might be questionable**

Ideally, one should have **a complete range of similarities**, including distantly-related examples linked through chains of close relationships

# Inferences from MSA

Some examples:

- Highly conserved regions - likely to be essential sites for structure/function, e.g. active site
- Regions rich in insertions/deletions - may correspond to loops/turns in proteins
- Build gene/protein families - use conserved regions to guide search
- Basis for phylogenetic analysis - infer evolutionary relationships between genes

## MSA of 8 fragments of immunoglobulin sequences

VTISCTGSSSNIGAG-NHVKWYQQLPG  
VTISCTGTSSNIGS--ITVNWYQQLPG  
LRLSCSSSGFIFSS--YAMYWVRQAPG  
LSLTCTVSGTSFDD--YYSTWVRQPPG  
PEVTCVVVDVSHEDPQVKFNWYVDG--  
ATLVCLISDFYPGA--VTVAWKADS--  
AALGCLVKDYFPEP--VTVSWNSG---  
VSLTCLVKGFYPSD--IAVEWESNG--

The alignment consists of 8 sequences. Conserved residues are highlighted by vertical bars: orange bars for positions 1, 2, and 3; a red bar for position 7; and a blue box for positions 14 through 18. Conserved regions are highlighted by a blue box for positions 14 through 18.

**Conserved residues, regions, patterns**

# Multiple Alignment: Evaluation

VTISCTGSSSNIG-AGNHVKWYQQLPG	VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIG--SITVNWYQQLPG	VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFS--SYAMYWVRQAPG	LRLSCS-SSGFIFSS-YAMYWVRQAPG
LSLTCTVSGTSFD--DYYSTWVRQPPG	LSLTCT-VSGTSFDD-YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNW--YVDG	PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPG--AVTVAW--KADS	ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPE--PVTVSW--NS-G	AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPS--DIAVEW--ESNG	VSLTCLVKGFYPSD--IAVEWESNG--

It is not enough to just get a multiple alignment;  
we need to score the alignment

# Multiple Alignment: Evaluation

A simple way to evaluate a multiple alignment is to evaluate the cost column by column

Sum of Pairs (SP)

$$= \sum_{i < j} D(S_i, S_j)$$

Using unit cost:  
mismatch costs 1,  
match 0, and  
indel costs 1

$$\text{ColumnCost} \begin{pmatrix} \text{L} \\ \text{L} \\ \text{A} \\ \text{P} \\ \text{G} \\ \text{S} \\ - \\ \text{G} \end{pmatrix} = ?$$

Summing the scores of all possible combinations of AA pairs in a column of MSA

# Multiple Alignment: Evaluation

$$\textit{ColumnCost} \begin{pmatrix} \text{L} \\ \text{L} \\ \text{A} \\ \text{P} \\ \text{G} \\ \text{S} \\ - \\ \text{G} \end{pmatrix} = 26.$$

Assumes a model for evolutionary change in which any of the sequence could be the ancestor of others



# SP Scoring Method

- There are problems with SP scoring system as illustrated in the example:

Sequence	Col. A	Col. B	Col. C
1	...N...	...N...	...N...
2	...N...	...N...	...N...
3	...N...	...N...	...N...
4	...N...	...N...	...C...
5	...N...	...C...	...C...
Score	60	24	9

(Using Blosom62):  
N-N: 6, N-C: - 3, C-C: 9

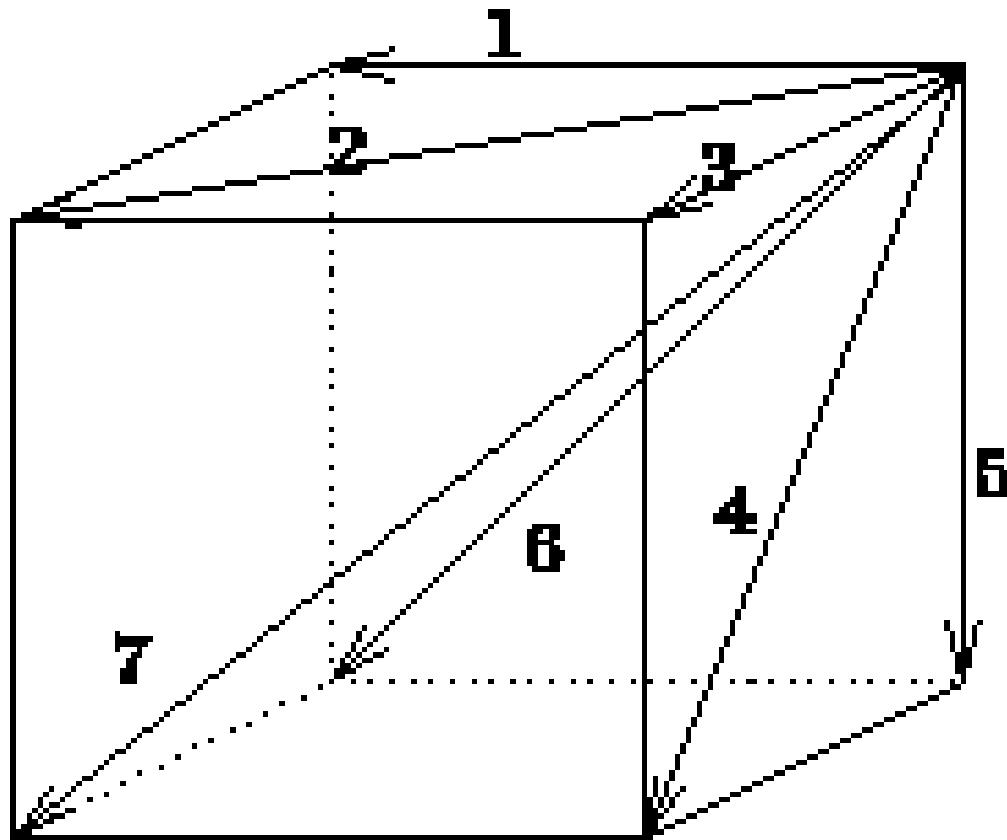
What's the problem?  
Score for N = 10 seq



# Multiple Alignment: DP

- Pair-wise alignment: DP involves a  $L \times L$  ( $L^2$ ) matrix
- Multiple alignment: DP involves an  $L^N$  matrix, an  $N$ -dimensional hyper-lattice,  $N$  - no. of sequences of length  $L$  each, to align simultaneously
- Computationally not feasible: for 5 sequences, each ~100 bp long,  $10^{10}$  matrix elements need to be computed
  - equivalent to a pairwise alignment of two 100,000 bp sequences

# The Recursive Relation



**current  
visit**

$i, j-1, k$   
 $i, j, k-1$   
 $i, j-1, k-1$   
 $i-1, j, k$   
 $i-1, j-1, k$   
 $i-1, j, k-1$   
 $i-1, j-1, k-1$

For 3 sequences, to assign a value to a node  $(i,j,k)$ , we need to consider 7 values; for 2 seqs, we needed only 3!

# Multiple Sequence Alignment

$\alpha_{i_1, i_2, \dots, i_N}$  - maximum score of an alignment up to the subsequences ending with  $x_{i_1}^1, x_{i_2}^2, \dots, x_{i_N}^N$

Recursive relation for multiple sequences:

$$\alpha_{i_1, i_2, \dots, i_N} = \max \left\{ \begin{array}{ll} \alpha_{i_1-1, i_2-1, \dots, i_N-1} & + S(x_{i_1}^1, x_{i_2}^2, \dots, x_{i_N}^N), \\ \alpha_{i_1, i_2-1, \dots, i_N-1} & + S(-, x_{i_2}^2, \dots, x_{i_N}^N), \\ \alpha_{i_1-1, i_2, i_3-1, \dots, i_N-1} & + S(x_{i_1}^1, -, \dots, x_{i_N}^N), \\ & \cdot \\ & \cdot \\ & \cdot \\ \alpha_{i_1-1, i_2-1, \dots, i_N} & + S(x_{i_1}^1, x_{i_2}^2, \dots, -), \\ \alpha_{i_1, i_2, i_3-1, \dots, i_N-1} & + S(-, -, \dots, x_{i_N}^N), \\ & \cdot \\ & \cdot \\ & \cdot \\ \alpha_{i_1, i_2-1, \dots, i_{N-1}-1, i_N} & + S(-, x_{i_2}^2, \dots, -), \\ & \cdot \\ & \cdot \\ & \cdot \end{array} \right.$$

# Multiple Sequence Alignment

To calculate each entry, need to maximize over all  $2^N - 1$  combinations of gaps in a column, excluding the case where all  $\Delta_k$  are zero.

Introducing the notation  $\Delta_i$  which is 0 or 1 and define the 'product'

$$\Delta_i \cdot x = \begin{cases} x & \text{if } \Delta_i = 1, \\ - & \text{if } \Delta_i = 0. \end{cases}$$

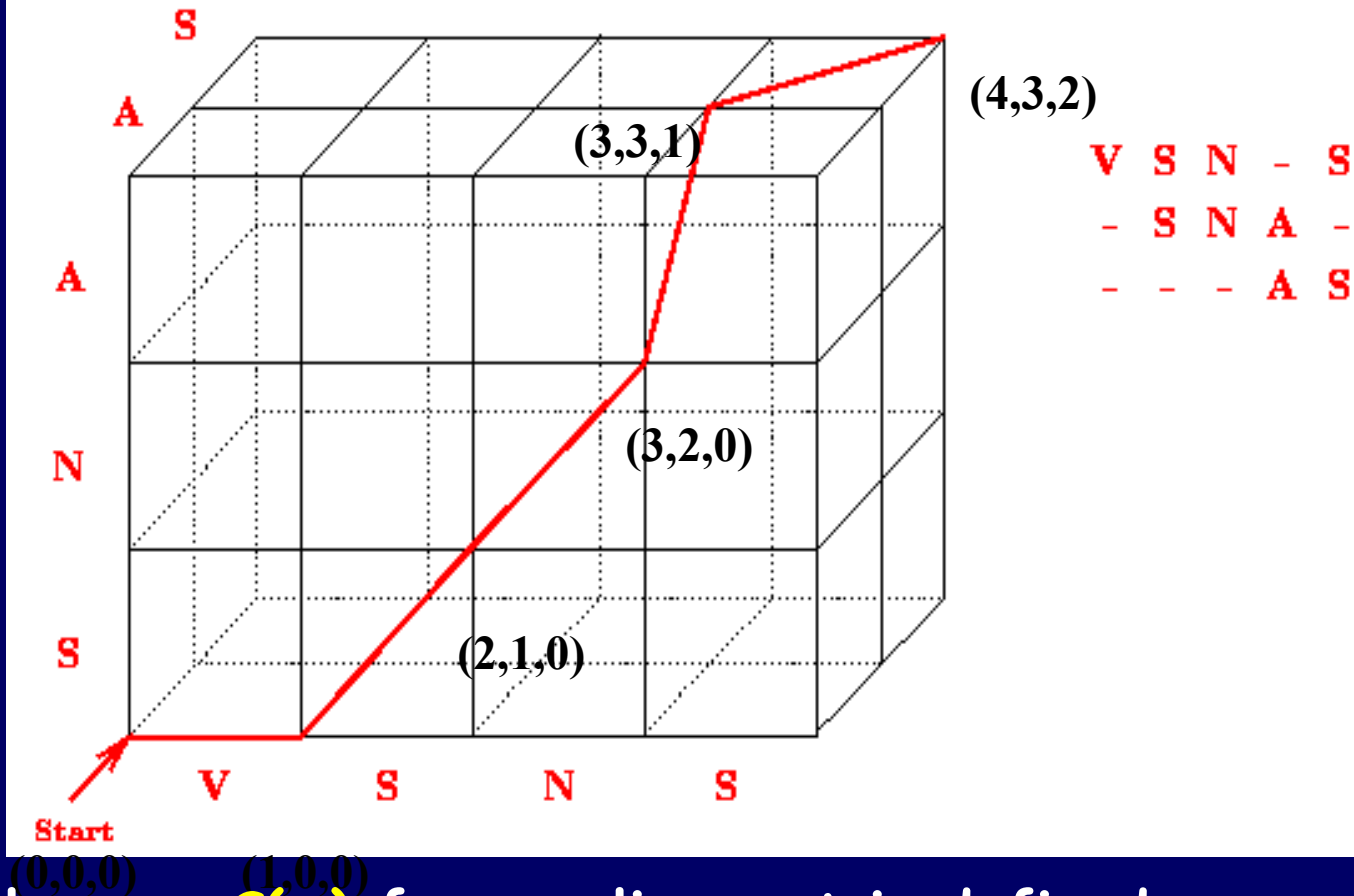
Recursion relation can now be written as

$$a_{i_1, i_2, \dots, i_N} = \max_{\Delta_1 + \dots + \Delta_N > 0} \left\{ a_{i_1 - \Delta_1, i_2 - \Delta_2, \dots, i_N - \Delta_N} + S(\Delta_1 \cdot x_{i_1}^1, \Delta_2 \cdot x_{i_2}^2, \dots, \Delta_N \cdot x_{i_N}^N) \right\}$$

For  $N = 3, 4, 5$ , and  $10$ ,  $2^N - 1 = ?$

# DP Hyperlattice: Example

Figure 1: Alignment Path for 3 Sequences.



Overall score  $S(m)$  for an alignment is defined as a sum of scores  $S(m_i)$  for each column  $i$ :  $S(m) = \sum_i S(m_i)$

# Time and Space complexity

Assuming all sequences of roughly same length  $L$ , memory complexity of the multi-dimensional DP algorithm is  $O(L^N)$  and time complexity is  $O(2^N L^N)$

- impractical for more than a few sequences.

For 6 sequences, each 100bp long, time taken will be  $2^6 \times 100^6 \times 10^{-9} = 64000$  seconds ( $\sim 18$  hrs)

Add 2 more sequences of same length and the no. is  $2.56 \times 10^9$  seconds (over 81 yrs)

Even worse is memory space requirement -  $10^{12}$  for 6 sequences!

# Progressive approach

- Align each sequence to every other pair-wise
- Compute distances between each aligned pair (e.g. no. of mismatches)
- Construct a phylogenetic tree
- Cluster closely related sequences
- Align closely related sequences first
- Gaps inserted in closely related sequences are propagated throughout

Progressive alignment - involves constructing a succession of pairwise alignments.

Tools: ClustalW, T-Coffee, MUSCLE

# Pairwise alignments

Sheep	STCVLSAYWKDLNNYH	Pig	STCVLSAYWRNELNNFH
Cattle	STCVLSAYWKDLNNYH	Rat	STCMLGTY-QD-LNKFH
Sheep	STCVLSAYWK-DLNNYH	Pig	STCVLSAYWRNELNNFH
Pig	STCVLSAYWRNELNNFH	Salmon	STCVLGKL-SQELHKLQ
Sheep	STCVLSAYWKDLNNYH	Human	STCMLGTYQDENKFH
Human	STCMLGTY-QDFNKFH	Rat	STCMLGTYQDLNKT
Sheep	STCVLSAYWKDLNNYH	Human	STCMLGTY-QDFNKT
Rat	STCMLGTY-QDLNKFH	Salmon	STCVLGKLSQELHKLQ
Sheep	STCVLSAYWKD-LNNYH	Rat	STCMLGTY-QDLNKT
Salmon	STCVLGKL-SQELHKLQ	Salmon	STCVLGKLSQELHKLQ
Pig	STCVLSAYWRNELNNEH		
Human	STCMLGTY-QD-ENKFH		



# Hierarchy of Addition

Sheep-Cattle	0	Pig-Rat	8
Sheep-Pig	4	Pig-Salmon	10
Sheep-Human	8	Human-Rat	1
Sheep-Rat	7	Human-Salmon	9
Sheep-Salmon	11	Rat-Salmon	8
Pig-Human	9		

- Align Sheep and Cattle first
- Align Human and Rat
- Align Pig to Sheep and Cattle
- Align these two clusters to each other
- ...
- Align Salmon to large alignment last

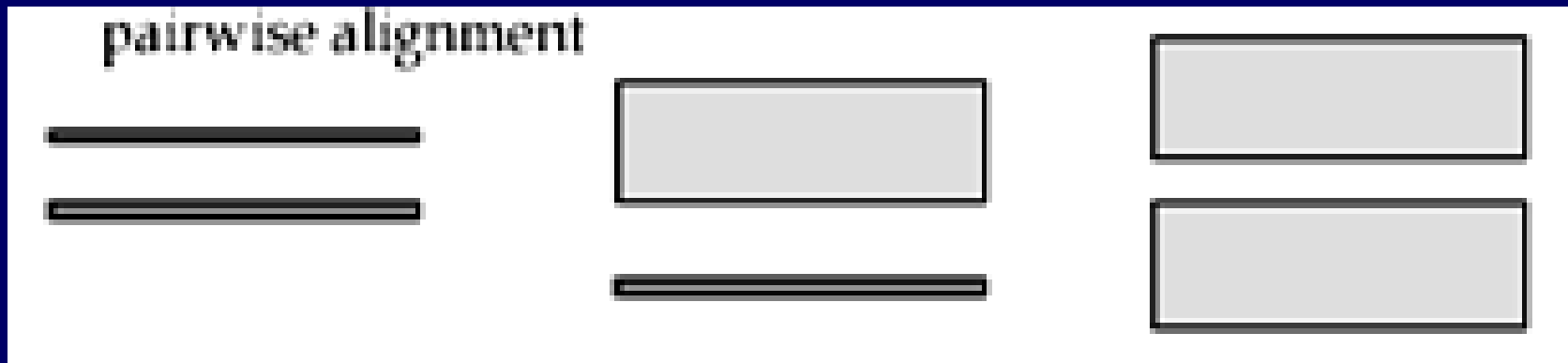
# Progressive alignment

step 1	Sheep	STCVLSAYWKDLNNYH
	Cattle	STCVLSAYWKDLNNYH
step 2	Human	STCMLGTYQDENKEFH
	Rat	STCMLGTYQDLNKEFH
step 3	Sheep	STCVLSAYWK-DLNNYH
	Cattle	STCVLSAYWK-DLNNYH
	Pig	STCVLSAYWRNELNNEFH
...	Salmon	STCVLGKLSQE-LHKLQ

# Aligning Alignments

All possible cases that arise in progressive alignment approach:

- Align two sequences to each other
- Align a sequence to an existing alignment
- Align two alignments to each other



Note - computationally this is always a PW

# Aligning Alignments

- Pairwise alignment of alignments is also called a profile alignment

Again, we can use DP

$$S(i,j) = \max[S(i-1,j-1)+m(i,j), S(i-1,j)+g, S(i,j-1)+g]$$

$m(i,j)$  – similarity score averaged over characters at that position, here  $i$  and  $j$  no longer stands for single sequences, but an average of several, when aligning alignments

$g$  – gap penalty

S T C V L S A YWKD-LNNYH

Sheep

S T C V L S A YWKD-LNNYH

Cattle

S T C V L S A YWRNELNNFH

Pig

SS

TT

CC

MM

LL

GG

TT

YY

QQ

DD

FL

NN

profile alignment

# Aligning Alignments

```
Alignment 1:  ATA
               CCA
               ↓
Alignment 2:  TCAFE
               TAT-E
               TATF-
               AGTFD
```

Score 1<sup>st</sup> column of 1<sup>st</sup> alignment against 2<sup>nd</sup> column in the other alignments using:

$$= 1/8 (\text{score}(A,C) + \text{score}(A,A) + \text{score}(A,A) + \text{score}(A,G) + \\ \text{score}(C,C) + \text{score}(C,A) + \text{score}(C,A) + \text{score}(C,G))$$

# Aligning Alignments

- Once sequences are aligned & gaps introduced, these are not altered – the alignment method is hierarchical
- ClustalW finds a local optimum as early alignment decisions are “locked in” by the “greedy” algorithm
- Early errors will be propagated and cause the final alignment to be worse

Sheep	STCVLSAYWK-DLNNYH
Cattle	STCVLSAYWK-DLNNYH
Pig	STCVLSAYWRNELNNFH

# ClustalX/ClustalW

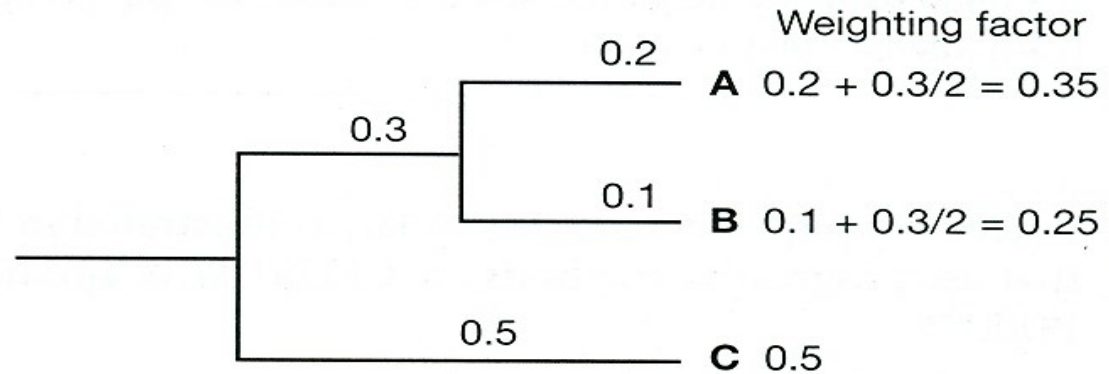
- **Performs pair-wise alignments of all the sequences**
  - k-tuple based alignment (fast/approx.), full DP (slow/accurate)
- **Uses alignment scores to produce a phylogenetic tree**
  - Genetic distance: no. of mismatches/no. of matches (positions against gaps not considered)
- **Aligns sequences sequentially, guided by the phylogenetic relationships indicated by tree.**
  - Sequence contributions are 'weighted' by their relationship on the predicted tree, **weights** based on the distance of each sequence from root



**Alignment scores between two positions in msa calculated using resulting weights as multiplication factors**

**Weights are normalized so that largest weight is 1**

### A. Calculation of sequence weights



### B. Use of sequence weights

Column in alignment 1

Sequence A (weight a) .....K.....

Sequence B (weight b) .....I.....

Column in alignment 2

Sequence C (weight c) .....L.....

Sequence D (weight d) .....V.....

Score for matching these two column in an msa =

$$[ a \times c \times \text{score}(K,L) + a \times d \times \text{score}(K,V) + b \times c \times \text{score}(I,L) + b \times d \times \text{score}(I,V) ] / 4$$

## Basic idea in Progressive Heuristic Approach:

- compute pairwise alignments and merge alignments consistently

Consider alignment of 3 sequences:

acg, cga, gac

Get optimal pairwise alignments:

a c g -	- a c g	c g a -
- c g a	g a c -	- g a c

1&2

a	c	g	-
-	c	g	a

1&3

-	a	c	g
g	a	c	-

2&3

c	g	a	-
-	g	a	c

Merge using  
alignments  
with 1<sup>st</sup> sequence

-acg-  
--cga  
gac--

Merge using  
alignments  
with 3<sup>rd</sup> sequence

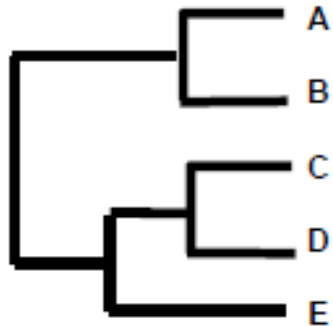
--acg  
cga--  
-gac-

Merge using  
alignments  
with 2<sup>nd</sup> sequence

acg--  
-cga-  
--gac

**Order of merging matters !**  
**Note once a gap, always a gap ...**

# Example



C PADKTNVKAANGKVG**A**HAGEYGA

D AADKTNVKAAWSKVGGHAGEYGA

A PEEKSAVTALWGKVN**V**DEYGG

B GEEKAAVLALWDKVNEEEYGG

C PADKTNVKAANG\_KVG**A**HAGEYGA

D AADKTNVKAAWS\_KVGGHAGEYGA

E AA\_\_TNVKTAWSSKVGGHAPA\_\_**A**

A PEEKSAV\_TALWG\_KVN\_\_VDEYGG

B GEEKAAV\_LALWD\_KVN\_\_EEYGG

C PADKTNVKA**A**\_WG\_KVG**A**HAGEYGA

D AADKTNVKA**A**\_WS\_KVGGHAGEYGA

E AA\_\_TNV**K**TA\_WSSKVGGHAPA\_\_**A**

Once a gap, always a gap

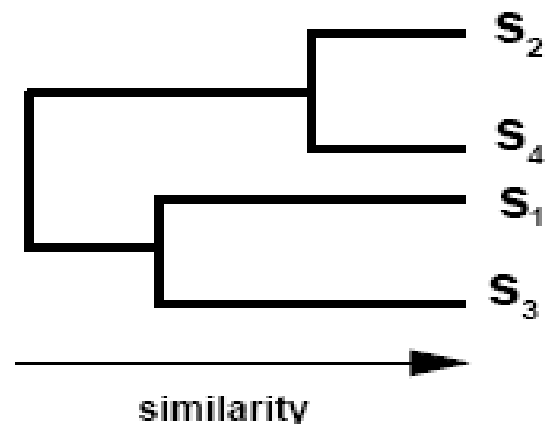
# Steps in Progressive Alignment

## (A) Pairwise Alignment

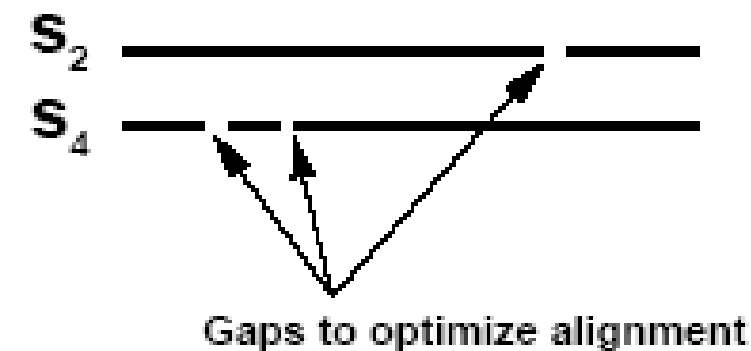
Example - 4 sequences  $S_1$   $S_2$   $S_3$   $S_4$

$S_1$  \_\_\_\_\_  
 $S_2$  \_\_\_\_\_  
 $S_3$  \_\_\_\_\_  
 $S_4$  \_\_\_\_\_

6 pairwise comparisons  
then cluster analysis



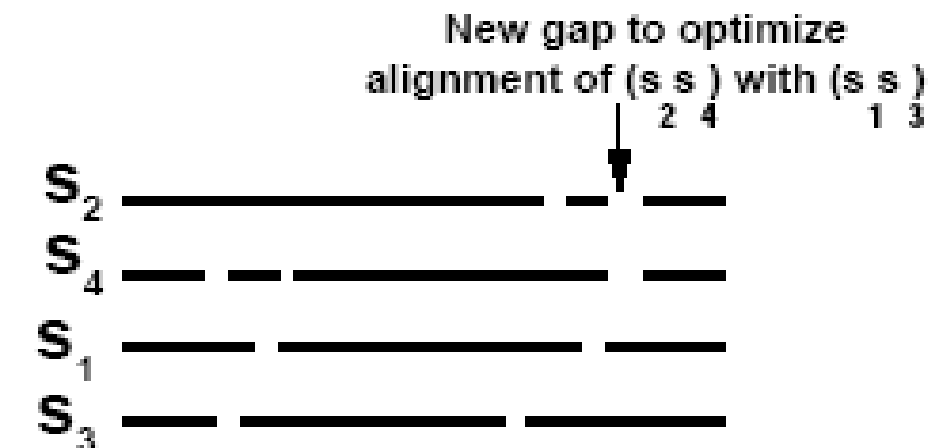
## (B) Multiple alignment following the tree from A



align most similar pair



align next most similar pair



align alignments - preserve gaps

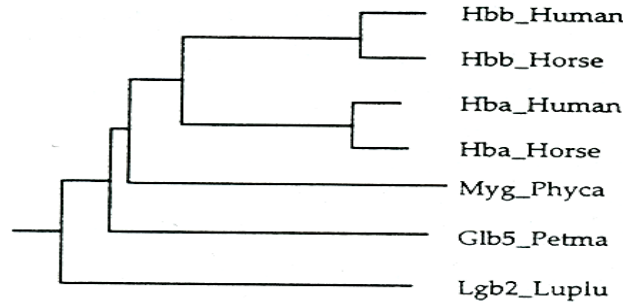
# ClustalW: Refinements

- Different BLOSUM matrices used based on similarity of sequences – to reflect evolutionary changes better
- Recall: BLOSUM80/ BLOSUM62 are based on sequences that are 80% / 62% identical, *i.e.* lower numbers for more distant sequences
- ClustalW calculates gaps in a way to place them between secondary structural elements
- Frequency of gaps next to each amino acid is based on the table provided by Pascarella and Argos.

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_Luplu	7	.87	.86	.86	.88	.93	.90
		1	2	3	4	5	6

Pairwise alignment:  
Calculate distance matrix

## MSA of 7 globins by CLUSTALW



Rooted Neighbor Joining  
tree (guide tree)

```

-----VHLTPEEKSAVTALWGKVNVDEVGGGEALGRLLVVYPWTQRFESFGDLST
-----VQLSGEKAQAVLALWDKVNVEEVGGGEALGRLLVVYPWTQRFDSFGDLSN
-----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLS--
-----VLSAADKTNVKAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLS--
-----VLSEGEWQLVLHVWAKVEALVAGHGQDILIRLFKSHPETLEKFDKFKHLKT
PIVDTGTSVAPLSAAEKTIRSAPVYSTYETSGVDILVKFFTSPTAAQEFFPKFKGLTT
-----GALTESQAALVKSSWEEFNANIPKHTHREFILVLEIAPAAKILFSFLKGTSE
          *               *               *               *

```

Progressive  
alignment:  
Align following  
the guide tree

```

PDAVMGNPKVKAHGKKV LGA FSDGLAHLD-----NLKGTTFATLSELHCDKLHVDPENFRL
PGAVMGNPKVKAHGKKV LHSFGEGVHHLD-----NLKGTFAALSELHCDKLHVDPENFRL
----HGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRLVDPVNFKL
----HGSAQVKAHGKKVGDALTLAVGHLD-----DLPGALSNSDLHAHKLRLVDPVNFKL
EAEMKASEDLKKHGVTVLTAALGAILKKKG----HHEAELKPLAQSHATKHKIPIKYLEF
ADQLKKSA D V R W H A E R I I N A V N D A V A S M D D T --E K M S M K L R D L S G K H A K S F Q V D P Q Y F K V
V P --Q N N P E L O A H A G K V F K L V Y E A A I Q L Q V T G V V V T D A T L K N L G S V H V S K G - V A D A H F P V
          *               *               *               *

```

```

LGNVLVCVLAHHEFGKEFTTPPVQAAAYQKVVAGVANALAHKYH-----
LGNVLVVVLARHFGKDFTPELQASQYQKVVAGVANALAHKYH-----
LSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTISKYR-----
LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTISKYR-----
ISEAIIHVLHSRHPGDFGADAGQAMNKALELFRKDIAAKYKELGYQG
LAAVIADTVAAG-----DAGFEKLMSMICILLRSAY-----
VKEAIIKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---

```

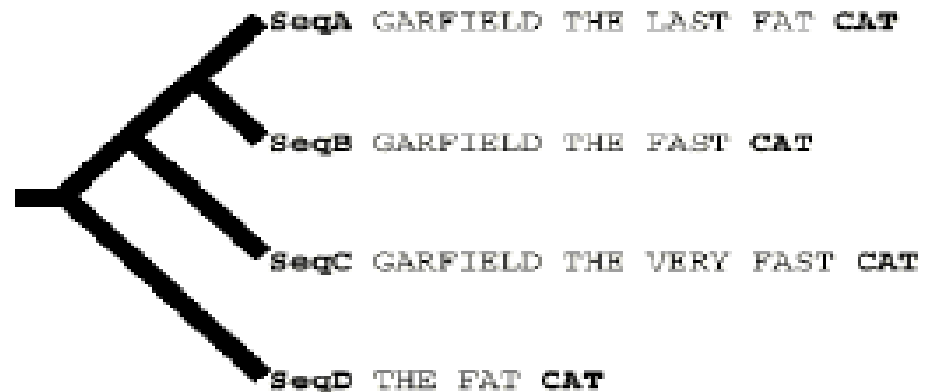
Known locations of 7  $\alpha$ -  
helices in the structure  
of this group shown in  
boxes



# Problems with Progressive Alignment

Gaps at the ends are penalized less, so CAT is aligned with FAT in sequence 2

The greedy approach results in efficiency of the algorithm at the cost of accuracy



CLUSTALW (Score=20, Gap=-1, Gap=0, M=1)

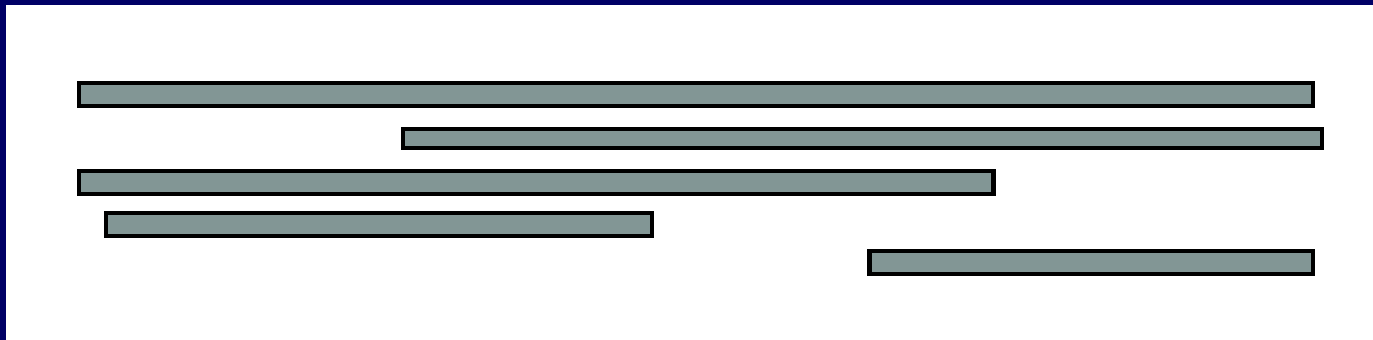
SeqA	GARFIELD	THE	LAST	FA-T	CAT
SeqB	GARFIELD	THE	FAST	CA-T	---
SeqC	GARFIELD	THE	VERY	FAST	CAT
SeqD	-----	THE	----	FA-T	CAT

CORRECT (Score=24)

SeqA	GARFIELD	THE	LAST	FA-T	CAT
SeqB	GARFIELD	THE	FAST	----	CAT
SeqC	GARFIELD	THE	VERY	FAST	CAT
SeqD	-----	THE	----	FA-T	CAT

# ClustalW Misapplied

ClustalW and other algorithms that include an initial pair-wise comparison step should not be used to align sequences that do not all share a common block.



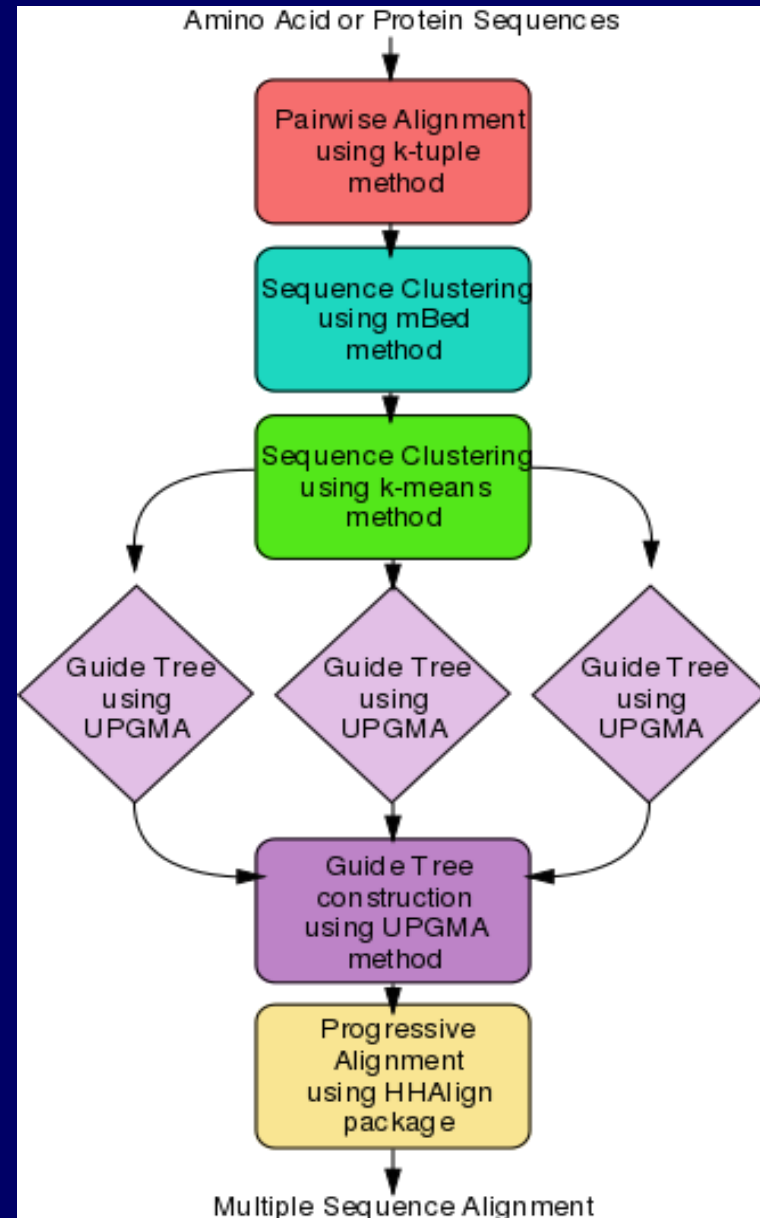
**Not suitable for sequence assembly!**

The latest version, called **Clustal Omega**, uses seeded guide trees and HMM profile-profile techniques to generate alignments

# ClustalW / ClustalΩ online

1. Produce pairwise alignment using k-tuple method
2. Sequences are clustered using mBed method, which calculates pairwise distance using sequence embedding.
3. This is followed by k-means clustering method.
4. Next, the guide tree is constructed using the UPGMA method.
5. Finally, MSA is produced using HHAAlign package from the HH-Suite, which uses two profile HMM's.

<https://www.ebi.ac.uk/Tools/msa/>



# Problems with Progressive Alignment

- Depends on the very first closely related sequences used for constructing the multiple alignment
- If these sequences align well, there will be few errors
- More distantly related these sequences are, more errors will get propagated through the alignment
- Using Bayesian methods such as Hidden Markov models (HMMs) may be useful for aligning more distantly related sequences.

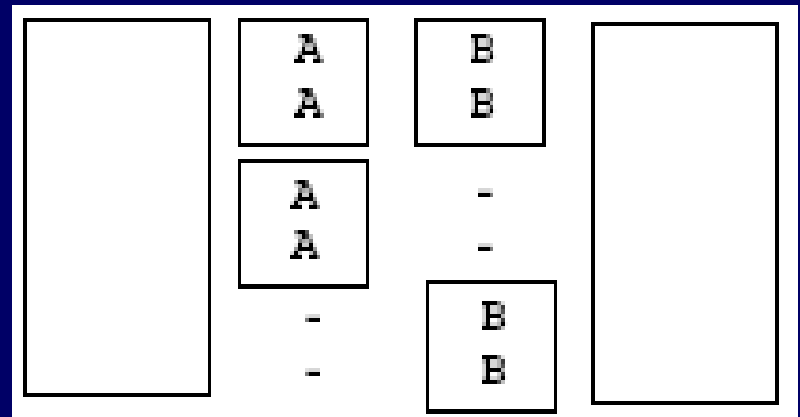
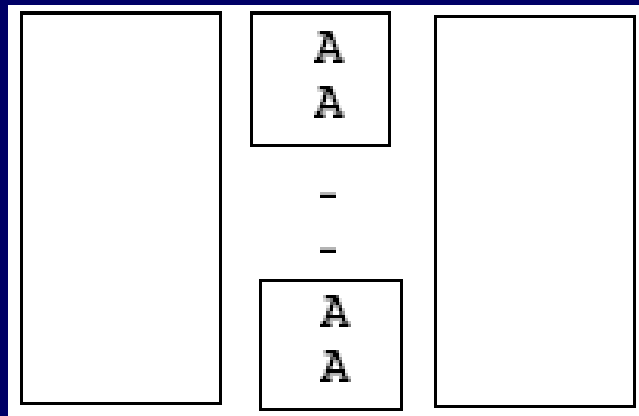
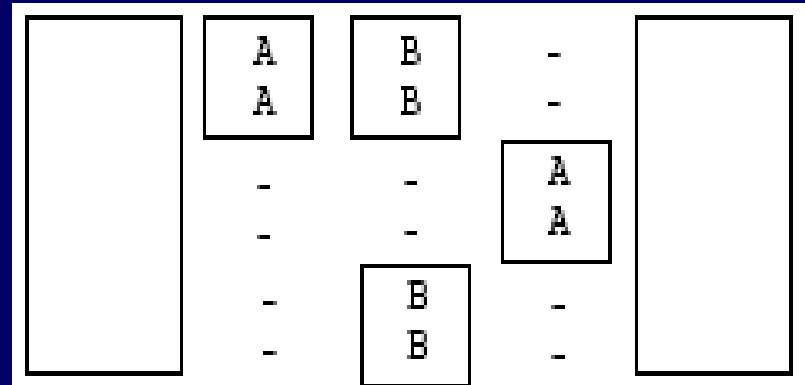
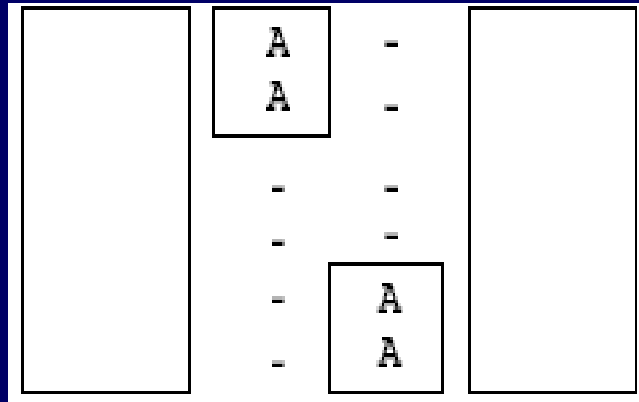
**Solution: Stochastic or Iterative Methods**

# Iterative Methods of MSA

- To correct for errors introduced by initial alignment, use iterative methods: re-align sub-groups of sequences and then align these sub-groups into a global alignment
- **Objective – to improve overall alignment score**
- Selection of groups may be based on phylogenetic tree, separation of one or two sequences from the rest, or a random selection of the groups.
- **Programs using iterative methods – MultiAlin, PRRP and DIALIGN**

# Examine alignments by eye

Some artifacts observed in the output of algorithms:



Always examine your alignment manually to see if it can be improved.

# Edit alignments manually

## Multiple sequence alignment tools

- **Viewers**
  - ClustalX, Jalview, Cinema, Sequence logos
- **Editors / annotation**
  - SeqVu, MACAW, BioEdit
- **BioEdit available at:**  
<https://bioedit.software.informer.com/7.2/>

# Work with proteins

- Twenty symbols to match as against four for DNA
- No noise resulting from the degeneracy of genetic code
- More sensitive scoring matrices
- Requires less of manual editing



# Choose genes judiciously

- When inferring phylogeny choose genes carefully
- For closely related organisms choose genes which mutate fast
- For distantly related species choose slowly mutating genes
- Compare orthologous genes between species and paralogous ones within an organism

# Summary

- Treat the output of multiple alignment programs as a first alignment
- Examine it by eye and edit it manually to improve it
- Get rid of low confidence or highly divergent regions
- Ensure you have started with a sensible evolutionary hypothesis

# Summary

How does one perform an MSA?

- By hand: too hard!
- Automated alignment: Fast, but doesn't necessarily produce the "correct" alignment

**Best approach = Automated alignment with  
manual editing**

# References

- David W. Mount, *Bioinformatics: Sequence and Genome Analysis*, CBS Publishers & distributors, New Delhi, and references therein.
- Thompson, J.D., Higgins, D.G. and Gibson. T.J. "CLUSTALW: improving the sensitivity of multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice" *Nuc. Acids Res.* 22, 4673-80 (1994).